

Statistical Inference via Convex Optimization

Anatoli Juditsky and Arkadi Nemirovski

StatOpt[•]LN[•]NS January 21, 2019 7x10

Contents

Li	st of	Figur	es	\mathbf{xi}
Pr	efac	e		xiii
A	ckno	wledgi	nents	xv
No	otati	onal c	onventions	1
Ał	oout	proof	s	4
		proor		-
1	Spa	rse Re	ecovery via ℓ_1 Minimization	5
	1.1	Comp	Circal Brancesses Brackless	5 F
		1.1.1	Signal Recovery Problem	5 6
		1.1.2 1 1 2	Signal Recovery: parametric and non-parametric cases	10
		1.1.0	Compressed Sensing via ℓ_1 minimization. Motivation	10
	19	Vəlidi	ty of sparse signal recovery via ℓ_1 minimization	10
	1.2	121	Validity of l_1 minimization in the noiseless case	12
		1.2.1	1.2.1.1 Notational convention	12
			1.2.1.2 s-Goodness	13
			1.2.1.3 Nullspace property	15^{-5}
		1.2.2	Imperfect ℓ_1 minimization	16
			1.2.2.1 Contrast matrices and quantifications of Nullspace	
			property	16
		1.2.3	Regular ℓ_1 recovery	18
		1.2.4	Penalized ℓ_1 recovery	18
		1.2.5	Discussion	19
	1.3	Verifia	ability and tractability issues	22
		1.3.1	Restricted Isometry Property and s-goodness of random ma-	
			trices	25
		1.3.2	Verifiable sufficient conditions for $\mathbf{Q}_q(s,\kappa)$	25
		1.3.3	Tractability of $\mathbf{Q}_{\infty}(s,\kappa)$	27
			1.3.3.1 Mutual Incoherence	28
			1.3.3.2 From RIP to conditions $\mathbf{Q}_q(\cdot,\kappa)$	29
			1.3.3.3 Limits of performance of verifiable sufficient condi-	20
	1 4	D	tions for goodness	29
	1.4	Exerc	Ises for Lecture 1	31 25
	1.0	1 5 1	Proofs of Theorem 1.3, 1.4	
		1.5.1	Proof of Theorem 1.5	36
		1.5.2	Proof of Proposition 1.7	38
		1.5.5 1.5.4	Proof of Propositions 1.8, 1.12	41
		155	Proof of Proposition 1 10	43

		1.5.6	Proof of Proposition 1.13	44
2	Hyj	oothesi	is Testing	46
	2.1	Prelin	inaries from Statistics: Hypotheses, Tests, Risks	46
		2.1.1	Hypothesis Testing Problem	46
		2.1.2	Tests	47
		2.1.3	Testing from repeated observations	47
			2.1.3.1 Stationary K-repeated observations	48
			2.1.3.2 Semi-stationary K-repeated observations	48
			2.1.3.3 Quasi-stationary K-repeated observations	49
		2.1.4	Risk of a simple test	50
		2.1.5	Two-point lower risk bound	51
	2.2	Hypot	hesis Testing via Euclidean Separation	54
		2.2.1	Situation	54
		2.2.2	Pairwise Hypothesis Testing via Euclidean Separation	55
			2.2.2.1 The simplest case	55
			2.2.2.2 Extension	56
			2.2.2.3 Further extensions: spherical families of distributions	58
		2.2.3	Euclidean Separation, Repeated Observations, and Majority	
			Tests	60
			2.2.3.1 Preliminaries: Repeated observations in "signal plus	
			noise" observation model	61
			2.2.3.2 Majority Test	61
		2.2.4	From Pairwise to Multiple Hypotheses Testing	64
			2.2.4.1 Situation	64
			2.2.4.2 Closeness relation and "up to closeness" risks	65
			2.2.4.3 Multiple Hypothesis Testing via pairwise tests	65
			2.2.4.4 Testing Multiple Hypotheses via Euclidean separa-	
			tion	67
	2.3	Detect	ors and Detector-Based Tests	70
		2.3.1	Detectors and their risks	70
		2.3.2	Detector-based tests	71
			2.3.2.1 Structural properties of risks	71
			2.3.2.2 Renormalization	71
			2.3.2.3 Detector-based testing from repeated observations	72
			2.3.2.4 Limits of performance of detector-based tests	75
	2.4	Simple	e observation schemes	77
		2.4.1	Simple observation schemes – Motivation	77
		2.4.2	Simple observation schemes – Definition	78
		2.4.3	Simple observation schemes – Examples	79
			2.4.3.1 Gaussian observation scheme	79
			2.4.3.2 Poisson observation scheme	80
			2.4.3.3 Discrete observation scheme	82
			2.4.3.4 Direct products of simple observation schemes	82
		2.4.4	Simple observation schemes – Main result	84
			2.4.4.1 Executive summary of convex-concave saddle point	
			problems	84
			2.4.4.2 Main Result	86
		2.4.5	Simple observation schemes – Examples of optimal detectors	91
			2.4.5.1 Gaussian o.s	92

iv

CONTENTS	S
----------	---

V

		2.4.5.2	Poisson o.s.	92
		2.4.5.3	Discrete o.s.	93
		2.4.5.4	K-th power of simple o.s	93
2.5	Testin	g multiple	e hypotheses	95
	2.5.1	Testing	unions	95
		2.5.1.1	Situation and goal	95
		2.5.1.2	The construction	96
	2.5.2	Testing	multiple hypotheses "up to closeness"	99
		2.5.2.1	Situation and goal	99
		2.5.2.2	"Building blocks" and construction	101
		2.5.2.3	Testing multiple hypotheses via repeated observa-	
			tions	102
		2.5.2.4	Consistency and near-optimality	103
	2.5.3	Illustrat	ion: Selecting the best among a family of estimates	105
		2.5.3.1	The construction	106
		2.5.3.2	A modification	108
		2.5.3.3	"Near-optimality"	110
		2.5.3.4	Numerical illustration	112
2.6	Seque	ntial Hyp	othesis Testing	113
	2.6.1	Motivati	ion: Election Polls	113
	2.6.2	Sequenti	ial hypothesis testing	116
	2.6.3	Conclud	ing remarks	121
2.7	Measu	rement D	Design in simple observation schemes	121
	2.7.1	Motivati	ion: Opinion Polls revisited	121
	2.7.2	Measure	ement Design: SetUp	123
	2.7.3	Formula	ting the MD problem	124
		2.7.3.1	Simple case, Discrete o.s.	126
		2.7.3.2	Simple case, Poisson o.s.	128
		2.7.3.3	Simple case, Gaussian o.s	129
2.8	Affine	detectors	beyond simple observation schemes	131
	2.8.1	Situation	n	132
		2.8.1.1	Preliminaries: Regular data and associated families	
			of distributions	132
		2.8.1.2	Basic examples of simple families of probability dis-	
			tributions	133
		2.8.1.3	Calculus of regular and simple families of probabil-	
			ity distributions	134
	2.8.2	Main res	sult	140
		2.8.2.1	Situation & Construction	140
		2.8.2.2	Main Result	140
	_	2.8.2.3	Illustration: sub-Gaussian and Gaussian cases	144
2.9	Beyon	d the sco	pe of affine detectors: lifting observations	147
	2.9.1	Motivati	ion	147
	2.9.2	Quadrat	ic lifting: Gaussian case	147
	2.9.3	Quadrat	tic lifting – does it help?	150
	2.9.4	Quadrat	ic lifting: sub-Gaussian case	153
	2.9.5	Recoveri	ing quadratic form of discrete distribution	155
	2.9.6	Generic	application: quadratically constrained hypotheses .	157
	-	2.9.6.1	Simple change detection	159
2.10	Exerc	ises for Le	ecture 2	166

		2.10.1	Two-point lower risk bound	167
		2.10.2	Around Euclidean Separation	167
		2.10.3	Hypothesis testing via ℓ_1 -separation	167
		2.10.4	Miscellaneous exercises	172
	2.11	Proofs		176
		2.11.1	Proof of Claim in Remark 2.10	176
		2.11.2	Proof of Proposition 2.8 in the case of quasi-stationary K -	
			repeated observations	177
		2.11.3	Proof of Proposition 2.40	180
		2.11.4	Proof of Proposition 2.46	181
		2.11.5	Proof of Proposition 2.49	186
	D (*			100
3	2 1	Fatime	g Functions via Hypotnesis Testing	100
	0.1	2 1 1	The problem	100
		$\begin{array}{c} 0.1.1 \\ 0.1.0 \end{array}$	The estimate	101
		0.1.2	The estimate	191
		3.1.3 2.1.4	Main result	192
		0.1.4 2.1.5	Wear-optimality	195
	2 9	5.1.5 Fatime	mustration	197
	3.2	291	Outline	199
		0.4.1 2.0.0	Estimating N convex functions: problem's setting	202
		3.2.2	Estimating N -convex functions. problem's setting \dots	202
		292	Bisogtion Estimate: Construction	203
		0.2.0	3.2.2.1 Proliminaries	204
		394	Building the Bisection estimate	204
		3.2.4	3.2.4.1 Control parameters	200
			3.2.4.2 Bisoction estimate: construction	200
		395	Bisoction estimate: Main result	$200 \\ 207$
		3.2.0	Illustration	201
		3.2.0 3.2.7	Estimating N convex functions: an alternative	200
		0.2.1	3.2.7.1 Numerical illustration	211 919
			3.2.7.2 Estimating dissipated power	212
	22	Estime	5.2.1.2 Estimating dissipated power	214
	0.0	221	Situation and goal	$210 \\ 217$
		332	Construction & Main results	217
		222	Estimation from repeated observations	210
		3.3.4	Application: Estimating linear form of sub-Gaussianity pa-	220
		0.0.1	rameters	222
			3 3 4 1 Consistency	222
			3 3 4 2 Direct product case	224
			3.3.4.3 Numerical illustration	224 225
	3 /	Estime	ating quadratic forms via quadratic lifting	$\frac{220}{297}$
	0.4	3/1	Estimating quadratic forms, Caussian case	$\frac{221}{227}$
		0.7.1	3 4 1 1 Preliminaries	227
			3.4.1.2 Estimating quadratic form: Situation & goal	221
			3.4.1.3 Construction & Result	220
			3414 Consistency	229 921
			3415 A modification	201 222
		319	Estimating quadratic form sub Caussian case	202 222
		0.4.4	Louina quadratic torm, sub-Gaussian case	404

vi

CONT	ΓENTS
------	-------

			3.4.2.1 Situation	232
			3.4.2.2 Construction & Result	233
			3.4.2.3 Numerical illustration, direct observations	237
			3.4.2.4 Numerical illustration, indirect observations	241
	3.5	Exerci	ises for Lecture 3	244
	3.6	Proofs	3	256
		3.6.1	Proof of Proposition 3.5	256
			3.6.1.1 Proof of Proposition 3.5.i	256
			3.6.1.2 Proof of Proposition 3.5.ii	258
		3.6.2	2-convexity of conditional quantile	259
		3.6.3	Proof of Proposition 3.15	262
	~ .			
4	Sig	nal Ree	covery from Gaussian Observations and Beyond	265 265
	4.1	Prelim	inaries: Executive Summary on Conic Programming	26
		4.1.1	Cones	267
		4.1.2	Conic problems and their duals	268
		4.1.3	Schur Complement Lemma	270
	4.2	Near-(Optimal Linear Estimation	270
		4.2.1	Situation and goal	270
			4.2.1.1 Ellitopes	270
			4.2.1.2 Estimates and their risks	272
			4.2.1.3 Main goal	27:
		4.2.2	Building linear estimate	272
			4.2.2.1 Illustration: Recovering temperature distribution.	27_{-}
		4.2.3	Near-optimality of \widehat{x}_{H_*}	276
			4.2.3.1 Relaxing the symmetry requirement	277
		4.2.4	Numerical illustration	278
		4.2.5	Byproduct on semidefinite relaxation	278
	4.3	From	ellitopes to spectratopes	279
		4.3.1	Spectratopes: definition and examples	280
			4.3.1.1 Examples of spectratopes	281
		4.3.2	Semidefinite relaxation on spectratopes	282
		4.3.3	Linear estimates beyond ellitopic signal sets and $\ \cdot\ _2$ -risk.	283
			4.3.3.1 Situation and goal	283
			$4.3.3.2 \text{Assumptions} \dots \dots \dots \dots \dots \dots \dots \dots \dots $	28_{-}
			4.3.3.3 Building linear estimate	286
			4.3.3.4 Upper-bounding $\ \cdot\ _{\mathcal{X},\ \cdot\ }$	280
			4.3.3.5 Upper-bounding $\Psi_{\Pi}(\cdot)$	28'
			4.3.3.6 Putting things together	289
			4.3.3.7 Illustration: covariance matrix estimation	289
			4.3.3.8 Estimation from repeated observations	294
			4.3.3.9 Near-optimality in Gaussian case	295
	4.4	Linear	estimates of stochastic signals	296
		4.4.1	Minimizing Euclidean risk	297
		4.4.2	$Minimizing \ \cdot \ \text{-risk} \dots \dots \dots \dots \dots \dots \dots \dots \dots $	298
	4.5	Linear	estimation under uncertain-but-bounded noise	300
		4.5.1	Uncertain-but-bounded noise	300
			4.5.1.1 Building linear estimate	301
			4.5.1.2 Near-optimality	301
			4.5.1.3 Nonlinear estimation	302

		4.5.1.4 Quantifying risk
	4.5.2	Mixed noise
4.6	Beyond	d the scope of linear estimation: polyhedral estimate 3
	4.6.1	Motivation
	4.6.2	Generic polyhedral estimate
	4.6.3	Specifying sets \mathcal{H}_{δ} for basic observation schemes 3
		4.6.3.1 Sub-Gaussian case
		4.6.3.2 Discrete case
		4.6.3.3 Poisson case
	4.6.4	Efficient upper-bounding of $\mathfrak{R}[H]$ and Contrast Design, I 3
		4.6.4.1 Assumptions
		4.6.4.2 Simple observation
		4.6.4.3 Specifying contrasts
		4.6.4.4 Illustration: Diagonal case
	4.6.5	Efficient upper-bounding of $\mathfrak{R}[H]$ and Contrast Design, II. 3
		4.6.5.1 Outline
		4.6.5.2 Specifying cones H
		4.6.5.3 Specifying functions <i>M</i>
		4.6.5.4 Putting things together 3
		4.6.5.5 Compatibility: basic examples and calculus 3
		4 6 5 6 Spectratopic Sub-Gaussian case 3
	466	Numerical illustration 3
47	Recove	pring signals from nonlinear observations by Stochastic Opti-
1.1	mizati	on ?
	1 17 1	Problem's setting 3
	4.7.2	Assumptions 3
	4.1.2	Main observation 3
	4.7.5	Numerical illustration
	4.7.4	"Single observation" case 3
18	4.7.0 Appon	div: Calculus of Ellitopos/Spectratopos
4.0	Fyorcia	see for Locture 4
4.9		Linear Estimates vs. Maximum Likelihood
	4.9.1	Moscurement Design in Signal Recovery
	4.9.4	Around somidefinite relevation
	4.3.3	Around Dropositions 4.4 and 4.14
	4.9.4	A 0 4 1 Optimizing linear estimates on converting linear
		4.5.4.1 Optimizing inleaf estimates on convex nuns of unions
		4042 Recovering nonlinear vector valued functions
		4.0.4.3 Suboptimal linear estimation
		4.0.4.4 Probabilities of large deviations in linear estimation
		4.5.4.4 r robabilities of large deviations in linear estimation
		40.45 Lincon processory under size 1 dependent poist
	4.0.5	4.9.4.0 Linear recovery under signal-dependent noise 3
	4.9.5	Signal recovery in Discrete and Poisson observation schemes 3
	4.9.0	Numerical lower-bounding minimax risk
	4.9.7	Around S-Lemma 4
4.10	4.9.8	Estimation by stochastic optimization
4.10	Proofs	
	4.10.1	Preliminaries
		4.10.1.1 Technical lemma
		4.10.1.2 Noncommutative Khintchine Inequality 4

viii

	4.10.1.3 Anderson's Lemma	414
4.10.2	Proof of Proposition 4.6	415
4.10.3	Proof of Proposition 4.8	417
4.10.4	Proof of Lemma 4.17	418
4.10.5	Proofs of Propositions 4.5, 4.16, 4.19	422
	$4.10.5.1 Proof of Proposition 4.16 \dots \dots \dots \dots \dots \dots \dots$	422
	$4.10.5.2 Proof of Proposition 4.5 \dots \dots \dots \dots \dots \dots$	428
4.10.6	Proofs of Propositions 4.18, 4.19, and justification of Remark	
	4.20	436
	4.10.6.1 Proof of Proposition 4.18 \ldots	436
	$4.10.6.2 \text{Proof of Proposition } 4.19 \dots \dots \dots \dots \dots \dots \dots \dots \dots $	436
	4.10.6.3 Justification of Remark 4.20	438
4.10.7	Proof of (4.96)	438
4.10.8	Proof of Lemma 4.28	439
4.10.9	Justification of (4.141)	440
Bibliography		441

Index

 $\mathbf{451}$

StatOpt[•]LN[•]NS January 21, 2019 7x10

List of Figures

1.1	Top: true 256×256 image; bottom: sparse in the wavelet basis approximations of the image. Wavelet basis is orthonormal, and a natural way to quantify near-sparsity of a signal is to look at the fraction of total energy (sum of squares of wavelet coefficients) stored in the leading coefficients; these are the "energy data" presented on the figure.	9
1.2	Singe-pixel camera	10
1.3	Regular and penalized ℓ_1 recovery of nearly s-sparse signals. Red circles: true time series, blue crosses: recovered time series (to make the plots readable, one per eight consecutive terms in the time series is shown). Problem's sizes are $m = 256$ and $n = 2m = 512$, noise level is $\sigma = 0.01$, deviation from s-sparsity is $ x - x^s _1 = 1$, contrast prime is $(U_{n_1} - \sqrt{n_1/2} + 1)$. In problem a period resource $\lambda = 256$	
	trast pair is $(H = \sqrt{n/mA}, \ \cdot\ _{\infty})$. In penalized recovery, $\lambda = 2s$, parameter ρ in regular recovery is set to ErfInv $(0.005/n)$.	23
1.4	Erroneous ℓ_1 recovery of 25-sparse signal, no observation noise. Magenta: true signal, blue: ℓ_1 recovery. Top: frequency domain, bot-	
	tom: time domain	30
2.1	"Gaussian Separation" (Example 2.5): Optimal test deciding on whether the mean of Gaussian r.v. belongs to the dark red (H_1) or to the dark blue (H_2) domains. Dark and light red: acceptance domain for H_1 . Dark	
2.2	and light blue: acceptance domain for H_2	53
	$\mathcal{N}(\mu, I_2)$, each stating that μ belongs to the polygon of specific color.	100
2.3	Signal (top, magenta) and its candidate estimates (top,blue). Bot-	119
2.4	3-candidate hypotheses in probabilistic simplex Δ_3 :	$113 \\ 117$
2.5	Frames from a "movie"	160
3.1	Bisection via Hypothesis Testing	200
3.2	A circuit (9 nodes, 16 arcs). Red: arc of interest; Green: arcs with	014
33	Histograms of recovery errors in experiments data over 1000 simu-	214
0.0	lations per experiment.	243
4.1	True distribution of temperature $U_* = B(x)$ at time $t_0 = 0.01$ (left) along with its recovery \hat{U} via optimal linear estimate (center) and	
	the "naive" recovery U (right)	276

StatOpt[•]LN[•]NS January 21, 2019 7x10

xii

CONTENTS

4.2	Recovery errors for near-optimal linear estimate (red) and the poly-	
	hedral estimates yielded by Proposition 4.29 (PolyI, blue) and by	
	construction from Section 4.6.4 (PolyII, cyan).	332
4.3	Performance of SAA recovery, $m = 100, m = 1$	346

Preface

WHEN SPEAKING about links between Statistics and Optimization, what comes to mind first is the indispensable role played by optimization algorithms in the "computational toolbox" of Statistics (think about the numerical implementation of the fundamental Maximum Likelihood method). However, on a second thought, we should conclude that whatever high this role could be, the fact that it comes to our mind first primarily reflects the weaknesses of Optimization rather than its strengths; were optimization algorithms used in Statistics as efficient and as reliable as, say, Linear Algebra techniques used there, nobody would think about special links between Statistics and Optimization, same as nobody usually thinks about special links between Statistics and Linear Algebra. When computational, rather than methodological, issues are concerned, we start to think about links with Optimization, Linear Algebra, Numerical Analysis, etc., only when computational tools offered to us by these disciplines do not work well and need the attention of experts in these disciplines.

The goal of Lectures is to present another type of links between Optimization and Statistics, those which have nothing in common with algorithms and numbercrunching. What we are speaking about, are the situations where Optimization theory (theory, not algorithms!) seems to be of methodological value in Statistics, acting as the source of statistical inferences with provably optimal, or nearly so, performance. In this context, we focus on utilizing Convex Programming theory, mainly due to its power, but also due to the desire to end up with inference routines reducing to solving convex optimization problems and thus implementable in a computationally efficient fashion. Thus, while we do not mention computational issues explicitly, we do remember that at the end of the day we need a number, and in this respect, intrinsically computationally friendly convex optimization models are the first choice.

The three topics we intend to consider are:

- 1. Sparsity-oriented Compressive Sensing. Here the role of Convex Optimization theory, by itself by far not negligible (it allows, e.g., to derive from "first principles" the necessary and sufficient conditions for the validity of ℓ_1 recovery) is relatively less important than in two other topics. Nevertheless, we believe that Compressive Sensing, due to its popularity and the fact that now it is one of the major "customers" of advanced convex optimization algorithms, is worthy of being considered.
- 2. Pairwise and Multiple Hypothesis Testing, including sequential tests, estimation of linear functionals, and some rudimentary design of experiments. This is the topic where, as of now, the approaches based on Convex Optimization theory were most successful.
- 3. Recovery of signals from noisy observations of their linear images.

 xiv

PREFACE

The exposition does *not* require prior knowledge of Statistics and Optimization; as far as these disciplines are concerned, all necessary for us facts and concepts are incorporated into the text. The actual prerequisites are elementary Calculus, Probability, and Linear Algebra and (last but by far not least) general mathematical culture.

Anatoli Juditsky & Arkadi Nemirovski Date details

Acknowledgments

WE ARE greatly indebted to Prof. H. Edwin Romeijn who initiated creating the Ph.D. course "Topics in Data Science," part of which are the lectures to follow. The second author gratefully acknowledges support from NSF Grant CCF-1523768 *Statistical Inference via Convex Optimization*; this research project is the source of basically all novel results presented in Lectures 2 - 4. Needless to say, responsibility for all drawbacks of Lectures is ours.

StatOpt[•]LN[•]NS January 21, 2019 7x10

Notational conventions

VECTORS AND MATRICES.

By default, all vectors are column ones; to write them down, we use "Matlab notation:" $\begin{bmatrix} 1\\2\\3 \end{bmatrix}$ is written as [1;2;3]. More generally, for vectors/matrices A, B, ..., Z of the same "width" [A; B; C; ...; D] is the matrix obtained by writing B beneath of A, C beneath of B, and so on. For vectors/matrices A, B, C, ..., Z of the same "height," [A, B, C, ..., Z] denotes the matrix obtained by writing B to the right of A, C to the right of B, and so on. Examples: for what in the "normal" notation is written down as $A = \begin{bmatrix} 1 & 2\\ 3 & 4 \end{bmatrix}, B = \begin{bmatrix} 5 & 6 \end{bmatrix}, C = \begin{bmatrix} 7\\8 \end{bmatrix}$, we have

$$[A;B] = \begin{bmatrix} 1 & 2\\ 3 & 4\\ 5 & 6 \end{bmatrix} = [1,2;3,4;5,6], \ [A,C] = \begin{bmatrix} 1 & 2 & 7\\ 3 & 4 & 8 \end{bmatrix} = [1,2,7;3,4,8].$$

Blanks in matrices replace (blocks of) zero entries. For example,

$$\left[\begin{array}{rrrr}1&&\\2&&\\3&4&5\end{array}\right] = \left[\begin{array}{rrrr}1&0&0\\2&0&0\\3&4&5\end{array}\right]$$

 $Diag\{A_1, A_2, ..., A_k\}$ stands for block-diagonal matrix with diagonal blocks A_1 , $A_2, ..., A_k$. For example,

$$\operatorname{Diag}\{1,2,3\} = \begin{bmatrix} 1 & & \\ & 2 & \\ & & 3 \end{bmatrix}, \operatorname{Diag}\{[1,2];[3;4]\} = \begin{bmatrix} 1 & 2 & \\ & & 3 \\ & & 4 \end{bmatrix}.$$

For an $m \times n$ matrix A, dg(A) is the diagonal of A – vector of dimension min[m, n] with entries A_{ii} , $1 \le i \le \min[m, n]$.

STANDARD LINEAR SPACES

in our course are \mathbf{R}^n (the space of *n*-dimensional column vectors), $\mathbf{R}^{m \times n}$ (the space of $m \times n$ real matrices), and \mathbf{S}^n (the space of $n \times n$ real symmetric matrices). All these linear spaces are equipped with the standard inner product:

$$\langle A, B \rangle = \sum_{i,j} A_{ij} B_{ij} = \operatorname{Tr}(AB^T) = \operatorname{Tr}(BA^T) = \operatorname{Tr}(A^T B) = \operatorname{Tr}(B^T A);$$

ACKNOWLEDGMENTS

in the case when A = a and B = b are column vectors, this simplifies to $\langle a, b \rangle = a^T b = b^T a$, and when A, B are symmetric, there is no need to write B^T in $\text{Tr}(AB^T)$.

Usually we denote vectors by lowercase, and matrices – by uppercase letters; sometimes, however, lowercase letters are used also for matrices.

Given a linear mapping $\mathcal{A}(x) : E_x \to E_y$, where E_x , E_y are standard linear spaces, one can define the *conjugate* mapping $\mathcal{A}^*(y) : E_y \to E_x$ via the identity

$$\langle \mathcal{A}(x), y \rangle = \langle x, \mathcal{A}^*(y) \rangle \ \forall (x \in E_x, y \in E_y).$$

One always has $(\mathcal{A}^*)^* = \mathcal{A}$. When $E_x = \mathbf{R}^n$, $E_y = \mathbf{R}^m$ and $\mathcal{A}(x) = Ax$, one has $\mathcal{A}^*(y) = A^T y$; when $E_x = \mathbf{R}^n$, $E_y = \mathbf{S}^m$, so that $\mathcal{A}(x) = \sum_{i=1}^n x_i A_i$, $A_i \in \mathbf{S}^m$, we have

$$\mathcal{A}^*(Y) = [\operatorname{Tr}(A_1Y); ...; \operatorname{Tr}(A_nY)].$$

 \mathbf{Z}^n is the set of *n*-dimensional integer vectors.

NORMS.

For $1 \le p \le \infty$ and for a vector $x = [x_1; ...; x_n] \in \mathbf{R}^n$, $||x||_p$ is the standard *p*-norm of x:

$$||x||_{p} = \begin{cases} (\sum_{i=1}^{n} |x_{i}|^{p})^{1/p} & , 1 \le p < \infty \\ \max_{i} |x_{i}| = \lim_{p' \to \infty} ||x||_{p'} & , p = \infty \end{cases},$$

Notation for various norms of matrices is specified when used.

STANDARD CONES.

 \mathbf{R}_{+} is the nonnegative ray on the real axis, \mathbf{R}_{+}^{n} stands for the *n*-dimensional nonnegative orthant – the cone comprised of all entrywise nonnegative vectors from \mathbf{R}^{n} , \mathbf{S}_{+}^{n} stands for the positive semidefinite cone in \mathbf{S}^{n} – the cone comprised of all positive semidefinite matrices from \mathbf{S}^{n} .

MISCELLANEOUS.

• For matrices A, B, relation $A \leq B$, or, equivalently, $B \geq A$, means that A, B are symmetric matrices of the same size such that B - A is positive semidefinite; we write $A \geq 0$ to express the fact that A is a symmetric positive semidefinite matrix. Strict version $A \succ B$ ($\Leftrightarrow B \prec A$) of $A \succeq B$ means that A - B is positive definite (and, as above, A and B are symmetric matrices of the same size).

• Linear Matrix Inequality (LMI, a.k.a. *semidefinite constraint*) in variables x is the constraint on x stating that a symmetric matrix affinely depending on x is positive semidefinite. When $x \in \mathbf{R}^n$, LMI reads

$$a_0 + \sum_i x_i a_i \succeq 0 \qquad \qquad [a_i \in \mathbf{S}^m, 0 \le i \le n]$$

• $\mathcal{N}(\mu, \Theta)$ stands for the Gaussian distribution with mean μ and covariance matrix Θ .

- For a probability distribution P,
- $\xi \sim P$ means that ξ is a random variable with distribution P. Sometimes we

ACKNOWLEDGMENTS

express the same fact by writing $\xi \sim p(\cdot)$, where p is the density of P taken w.r.t. some reference measure (the latter always is fixed by the context);

• $\mathbf{E}_{\xi \sim P}\{f(\xi)\}$ is the expectation of $f(\xi)$, $\xi \sim P$; when P is clear from the context, this notation can be shortened to $\mathbf{E}_{\xi}\{f(\xi)\}$, or $\mathbf{E}_{P}\{f(\xi)\}$, or even $\mathbf{E}\{f(\xi)\}$. Similarly, $\operatorname{Prob}_{\xi \sim P}\{...\}$, $\operatorname{Prob}_{\xi}\{...\}$, $\operatorname{Prob}_{P}\{...\}$ denote the P-probability of the event specified inside the braces.

• O(1)'s stand for positive *absolute* constants – positive reals which we do not want or are too lazy to write down explicitly, like in $\sin(x) \leq O(1)|x|$.

• $\int_{\Omega} f(\xi) \Pi(d\xi)$ stands for the integral, taken w.r.t. measure Π over domain Ω , of function f.

About proofs

Lecture Notes are basically self-contained in terms of proofs of the statements to follow. Simple proofs usually are placed immediately after the corresponding statements; more technical proofs are transferred to dedicated sections titled "Proof of ...," and this is where a reader should look for "missing" proofs.

Lecture One

Sparse Recovery via ℓ_1 Minimization

In this lecture, we overview basic results of Compressed Sensing – a relatively new and extremely rapidly developing area in Signal Processing dealing with recovering signals (vectors x from some \mathbb{R}^n) from their noisy observations $Ax + \eta$ (A is a given $m \times n$ sensing matrix, η is observation noise) in the case when the number of observations m is much smaller than the signal's dimension n, but is essentially larger than the "true" dimension – the number of nonzero entries – in the signals. This setup leads to extremely deep, elegant and highly innovative theory and possesses quite significant applied potential. It should be added that along with the plain sparsity (small number of nonzero entries), Compressed Sensing deals with other types of "low-dimensional structure" hidden in high-dimensional signals, most notably, with the case of low rank matrix recovery, when signal is a matrix, and sparse signals are matrices with low ranks, and the case of block sparsity, where signal is a block vector, and sparsity means that only small number of blocks are nonzero. In our presentation, we do not consider these extensions of the simplest sparsity paradigm.

1.1 COMPRESSED SENSING: WHAT IT IS ABOUT?

1.1.1 Signal Recovery Problem

One of the basic problems in Signal Processing is the problem of recovering a signal $x \in \mathbf{R}^n$ from noisy observations

$$y = Ax + \eta \tag{1.1}$$

of the affine image of the signal under a given sensing mapping $x \mapsto Ax : \mathbf{R}^n \to \mathbf{R}^m$; in (1.1), η is the observation error. Matrix A in (1.1) is called sensing matrix.

Recovery problem of outlined types arise in many applications, including, but $by \ far$ not reducing to,

• communications, where x is the signal sent by transmitters, y is the signal recorded by receivers, A represents the communication channel (reflecting, e.g., dependencies of decays in signals' amplitude on the transmitter-receiver distances); η here typically is modeled as the standard (zero mean, unit covariance

matrix) m-dimensional Gaussian noise¹;

- *image reconstruction*, where the signal x is an image a 2D array in the usual photography, or a 3D array in Tomography, and y is data acquired by the imaging device. Here η in many cases (although not always) can again be modeled as the standard Gaussian noise;
- linear regression arising in a wide range of applications. In linear regression, one is given m pairs "input $a^i \in \mathbf{R}^n$ " to a "black box" output $y_i \in \mathbf{R}$ of the black box." Sometimes we have reasons to believe that the output is a corrupted by noise version of the "existing in the nature," but unobservable, ideal output" $y_i^* = x^T a^i$ which is just a linear function of the input (this is called "linear regression model," with inputs a^i called "regressors"). Our goal is to convert actual observations $(a^i, y_i), 1 \leq i \leq m$, into estimates of the unknown "true" vector of parameters x. Denoting by A the matrix with the rows $[a^i]^T$ and assembling individual observations y_i into a single observation $y = [y_1; ...; y_m] \in \mathbf{R}^m$, we arrive at the problem of recovering vector x from noisy observations of Ax. Here again the most popular model for η is the standard Gaussian noise.

1.1.2 Signal Recovery: parametric and non-parametric cases

Recovering signal x from observation y would be easy if there were no observation noise $(\eta = 0)$ and the rank of matrix A were equal to the dimension n of signals. In this case, which can take place only when $m \ge n$ ("more observations that unknown parameters"), and is typical in this range of sizes m, n, the desired x would be the unique solution to the system of linear equation, and to find x would be a simple problem of Linear Algebra. Aside of this trivial "enough observations, no noise" case, people over the years looked at the following two versions of the recovery problem:

Parametric case: $m \gg n$, η is nontrivial noise with zero mean, say, standard Gaussian one. This is the classical statistical setup considered in thousands of papers, with the emphasis on how to use the numerous observations we have at our disposal in order to suppress in the recovery, to the extent possible, the influence of observation noise.

Nonparametric case: $m \ll n$. Literally treated, this case seems to be senseless: when the number of observations is less that the number of unknown parameters, even in the no-noise case we arrive at the necessity to solve an undetermined (less equations than unknowns) system of linear equations. Linear Algebra says that if solvable, the system has infinitely many solutions; moreover, the solution set (an affine subspace of positive dimension) is unbounded, meaning that the solutions are in no sense close to each other. Typical way to make the case of $m \ll n$ meaningful is to add to the observations (1.1) some a priori information on the signal. In

$$Q^{-1}y = [Q^{-1}A]x + \xi$$

6

¹The "physical" noise usually indeed is Gaussian with zero mean, but its covariance matrix not necessarily is the unit matrix. Note, however, that a zero mean Gaussian noise η always can be represented as $Q\xi$ with standard Gaussian ξ ; assuming Q nonsingular (which indeed is so when the covariance matrix of η is positive definite), we can rewrite (1.1) equivalently as

and treat $Q^{-1}y$ and $Q^{-1}A$ as our new observation and new sensing matrix; new observation noise ξ is indeed standard. Thus, in the case of Gaussian zero mean observation noise, to assume the noise standard Gaussian is the same as to assume that its covariance matrix is known.

SPARSE RECOVERY VIA ℓ_1 MINIMIZATION

7

traditional Nonparametric Statistics this additional information is summarized in a given to us in advance bounded convex set $X \subset \mathbf{R}^n$ known to contain the true signal x. This set usually is such that every signal $x \in X$ can be approximated by a linear combination of s = 1, 2, ..., n of vectors from properly selected and known to us in advance orthonormal basis ("dictionary" in the slang of signal processing) within accuracy $\delta(s)$, where $\delta(s)$ is a known in advance function approaching 0 as $s \to \infty$. In this situation, with appropriate A (e.g., just the unit matrix, as in denoising problem), we can select somehow $s \ll m$ and try to recover x as if it were a vector from the linear span E_s of the first s vectors of the outlined basis. In the "ideal case" $x \in E_s$, recovering x in fact reduces to the case where the dimension of the signal is $s \ll m$ rather than $n \gg m$, and we arrive at the well-studied situation of recovering signal of low (as compared to the number of observations) dimension. In the "realistic case" of $x \delta(s)$ -close to E_s , deviation of x from E_s results in additional component in the recovery error ("bias"); a typical result of traditional Nonparametric Statistics quantifies the resulting error and minimizes it in s. Of course, this outline of traditional statistical approach to "nonparametric" (with $n \gg m$) recovery problems is extremely sketchy, but it captures the most important in our context fact: with the traditional approach to nonparametric signal recovery, one assumes that after representing the signals by vectors of their coefficients in properly selected orthonormal basis, the *n*-dimensional signal to be recovered can be well approximated by s-sparse (at most s nonzero entries) signal, with $s \ll n$, and this sparse approximation can be obtained by zeroing out all but the first s entries in the signal vector.

The just formulated assumption indeed takes place for signals obtained by discretization of *smooth* uni- and multivariate functions, and this class of signals for several decades was the main, if not the only, focus of Nonparametric Statistics.

To the best of our knowledge, developments in the traditional Nonparametric Statistics had nearly nothing to do with Convex Optimization.

Compressed Sensing. The situation changed dramatically around Year 2000 as a consequence of the breakthroughs due to D. Donoho, T. Tao, J. Romberg, E. Candes, J. Fuchs and several other researchers; as a result of these breakthroughs, an extremely popular and completely novel area of research, called *Compressed Sensing*, emerged.

In the Compressed Sensing (CS) setup of the Signal Recovery problem, same as in the traditional Nonparametric Statistics, is assumed that after passing to an appropriate basis, the signal to be recovered is s-sparse (has $\leq s$ nonzero entries), or is well approximated by s-sparse signal. The difference with the traditional approach is that now we assume *nothing* on the location of the nonzero entries. Thus, the a priori information on the signal x both in the traditional and in the CS settings is summarized in a set X known to contain the signal x we want to recover. The difference is, that in the traditional setting, X is a bounded convex and "nice" (well approximated by its low-dimensional cross-sections) set, while in CS this set is, computationally speaking, a "monster:" already in the simplest case of recovering *exactly s-sparse* signals, X is the union of all s-dimensional coordinate planes, which is a heavily combinatorial entity.

Note that in many applications we indeed can be sure that the true vector of parameters θ^* is sparse. Consider, e.g., the following story about signal detection. There are n locations where signal transmitters could be placed, and m locations with the receivers. The contribution of a signal of unit

LECTURE 1

magnitude originating in location j to the signal measured by receiver i is a known quantity A_{ij} , and signals originating in different locations merely sum up in the receivers; thus, if x is the n-dimensional vector with entries x_j representing the magnitudes of signals transmitted in locations j = 1, 2, ..., n, then the m-dimensional vector y of measurements of the m receivers is y = $Ax + \eta$, where η is the observation noise. Given y, we intend to recover x.

Now, if the receivers are hydrophones registering noises emitted by submarines in certain part of Atlantic, tentative positions of submarines being discretized with resolution 500 m, the dimension of the vector x (the number of points in the discretization grid) will be in the range of tens of thousands, if not tens of millions. At the same time, the total number of submarines (i.e., nonzero entries in x) can be safely upper-bounded by 50, if not by 20.

In order to see sparsity on our everyday life, look at the 256×256 image on the top of Figure 1.1. The image can be thought of as a $256^2 = 65536$ -dimensional vector comprised of pixels' intensities in gray scale, and there is no much sparsity in this vector. However, when representing the image in the *wavelet basis*, whatever it means, we get a "nearly sparse" vector of wavelet coefficients (this is true for typical "non-pathological" images). On the bottom of Figure 1.1 we see what happens when we zero out all but a percentage of the largest in magnitude wavelet coefficients and replace the true image by its sparse, in the wavelet basis, approximations.

Our visual illustration along with numerous similar examples show the "everyday presence" of sparsity and the possibility to utilize it when compressing signals. The difficulty, however, is that simple compression – compute the coefficients of the signal in an appropriate basis and then keep, say, 10% of the largest in magnitude coefficients – requires to start with digitalizing the signal – representing it as an array of all its coefficients in some orthonormal basis. These coefficients are inner products of the signal with vectors of the basis; for a "physical" signal, like speech or image, these inner products are computed by analogous devices, with subsequent discretization of the results. After the measurements are discretized, processing the signal (denoising, compression, storing, etc., etc.) can be fully computerized. The major potential (to some extent, already actual) advantage of Compressed Sensing is in the possibility to reduce the "analogous effort" in the outlined process: instead of computing analogously n linear forms of n-dimensional signal x (its coefficients in a basis), we use analogous device to compute $m \ll n$ other linear forms of the signal and then use signal's sparsity in a known to us basis in order to recover the signal reasonably well from these m observations.

In our "picture illustration" this technology would work (in fact, works - it is called "single pixel camera," see Figure 1.2) as follows: in reality, the digital 256×256 image on the top of Figure 1.1 was obtained by analogous device – a digital camera which gets on input analogous signal (light of varying along the field of view intensity caught by camera's lenses) and discretizes lights's intensity in every pixel to get the digitalized image. We then can compute the wavelet coefficients of the digitalized image, compress its representation by keeping, say, just 10% of leading coefficients, etc., etc., but "the damage is already done" – we have already spent our analogous resources to get the entire digitalized image. The technology utilizing Compressed Sensing would work as follows: instead of measuring and discretizing light intensity in every one of the 65,536 pixels, we compute analogously the integral, taken over the field of view, of the product of light intensity and an analogously generated "mask," and do it for, say, 20,000 different masks, thus ob-



Figure 1.1: Top: true 256×256 image; bottom: sparse in the wavelet basis approximations of the image. Wavelet basis is orthonormal, and a natural way to quantify near-sparsity of a signal is to look at the fraction of total energy (sum of squares of wavelet coefficients) stored in the leading coefficients; these are the "energy data" presented on the figure.



Figure 1.2: Singe-pixel camera

taining measurements of 20,000 linear forms of our 65,536-dimensional signal. Next we utilize, via the Compressed Sensing machinery, signal's sparsity in the wavelet basis in order to recover the signal from these 20,000 measurements. With this approach, we reduce the "analogous component" of signal processing effort, at the price of increasing the "computerized component" of the effort (instead of ready-touse digitalized image directly given by 65,536 analogous measurements, we need to recover the image by applying computationally not so trivial decoding algorithms to our 20,000 "indirect" measurements). When taking pictures by your camera or ipad, the game is not worth the candle – analogous component of taking usual pictures is cheap enough, and decreasing it at the price of nontrivial decoding of the digitalized measurements would be counter-productive. There are, however, important applications where the advantages stemming from reduced "analogous effort" overweight significantly the drawbacks caused by the necessity to use nontrivial computerized decoding.

1.1.3 Compressed Sensing via ℓ_1 minimization: Motivation

1.1.3.1 Preliminaries

In principle there is nothing surprising in the fact that under reasonable assumption on $m \times n$ sensing matrix A we may hope to recover from noisy observations of Axan s-sparse, with $s \ll m$, signal x. Indeed, assume for the sake of simplicity that there are no observation errors, and let $\operatorname{Col}_j[A]$ be j-th column in A. If we knew the locations $j_1 < j_2 < \ldots < j_s$ of the nonzero entries in x, identifying x could be reduced to solving system of linear equations $\sum_{\ell=1}^{s} x_{i_\ell} \operatorname{Col}_{j_\ell}[A] = y$ with m equations and $s \ll m$ unknowns; assuming every s columns in A linearly independent (a quite unrestrictive assumption on a matrix with $m \ge s$ rows), the

10

SPARSE RECOVERY VIA ℓ_1 MINIMIZATION

solution to the above system is unique, and is exactly the signal we are looking for. Of course, the assumption that we know the locations of nonzeros in x makes the recovery problem completely trivial. However, it suggests the following course of actions: given noiseless observation y = Ax of an s-sparse signal x, let us solve the combinatorial optimization problem

$$\min\{\|z\|_0 : Az = y\},\tag{1.2}$$

where $||z||_0$ is the number of nonzero entries in z. Clearly, the problem has a solution with the value of the objective at most s. Moreover, it is immediately seen (verify it!) that if every 2s columns in A are linearly independent (which again is a very unrestrictive assumption on the matrix A provided that $m \ge 2s$), then the true signal x is the unique optimal solution to (1.2).

What was said so far can be extended to the case of noisy observations and "nearly *s*-sparse" signals *x*. For example, assuming that the observation error is "uncertainbut-bounded," specifically some known norm $\|\cdot\|$ of this error does not exceed a given $\epsilon > 0$, and that the true signal is *exactly s*-sparse (think how to relax this to "near *s*-sparsity"), we could solve the combinatorial optimization problem

$$\min\{\|z\|_0: \|Az - y\| \le \epsilon\}.$$
(1.3)

Assuming that every $m \times 2s$ submatrix \overline{A} of A is not just with linearly independent columns (i.e., with trivial kernel), but is reasonably well conditioned:

$$\|\bar{A}w\| \ge C^{-1} \|w\|_2$$

for all (2s)-dimensional vectors w, with some constant C, it is immediately seen that the true signal x underlying observation and the optimal solution \hat{x} of (1.3) are close to each other within accuracy of order of ϵ : $||x - \hat{x}||_2 \leq 2C\epsilon$; it is easily seen that the resulting error bound is basically as good as it could be.

We see that the difficulties with recovering sparse signals stem not from the lack of information, they are of purely computational nature: (1.2) is a disastrously difficult combinatorial problem, and the only known way to process it is by "brute force" search through all guesses on where the nonzeros in x are located – by inspecting first the only option that there are no nonzeros in x at all, then by inspecting n options that there is only one nonzero, for every one of n locations of this nonzero, then n(n-1)/2 options that there are exactly two nonzeros, etc., etc. until the current option will result in a solvable system of linear equations Az = y in variables z with entries restricted to vanish outside the locations prescribed by the option under consideration. Running time of this "brute force" search, beyond the range of small values of s and n (by far too small to be of any applied interest) is by many orders of magnitude larger than what we can afford to ourselves in reality².

A partial remedy is as follows. Well, if we do not know how to minimize under linear constraints, as in (1.2), the "bad" objective $||z||_0$, let us "approximate" this

²When s = 5 and n = 100, a sharp upper bound on the number of linear systems we should process before termination in the "brute force" algorithm is $\approx 7.53e7$ — much, but perhaps doable. When n = 200 and s = 20, the number of systems to be processed jumps to $\approx 1.61e27$, which is by many orders of magnitude beyond our "computational grasp"; we would be unable to carry out that many computations even if the fate of the mankind were at stake. And from the perspective of Compressed Sensing, n = 200 still is a completely toy size, by 3-4 orders of magnitude less than we would like to handle.

LECTURE 1

objective with one which we do know how to minimize. The true objective is separable: $||z|| = \sum_{i=1}^{n} \xi(z_j)$, where $\xi(s)$ is the function on the axis equal to 0 at the origin and equal to 1 otherwise. As a matter of fact, the separable functions which we do know how to minimize under linear constraints are sums of *convex* functions of $z_1, ..., z_n$. The most natural candidate to the role of *convex* approximation of $\xi(s)$ is |s|; with this approximation, (1.2) converts into the ℓ_1 minimization problem

$$\min_{z} \left\{ \|z\|_{1} := \sum_{i=1}^{n} |z_{j}| : Az = y \right\},$$
(1.4)

and (1.3) becomes the convex optimization problem

$$\min_{z} \left\{ \|z\|_{1} := \sum_{i=1}^{n} |z_{j}| : \|Az - y\| \le \epsilon \right\}.$$
 (1.5)

Both problems are efficiently solvable, which is nice; the question, however, is how relevant these problems are in our context – whether it is true that they do recover the "true" s-sparse signals in the noiseless case, or "nearly recover" these signals when the observation error is small. Since we want to be able to handle whatever s-sparse signals, the validity of ℓ_1 recovery – it ability to recover well every s-sparse signal – depends solely on the sensing matrix A. Our current goal is to understand what are "good" in this respect sensing matrices.

1.2 VALIDITY OF SPARSE SIGNAL RECOVERY VIA ℓ_1 MINIMIZATION

What follows is based on the standard basic results of Compressed Sensing theory originating from [35, 34, 36, 37, 32, 33, 48, 46, 47, 60, 61] and augmented by the results of $[85]^3$.

1.2.1 Validity of ℓ_1 minimization in the noiseless case

The minimal requirement on sensing matrix A which makes ℓ_1 minimization valid is to guarantee the correct recovery of *exactly* s-sparse signals in the *noiseless* case, and we start with investigating this property.

1.2.1.1 Notational convention

From now on, for a vector $x \in \mathbf{R}^n$

• $I_x = \{j : x_j \neq 0\}$ stands for the *support* of x; we also set

$$I_x^+ = \{j : x_j > 0\}, \ I_x^- = \{j : x_j < 0\} \qquad [\Rightarrow I_x = I_x^+ \cup I_x^-]$$

• for a subset I of the index set $\{1, ..., n\}$, x_I stands for the vector obtained from

 $^{^{3}}$ in fact, in the latter source, an extension of the sparsity, the so called block sparsity, is considered; in what follows, we restrict the results of [85] to the case of plain sparsity.

SPARSE RECOVERY VIA ℓ_1 MINIMIZATION

x by zeroing out entries with indexes not in I, and I^{o} for the complement of I:

$$I^{o} = \{ i \in \{1, ..., n\} : i \notin I \};\$$

- for $s \leq n$, x^s stands for the vector obtained from x by zeroing our all but the s largest in magnitude entries⁴ Note that x^s is the best s-sparse approximation of x in any one of the ℓ_p norms, $1 \leq p \leq \infty$;
- for $s \leq n$ and $p \in [1, \infty]$, we set

$$||x||_{s,p} = ||x^s||_p;$$

note that $\|\cdot\|_{s,p}$ is a norm (why?).

1.2.1.2 s-Goodness

Definition of *s*-goodness. Let us say that an $m \times n$ sensing matrix *A* is *s*-good, if whenever the true signal *x* underlying *noiseless* observations is *s*-sparse, this signal will be recovered *exactly* by ℓ_1 minimization. In other words, *A* is *s*-good, if whenever in *y* in (1.4) is of the form y = Ax with *s*-sparse *x*, *x* is the unique optimal solution to (1.4).

Nullspace property. There is a simply-looking necessary and sufficient condition for a sensing matrix A to be s-good – the nullspace property. After this property is guessed, it is easy to see that it indeed is necessary and sufficient for s-goodness; we, however, prefer to derive this condition from the "first principles," which can be easily done via Convex Optimization; thus, in the case in question, same as in many other cases, there is no necessity to be smart to arrive at the truth via "lucky guess," it suffices to be knowledgeable and use the standard tools.

Let us start with necessary and sufficient condition for A to be such that whenever x is s-sparse, x is an optimal solution (perhaps, not the unique one) of the optimization problem

$$\min\{\|z\|_1 : Az = Ax\}, \tag{(*)}$$

let us call the latter property of A weak s-goodness. Our first observation is as follows:

Proposition 1.1. A is weakly s-good if and only if the following condition holds true: whenever I is a subset of $\{1, ..., n\}$ of cardinality $\leq s$, we have

$$\forall w \in \operatorname{Ker} A : \|w_I\|_1 \le \|w_{\bar{I}}\|_1 \tag{1.6}$$

Proof is immediate. In one direction: Assume A is weakly s-good, and let us verify (1.6). Let I be an s-element subset of $\{1, ..., n\}$, and x be s-sparse vector with support I. Since A is weakly s-good, x is an optimal solution to (*). Rewriting the

⁴note that in general x^s is not uniquely defined by x and s, since the s-th largest among the magnitudes of entries in x can be achieved at several entries. In our context, it does not matter how the ties of this type are resolved; for the sake of definiteness, we can assume that when ordering the entries in x according to their magnitudes, from the largest to the smallest, entries of equal magnitude are ordered in the order of their indexes.

LECTURE 1

latter problem in the form of LP, that is, as

$$\min_{z,t} \{ \sum_{j} t_j : t_j + z_j \ge 0, t_j - z_j \ge 0, Az = Ax \},\$$

and invoking LP optimality conditions, the necessary and sufficient condition for z = x to be the z-component of an optimal solution is the existence of λ_j^+ , λ_j^- , $\mu \in \mathbf{R}^m$ (Lagrange multipliers for the constraints $t_j - z_j \ge 0$, $t_j + z_j \ge 0$, and Az = Ax, respectively) such that

From (c, d), we have $\lambda_j^+ = 1, \lambda_j^- = 0$ for $j \in I_x^+$ and $\lambda_j^+ = 0, \lambda_j^- = 1$ for $j \in I_x^-$. From (a) and nonnegativity of λ_j^{\pm} it follows that for $j \notin I_x$ we should have $-1 \leq \lambda_j^+ - \lambda_j^- \leq 1$. With this in mind, the above optimality conditions admit eliminating λ 's and reduce to the following conclusion:

(!) x is an optimal solution to (*) if and only if there exists vector $\mu \in \mathbf{R}^m$ such that j-th entry of $A^T \mu$ is -1, if $x_j > 0$, +1, if $x_j < 0$, and a real from [-1,1], if $x_j = 0$.

Now let $w \in \text{Ker } A$ be a vector with the same signs of entries $w_i, i \in I$, as these of the entries in x. Then

$$\begin{array}{l} 0 = \mu^T A w = [A^T \mu]^T w = \sum_j [A^T \mu]_j w_j \\ \Rightarrow \sum_{j \in I_x} |w_j| = \sum_{j \in I_x} [A^T \mu]_j w_j = -\sum_{j \notin I_x} [A^T \mu]_j w_j \le \sum_{j \notin I_x} |w_j| \end{array}$$

(we have used the fact that $[A^T \mu]_j = \operatorname{sign} x_j = \operatorname{sign} w_j$ for $j \in I_x$ and $|[A^T \mu]_j| \leq 1$ for all j). Since I can be an arbitrary *s*-element subset of $\{1, ..., n\}$ and the pattern of signs of an *s*-sparse vector x supported on I can be arbitrary, (1.6) holds true.

Now let us assume that (1.6) holds true, and let us prove that A is weakly ssparse. Assume the opposite; then for some s-sparse x, x is not an optimal solution to (*), meaning that system (1.7) of linear constraints in variables λ^{\pm} , μ has no solution. Applying Theorem on Alternative ([11, Theorem 1.2.1]), we can assign the constraints (a) - (f) in (1.7) with respective vectors of weights $w_a, ..., w_e, w_f$, with the weights w_e, w_f of inequality constraints (e), (f) being nonnegative, such that multiplying the constraints by the weights and summing up the results, we get as a consequence of (1.7) a contradictory inequality – one with no solutions at all. This contradictory consequence of (1.7) is the linear inequality in variables λ^{\pm}, μ :

$$[w_a + w_b + G_+ w_c + w_e]^T \lambda^+ + [w_a - w_b + G_- w_d + w_f]^T \lambda^- + w_b^T A^T \mu \ge \sum_j (w_a)_j, \quad (**)$$

where G_+ , G_- are diagonal matrices with *j*-th diagonal entry equal to $|x_j| - x_j$ (G_+) and $|x_j| + x_j$ (G_-). Thus, we can find $w_a, ..., w_f$ with nonnegative w_e and

14

SPARSE RECOVERY VIA ℓ_1 MINIMIZATION

 w_f such that

$$w_a + w_b + G_+ w_c + w_e = 0, \ w_a - w_b + G_- w_d + w_f = 0, \ Aw_b = 0, \ \sum_j (w_a)_j > 0.$$

or, equivalently, there exist w_a, w_b, w_c, w_d such that

$$\begin{array}{ll} (p) & w_{a}+w_{b}+\underbrace{G_{+}w_{c}}_{g}\leq 0,\\ (q) & w_{a}-w_{b}+\underbrace{G_{-}w_{d}}_{h}\leq 0,\\ (r) & Aw_{b}=0,\\ (s) & \sum_{j}(w_{a})_{j}>0. \end{array}$$

Now note that when $j \in I_x^+$, we have $g_j = 0$ and thus (p) says that $|[w_b]_j| \ge [w_a]_j$, and when $j \in I_x^-$, we have $h_j = 0$ and thus (q) says that $|[w_b]_j| \ge [w_a]_j$. And when $j \notin I_x := I_x^+ \cup I_x^-$, (p) and (q) say that $[w_a]_j \le -|[w_b]_j|$. With this in mind, (s)implies that $-\sum_{j \notin I_x} |[w_b]_j| + \sum_{j \in I_x} |[w_b]_j| \ge \sum_j [w_a]_j > 0$. Thus, assuming that A is not weakly s-good, we have found a set I_x of indexes of cardinality $\le s$ and a vector $w_b \in \text{Ker } A$ (see (r)) such that $\sum_{j \in I_x} |[w_b]_j| > \sum_{j \notin I_x} |[w_b]_j|$, contradicting the condition (1.6). \Box

1.2.1.3 Nullspace property

We have established necessary and sufficient condition for A to be weakly s-good; it states that $||w_I||_1$ should be $\leq ||w_{I^o}||_1$ for all $w \in \text{Ker } A$ and all I of cardinality s. It may happen that this inequality holds true as equality, for some nonzero $w \in \text{Ker } A$:

$$\exists (w \in \text{Ker } A \setminus \{0\}, I, \text{Card}(I) \leq s) : ||w_I||_1 = ||w_{I^o}||_1.$$

In this case matrix A clearly is not s-good, since the s-sparse signal $x = w_I$ is not the unique optimal solution to (*) – the vector $-w_{I^o}$ is a different feasible solution to the same problem and with the same value of the objective. We conclude that for A to be s-good, a necessary condition is for the inequality in (1.6) to be strict whenever $w \in \text{Ker } A$ is nonzero. By the standard compactness arguments, the latter condition means the existence of $\gamma \in (0, 1)$ such that

$$\forall (w \in \operatorname{Ker} A, I, \operatorname{Card}(I) \leq s) : \|w_I\|_1 \leq \gamma \|w_{I^o}\|_1,$$

or, which is the same, existence of $\kappa \in (0, 1/2)$ such that

$$\forall (w \in \operatorname{Ker} A, I, \operatorname{Card}(I) \le s) : \|w_I\|_1 \le \kappa \|w\|_1.$$

Finally, the supremum of $||w_I||_1$ over I of cardinality s is what we have defined the norm $||w||_{1,s}$ (the sum of s largest magnitudes of entries) of w, so that the condition we are processing finally can be formulated as

$$\exists \kappa \in (0, 1/2) : \|w\|_{1,s} \le \kappa \|w\|_1 \ \forall w \in \operatorname{Ker} A.$$
(1.8)

LECTURE 1

The resulting *nullspace condition* in fact is necessary *and sufficient* for A to be s-good:

Proposition 1.2. Condition (1.8) is necessary and sufficient for A to be s-good.

Proof. We have already seen that the nullspace condition is necessary for sgoodness. To verify sufficiency, let A satisfy nullspace condition, and let us prove that A is s-good. Indeed, let x be an s-sparse vector. By Proposition 1.1, A is weakly s-good, so that x is an optimal solution to (*), and the only thing we need to prove is that if y is another optimal solution to (*), then y = x. Assuming y optimal for (*), let I be the support of x. Setting w = y - x, we have

$$Aw = 0 \& \underbrace{\|y_I - x\|_1}_{\|w_I\|_1} \le \kappa \left[\|y_I - x\|_1 + \|y_{I^o} - x_{I^o}\|_1\right] = \kappa \left[\|w_I\|_1 + \|y_{I^o}\|_1\right],$$

whence

$$(1-\kappa)\|w_I\|_1 \le \kappa \|y_{I^o}\|_1 = \kappa (\|y\|_1 - \|y_I\|_1).$$

Since $||w_I||_1 = ||y_I - x||_1 \ge ||x||_1 - ||y_I||_1$, we arrive at

$$(1-\kappa)(\|x\|_1 - \|y_I\|_1) \le \kappa(\|y\|_1 - \|y_I\|_1),$$

which, due to $||x||_1 = ||y||_1$ (since x and y are optimal solutions of (*)) and $\kappa < 1/2$, boils down to

$$[\|y\|_1 =]\|x\|_1 \le \|y_I\|_1,$$

implying, due to $||x||_1 = ||y||_1$, that $y_I = y$, that is, y is supported on the support I of x. In other words, w = y - x is supported on s-element set, and since Aw = 0, nullspace property implies that y = x.

1.2.2 Imperfect ℓ_1 minimization

We have found a necessary and sufficient condition for ℓ_1 minimization to recover exactly s-sparse signals in the noiseless case. "In reality," both these assumptions typically are violated: instead of s-sparse signals, we should speak about "nearly s-sparse ones," quantifying the deviation from sparsity by the distance from the signal x underlying observations to its best s-sparse approximation x^s . Similarly, we should allow for nonzero observation noise. With noisy observations and/or imperfect sparsity, we cannot hope to recover signal exactly; all we may hope for, is to recover it with some error depending on the level of observation noise and "deviation from s-sparsity" and tending to zero as these level and deviation tend to 0. We are about to quantify the Nullspace property to allow for instructive "error analysis."

1.2.2.1 Contrast matrices and quantifications of Nullspace property

By itself, Nullspace property says something about the signals from the kernel of the sensing matrix. We can reformulate it equivalently to say something important about *all* signals. Namely, observe that given sparsity s and $\kappa \in (0, 1/2)$, the Nullspace property

$$\|w\|_{s,1} \le \kappa \|w\|_1 \ \forall w \in \operatorname{Ker} A \tag{1.9}$$

16

SPARSE RECOVERY VIA ℓ_1 MINIMIZATION

is satisfied if and only if for a properly selected constant C one has

$$\|w\|_{s,1} \le C \|Aw\|_2 + \kappa \|w\|_1 \,\forall w. \tag{1.10}$$

Indeed, (1.10) clearly implies (1.9); to get the inverse implication, note that for every h orthogonal to Ker A it holds

$$\|Ah\|_2 \ge \sigma \|h\|_2,$$

where $\sigma > 0$ is the minimal positive singular value of A. Now, given $w \in \mathbf{R}^n$, we can decompose w into the sum of $\bar{w} \in \operatorname{Ker} A$ and $h \in (\operatorname{Ker} A)^{\perp}$, so that

$$\begin{aligned} \|w\|_{s,1} &\leq \|\bar{w}\|_{s,1} + \|h\|_{s,1} \leq \kappa \|\bar{w}\|_{1} + \sqrt{s}\|h\|_{s,2} \leq \kappa \|w\|_{1} + \|h\|_{1} + \sqrt{s}\|h\|_{2} \\ &\leq \kappa \|w\|_{1} + [\kappa\sqrt{n} + \sqrt{s}]\|h\|_{2} \leq \underbrace{\sigma^{-1}[\kappa\sqrt{n} + \sqrt{s}]}_{C} \underbrace{\|Ah\|_{2}}_{=\|Aw\|_{2}} + \kappa \|w\|_{1}, \end{aligned}$$

as required in (1.10).

Condition $\mathbf{Q}_1(s,\kappa)$. For our purposes, it is convenient to present the condition (1.10) in the following flexible form:

$$\|w\|_{s,1} \le s \|H^T A w\| + \kappa \|w\|_1, \tag{1.11}$$

where H is an $m \times N$ contrast matrix and $\|\cdot\|$ is some norm on \mathbb{R}^N . Whenever a pair $(H, \|\cdot\|)$, called contrast pair, satisfies (1.11), we say that $(H, \|\cdot\|)$ satisfies condition $\mathbb{Q}_1(s, \kappa)$. From what we have seen, If A possesses Nullspace property with some sparsity level s and some $\kappa \in (0, 1/2)$, then there are many ways to select pairs $(H, \|\cdot\|)$ satisfying $\mathbb{Q}_1(s, \kappa)$, e.g., to take $H = CI_m$ with appropriately large C and $\|\cdot\| = \|\cdot\|_2$.

Conditions $\mathbf{Q}_q(s,\kappa)$. As we shall see in a while, it makes sense to embed the condition $\mathbf{Q}_1(s,\kappa)$ into a parametric family of conditions $\mathbf{Q}_q(s,\kappa)$, where the parameter q runs through $[1,\infty]$. Specifically,

Given $m \times n$ sensing matrix A, sparsity level $s \leq n$ and $\kappa \in (0, 1/2)$, we say that $m \times N$ matrix H and a norm $\|\cdot\|$ on \mathbb{R}^N satisfy condition $\mathbb{Q}_q(s, \kappa)$, if

$$\|w\|_{s,q} \le s^{\frac{1}{q}} \|H^T A w\| + \kappa s^{\frac{1}{q}-1} \|w\|_1 \,\forall w \in \mathbf{R}^n.$$
(1.12)

Let us make two immediate observations on relations between the conditions:

A. When a pair $(H, \|\cdot\|)$ satisfies condition $\mathbf{Q}_q(s, \kappa)$, the pair satisfies also all conditions $\mathbf{Q}_{q'}(s, \kappa)$ with $1 \leq q' \leq q$.

Indeed in the situation in question for $1 \le q' \le q$ it holds

$$\begin{aligned} \|w\|_{s,q'} &\leq s^{\frac{1}{q'} - \frac{1}{q}} \|w\|_{q,s} \leq s^{\frac{1}{q'} - \frac{1}{q}} \left[s^{\frac{1}{q}} \|H^T A w\| + \kappa s^{\frac{1}{q} - 1} \|w\|_1 \right] \\ &= s^{\frac{1}{q'}} \|H^T A w\| + \kappa s^{\frac{1}{q'} - 1} \|w\|_1, \end{aligned}$$

where the first inequality is the standard inequality between ℓ_p -norms of the s-dimensional vector w^s .

B. When a pair $(H, \|\cdot\|)$ satisfies condition $\mathbf{Q}_q(s, \kappa)$ and $1 \leq s' \leq s$, the pair $((s/s')^{\frac{1}{q}}H, \|\cdot\|)$ satisfies the condition $\mathbf{Q}_q(s', \kappa)$.

LECTURE 1

Indeed, in the situation in question we clearly have for $1 \le s' \le s$:

$$||w||_{s',q} \le ||w||_{s,q} \le (s')^{\frac{1}{q}} || \left[(s/s')^{\frac{1}{q}} H \right] Aw|| + \kappa \underbrace{s^{\frac{1}{q}-1}}_{\le (s')^{\frac{1}{q}-1}} ||w||_1.$$

1.2.3 Regular ℓ_1 recovery

Given observation scheme (1.1) with $m \times n$ sensing matrix A, we define the *regular* ℓ_1 recovery of x via observation y as

$$\widehat{x}_{reg}(y) \in \underset{u}{\operatorname{Argmin}} \left\{ \|u\|_1 : \|H^T(Au - y)\| \le \rho \right\},$$
(1.13)

where the *contrast matrix* $H \in \mathbf{R}^{m \times N}$, the norm $\|\cdot\|$ on \mathbf{R}^N and $\rho > 0$ are parameters of the construction.

The role of **Q**-conditions we have introduced is clear from the following

Theorem 1.3. Let s be a positive integer, $q \in [1, \infty]$ and $\kappa \in (0, 1/2)$. Assume that the pair $(H, \|\cdot\|)$ satisfies the condition $\mathbf{Q}_q(s, \kappa)$ associated with A, and let

$$\Xi_{\rho} = \{\eta : \|H^T \eta\| \le \rho\}.$$
(1.14)

Then for all $x \in \mathbf{R}^n$ and $\eta \in \Xi_{\rho}$ one has

$$\|\widehat{x}_{reg}(Ax+\eta) - x\|_p \le \frac{4(2s)^{\frac{1}{p}}}{1-2\kappa} \left[\rho + \frac{\|x - x^s\|_1}{2s}\right], \ 1 \le p \le q.$$
(1.15)

The above result can be slightly strengthened by replacing the assumption that $(H, \|\cdot\|)$ satisfies $\mathbf{Q}_q(s, \kappa)$ with some $\kappa < 1/2$, with a weaker, by observation \mathbf{A} from Section 1.2.2.1, assumption that $(H, \|\cdot\|)$ satisfies $\mathbf{Q}_1(s, \varkappa)$ with $\varkappa < 1/2$ and satisfies $\mathbf{Q}_q(s, \kappa)$ with some (perhaps large) κ :

Theorem 1.4. Given A, integer s > 0 and $q \in [1, \infty]$, assume that $(H, \|\cdot\|)$ satisfies the condition $\mathbf{Q}_1(s, \varkappa)$ with $\varkappa < 1/2$ and the condition $\mathbf{Q}_q(s, \kappa)$ with some $\kappa \geq \varkappa$, and let Ξ_{ρ} be given by (1.14). Then for all $x \in \mathbf{R}^n$ and $\eta \in \Xi_{\rho}$ it holds:

$$\|\widehat{x}_{reg}(Ax+\eta) - x\|_{p} \le \frac{4(2s)^{\frac{1}{p}} [1+\kappa-\varkappa]^{\frac{q(p-1)}{p(q-1)}}}{1-2\varkappa} \left[\rho + \frac{\|x-x^{s}\|_{1}}{2s}\right], \ 1 \le p \le q.$$
(1.16)

Before commenting on the above results, let us present their alternative versions.

1.2.4 Penalized ℓ_1 recovery

Penalized ℓ_1 recovery of signal x from its observation (1.1) is

$$\widehat{x}_{pen}(y) \in \underset{u}{\operatorname{Argmin}} \left\{ \|u\|_1 + \lambda \|H^T(Au - y)\| \right\},$$
(1.17)

where $H \in \mathbf{R}^{m \times N}$, a norm $\|\cdot\|$ on \mathbf{R}^N and a positive real λ are parameters of the construction.

Theorem 1.5. Given A, positive integer s, and $q \in [1, \infty]$, assume that $(H, \|\cdot\|)$ satisfies the conditions $\mathbf{Q}_q(s, \kappa)$ and $\mathbf{Q}_1(s, \varkappa)$ with $\varkappa < 1/2$ and $\kappa \ge \varkappa$. Then

18

SPARSE RECOVERY VIA ℓ_1 MINIMIZATION

(i) Let
$$\lambda \geq 2s$$
. Then for all $x \in \mathbf{R}^n$, $y \in \mathbf{R}^m$ it holds:

$$\widehat{x}_{pen}(y) - x \|_p \leq \frac{4\lambda^{\frac{1}{p}}}{1-2\varkappa} \left[1 + \frac{\kappa\lambda}{2s} - \varkappa\right]^{\frac{q(p-1)}{p(q-1)}} \left[\|H^T(Ax-y)\| + \frac{\|x-x^s\|_1}{2s} \right], \ 1 \leq p \leq q.$$
(1.18)

In particular, with $\lambda = 2s$ we have:

$$\|\widehat{x}_{pen}(y) - x\|_{p} \leq \frac{4(2s)^{\frac{1}{p}}}{1-2\varkappa} \left[1 + \kappa - \varkappa\right]^{\frac{q(p-1)}{p(q-1)}} \left[\|H^{T}(Ax - y)\| + \frac{\|x - x^{s}\|_{1}}{2s}\right].$$
(1.19)

(ii) Let $\rho \geq 0$, and let Ξ_{ρ} be given by (1.14). Then for all $x \in \mathbf{R}^n$ and all $\eta \in \Xi_{\rho}$ one has:

$$\begin{split} \lambda &\geq 2s \quad \Rightarrow \\ & \|\widehat{x}_{pen}(Ax+\eta) - x\|_{p} \leq \frac{4\lambda^{\frac{1}{p}}}{1-2\varkappa} \left[1 + \frac{\kappa\lambda}{2s} - \varkappa\right]^{\frac{q(p-1)}{p(q-1)}} \left[\rho + \frac{\|x-x^{s}\|_{1}}{2s}\right], \ 1 \leq p \leq q; \\ \lambda &= 2s \quad \Rightarrow \\ & \|\widehat{x}_{pen}(Ax+\eta) - x\|_{p} \leq \frac{4(2s)^{\frac{1}{p}}}{1-2\varkappa} \left[1 + \kappa - \varkappa\right]^{\frac{q(p-1)}{p(q-1)}} \left[\rho + \frac{\|x-x^{s}\|_{1}}{2s}\right], \ 1 \leq p \leq q. \end{split}$$
(1.20)

1.2.5 Discussion

Some remarks are in order.

A. Qualitatively speaking, Theorems 1.3, 1.4, 1.5 say the same: under **Q**conditions, the regular, resp., penalized recoveries are capable to reproduce the true signal *exactly* when there is no observation noise and the signal is *s*-sparse; in the presence of observation error η and imperfect sparsity, the signal is recovered within the error which can be upper-bounded by the sum of two terms, one proportional to the magnitude of observation noise and one proportional to the deviation $\|x-x^s\|_1$ of the signal from *s*-sparse ones. In the penalized recovery, the observation error is measured in the scale given by the contrast matrix and the norm $\|\cdot\|$ - as $\|H^T\eta\|$, and in the regular one – by an a priori upper bound ρ on $\|H^T\eta\|$ — when $\rho \geq \|H^T\eta\|$, η belongs to Ξ_{ρ} and thus the bounds (1.15), (1.16) are applicable to the actual observation error η . Clearly, in qualitative terms error bound of this type is the best we may hope for. Now let us look at the quantitative aspect. Assume that in the regular recovery we use $\rho \approx \|H^T\eta\|$, and in the penalized one use $\lambda = 2s$. In this case, error bounds (1.15), (1.16), (1.20), up to factors *C* depending solely on \varkappa and κ , are the same, specifically,

$$\|\widehat{x} - x\|_p \le Cs^{1/p}[\|H^T\eta\| + \|x - x^s\|_1/s], \ 1 \le p \le q.$$
(!)

Is this error bound bad or good? The answer depends on many factors, including on how well we select H and $\|\cdot\|$. To get a kind of orientation, consider the trivial case of *direct* observations, where matrix A is square and, moreover, is proportional to the unit matrix: $A = \alpha I$; assume in addition that x is exactly *s*-sparse. In this case, the simplest way to ensure condition $\mathbf{Q}_q(s,\kappa)$, even with $\kappa = 0$, is to take $\|\cdot\| = \|\cdot\|_{s,q}$ and $H = s^{-1/q} \alpha^{-1} I$, so that (!) becomes

$$\|\widehat{x} - x\|_p \le C\alpha^{-1} s^{1/p - 1/q} \|\eta\|_{s,q}, \ 1 \le p \le q.$$
(!!)

As far as the dependence of the bound on the magnitude $\|\eta\|_{s,q}$ of the observation noise is concerned, this dependence is as good as it can be – even if we knew in advance the positions of the *s* largest in magnitude entries of *x*, we would be unable to recover *x* is *q*-norm with error $\leq \alpha^{-1} \|\eta\|_{s,q}$ (why?); in addition, with the equal to each other *s* largest magnitudes of entries in η , the $\|\cdot\|_p$ -norm of the recovery error clearly cannot be guaranteed to be less than $\alpha^{-1} \|\eta\|_{s,p} = \alpha^{-1} s^{1/p-1/q} \|\eta\|_{s,q}$. Thus, at least for *s*-sparse signals *x*, our error bound is, basically, the best one can get already in the "ideal" case of direct observations.

B. Given that $(H, \|\cdot\|)$ obeys $\mathbf{Q}_1(s, \varkappa)$ with some $\varkappa < 1/2$, the larger is q such that the pair $(H, \|\cdot\|)$ obeys the condition $\mathbf{Q}_q(s, \kappa)$ with a given $\kappa \ge \varkappa$ (κ can be $\ge 1/2$) and s, the larger is the range $p \le q$ of values of p where the error bounds (1.16), (1.20) are applicable. This is in full accordance with the fact that if a pair $(H, \|\cdot\|)$ obeys condition $\mathbf{Q}_q(s, \kappa)$, it obeys also all conditions $\mathbf{Q}_{q'}(s, \kappa)$ with $1 \le q' \le q$ (item **A** in Section 1.2.2.1).

C. Flexibility offered by contrast matrix H and norm $\|\cdot\|$ allows to adjust, to some extent, the recovery to the "geometry of observation errors." For example, when η is "uncertain but bounded," say, all we know is that $\|\eta\|_2 \leq \delta$ with some given δ , all what matters (on the top of the requirement for $(H, \|\cdot\|)$ to obey **Q**-conditions) is how large could be $\|H^T\eta\|$ when $\|\eta\|_2 \leq \delta$. In particular, when $\|\cdot\| = \|\cdot\|_2$, the error bound "is governed" by the spectral norm of H; consequently, if we have a technique allowing to design H such that $(H, \|\cdot\|_2)$ obeys **Q**-condition(s) with given parameters, it makes sense to look for design with as small spectral norm of H as possible. In contrast to this, in the most interesting for applications case of Gaussian noise:

$$y = Ax + \eta, \ \eta \sim \mathcal{N}(0, \sigma^2 I_m) \tag{1.21}$$

looking at the spectral norm of H, with $\|\cdot\|_2$ in the role of $\|\cdot\|$, is counter-productive, since a typical realization of η is of Euclidean norm of order of $\sqrt{m\sigma}$ and thus is quite large when m is large. In this case to quantify "the magnitude" of $H^T\eta$ by the product of the spectral norm of H and the Euclidean norm of η is completely misleading – in typical cases, this product will grow rapidly with the number of observations m, completely ignoring the fact that η is random with zero mean⁵. What is much better suited for the case of Gaussian noise, is $\|\cdot\|_{\infty}$ norm in the role of $\|\cdot\|$ and the norm "the maximum of $\|\cdot\|_2$ -norms of the columns in H," let it be denoted by $\|H\|_{1,2}$, of H. Indeed, with $\eta \sim \mathcal{N}(0, \sigma^2 I_m)$, the entries in $H^T\eta$ are Gaussian with zero mean and variance bounded by $\sigma^2 \|H\|_{1,2}^2$, so that $\|H^T\eta\|_{\infty}$ is the maximum of magnitudes of N zero mean Gaussian random variables with standard deviations bounded by $\sigma \|H\|_{1,2}^2$. As a result,

$$\operatorname{Prob}\{\|H^{T}\eta\|_{\infty} \ge \rho\} \le N\operatorname{Erf}(\rho/\sigma)\|H\|_{1,2} \le N\operatorname{e}^{\frac{-\rho^{2}}{2\sigma^{2}}}\|H\|_{1,2}, \qquad (1.22)$$

where

$$\operatorname{Erf}(s) = \frac{1}{\sqrt{2\pi}} \int_{s}^{\infty} e^{-t^{2}/2} dt$$

is the error function. It follows that the typical values of $||H^T\eta||_{\infty}$, $\eta \sim \mathcal{N}(0, \sigma^2 I_m)$ are of order of at most $\sigma \sqrt{\ln(N)} ||H||_{1,2}$; typically, N = O(m), so that with σ and

⁵the simplest way to see the difference is to look at a particular entry $h^T \eta$ in $H^T \eta$. Operating with spectral norms, we upper-bound this entry by $\|h\|_2 \|\eta\|_2$, and the second factor for $\eta \sim \mathcal{N}(0, \sigma^2 I_m)$ is typically as large as $\sigma \sqrt{m}$, in sharp contrast to the fact that typical values of $h^T \eta$ are of order of σ , completely independently of what m is!
$||H||_{1,2}$ given, typical values $||H^T\eta||_{\infty}$ are nearly independent of m. The bottom line is that ℓ_1 minimization is capable to handle large-scale Gaussian observation noise incomparably better than "uncertain-but-bounded" observation noise of similar magnitude (measured in Euclidean norm).

D. As far as comparison of regular and penalized ℓ_1 recoveries with the same pair $(H, \|\cdot\|)$ is concerned, the situation is as follows. Assume for the sake of simplicity that $(H, \|\cdot\|)$ satisfies $\mathbf{Q}_q(s, \kappa)$ with some s and some $\kappa < 1/2$, and let the observation error be random. Given $\epsilon \in (0, 1)$, let

$$\rho_{\epsilon}[H, \|\cdot\|] = \min\left\{\rho : \operatorname{Prob}\left\{\eta : \|H^{T}\eta\| \le \rho\right\} \ge 1 - \epsilon\right\}; \quad (1.23)$$

this is nothing but the smallest ρ such that

$$\operatorname{Prob}\{\eta \in \Xi_{\rho}\} \ge 1 - \epsilon \tag{1.24}$$

(see (1.14)) and thus – the smallest ρ for which the error bound (1.15) for the regular ℓ_1 recovery holds true with probability $1 - \epsilon$ (or at least the smallest ρ for which the latter claim is supported by Theorem 1.3). With $\rho = \rho_{\epsilon}[H, \|\cdot\|]$, the regular ℓ_1 recovery guarantees (and that is the best guarantee one can extract from Theorem 1.3) that

(#) For some set Ξ , $\operatorname{Prob}\{\eta \in \Xi\} \geq 1 - \epsilon$, of "good" realizations of $\eta \sim \mathcal{N}(0, \sigma^2 I_m)$, one has

$$\|\widehat{x}(Ax+\eta) - x\|_p \le \frac{4(2s)^{\frac{1}{p}}}{1-2\kappa} \left[\rho_{\epsilon}[H, \|\cdot\|] + \frac{\|x-x^s\|_1}{2s}\right], \ 1 \le p \le q, \quad (1.25)$$

whenever $x \in \mathbf{R}^n$ and $\eta \in \Xi_{\rho}$.

The error bound (1.19) (where we set $\varkappa = \kappa$) says that (#) holds true for the penalized ℓ_1 recovery with $\lambda = 2s$. The latter observation suggests that the penalized ℓ_1 recovery associated with $(H, \|\cdot\|)$ and $\lambda = 2s$ is better than its regular counterpart, the reason being twofold. First, in order to ensure (#) with the regular recovery, the "built in" parameter ρ of this recovery should be set to $\rho_{\epsilon}[H, \|\cdot\|]$, and the latter quantity not always is easy to identify. In contrast to this, the construction of penalized ℓ_1 recovery is completely independent of a priori assumptions on the structure of observation errors, while automatically ensuring (#) for the error model we use. Second, and more importantly, for the penalized recovery the bound (1.25) is no more than the "worst, with confidence $1 - \epsilon$, case," while the typical values of the quantity $||H^T\eta||$ which indeed participates in the error bound (1.18) are essentially smaller than $\rho_{\epsilon}[H, \|\cdot\|]$. Numerical experience fully supports the above suggestion: the difference in observed performance of the two routines in question, although not dramatic, is definitely in favor of the penalized recovery. The only potential disadvantage of the latter routine is that the penalty parameter λ should be tuned to the level s of sparsity we aim at, while the regular recovery is free of any guess of this type. Of course, the "tuning" is rather loose – all we need (and experiments show that we indeed need this) is the relation $\lambda \geq 2s$, so that a rough upper bound on s will do; note, however, that bound (1.18) deteriorates as λ grows.

22

LECTURE 1

Finally, we remark that when H is $m \times N$ and $\eta \sim \mathcal{N}(0, \sigma^2 I_m)$, we have

$$\rho_{\epsilon}[H, \|\cdot\|_{\infty}] \leq \operatorname{ErfInv}(\epsilon/N) \|H\|_{1,2} \leq \sqrt{2\ln(N/\epsilon)} \|H\|_{1,2}$$

(see 1.22)); here $\operatorname{ErfInv}(\delta)$ is the inverse error function:

$$\operatorname{Erf}(\operatorname{ErfInv}(\delta)) = \delta, \ 0 < \delta < 1.$$

How it works. Here we present a small numerical illustration. We observe in Gaussian noise m = n/2 randomly selected terms in *n*-element "time series" $z = (z_1, ..., z_n)$ and want to recover this series under the assumption that the series is "nearly *s*-sparse in frequency domain," that is, that

$$z = Fx$$
 with $||x - x^s||_1 \le \delta$,

where F is the matrix of $n \times n$ Inverse Discrete Cosine Transform, x^s is the vector obtained from x by zeroing out all but s largest in magnitude entries, and δ upperbounds the distance from x to s-sparse signals. Denoting by A the $m \times n$ submatrix of F corresponding to the time instants t where z_t is observed, our observation scheme becomes

$$y = Ax + \sigma\xi,$$

where ξ is the standard Gaussian noise. After the signal in frequency domain, that is, x, is recovered by ℓ_1 minimization, let the recovery be \hat{x} , we recover the signal in the time domain as $\hat{z} = F\hat{x}$. On Figure 1.3, we present four test signals, of different (near) sparsity, along with their regular and penalized ℓ_1 recoveries. The data on Figure 1.3 clearly show how the quality of ℓ_1 recovery deteriorates as the number sof "essential nonzeros" of the signal in the frequency domain grows. It is seen also that the penalized recovery meaningfully outperforms the regular one in the range of sparsities up to 64.

1.3 VERIFIABILITY AND TRACTABILITY ISSUES

Good news on ℓ_1 recovery stated in Theorems 1.3, 1.4, 1.5 are "conditional" – we assume that we are smart enough to point out a pair $(H, \|\cdot\|)$ satisfying condition $\mathbf{Q}_1(s, \varkappa)$ with $\varkappa < 1/2$ (and condition $\mathbf{Q}_q(s, \kappa)$ with a "moderate" \varkappa^6). The related issues are twofold:

- 1. First, we do not know in which range of s, m, n these conditions, or even the weaker than $\mathbf{Q}_1(s, \varkappa), \varkappa < 1/2$, Nullspace property can be satisfied; and without the Nullspace property, ℓ_1 minimization becomes useless, at least when we want to guarantee its validity whatever be s-sparse signal we want to recover;
- 2. Second, it is unclear how to verify whether a given sensing matrix A satisfies the Nullspace property for a given s, or a given pair $(H, \|\cdot\|)$ satisfies the condition

 $^{{}^{6}\}mathbf{Q}_{q}(s,\kappa)$ always is satisfied with "large" κ , namely, $\kappa = s$, but this large value of κ is of no interest: the associated bounds on *p*-norms of recovery error are straightforward consequences of the bounds on $\|\cdot\|_{1}$ -norm of this error yielded by the condition $\mathbf{Q}_{1}(s,\varkappa)$.



Figure 1.3: Regular and penalized ℓ_1 recovery of nearly *s*-sparse signals. Red circles: true time series, blue crosses: recovered time series (to make the plots readable, one per eight consecutive terms in the time series is shown). Problem's sizes are m = 256 and n = 2m = 512, noise level is $\sigma = 0.01$, deviation from *s*-sparsity is $||x - x^s||_1 = 1$, contrast pair is $(H = \sqrt{n/m}A, || \cdot ||_{\infty})$. In penalized recovery, $\lambda = 2s$, parameter ρ in regular recovery is set to $\operatorname{ErfInv}(0.005/n)$.

24

LECTURE 1

$\mathbf{Q}_q(s,\kappa)$ with given parameters.

What is known on these crucial issues, can be outlined as follows.

1. It is known that for given m, n with $m \ll n$ (say, $m/n \leq 1/2$), there exist $m \times n$ sensing matrices which are s-good for the values of s "nearly as large as m", specifically, for $s \leq O(1) \frac{m}{\ln(n/m)}$ ⁷. Moreover, there are natural families of matrices where this level of goodness "is a rule." E.g., when drawing an $m \times n$ matrix at random from the Gaussian or the ± 1 distributions (i.e., filling the matrix with independent realizations of a random variable which is either a standard (zero mean, unit variance) Gaussian one, or takes values ± 1 with probabilities 0.5), the result will be s-good, for the outlined value of s, with probability approaching 1 as m and n grow. All this remains true when instead of speaking about matrices A for which it is easy to point out a pair $(H, \|\cdot\|)$ satisfying the condition $\mathbf{Q}_2(s, \varkappa)$ with, say, $\varkappa = 1/4$.

The above results can be considered as a good news. A bad news is, that we do *not* know how to check efficiently, given an *s* and a sensing matrix *A*, that the matrix is *s*-good, same as we do not know how to check that *A* admits good (i.e., satisfying $\mathbf{Q}_1(s, \varkappa)$ with $\varkappa < 1/2$) pairs $(H, \|\cdot\|)$. Even worse: we do not know an efficient recipe allowing to build, given *m*, an $m \times 2m$ matrix A^m which is provably *s*-good for *s* larger than $O(1)\sqrt{m}$, which is a much smaller "level of goodness" then the one promised by theory for randomly generated matrices⁸. The "common life" analogy of this pitiful situation would be as follows: you know that 90% of bricks in your wall are made of gold, and at the same time, you do not know how to tell a golden brick from a usual one.⁹

2. There exist verifiable sufficient conditions for s-goodness of a sensing matrix,

⁷From now on, O(1)'s denote positive *absolute constants* – appropriately chosen numbers like 0.5, or 1, or perhaps 100,000. We could, in principle, replace all O(1)'s by specific numbers; following the standard mathematical practice, we do not do it, partly from laziness, partly because the particular values of these numbers in our context are irrelevant.

⁸Note that the naive algorithm "generate $m \times 2m$ matrices at random until an *s*-good, with *s* promised by the theory, matrix is generated" is *not* an efficient recipe, since we do not know how to check *s*-goodness efficiently.

⁹This phenomenon is met in many other situations. E.g., in 1938 Claude Shannon (1916-2001), "the father of Information Theory," made (in his M.Sc. Thesis!) a fundamental discovery as follows. Consider a Boolean function of n Boolean variables (i.e., both the function and the variables take values 0 and 1 only); as it is easily seen there are 2^{2^n} function of this type, and every one of them can be computed by a dedicated circuit comprised of "switches" implementing just 3 basic operations AND, OR and NOT (like computing a polynomial can be carried out on a circuit with nodes implementing just two basic operation: addition of reals and their multiplication). The discovery of Shannon was that every Boolean function of n variables can be computed on a circuit with no more than $Cn^{-1}2^n$ switches, where C is an appropriate absolute constant. Moreover, Shannon proved that "nearly all" Boolean functions of n variables require circuits with at least $cn^{-1}2^n$ switches, c being another absolute constant; "nearly all" in this context means that the fraction of "easy to compute" functions (i.e., those computable by circuits with less than $cn^{-1}2^n$ switches) among all Boolean functions of n variables goes to 0 as n goes to ∞ . Now, computing Boolean functions by circuits comprised of switches was an important technical task already in 1938; its role in our today life can hardly be overestimated — the outlined computation is nothing but what is going on in a computer. Given this observation, it is not surprising that the Shannon discovery of 1938 was the subject of countless refinements, extensions, modifications, etc., etc. What is still missing, is a *single individual example* of a "difficult to compute" Boolean function: as a matter of fact, all multivariate Boolean functions $f(x_1, ..., x_n)$ people managed to describe explicitly are computable by circuits with just *linear* in n number of switches!

same as verifiable sufficient conditions for a pair $(H, \|\cdot\|)$ to satisfy condition $\mathbf{Q}_q(s,\kappa)$. A bad news that when $m \ll n$, these verifiable sufficient conditions can be satisfied only when $s \leq O(1)\sqrt{m}$ – once again, in a much more narrow range of values of s than the one where typical randomly selected sensing matrices are s-good. In fact, $s = O(\sqrt{m})$ is the best known so far sparsity level for which we know individual s-good $m \times n$ sensing matrices with $m \leq n/2$.

1.3.1 Restricted Isometry Property and *s*-goodness of random matrices

There are several sufficient conditions for s-goodness, equally difficult to verify, but provably satisfied for typical random sensing matrices. The best known of them is the *Restricted Isometry Property* (RIP) defined as follows:

Definition 1.6. Let k be an integer and $\delta \in (0,1)$. We say that an $m \times n$ sensing matrix A possesses the Restricted Isometry Property with parameters δ and k, $\operatorname{RIP}(\delta, k)$, if for every k-sparse $x \in \mathbf{R}^n$ one has

$$(1-\delta)\|x\|_{2}^{2} \le \|Ax\|_{2}^{2} \le (1+\delta)\|x\|_{2}^{2}.$$
(1.26)

It turns out that for natural ensembles of random $m \times n$ matrices, a typical matrix from the ensemble satisfies $\operatorname{RIP}(\delta, k)$ with small δ and k "nearly as large as m," and that $\operatorname{RIP}(\frac{1}{6}, 2s)$ implies Nullspace condition, and more. The simplest versions of the corresponding results are as follows.

Proposition 1.7. Given $\delta \in (0, \frac{1}{5}]$, with properly selected positive $c = c(\delta)$, $d = d(\delta)$, $f = f(\delta)$ for all $m \leq n$ and all positive integers k such that

$$k \le \frac{m}{c\ln(n/m) + d} \tag{1.27}$$

the probability for a random $m \times n$ matrix A with independent $\mathcal{N}(0, \frac{1}{m})$ entries to satisfy $\operatorname{RIP}(\delta, k)$ is at least $1 - \exp\{-fm\}$.

Proposition 1.8. Let $A \in \mathbb{R}^{m \times n}$ satisfy $\operatorname{RIP}(\delta, 2s)$ for some $\delta < 1/3$ and positive integer s. Then

(i) The pair $\left(H = \frac{s^{-1/2}}{\sqrt{1-\delta}}I_m, \|\cdot\|_2\right)$ satisfies the condition $\mathbf{Q}_2\left(s, \frac{\delta}{1-\delta}\right)$ associated with A;

(ii) The pair $(H = \frac{1}{1-\delta}A, \|\cdot\|_{\infty})$ satisfies the condition $\mathbf{Q}_2\left(s, \frac{\delta}{1-\delta}\right)$ associated with A.

1.3.2 Verifiable sufficient conditions for $\mathbf{Q}_q(s,\kappa)$

When speaking about verifiable sufficient conditions for a pair $(H, \|\cdot\|)$ to satisfy $\mathbf{Q}_q(s, \kappa)$, it is convenient to restrict ourselves with the case when H, same as A, is an $m \times n$ matrix, and $\|\cdot\| = \|\cdot\|_{\infty}$.

Proposition 1.9. Let A be an $m \times n$ sensing matrix, and $s \leq n$ be a sparsity level.

Given $m \times n$ matrix H and $q \in [1, \infty]$, let us set

$$\nu_{s,q}[H] = \max_{j \le n} \|\text{Col}_j[I - H^T A]\|_{s,q},$$
(1.28)

where $\operatorname{Col}_{j}[C]$ is *j*-th column of matrix C. Then

$$||w||_{s,q} \le s^{1/q} ||H^T A w||_{\infty} + \nu_{s,q} [H] ||w||_1 \ \forall w \in \mathbf{R}^n,$$
(1.29)

implying that the pair $(H, \|\cdot\|_{\infty})$ satisfies the condition $\mathbf{Q}_q(s, s^{1-\frac{1}{q}}\nu_{s,q}[H])$.

Proof is immediate. Setting $V = I - H^T A$, we have

$$\begin{aligned} \|w\|_{s,q} &= \|[H^T A + V]w\|_{s,q} \le \|H^T A w\|_{s,q} + \|Vw\|_{s,q} \\ &\le s^{1/q} \|H^T A w\|_{\infty} + \sum_{j} |w_j| \|\operatorname{Col}_j[V]\|_{s,q} \le s^{1/q} \|H^T A\|_{\infty} + \nu_{s,q}[H] \|w\|_1. \end{aligned}$$

Observe that the function $\nu_{s,q}[H]$ is an efficiently computable convex function of H, so that the set

$$\mathcal{H}_{s,q}^{\kappa} = \{ H \in \mathbf{R}^{m \times n} : \nu_{s,q}[H] \le s^{\frac{1}{q} - 1} \kappa \}$$
(1.30)

is a computationally tractable convex set. When this set is nonempty for some $\kappa < 1/2$, every point H in this set is a contrast matrix such that $(H, \|\cdot\|_{\infty})$ satisfies the condition $\mathbf{Q}_a(s,\kappa)$, that is, we can find contrast matrices making ℓ_1 minimization valid. Moreover, we can *design* contrast matrix, e.g., by minimizing over $\mathcal{H}_{s,q}^{\kappa}$ the function $\|H\|_{1,2}$, thus optimizing the sensitivity of the corresponding ℓ_1 recoveries to Gaussian observation noise, see items \mathbf{C} , \mathbf{D} in Section 1.2.5.

Explanation. The sufficient condition for s-goodness of A stated in Proposition 1.9 looks as coming out of thin air; in fact it is a particular case of a simple and general construction as follows. Let f(x) be a real-valued convex function on \mathbb{R}^n , and $X \subset \mathbb{R}^n$ be a nonempty bounded polytope represented as

$$X = \{x \in \text{Conv}\{g_1, ..., g_N\} : Ax = 0\},\$$

where $\operatorname{Conv}\{g_1, ..., g_N\} = \{\sum_i \lambda_i g_i : \lambda \ge 0, \sum_i \lambda_i = 1\}$ is the convex hull of $g_1, ..., g_N$. Our goal is to upper-bound the maximum $\operatorname{Opt} = \max_{x \in X} f(x)$; this is a meaningful problem, since precise maximizing a convex function over a polytope typically is a computationally intractable task. Let us act as follows: clearly, for a whatever matrix H of the same sizes as A we have $\max_{x \in X} f(x) = \max_{x \in X} f([I - H^T A]x)$, since on X we have $[I - H^A]x = x$. As a result,

Opt :=
$$\max_{x \in X} f(x) = \max_{x \in X} f([I - H^T A]x)$$

 $\leq \max_{x \in \operatorname{Conv}\{g_1, \dots, g_N\}} f([I - H^T A]x)$
 $= \max_{j \leq N} f([I - H^T A]g_j).$

We get a parametric, the parameter being H, upper bound on Opt, namely, the bound $\max_{j \leq N} f([I - H^T A]g_j)$. This parametric bound is convex in the parameter H, and thus is well suited for minimization over this parameter.

The result of Proposition 1.9 is inspired by this construction as applied to the

nullspace property: given $m \times n$ sensing matrix A and setting

$$X = \{x \in \mathbf{R}^n : \|x\|_1 \le 1, Ax = 0\} = \{x \in \text{Conv}\{\pm e_1, ..., \pm e_n\} : Ax = 0\}$$

 $(e_i \text{ are the basic orths in } \mathbf{R}^n)$, A is s-good if and only if

$$\mathrm{Opt}_s := \max_{x \in X} \{ f(x) := \|x\|_{s,1} \} < 1/2;$$

A verifiable sufficient condition for this yielded by the above construction is the existence of $m \times n$ matrix H such that

$$\max_{j \le n} \max[f([I_n - H^T A]e_j), f(-[I_n - H^T A]e_j)] < 1/2,$$

or, which is the same,

$$\max_{i} \|\operatorname{Col}_{j}[I_{n} - H^{T}A]\|_{s,1} < 1/2,$$

bringing to our attention the matrix $I - H^T A$ with varying H and the idea to express sufficient conditions for s-goodness and related properties in terms of this matrix.

1.3.3 Tractability of $\mathbf{Q}_{\infty}(s,\kappa)$

As we have already mentioned, the conditions $\mathbf{Q}_q(s,\kappa)$ are intractable, in the sense that we do not know how to verify whether a given pair $(H, \|\cdot\|)$ satisfies the condition. Surprisingly, this is *not* the case with the strongest of these conditions, the one with $q = \infty$. Specifically,

Proposition 1.10. Let A be an $m \times n$ sensing matrix, s be a sparsity level, and $\kappa \geq 0$. Then whenever a pair $(\bar{H}, \|\cdot\|)$ satisfies the condition $\mathbf{Q}_{\infty}(s, \kappa)$, there exists an $m \times n$ matrix H such that

$$\|\operatorname{Col}_{j}[I_{n} - H^{T}A]\|_{s,\infty} = \|\operatorname{Col}_{j}[I_{n} - H^{T}A]\|_{\infty} \le s^{-1}\kappa, \ 1 \le j \le n,$$

(so that $(H, \|\cdot\|_{\infty})$ satisfies $\mathbf{Q}_{\infty}(s, \kappa)$ by Proposition 1.9) and, in addition,

$$\|H^T\eta\|_{\infty} \le \|\bar{H}^T\eta\| \ \forall \eta \in \mathbf{R}^m.$$
(1.31)

In addition, $m \times n$ contrast matrix H such that the pair $(H, \|\cdot\|_{\infty})$ satisfies the condition $\mathbf{Q}_{\infty}(s, \kappa)$ with as small κ as possible can be found as follows: we consider n LP programs

$$Opt_{i} = \min_{\nu, h} \left\{ \nu : \|A^{T}h - e^{i}\|_{\infty} \le \nu \right\}, \qquad (\#_{i})$$

where e^i is *i*-th basic orth in \mathbb{R}^n , find optimal solutions Opt_i , h_i to these problems, and make h_i , i = 1, ..., n, the columns of H; the corresponding value of κ is

$$\kappa_* = s \max_i \operatorname{Opt}_i.$$

Besides this, there exists a transparent alternative description of the quantities Opt_i

28

LECTURE 1

(and thus – of κ_*); specifically,

$$\operatorname{Opt}_{i} = \max\left\{x_{i} : \|x\|_{1} \le 1, Ax = 0\right\}.$$
(1.32)

Looking at (1.31) and error bounds in Theorems 1.3, 1.4, 1.5, Proposition 1.10 says that

As far as the condition $\mathbf{Q}_{\infty}(s,\kappa)$ is concerned, we lose nothing when restricting ourselves with pairs $(H \in \mathbf{R}^{m \times n}, \|\cdot\|_{\infty})$ and contrast matrices Hsatisfying the condition

$$|[I_n - H^T A]_{ij}| \le s^{-1}\kappa \tag{1.33}$$

implying that $(H, \|\cdot\|_{\infty})$ satisfies $\mathbf{Q}_{\infty}(s, \kappa)$.

A good news is that (1.33) is an explicit convex constraint on H (in fact, even on H and κ), so that we can solve the *design problems*, where we want to optimize a convex function of H under the requirement that $(H, \|\cdot\|_{\infty})$ satisfies the condition $\mathbf{Q}_{\infty}(s, \kappa)$ (and, perhaps, additional convex constraints on H and κ).

1.3.3.1 Mutual Incoherence

The simplest (and up to some point in time, the only) verifiable sufficient condition for s-goodness of a sensing matrix A is expressed in terms of mutual incoherence of A defined as

$$\mu(A) = \max_{i \neq j} \frac{|\operatorname{Col}_{i}^{T}[A]\operatorname{Col}_{j}[A]|}{||\operatorname{Col}_{i}[A]||_{2}^{2}};$$
(1.34)

this quantity is well defined whenever A has no zero columns (otherwise A is not even 1-good). Note that when A is normalized to have all columns of equal $\|\cdot\|_2$ -lengths¹⁰, $\mu(A)$ is small when the directions of distinct columns in A are nearly orthogonal. The standard related result is that

Whenever A and a positive integer s are such that $\frac{2\mu(A)}{1+\mu(A)} < \frac{1}{s}$, A is s-good.

It is immediately seen that the latter condition is weaker than what we can get with the aid of (1.33):

Proposition 1.11. Let A be an $m \times n$ matrix, and let the columns in $m \times n$ matrix H be given by

$$\operatorname{Col}_{j}(H) = \frac{1}{(1+\mu(A)) \|\operatorname{Col}_{j}(A)\|_{2}^{2}} \operatorname{Col}_{j}(A), \ 1 \le j \le n.$$

Then

$$|[I_m - H^T A]_{ij}| \le \frac{\mu(A)}{1 + \mu(A)} \,\forall i, j.$$
(1.35)

In particular, when $\frac{2\mu(A)}{1+\mu(A)} < \frac{1}{s}$, A is s-good.

¹⁰as far as ℓ_1 minimization is concerned, this normalization is non-restrictive: we always can enforce it by diagonal scaling of the signal underlying observations (1.1), and ℓ_1 minimization in scaled variables is the same as weighted ℓ_1 minimization in original variables.

Proof. With *H* as above, the diagonal entries in $I - H^T A$ are equal to $1 - \frac{1}{1+\mu(A)} = \frac{\mu(A)}{1+\mu(A)}$, while by definition of mutual incoherence the magnitudes of the off-diagonal entries in $I - H^T A$ are $\leq \frac{\mu(A)}{1+\mu(A)}$ as well, implying (1.35). The "in particular" claim is given by (1.35) combined with Proposition 1.9.

1.3.3.2 From RIP to conditions $\mathbf{Q}_q(\cdot,\kappa)$

It turns out that when A is $\operatorname{RIP}(\delta, k)$ and $q \ge 2$, it is easy to point out pairs $(H, \|\cdot\|)$ satisfying $\mathbf{Q}_q(t, \kappa)$ with a desired $\kappa > 0$ and properly selected t:

Proposition 1.12. Let A be an $m \times n$ sensing matrix satisfying RIP $(\delta, 2s)$ with some s and some $\delta \in (0, 1)$, and let $q \in [2, \infty]$ and $\kappa > 0$ be given. Then

(i) Whenever a positive integer t satisfies

$$t \le \min\left[\left[\frac{\kappa\sqrt{1-\delta}}{\delta\sqrt{1+\delta}}\right]^{\frac{q}{q-1}}, s^{\frac{q-2}{q-1}}\right]s^{\frac{q}{2(q-1)}},$$
(1.36)

the pair $(H = \frac{t^{-1/2}}{\sqrt{1-\delta}}I_m, \|\cdot\|_2)$ satisfies $\mathbf{Q}_q(t, \kappa);$ (ii) Whenever a positive integer t satisfies

$$t \le \min\left[\left[\frac{\kappa(1-\delta)}{\delta}\right]^{\frac{q}{q-1}}, s^{\frac{q-2}{2q-2}}\right]s^{\frac{q}{2(q-1)}},$$
(1.37)

the pair $(H = \frac{s^{\frac{1}{2}}t^{-\frac{1}{q}}}{1-\delta}A, \|\cdot\|_{\infty})$ satisfies $\mathbf{Q}_q(t,\kappa)$.

The most important consequence of Proposition 1.12 deals with the case of $q = \infty$ and states that when s-goodness of a sensing matrix A can be ensured by difficult to verify condition RIP $(\delta, 2s)$ with, say, $\delta = 0.2$, the somehow worse level of sparsity, $t = O(1)\sqrt{s}$ with properly selected absolute constant O(1) can be certified via condition $\mathbf{Q}_{\infty}(t, \frac{1}{3})$ – there exists pair $(H, \|\cdot\|_{\infty})$ satisfying this condition. The point is that by Proposition 1.10, if the condition $\mathbf{Q}_{\infty}(t, \frac{1}{3})$ can at all be satisfied, a pair $(H, \|\cdot\|_{\infty})$ satisfying this condition can be found efficiently.

Unfortunately, the significant "dropdown" in the level of sparsity when passing from unverifiable RIP to verifiable \mathbf{Q}_{∞} is inevitable; this bad news is what is on our agenda now.

1.3.3.3 Limits of performance of verifiable sufficient conditions for goodness

Proposition 1.13. Let A be an $m \times n$ sensing matrix which is "essentially non-square," specifically, such that $2m \leq n$, and let $q \in [1, \infty]$. Whenever a positive integer s and an $m \times n$ matrix H are linked by the relation

$$\|\operatorname{Col}_{j}[I_{n} - H^{T}A]\|_{s,q} < \frac{1}{2}s^{\frac{1}{q}-1}, \ 1 \le j \le n,$$
(1.38)

one has

$$s \le \sqrt{2m}.\tag{1.39}$$

As a result, sufficient condition for the validity of $\mathbf{Q}_q(s,\kappa)$ with $\kappa < 1/2$ from Proposition 1.9 can never be satisfied when $s > \sqrt{2m}$. Similarly, the verifiable



Figure 1.4: Erroneous ℓ_1 recovery of 25-sparse signal, no observation noise. Magenta: true signal, blue: ℓ_1 recovery. Top: frequency domain, bottom: time domain.

sufficient condition $\mathbf{Q}_{\infty}(s,\kappa)$, $\kappa < 1/2$ for s-goodness of A cannot be satisfied when $s > \sqrt{2m}$.

We see that unless A is "nearly square," our (same as all other known to us) verifiable sufficient conditions for s-goodness are unable to justify this property for "large" s. This unpleasant fact is in full accordance with the already mentioned fact that no individual provably s-good "essentially nonsquare" $m \times n$ matrices with $s \ge O(1)\sqrt{m}$ are known.

Matrices for which our verifiable sufficient conditions do establish s-goodness with $s \leq O(1)\sqrt{m}$ do exist.

How it works: Numerical illustration. Let us apply our machinery to the 256×512 randomly selected submatrix A of the matrix of 512×512 Inverse Discrete Cosine Transform which we used in experiments reported on Figure 1.3. These experiments exhibit nice performance of ℓ_1 minimization when recovering sparse (even nearly sparse) signals with as much as 64 nonzeros. In fact, the level of goodness of A is at most 24, as is witnessed by Figure 1.4.

In order to upper-bound the level of goodness of a matrix A, one can try to maximize the convex function $||w||_{s,1}$ over the set $W = \{w : Aw = 0, ||w||_1 \leq 1\}$; if, for a given s, the maximum of $||\cdot||_{s,1}$ over W is $\geq 1/2$, the matrix is not s-good – it does not possess the Nullspace property. Now, while global maximization of the convex function $||w||_{s,1}$ over W is difficult, we can try to find suboptimal solutions as follows: let us start with a vector $w_1 \in W$ of $||\cdot||_1$ -norm 1, and let u^1 be obtained from w_1 by replacing the s largest in magnitude entries in w_1 by the signs of these entries and zeroing out all other entries, so that $w_1^T u^1 = ||w_1||_{s,1}$. After u^1 is found, let us solve the LO program $\max_w \{[u^1]^T w : w \in W\}$. w_1 is a feasible solution to this problem, so that for the optimal solution w_2 to it we have $[u^1]^T w_2 \geq [u^1]^T w_1 =$

 $||w_1||_{s,1}$; this inequality, by virtue of what u^1 is, implies that $||w_2||_{s,1} \ge ||w_1||_{s,1}$ and by construction $w_2 \in W$. We now can iterate the construction, with w_2 in the role of w_1 , to get $w_3 \in W$ with $||w_3||_{s,1} \ge ||w_2||_{s,1}$; proceeding in this way, we generate a sequence of points from W with monotonically increasing value of the objective $|| \cdot ||_{s,1}$ we want to maximize. Usually, people terminate this recurrence either when the achieved value of the objective becomes $\ge 1/2$ (then we know for sure that A is not s-good, and can proceed to investigating s-goodness for a smaller value of s) or when the recurrence becomes stuck – the observed progress in the objective falls below a given threshold, say, 1.e-6; after it happens we can restart this process from a new randomly selected in W starting point, after getting stuck, restart again, etc., etc., until exhausting our time budget. The output of the process is the best – with the largest $|| \cdot ||_{s,1}$ – of the points from W we have generated. Applying this approach to the matrix A in question, in a couple of minutes it turns out that the matrix is at most 24-good.

One can ask how happens that experiments with recovering 64-sparse signals went fine, when in fact some 25-sparse signals cannot be recovered by ℓ_1 minimization even in the ideal noiseless case. The answer is simple: in our experiments, we dealt with *randomly selected* signals, and, as it typically is the case, randomly selected data are much nicer, whatever be the purpose of a numerical experiment, that the worst-case data.

It is interesting to understand also which goodness we can certify with our verifiable sufficient conditions. Computation shows that the fully verifiable (and strongest in our scale of sufficient conditions for s-goodness) condition $\mathbf{Q}_{\infty}(s, \varkappa)$ can be satisfied with $\varkappa < 1/2$ when s is as large as 7 and $\varkappa = 0.4887$, and cannot be satisfied with $\varkappa < 1/2$ when s = 8. As about Mutual Incoherence, it can justify just 3-goodness, no more. We hardly could be happy with the resulting bounds – goodness at least 7 and at most 24; however, it could be worse...

1.4 EXERCISES FOR LECTURE 1

Exercise 1.14. k-th Hadamard matrix, \mathcal{H}_k (here k is nonnegative integer) is the $n_k \times n_k$ matrix, $n_k = 2^k$, given by the recurrence

$$\mathcal{H}_{0} = [1]; \mathcal{H}_{k+1} = \left[\begin{array}{c|c} \mathcal{H}_{k} & \mathcal{H}_{k} \\ \hline \mathcal{H}_{k} & -\mathcal{H}_{k} \end{array} \right]$$
(1.40)

In the sequel, we assume that k > 0. Now goes the exercise:

- 1. Check that \mathcal{H}_k is symmetric matrix with entries ± 1 , and columns of the matrix are mutually orthogonal, so that $\mathcal{H}_k/\sqrt{n_k}$ is an orthogonal matrix.
- 2. Check that when k > 0, \mathcal{H}_k has just two distinct eigenvalues, $\sqrt{n_k}$ and $-\sqrt{n_k}$, each of multiplicity $m_k := 2^{k-1} = n_k/2$.
- 3. Prove that whenever f is an eigenvector of \mathcal{H}_k , one has

$$\|f\|_{\infty} \le \|f\|_1 / \sqrt{n_k}$$

Derive from this observation the conclusion as follows:

Let $a_1, ..., a_{m_k} \in \mathbf{R}^{n_k}$ be orthogonal to each other unit vectors which are

eigenvectors of \mathcal{H}_k with eigenvalues $\sqrt{n_k}$ (by the above, the dimension of the eigenspace of \mathcal{H}_k associated with the eigenvalue $\sqrt{n_k}$ is m_k , so that the required $a_1, ..., a_{m_k}$ do exist), and let A be the $m_k \times n_k$ matrix with the rows $a_1^T, ..., a_{m_k}^T$. For every $x \in \text{Ker } A$ it holds

$$\|x\|_{\infty} \le \frac{1}{\sqrt{n_k}} \|x\|_1$$

whence A satisfies the nullspace property whenever the sparsity s satisfies $2s < \sqrt{n_k} = \sqrt{2m_k}$. Moreover, there exists (and can be found efficiently) an $m_k \times n_k$ contrast matrix $H = H_k$ such that for every $s < \frac{1}{2}\sqrt{n_k}$, the pair $(H_k, \|\cdot\|_{\infty})$ satisfies the associated with A condition $\mathbf{Q}_{\infty}(s, \kappa_s = \frac{s}{\sqrt{n_k}})$,

and the $\|\cdot\|_2$ -norms of columns of H_k do not exceed $\sqrt{2\frac{\sqrt{n_k}+1}{\sqrt{n_k}}}$.

Note that the above conclusion yields a sequence of individual $(m_k = 2^{k-1}) \times (n_k = 2^k)$ sensing matrices, k = 1, 2, ..., with "size ratio" $n_k/m_k = 2$, which make an efficiently verifiable condition for s-goodness, say, $\mathbf{Q}_{\infty}(s, \frac{1}{3})$ satisfiable in basically the entire range of values of s allowed by Proposition 1.13. It would be interesting to get similar "fully constructive" results for other size ratios, like m: n = 1: 4, m: n = 1: 8, etc.

Exercise 1.15. [Follow-up to Exercise 1.14] Exercise 1.14 provides us with an explicitly given $(m = 512) \times (n = 1024)$ sensing matrix \overline{A} such that the efficiently verifiable condition $\mathbf{Q}_{\infty}(15, \frac{15}{32})$ is satisfiable; in particular, \overline{A} is 15-good. With all we know about limits of performance of verifiable sufficient conditions for goodness, how should we evaluate this specific sensing matrix? Could we point out a sensing matrix of the same size which is provably *s*-good for a larger (or "much larger") than 15 value of *s*?

We do not know the answer, and you are requested to explore some possibilities, including (but not reducing to – you are welcome to investigate more options!) the following ones.

- 1. Generate at random a sample of $m \times n$ sensing matrices A, compute their mutual incoherences and look how large goodness levels they justify. What happens when the matrices are the Gaussian (independent $\mathcal{N}(0, 1)$ entries) and the Rademacher ones (independent entries taking values ± 1 with probabilities 1/2)?
- 2. Generate at random a sample of $m \times n$ matrices with independent $\mathcal{N}(0, 1/m)$ entries. Proposition 1.7 suggests that a sampled matrix A has good chances to satisfy RIP (δ, k) with some $\delta < 1/3$ and some k, and thus to be s-good (and even more than this, see Proposition 1.8) for every $s \leq k/2$. Of course, given A we cannot check whether the matrix indeed satisfies RIP (δ, k) with given δ, k ; what we can try to do is to certify that RIP (δ, k) does <u>not</u> take place. To this end, it suffices to select at random, say, 200 $m \times k$ submatrices \tilde{A} of A and compute the eigenvalues of $\tilde{A}^T \tilde{A}$; if A possesses RIP (δ, k) , all these eigenvalues should belong to the segment $[1 - \delta, 1 + \delta]$, and if in reality this does not happen, A definitely is not RIP (δ, k) .

Exercise 1.16. Let us start with preamble. Consider a finite Abelian group; the only thing which matters for us is that such a group G is specified by a collection of a

 $k \ge 1$ of positive integers $\nu_1, ..., \nu_k$ and is comprised of all collections $\omega = (\omega_1, ..., \omega_k)$ where every ω_i is an integer from the range $\{0, 1, ..., \nu_k - 1\}$; the group operation, denoted by \oplus , is

$$(\omega_1, ..., \omega_k) \oplus (\omega'_1, ..., \omega'_k) = ((\omega_1 + \omega'_1) \operatorname{mod} \nu_1, ..., (\omega_k + \omega'_k) \operatorname{mod} \nu_k),$$

where $a \mod b$ is the remainder, taking values in $\{0, 1, ..., b-1\}$, in the division of an integer a by positive integer b; say, $5 \mod 3 = 2$, and $6 \mod 3 = 0$. Clearly, the cardinality of the above group G is $n_k = \nu_1 \nu_2 ... \nu_k$. A *character* of group G is a homomorphism acting from G into the multiplicative group of complex numbers of modulus 1, or, in simple words, a complex-valued function $\chi(\omega)$ on G such that $|\chi(\omega)| = 1$ for all $\omega \in G$ and $\chi(\omega \oplus \omega') = \chi(\omega)\chi(\omega')$ for all $\omega, \omega' \in G$. Note that characters themselves form a group w.r.t. pointwise multiplication; clearly, all characters of our G are functions of the form

$$\chi((\omega_1, ..., \omega_k)) = \mu_1^{\omega_1} ... \mu_k^{\omega_k}$$

where μ_i are restricted to be roots of degree ν_i from 1: $\mu_i^{\nu_i} = 1$. It is immediately seen that the group G_* of characters of G is of the same cardinality $n_k = \nu_1 \dots \nu_k$ as G. We can associate with G the matrix \mathcal{F} of size $n_k \times n_k$; the columns in the matrix are indexed by the elements ω of G, the rows – by the characters $\chi \in G_*$ of G, and the element in cell (χ, ω) is $\chi(\omega)$. The standard example here corresponds to k = 1, in which case \mathcal{F} clearly is the $\nu_1 \times \nu_1$ matrix of Discrete Fourier Transform.

Now goes the exercise:

- 1. Verify that the above \mathcal{F} is, up to factor $\sqrt{n_k}$, a unitary matrix: denoting by \overline{a} the complex conjugate of a complex number a, $\sum_{\omega \in G} \chi(\omega) \overline{\chi'}(\omega)$ is n_k or 0 depending on whether $\chi = \chi'$ or $\chi \neq \chi'$.
- 2. Let $\bar{\omega}, \bar{\omega}'$ be two elements of G. Prove that there exists a permutation Π of elements of G which maps $\bar{\omega}$ into $\bar{\omega}'$ and is such that

$$\operatorname{Col}_{\Pi(\omega)}[\mathcal{F}] = D\operatorname{Col}_{\omega}[\mathcal{F}] \; \forall \omega \in G,$$

where D is diagonal matrix with diagonal entries $\chi(\bar{\omega}')/\chi(\bar{\omega}), \chi \in G_*$.

- 3. Consider the special case of the above construction where $\nu_1 = \nu_2 = ... = \nu_k = 2$. Verify that in this case \mathcal{F} , up to permutation of rows and permutation of columns (these permutations depend on how we assign the elements of G and of G_* their serial numbers) is exactly the Hadamard matrix \mathcal{H}_k .
- 4. Extract from the above the following fact: let m, k be positive integers such that $m \leq n_k := 2^k$, and let sensing matrix A be obtained from \mathcal{H}_k by selecting m distinct rows. Assume we want to find an $m \times n_k$ contrast matrix H such that the pair $(H, \|\cdot\|_{\infty})$ satisfies the condition $\mathbf{Q}_{\infty}(s, \kappa)$ with as small κ as possible; by Proposition 1.10, to this end we should solve n LP programs

$$\operatorname{Opt}_i = \min_h \|e^i - A^T h\|_{\infty},$$

where e^i is *i*-th basic orth in \mathbb{R}^n . Prove that with A coming from \mathcal{H}_k , all these problems have the same optimal value, and optimal solutions to all of the problems are readily given by the optimal solution to just one of them.

Exercise 1.17. Proposition 1.13 states that the verifiable condition $\mathbf{Q}_{\infty}(s,\kappa)$ can

certify s-goodness of "essentially nonsquare" (with $m \le n/2$) $m \times n$ sensing matrix A only when s is small as compared to m, namely, $s \le \sqrt{2m}$. The exercise to follow is aimed at investigating what happens when $m \times n$ "low" (with m < n) sensing matrix A is "nearly square", meaning that $m^o = n - m$ is small as compared to n. Specifically, you should prove that for properly selected individual $(n - m^o) \times n$ matrices A the condition $\mathbf{Q}_{\infty}(s,\kappa)$ with $\kappa < 1/2$ is satisfiable when s is as large as $O(1)n/\sqrt{m^o}$.

1. Let $n = 2^k p$ with positive integer p and integer $k \ge 1$, and let $m^o = 2^{k-1}$. Given $2m^o$ -dimensional vector u, let u^+ be n-dimensional vector built as follows: we split indexes from $\{1, ..., n = 2^k p\}$ into 2^k consecutive groups $I_1, ..., I_{2^k}, p$ elements per group, and all entries of u^+ with indexes from I_i are equal to *i*-th entry, u_i , of vector u. Now let U be the linear subspace in \mathbf{R}^{2^k} comprised of all eigenvectors, with eigenvalue $\sqrt{2^k}$, of the Hadamard matrix \mathcal{H}_k , see Exercise 1.14, so that the dimension of U is $2^{k-1} = m^o$, and let L be given by

$$L = \{u^+ : u \in U\} \subset \mathbf{R}^n$$

Clearly, L is a linear subspace in \mathbf{R}^n of dimension m^o . Prove that

$$\forall x \in L : \|x\|_{\infty} \le \frac{\sqrt{2m^o}}{n} \|x\|_1.$$

Conclude that if A is $(n - m^o) \times n$ sensing matrix with Ker A = L, then the verifiable sufficient condition $\mathbf{Q}_{\infty}(s,\kappa)$ does certify s-goodness of A whenever

$$1 \le s < \frac{n}{2\sqrt{2m^o}}.$$

2. Let L be m^{o} -dimensional subspace in \mathbb{R}^{n} . Prove that L contains nonzero vector x with

$$\|x\|_{\infty} \ge \frac{\sqrt{m^o}}{n} \|x\|_1,$$

so that the condition $\mathbf{Q}_{\infty}(s,\kappa)$ cannot certify s-goodness of $(n-m^o) \times n$ sensing matrix A whenever $s > O(1)n/\sqrt{m^o}$, for properly selected absolute constant O(1).

Exercise 1.18. Utilize the results of Exercise 1.16 in a numerical experiment as follows.

- select n as an integer power 2^k of 2, say, set $n = 2^{10} = 1024$
- select a "representative" sequence M of values of $m, 1 \le m < n$, including values of m close to n and "much smaller" than n, say, use

 $M = \{2, 5, 8, 16, 32, 64, 128, 256, 512, 7, 896, 960, 992, 1008, 1016, 1020, 1022, 1023\}$

- for every $m \in M$,
 - generate at random an $m \times n$ submatrix A of the $n \times n$ Hadamard matrix \mathcal{H}_k and utilize the result of item 4 of Exercise 1.16 in order to find the largest s such that s-goodness of A can be certified via the condition $\mathbf{Q}_{\infty}(\cdot, \cdot)$; call s(m) the resulting value of s.

- generate a moderate sample of Gaussian $m \times n$ sensing matrices A_i with independent $\mathcal{N}(0, 1/m)$ entries and use the construction from Exercise 1.15 to upper-bound the largest *s* for which a matrix from the sample satisfies RIP(1/3, 2s); call $\hat{s}(m)$ the largest, over your A_i 's, of the resulting upper bounds.

The goal of the exercise is to compare the computed values of s(m) and $\hat{s}(m)$; in other words, we again want to understand how "theoretically perfect" RIP compares to "conservative restricted scope" condition \mathbf{Q}_{∞} .

1.5 PROOFS

1.5.1 **Proofs of Theorem 1.3, 1.4**

All we need is to prove Theorem 1.4, since Theorem 1.3 is the particular case $\varkappa = \kappa < 1/2$ of Theorem 1.4.

Let us fix $x \in \mathbf{R}^n$ and $\eta \in \Xi_\rho$, and let us set $\hat{x} = \hat{x}_{reg}(Ax + \eta)$. Let also $I \subset \{1, ..., n\}$ be the set of indexes of the *s* largest in magnitude entries in *x*, I^o be the complement of *I* in $\{1, ..., n\}$, and let for $w \in \mathbf{R}^n$, w_I and w_{I^o} be the vectors obtained from *w* by zeroing entries with indexes $j \notin I$ and $j \notin I^o$, respectively, and keeping the remaining entries intact. Finally, let $z = \hat{x} - x$. **1**⁰. By the definition of Ξ_ρ and due to $\eta \in \Xi_\rho$ we have

$$\|H^T([Ax+\eta] - Ax)\| \le \rho, \tag{1.41}$$

so that x is a feasible solution to the optimization problem specifying \hat{x} , whence $\|\hat{x}\|_1 \leq \|x\|_1$. We therefore have

$$\begin{aligned} \|\widehat{x}_{I^{o}}\|_{1} &= \|\widehat{x}\|_{1} - \|\widehat{x}_{I}\|_{1} \le \|x\|_{1} - \|\widehat{x}_{I}\|_{1} = \|x_{I}\|_{1} + \|x_{I^{o}}\|_{1} - \|\widehat{x}_{I}\|_{1} \\ &\le \|z_{I}\|_{1} + \|x_{I^{o}}\|_{1}, \end{aligned}$$
(1.42)

and therefore

$$||z_{I^o}||_1 \le ||\widehat{x}_{I^o}||_1 + ||x_{I^o}||_1 \le ||z_I||_1 + 2||x_{I^o}||_1.$$

It follows that

$$||z||_1 = ||z_I||_1 + ||z_{I^o}||_1 \le 2||z_I||_1 + 2||x_{I^o}||_1.$$
(1.43)

Further, by definition of \hat{x} we have $||H^T([Ax+\eta] - A\hat{x})|| \le \rho$, which combines with (1.41) to imply that

$$\|H^T A(\hat{x} - x)\| \le 2\rho.$$
(1.44)

2⁰. Since $(H, \|\cdot\|)$ satisfies $\mathbf{Q}_1(s, \varkappa)$, we have

$$||z||_{s,1} \le s ||H^T A z|| + \varkappa ||z||_1.$$

By (1.44), it follows that $||z||_{s,1} \leq 2s\rho + \varkappa ||z||_1$, which combines with the evident inequality $||z_I|| \leq ||z||_{s,1}$ (recall that $\operatorname{Card}(I) = s$) and with (1.43) to imply that

$$||z_I||_1 \le 2s\rho + \varkappa ||z||_1 \le 2s\rho + 2\varkappa ||z_I||_1 + 2\varkappa ||x_{I^o}||_1,$$

whence

$$||z_I||_1 \le \frac{2s\rho + 2\varkappa ||x_{I^o}||_1}{1 - 2\varkappa}.$$

Invoking (1.43), we conclude that

$$\|z\|_{1} \leq \frac{4s}{1-2\varkappa} \left[\rho + \frac{\|x_{I^{o}}\|_{1}}{2s}\right].$$
(1.45)

3⁰. Since $(H, \|\cdot\|)$ satisfies $\mathbf{Q}_q(s, \kappa)$, we have

$$||z||_{s,q} \le s^{\frac{1}{q}} ||H^T A z|| + \kappa s^{\frac{1}{q}-1} ||z||_1,$$

which combines with (1.45) and (1.44) to imply that

$$\|z\|_{s,q} \le s^{\frac{1}{q}} 2\rho + \kappa s^{\frac{1}{q}} \frac{4\rho + 2s^{-1} \|x_{I^o}\|_1}{1 - 2\varkappa} \le \frac{4s^{\frac{1}{q}} [1 + \kappa - \varkappa]}{1 - 2\varkappa} \left[\rho + \frac{\|x_o\|_1}{2s}\right]$$
(1.46)

(we have taken into account that $\varkappa < 1/2$ and $\kappa \ge \varkappa$). Let θ be the (s + 1)-st largest magnitude of entries in z, and let $w = z - z^s$. Now (1.46) implies that

$$\theta \le \|z\|_{s,q} s^{-\frac{1}{q}} \le \frac{4[1+\kappa-\varkappa]}{1-2\varkappa} \left[\rho + \frac{\|x_{I^o}\|_1}{2s}\right].$$

Hence invoking (1.45) we have

$$\begin{split} \|w\|_{q} &\leq \|w\|_{\infty}^{\frac{q-1}{q}} \|w\|_{1}^{\frac{1}{q}} \leq \theta^{\frac{q-1}{q}} \|z\|_{1}^{\frac{1}{q}} \\ &\leq \theta^{\frac{q-1}{q}} \frac{(4s)^{\frac{1}{q}}}{[1-2\varkappa]^{\frac{1}{q}}} \left[\rho + \frac{\|x_{I^{o}}\|_{1}}{2s}\right]^{\frac{1}{q}} \\ &\leq \frac{4s^{\frac{1}{q}}[1+\kappa-\varkappa]^{\frac{q-1}{q}}}{1-2\varkappa} \left[\rho + \frac{\|x_{I^{o}}\|_{1}}{2s}\right]. \end{split}$$

Taking into account (1.46) and the fact that the supports of z^s and w do not intersect, we get

$$\begin{aligned} \|z\|_{q} &\leq 2^{\frac{1}{q}} \max[\|z^{s}\|_{q}, \|w\|_{q}] = 2^{\frac{1}{q}} \max[\|z\|_{s,q}, \|w\|_{q}] \\ &\leq \frac{4(2s)^{\frac{1}{q}}[1+\kappa-\varkappa]}{1-2\varkappa} \left[\rho + \frac{\|x_{I^{o}}\|_{1}}{2s}\right]. \end{aligned}$$

This bound combines with (1.45), the Hölder inequality and the relation $||x_{I^o}||_1 = ||x - x^s||_1$ to imply (1.16).

1.5.2 Proof of Theorem 1.5

Let us prove (i). Let us fix $x \in \mathbf{R}^n$ and η , and let us set $\hat{x} = \hat{x}_{pen}(Ax + \eta)$. Let also $I \subset \{1, ..., K\}$ be the set of indexes of the *s* largest in magnitude entries in *x*, I^o be the complement of *I* in $\{1, ..., n\}$, and for $w \in \mathbf{R}^n$ let w_I , w_{I^o} be the vectors obtained from *w* by zeroing out all entries with indexes not in *I*, respectively, not in I^o . Finally, let $z = \hat{x} - x$ and $\nu = ||H^T \eta||$. **1**⁰. We have

 $\|\widehat{x}\|_1 + \lambda \|H^T (A\widehat{x} - Ax - \eta)\| \le \|x\|_1 + \lambda \|H^T \eta\|$

and

$$||H^{T}(A\widehat{x} - Ax - \eta)|| = ||H^{T}(Az - \eta)|| \ge ||H^{T}Az|| - ||H^{T}\eta||,$$

whence

$$\|\widehat{x}\|_{1} + \lambda \|H^{T}Az\| \le \|x\|_{1} + 2\lambda \|H^{T}\eta\| = \|x\|_{1} + 2\lambda\nu.$$
(1.47)

We have

$$\begin{aligned} \|\widehat{x}\|_{1} &= \|x+z\|_{1} = \|x_{I}+z_{I}\|_{1} + \|x_{I^{o}}+z_{I^{o}}\|_{1} \\ &\geq \|x_{I}\|_{1} - \|z_{I}\|_{1} + \|z_{I^{o}}\|_{1} - \|x_{I^{o}}\|_{1}, \end{aligned}$$

which combines with (1.47) to imply that

$$||x_I||_1 - ||z_I||_1 + ||z_{I^o}||_1 - ||x_{I^o}||_1 + \lambda ||H^T A z|| \le ||x||_1 + 2\lambda\nu,$$

or, which is the same,

$$||z_{I^o}||_1 - ||z_I||_1 + \lambda ||H^T A z|| \le 2||x_{I^o}||_1 + 2\lambda\nu.$$
(1.48)

Since $(H, \|\cdot\|)$ satisfies $\mathbf{Q}_1(s, \varkappa)$, we have

$$||z_I||_1 \le ||z||_{s,1} \le s ||H^T A z|| + \varkappa ||z||_1,$$

so that

$$(1 - \varkappa) \|z_I\|_1 - \varkappa \|z_{I^o}\|_1 - s \|H^T A z\| \le 0.$$
(1.49)

Taking weighted sum of (1.48) and (1.49), the weights being 1, 2, respectively, we get

$$(1-2\varkappa) [\|z_I\|_1 + \|z_{I^o}\|_1] + (\lambda - 2s) \|H^T A z\| \le 2 \|x_{I^o}\|_1 + 2\lambda\nu,$$

that is (since $\lambda \geq 2s$),

$$\|z\|_{1} \leq \frac{2\lambda\nu + 2\|x_{I^{o}}\|_{1}}{1 - 2\varkappa} \leq \frac{2\lambda}{1 - 2\varkappa} \left[\nu + \frac{\|x_{I^{o}}\|_{1}}{2s}\right].$$
 (1.50)

Further, by (1.47) we have

$$\lambda \|H^T A z\| \le \|x\|_1 - \|\widehat{x}\|_1 + 2\lambda\nu \le \|z\|_1 + 2\lambda\nu,$$

which combines with (1.50) to imply that

$$\lambda \| H A^T z \| \le \frac{2\lambda\nu + 2\|x_{I^o}\|_1}{1 - 2\varkappa} + 2\lambda\nu = \frac{2\lambda\nu(2 - 2\varkappa) + 2\|x_{I^o}\|_1}{1 - 2\varkappa}.$$
 (1.51)

From $\mathbf{Q}_q(s,\kappa)$ it follows that

$$||z||_{s,q} \le s^{\frac{1}{q}} ||H^T A z|| + \kappa s^{\frac{1}{q}-1} ||z||_1,$$

which combines with (1.51) and (1.50) to imply that

(recall that $\lambda \geq 2s$, $\kappa \geq \varkappa$, and $\varkappa < 1/2$). It remains to repeat the reasoning following (1.46) in item 3⁰ of the proof of Theorem 1.4. Specifically, denoting by θ

the (s+1)-st largest magnitude of entries in z, (1.52) implies that

$$\theta \le s^{-1/q} \|z\|_{s,q} \le \frac{4}{1 - 2\varkappa} [1 + \kappa \frac{\lambda}{2s} - \varkappa] \left[\nu + \frac{\|x_{I^o}\|_1}{2s}\right], \tag{1.53}$$

so that for the vector $w = z - z^s$ one has

$$\|w\|_{q} \leq \theta^{1-\frac{1}{q}} \|w\|_{1}^{\frac{1}{q}} \leq \frac{4(\lambda/2)^{\frac{1}{q}}}{1-2\varkappa} \left[1+\kappa\frac{\lambda}{2s}-\varkappa\right]^{\frac{q-1}{q}} \left[\nu+\frac{\|x_{I^{o}}\|_{1}}{2s}\right]$$

(we have used (1.53), (1.50) and the fact that $\lambda \geq 2s$). Hence, taking into account that z^s and w have non-intersecting supports,

$$\begin{aligned} \|z\|_{q} &\leq 2^{\frac{1}{q}} \max[\|z^{s}\|_{q}, \|w\|_{q}] = 2^{\frac{1}{q}} \max[\|z\|_{s,q}, \|w\|_{q}] \\ &\leq \frac{4\lambda^{\frac{1}{q}}}{1-2\varkappa} \left[1 + \kappa \frac{\lambda}{2s} - \varkappa\right] \left[\nu + \frac{\|x_{I^{o}}\|_{1}}{2s}\right] \end{aligned}$$

(we have used (1.52) along with $\lambda \geq 2s$ and $\kappa \geq \varkappa$). This combines with (1.50) and Hölder inequality to imply (1.18). All remaining claims of Theorem 1.5 are immediate corollaries of (1.18).

1.5.3 Proof of Proposition 1.7

1⁰. Assuming $k \leq m$ and selecting a set I of k distinct from each other indexes from $\{1, ..., n\}$, consider an $m \times k$ submatrix A_I of A comprised of columns with indexes from I, and let u be a unit vector in \mathbf{R}^k . The entries in the vector $m^{1/2}A_I u$ are independent $\mathcal{N}(0, 1)$ random variables, so that for the random variable $\zeta_u = \sum_{i=1}^{m} (m^{1/2}A_I u)_i^2$ and $\gamma \in (-1/2, 1/2)$ it holds (in what follows, expectations and probabilities are taken w.r.t. our ensemble of random A's)

$$\ln\left(\mathbf{E}\{\exp\{\gamma\zeta\}\}\right) = m\ln\left(\frac{1}{\sqrt{2\pi}}\int e^{\gamma t^2 - \frac{1}{2}t^2}ds\right) = -\frac{m}{2}\ln(1-2\gamma).$$

Given $\alpha \in (0, 0.1]$ and selecting γ in such a way that $1 - 2\gamma = \frac{1}{1+\alpha}$, we get $0 < \gamma < 1/2$ and therefore

$$\begin{aligned} \operatorname{Prob}\{\zeta_u > m(1+\alpha)\} &\leq \mathbf{E}\{\exp\{\gamma\zeta_u\}\}\exp\{-m\gamma(1+\alpha)\}\\ &= \exp\{-\frac{m}{2}\ln(1-2\gamma) - m\gamma(1+\alpha)\}\\ &= \exp\{\frac{m}{2}\left[\ln(1+\alpha) - \alpha\right]\} \leq \exp\{-\frac{m}{5}\alpha^2\},\end{aligned}$$

and similarly, selecting γ in such a way that $1 - 2\gamma = \frac{1}{1-\alpha}$, we get $-1/2 < \gamma < 0$ and therefore

$$\begin{aligned} \operatorname{Prob}\{\zeta_u < m(1-\alpha)\} &\leq \mathbf{E}\{\exp\{\gamma\zeta_u\}\}\exp\{-m\gamma(1-\alpha)\}\\ &= \exp\{-\frac{m}{2}\ln(1-2\gamma) - m\gamma(1-\alpha)\}\\ &= \exp\{\frac{m}{2}\left[\ln(1-\alpha) + \alpha\right]\} \leq \exp\{-\frac{m}{5}\alpha^2\},\end{aligned}$$

and we end up with

$$u \in \mathbf{R}^{k}, \|u\|_{2} = 1 \Rightarrow \begin{cases} \operatorname{Prob}\{A : \|A_{I}u\|_{2}^{2} > 1 + \alpha\} \leq \exp\{-\frac{m}{5}\alpha^{2}\} \\ \operatorname{Prob}\{A : \|A_{I}u\|_{2}^{2} < 1 - \alpha\} \leq \exp\{-\frac{m}{5}\alpha^{2}\} \end{cases}$$
(1.54)

2⁰. Same as above, let $\alpha \in (0, 0.1]$, let

$$M = 1 + 2\alpha, \epsilon = \frac{\alpha}{2(1+2\alpha)},$$

and let us build an ϵ -net on the unit sphere S in \mathbb{R}^k as follows. We start with a point $u_1 \in S$; after $\{u_1, ..., u_t\} \subset S$ is already built, we check whether there is a point in S at the $\|\cdot\|_2$ -distance from all points of the set $> \epsilon$. If it is the case, we add such a point to the net built so far and proceed with building the net, otherwise we terminate with the net $\{u_1, ..., u_t\}$. By compactness of S and due to $\epsilon > 0$, this process eventually terminates; upon termination, we have at our disposal collection $\{u_1, ..., u_N\}$ of unit vectors such that every two of them are at the $\|\cdot\|_2$ -distance $> \epsilon$ from each other, and every point from S is at the distance at most ϵ from some point of the collection. We claim that the cardinality N of the resulting set can be bounded as

$$N \le \left[\frac{2+\epsilon}{\epsilon}\right]^k = \left[\frac{4+9\alpha}{\alpha}\right]^k \le \left(\frac{5}{\alpha}\right)^k.$$
(1.55)

Indeed, the interiors of the $\|\cdot\|_2$ -balls of radius $\epsilon/2$ centered at the points $u_1, ..., u_N$ are mutually disjoint, and their union is contained in the $\|\cdot\|_2$ -ball of radius $1 + \epsilon/2$ centered at the origin; comparing the volume of the union and the one of the ball, we arrive at (1.55).

3⁰. Consider event *E* comprised of all realizations of *A* such that for all *k*-element subsets *I* of $\{1, ..., n\}$ and all $t \leq n$ it holds

$$1 - \alpha \le \|A_I u_t\|_2^2 \le 1 + \alpha. \tag{1.56}$$

By (1.54) and the union bound,

$$\operatorname{Prob}\{A \notin E\} \le 2N \left(\begin{array}{c} n\\ k \end{array}\right) \exp\{-\frac{m}{5}\alpha^2\}.$$
(1.57)

We claim that

$$A \in E \Rightarrow (1 - 2\alpha) \le ||A_I u||_2^2 \le 1 + 2\alpha \ \forall \left(\begin{array}{c} I \subset \{1, ..., n\} : \operatorname{Card}(I) = k\\ u \in \mathbf{R}^k : ||u||_2 = 1 \end{array}\right).$$
(1.58)

Indeed, let $A \in E$, let us fix $I \in \{1, ..., n\}$, $\operatorname{Card}(I) = k$, and let M be the maximal value of the quadratic form $f(u) = u^T A_I^T A_I u$ on the unit $\|\cdot\|_2$ -ball B, centered at the origin, in \mathbf{R}^k . In this ball, f is Lipschitz continuous with constant 2M w.r.t. $\|\cdot\|_2$; denoting by \bar{u} a maximizer of the form on B, we lose nothing when assuming that \bar{u} is a unit vector. Now let u_s be the point of our net which is at the $\|\cdot\|_2$ -distance from \bar{u} at most ϵ . We have

$$M = f(\bar{u}) \le f(u_s) + 2M\epsilon \le 1 + \alpha + 2M\epsilon,$$

whence

$$M \le \frac{1+\alpha}{1-2\epsilon} = 1+2\alpha,$$

implying the right inequality in (1.58). Now let u be unit vector in \mathbf{R}^k , and u_s be a point in the net at the $\|\cdot\|$ -distance $\leq \epsilon$ from u. We have

$$f(u) \ge f(u_s) - 2M\epsilon \ge 1 - \alpha - 2\frac{1+\alpha}{1-2\epsilon}\epsilon = 1 - 2\alpha,$$

justifying the first inequality in (1.58).

The bottom line is:

$$\delta \in (0, 0.2], 1 \le k \le n$$

$$\Rightarrow \operatorname{Prob}\{A : A \text{ does not satisfy RIP}(\delta, k)\} \le \underbrace{2\left(\frac{10}{\delta}\right)^k}_{\le \left(\frac{20}{\delta}\right)^k} \binom{n}{k} \exp\{-\frac{m\delta^2}{20}\}.$$
(1.59)

Indeed, setting $\alpha = \delta/2$, we have seen that whenever $A \notin E$, we have $(1 - \delta) \leq ||Au||_2^2 \leq (1 + \delta)$ for all unit k-sparse u, which is nothing but $\operatorname{RIP}(\delta, k)$; with this in mind, (1.59) follows from (1.57) and (1.55).

4⁰. It remains to verify that with properly selected, depending solely on δ , positive quantities c, d, f, for every $k \geq 1$ satisfying (1.27) the right hand side in (1.59) is at most $\exp\{-fm\}$. Passing to logarithms, our goal is to ensure the relation

$$G := a(\delta)m - b(\delta)k - \ln \binom{n}{k} \ge mf(\delta) > 0$$
$$\left[a(\delta) = \frac{\delta^2}{20}, b(\delta) = \ln \left(\frac{20}{\delta}\right)\right]$$
(1.60)

provided that $k \ge 1$ satisfies (1.27).

Let k satisfy (1.27) with some c, d to be specified later, and let y = k/m. Assuming $d \ge 3$, we have $0 \le y \le 1/3$. Now, it is well known that

$$C := \ln \left(\begin{array}{c} n \\ k \end{array} \right) \le n \left[\frac{k}{n} \ln(\frac{n}{k}) + \frac{n-k}{n} \ln(\frac{n}{n-k}) \right],$$

whence

$$C \le n \left[\frac{m}{n} y \ln(\frac{n}{my}) + \frac{n-k}{n} \underbrace{\ln(1 + \frac{k}{n-k})}_{\le \frac{k}{n-k}}\right]$$
$$\le n \left[\frac{m}{n} y \ln(\frac{n}{my}) + \frac{k}{n}\right] = m \left[y \ln(\frac{n}{my}) + y\right] \le 2my \ln(\frac{n}{my})$$

(recall that $n \ge m$ and $y \le 1/3$). It follows that

$$G = a(\delta)m - b(\delta)k - C \ge a(\delta)m - b(\delta)ym - 2my\ln(\frac{n}{my})$$

=
$$m\underbrace{\left[a(\delta) - b(\delta)y - 2y\ln(\frac{n}{m}) - 2y\ln(\frac{1}{y})\right]}_{H},$$

and all we need is to select c, d in such a way that (1.27) would imply that $H \ge f$ with some positive $f = f(\delta)$. This is immediate: we can find $u(\delta) > 0$ such that when $0 \le y \le u(\delta)$, we have $2y \ln(1/y) + b(\delta)y \le \frac{1}{3}a(\delta)$; selecting $d(\delta) \ge 3$ large

enough, (1.27) would imply $y \leq u(\delta)$, and thus would imply

$$H \ge \frac{2}{3}a(\delta) - 2y\ln(\frac{n}{m})$$

Now we can select $c(\delta)$ large enough for (1.27) to ensure that $2y \ln(\frac{n}{m}) \leq \frac{1}{3}a(\delta)$. With just specified c, d, (1.27) implies that $H \geq \frac{1}{3}a(\delta)$, and we can take the latter quantity as $f(\delta)$.

1.5.4 Proof of Propositions 1.8, 1.12

Let us prove Proposition 1.12; as a byproduct of our reasoning, we shall prove Proposition 1.8 as well.

Let $x \in \mathbf{R}^n$, and let $x^1, ..., x^q$ be obtained from x by the following construction: x^1 is obtained from x by zeroing all but the s largest in magnitude entries; x^2 is obtained by the same procedure applied to $x - x^1$, $x^3 - by$ the same procedure applied to $x - x^1 - x^2$, and so on; the process is terminated at the first step q when it happens that $x = x^1 + ... + x^q$. Note that for $j \ge 2$ we have $||x^j||_{\infty} \le s^{-1} ||x^{j-1}||_1$ and $||x^j||_1 \le ||x^{j-1}||_1$, whence also $||x^j||_2 \le \sqrt{||x^j||_{\infty} ||x^j||_1} \le s^{-1/2} ||x^{j-1}||_1$. It is easily seen that if A is RIP $(\delta, 2s)$, then for every two s-sparse vectors u, v with non-overlapping supports we have

$$v^{T} A^{T} A u | \le \delta \|u\|_{2} \|v\|_{2}.$$
(*)

Indeed, for s-sparse u, v, let I be the index set of cardinality $\leq 2s$ containing the supports of u and v, so that, denoting by A_I the submatrix of A comprised of columns with indexes from I, we have $v^T A^T A u = v_I^T [A_I^T A_I] u_I$. By RIP, the eigenvalues $\lambda_i = 1 + \mu_i$ of the symmetric matrix $Q = A_I^T A_I$ are in-between $1 - \delta$ and $1 + \delta$; representing u_I and v_I by vectors w, z of their coordinates in the orthonormal eigenbasis of Q, we get $|v^T A^T A u| = |\sum_i \lambda_i w_i z_i| = |\sum_i w_i z_i + \sum_i \mu_i w_i z_i| \le |w^T z| + \delta ||w||_2 ||z||_2$. It remains to note that $w^T z = u_I^T v_I = 0$ and $||w||_2 = ||u||_2$, $||z||_2 = ||v||_2$.

(i): We have

$$\begin{aligned} \|Ax^{1}\|_{2}\|Ax\|_{2} &\geq [x^{1}]^{T}A^{T}Ax = \|Ax^{1}\|_{2}^{2} - \sum_{j=2}^{q} [x^{1}]^{T}A^{T}Ax^{j} \\ &\geq \|Ax^{1}\|_{2}^{2} - \delta \sum_{j=2}^{q} \|x^{1}\|_{2} \|x^{j}\|_{2} \ [by \ (*)] \\ &\geq \|Ax^{1}\|_{2}^{2} - \delta s^{-1/2} \|x^{1}\|_{2} \sum_{j=2}^{q} \|x^{j-1}\|_{1} \geq \|Ax^{1}\|_{2}^{2} - \delta s^{-1/2} \|x^{1}\|_{2} \|x\|_{1} \\ &\Rightarrow \|Ax^{1}\|_{2}^{2} \leq \|Ax^{1}\|_{2} \|Ax\|_{2} + \delta s^{-1/2} \|x^{1}\|_{2} \|x\|_{1} \\ &\Rightarrow \|x^{1}\|_{2} = \frac{\|x^{1}\|_{2}}{\|Ax^{1}\|_{2}^{2}} \|Ax^{1}\|_{2}^{2} \leq \frac{\|x^{1}\|_{2}}{\|Ax^{1}\|_{2}} \|Ax\|_{2} + \delta s^{-1/2} \left(\frac{\|x^{1}\|_{2}}{\|Ax^{1}\|_{2}}\right)^{2} \|x\|_{1} \\ &\Rightarrow \|x\|_{s,2} = \|x^{1}\|_{2} \leq \frac{1}{\sqrt{1-\delta}} \|Ax\|_{2} + \frac{\delta s^{-1/2}}{1-\delta} \|x\|_{1} \ [by \ \text{RIP}(\delta, 2s)] \end{aligned}$$

and we see that the pair $\left(H = \frac{s^{-1/2}}{\sqrt{1-\delta}}I_m, \|\cdot\|_2\right)$ satisfies $\mathbf{Q}_2(s, \frac{\delta}{1-\delta})$, as claimed in (i).

In addition, the relation after the first \Rightarrow implies that

$$||Ax^{1}||_{2} \le ||Ax||_{2} + \delta s^{-1/2} \left[\frac{||x_{1}||_{2}}{||Ax^{1}||_{2}} \right] ||x||_{1}.$$

By RIP, the left hand side in this inequality is $\geq ||x^1||_2 \sqrt{1-\delta}$, while the ratio of norms in the right hand side is $\leq \frac{1}{\sqrt{1-\delta}}$, so that

$$\|x\|_{s,2} = \|x^1\|_2 \le \frac{1}{\sqrt{1-\delta}} \|Ax\|_2 + \frac{\delta s^{-1/2}}{1-\delta} \|x\|_1,$$

implying Proposition 1.8.i. Moreover, when $q \ge 2$, $\kappa > 0$ and integer $t \ge 1$ satisfy $t \le s$ and $\kappa t^{1/q-1} \ge \frac{\delta s^{-1/2}}{1-\delta}$, we have

$$||x||_{t,q} \le ||x||_{s,q} \le ||x||_{s,2} \le \frac{1}{\sqrt{1-\delta}} ||Ax||_2 + \kappa t^{1/q-1} ||x||_1,$$

or, equivalently,

$$1 \le t \le \min\left[\left[\frac{\kappa(1-\delta)}{\delta}\right]^{\frac{q}{q-1}}, s^{\frac{q-2}{2q-2}}\right] s^{\frac{q}{2(q-1)}} \Rightarrow \quad (H = \frac{t^{-1/2}}{\sqrt{1-\delta}} I_m, \|\cdot\|_2) \text{ satisfies } \mathbf{Q}_q(t, \kappa),$$

as required in item (i) of Proposition 1.12.

(ii): We have

 $\Rightarrow \begin{array}{l} (1-\delta) \|x^{1}\|_{2}^{2} \leq \|x^{1}\|_{1} \|A^{T}Ax\|_{\infty} + \delta s^{-1/2} \|x^{1}\|_{2} \|x\|_{1} \ [by \ \text{RIP}(\delta, 2s)] \\ \leq s^{1/2} \|x^{1}\|_{2} \|A^{T}Ax\|_{\infty} + \delta s^{-1/2} \|x^{1}\|_{2} \|x\|_{1} \\ \Rightarrow \ \|x\|_{s,2} = \|x^{1}\|_{2} \leq \frac{s^{1/2}}{1-\delta} \|A^{T}Ax\|_{\infty} + \frac{\delta}{1-\delta} s^{-1/2} \|x\|_{1}, \end{array}$

and we see that the pair $\left(H = \frac{1}{1-\delta}A, \|\cdot\|_{\infty}\right)$ satisfies the condition $\mathbf{Q}_2\left(s, \frac{\delta}{1-\delta}\right)$, as required in Proposition 1.8.ii.

In addition, the inequality after the second \Rightarrow implies that

$$\|x^{1}\|_{2} \leq \frac{1}{1-\delta} \left[s^{1/2} \|A^{T}Ax\|_{\infty} + \delta s^{-1/2} \|x\|_{1} \right],$$

Consequently, when $q \ge 2$, $\kappa > 0$ and integer $t \ge 1$ satisfy $t \le s$ and $\kappa t^{1/q-1} \ge \frac{\delta}{1-\delta}s^{-1/2}$, we have

$$\|x\|_{t,q} \le \|x\|_{s,q} \le \|x\|_{s,2} \le \frac{1}{1-\delta} s^{1/2} \|A^T A x\|_{\infty} + \kappa t^{1/q-1} \|x\|_1,$$

or, equivalently,

$$1 \le t \le \min\left[\left[\frac{\kappa(1-\delta)}{\delta}\right]^{\frac{q}{q-1}}, s^{\frac{q-2}{2q-2}}\right]s^{\frac{q}{2(q-1)}} \Rightarrow (H = \frac{s^{\frac{1}{2}}t^{-\frac{1}{q}}}{1-\delta}A, \|\cdot\|_{\infty}) \text{ satisfies } \mathbf{Q}_q(t,\kappa),$$

as required in item (ii) of Proposition 1.12.

1.5.5 Proof of Proposition 1.10

(i): Let $\overline{H} \in \mathbf{R}^{m \times N}$ and $\|\cdot\|$ satisfy $\mathbf{Q}_{\infty}(s,\kappa)$. Then for every $k \leq n$ we have

$$|x_k| \le \|\bar{H}^T A x\| + s^{-1} \kappa \|x\|_1,$$

or, which is the same by homogeneity,

$$\min_{x} \left\{ \|\bar{H}^{T}Ax\| - x_{k} : \|x\|_{1} \le 1 \right\} \ge -s^{-1}\kappa.$$

In other words, the optimal value Opt_k of the conic optimization problem¹¹

$$Opt_k = \min_{x,t} \left\{ t - [e^k]^T x : \|\bar{H}^T A x\| \le t, \|x\|_1 \le 1 \right\},\$$

where $e^k \in \mathbf{R}^n$ is k-th basic orth, is $\geq -s^{-1}\kappa$. Since the problem clearly is strictly feasible, this is the same as to say that the dual problem

$$\max_{\mu \in \mathbf{R}, g \in \mathbf{R}^{n}, \eta \in \mathbf{R}^{N}} \left\{ -\mu : A^{T} \bar{H} \eta + g = e^{k}, \|g\|_{\infty} \le \mu, \|\eta\|_{*} \le 1 \right\},\$$

where $\|\cdot\|_*$ is the norm conjugate to $\|\cdot\|$:

$$\|u\|_* = \max_{\|h\| \le 1} h^T u$$

has a feasible solution with the value of the objective $\geq -s^{-1}\kappa$. It follows that there exists $\eta = \eta^k$ and $g = g^k$ such that

$$\begin{aligned} (a) &: e^{k} = A^{T} h^{k} + g^{k}, \\ (b) &: h^{k} := \bar{H} \eta^{k}, \|\eta^{k}\|_{*} \leq 1, \\ (c) &: \|g^{k}\|_{\infty} \leq s^{-1} \kappa. \end{aligned}$$
 (1.61)

Denoting $H = [h^1, ..., h^n], V = I - H^T A$, we get

$$\operatorname{Col}_k[V^T] = e^k - A^T h^k = g^k,$$

implying that $\|\text{Col}_k[V^T]\|_{\infty} \le s^{-1}\kappa$. Since the latter inequality it is true for all $k \le n$, we conclude that

$$\|\operatorname{Col}_k[V]\|_{s,\infty} = \|\operatorname{Col}_k[V]\|_{\infty} \le s^{-1}\kappa, \ 1 \le k \le n,$$

whence, by Proposition 1.9, $(H, \|\cdot\|_{\infty})$ satisfies $\mathbf{Q}_{\infty}(s, \kappa)$. Moreover, for every $\eta \in \mathbf{R}^m$ and every $k \leq n$ we have, in view of (b) and (c),

$$|[h^{k}]^{T}\eta| = |[\eta^{k}]^{T}\bar{H}^{T}\eta| \le ||\eta^{k}||_{*}||\bar{H}^{T}\eta||_{*}$$

whence $||H^T\eta||_{\infty} \leq ||\bar{H}^T\eta||$.

Now let us prove the "In addition" part of Proposition. Let $H = [h_1, ..., h_n]$ be the contrast matrix specified in this part. We have

$$|[I_m - H^T A]_{ij}| = |[[e^i]^T - h_i^T A]_j| \le ||[e^i]^T - h_i^T A||_{\infty} = ||e^i - A^T h_i||_{\infty} \le \operatorname{Opt}_i,$$

¹¹For summary on conic programming, see Section 4.1.

implying by Proposition 1.9 that $(H, \|\cdot\|_{\infty})$ does satisfy the condition $\mathbf{Q}_{\infty}(s, \kappa_*)$ with $\kappa_* = s \max_i \operatorname{Opt}_i$. Now assume that there exists a matrix H' which, taken along with some norm $\|\cdot\|$, satisfies the condition $\mathbf{Q}_{\infty}(s,\kappa)$ with $\kappa < \kappa_*$, and let us lead this assumption to a contradiction. By the already proved first part of Proposition 1.10, our assumption implies that there exists $m \times n$ matrix $\overline{H} = [\overline{h}_1, ..., \overline{h}_n]$ such that $\|\operatorname{Col}_j[I_n - \overline{H}^T A]\|_{\infty} \leq s^{-1}\kappa$ for all $j \leq n$, implying that $\|[e^i]^T - \overline{h}_i^T A]_j| \leq s^{-1}\kappa$ for all i and j, or, which is the same, $\|e^i - A^T \overline{h}_i\|_{\infty} \leq s^{-1}\kappa$ for all i. Due to the origin of Opt_i , we have $\operatorname{Opt}_i \leq \|e^i - A^T \overline{h}_i\|_{\infty}$ for all i, and we arrive at $s^{-1}\kappa_* = \max_i \operatorname{Opt}_i \leq s^{-1}\kappa$, that is, $\kappa_* \leq \kappa$, which is a desired contradiction.

It remains to prove (1.32), which is just an exercise on LP duality: denoting by **e** *n*-dimensional all-ones vector, we have

$$\begin{aligned} \operatorname{Opt}_{i} &:= \min_{h} \|e^{i} - A^{T}h\|_{\infty} = \min_{h,t} \left\{ t : e^{i} - A^{T}h \leq t\mathbf{e}, A^{T}h - e^{i} \leq t\mathbf{e} \right\} \\ &= \max_{\lambda,\mu} \left\{ \lambda_{i} - \mu_{i} : \lambda, \mu \geq 0, A[\lambda - \mu] = 0, \sum_{i} \lambda_{i} + \sum_{i} \mu_{i} = 1 \right\} \\ & \text{[LP duality]} \\ &= \max_{x:=\lambda-\mu} \left\{ x_{i} : Ax = 0, \|x\|_{1} \leq 1 \right\} \end{aligned}$$

where the concluding equality follows from the fact that vectors x representable as $\lambda - \mu$ with $\lambda, \mu \ge 0$ satisfying $\|\lambda\|_1 + \|\mu\|_1 = 1$ are exactly vectors x with $\|x\|_1 \le 1$. \Box

1.5.6 Proof of Proposition 1.13

Let *H* satisfy (1.38). Since $||v||_{s,1} \leq s^{1-1/q} ||v||_{s,q}$, it follows that *H* satisfies for some $\alpha < 1/2$ the condition

$$\|\operatorname{Col}_{j}[I_{n} - H^{T}A]\|_{s,1} \le \alpha, \ 1 \le j \le n.$$
(1.62)

whence, as we know,

$$||x||_{s,1} \le s ||H^T A x||_{\infty} + \alpha ||x||_1 \,\forall x \in \mathbf{R}^n$$

It follows that $s \leq m$, since otherwise there exists a nonzero s-sparse vector x with Ax = 0; for this x, the inequality above cannot hold true.

Let us set $\bar{n} = 2m$, so that $\bar{n} \leq n$, and let H and A be the $m \times \bar{n}$ matrices comprised of the first 2m columns of H, respectively, A. Relation (1.62) implies that the matrix $V = I_{\bar{n}} - \bar{H}^T \bar{A}$ satisfies

$$\|\operatorname{Col}_{j}[V]\|_{s,1} \le \alpha < 1/2, 1 \le j \le \bar{n}.$$
(1.63)

Now, $V = I_{\bar{n}} - \bar{H}^T \bar{A}$, and since the rank of $\bar{H}^T \bar{A}$ is $\leq m$, at least $\bar{n} - m$ singular values of V are ≥ 1 , and therefore the squared Frobenius norm $||V||_F^2$ of V is at least $\bar{n} - m$. On the other hand, we can upper-bound this squared norm as follows. Observe that for every \bar{n} -dimensional vector f one has

$$\|f\|_{2}^{2} \leq \max\left[\frac{\bar{n}}{s^{2}}, 1\right] \|f\|_{s,1}^{2}.$$
(1.64)

Indeed, by homogeneity it suffices to verify the inequality when $||f||_{s,1} = 1$; besides, we can assume w.l.o.g. that the entries in f are nonnegative, and that $f_1 \ge f_2 \ge \ldots \ge f_{\bar{n}}$. We have $f_s \le ||f||_{s,1}/s = \frac{1}{s}$; in addition,

$$\sum_{j=1}^{s} f_j^2 \le f_s^2 + \max_t \{ \sum_{j=1}^{s-1} t_s^2 : t_j \ge f_s, j \le s-1, \sum_{j=1}^{s-1} t_j = 1-f_s \}.$$

The maximum in the right hand side is the maximum of a convex function over a bounded polytope; it is achieved at an extreme point, that is, at a point where one of the t_j is equal to $1 - (s - 1)f_s$, and all remaining t_j are equal to f_s . As a result,

$$\sum_{j} f_j^2 \le \left[(1 - (s - 1)f_s)^2 + (s - 1)f_s^2 \right] + (\bar{n} - s)f_s^2 \le (1 - (s - 1)f_s)^2 + (\bar{n} - 1)f_s^2.$$

The right hand side in the latter inequality is convex in f_s and thus achieves its maximum over the range [0, 1/s] of allowed values of f_s at an endpoint, yielding $\sum_j f_j^2 \leq \max[1, \bar{n}/s^2]$, as claimed.

Applying (1.64) to the columns of V and recalling that $\bar{n} = 2m$, we get

$$\|V\|_F^2 = \sum_{j=1}^{2m} \|\operatorname{Col}_j[V]\|_2^2 \le \max[1, 2m/s^2] \sum_{j=1}^{2m} \|\operatorname{Col}_j[V]\|_{s,1}^2 \le 2\alpha m \max[1, 2m/s^2].$$

The left hand side in this inequality, as we remember, is $\geq \bar{n} - m = m$, and we arrive at

$$m \le 2\alpha m \max[1, 2m/s^2].$$

Since $\alpha < 1/2$, this inequality implies $2m/s^2 \ge 1$, whence $s \le \sqrt{2m}$.

It remains to prove that when $m \leq n/2$, the condition $\mathbf{Q}_{\infty}(s,\kappa)$ with $\kappa < 1/2$ can be satisfied only when $s \leq \sqrt{2m}$. This is immediate: by Proposition 1.10, assuming $\mathbf{Q}_{\infty}(s,\kappa)$ satisfiable, there exists $m \times n$ contrast matrix H such that $|[I_n - H^T A]_{ij}| \leq \kappa/s$ for all i, j, which, by the already proved part of Proposition 1.13, is impossible when $s > \sqrt{2m}$.

Lecture Two

Hypothesis Testing

Disclaimer for experts. In what follows, we allow for "general" probability and observation spaces, general probability distributions, etc., which, formally, would make it necessary to address the related measurability issues. In order to streamline our exposition, and taking into account that we do not expect from our target audience to be experts in formal nuances of the measure theory, we decided to omit in the text comments (always self-evident for an expert) on measurability and replace them with a "disclaimer" as follows:

Below, unless the opposite is explicitly stated,

- all probability and observation spaces are Polish (complete separable metric) spaces equipped by σ -algebras of Borel sets;
- all random variables (i.e., functions from a probability space to some other space) take values in Polish spaces; these variables, same as other functions we deal with, are Borel;
- all probability distributions we are dealing with are σ-additive Borel measures on the respective probability spaces; the same is true for all reference measures and probability densities taken w.r.t. these measures.

When an entity (a random variable, or a probability density, or a function, say, a test) is part of the data, the Borel property is a default assumption; e.g., the sentence "Let random variable η be a deterministic transformation of random variable ξ " should be read as "let $\eta = f(\xi)$ for some Borel function f", and the sentence "Consider test \mathcal{T} deciding on hypotheses H_1, \ldots, H_L via observation $\omega \in \Omega$ " should be read as "Consider a Borel function \mathcal{T} on Polish space Ω , the values of the function being subsets of the set $\{1, \ldots, L\}$." When an entity is built by us rather than being part of the data, the Borel property is (always straightforwardly verifiable) property of the construction. For example, the statement "The test \mathcal{T} given by... is such that..." should be read as "The test \mathcal{T} given by... is a Borel function of observations and is such that..."

On several occasions, we still use the word "Borel;" those not acquainted with the notion are welcome to just ignore this word.

2.1 PRELIMINARIES FROM STATISTICS: HYPOTHESES, TESTS, RISKS

2.1.1 Hypothesis Testing Problem

Hypothesis Testing is one of the most basic problems of Statistics. Informally, this is the problem where one is given an *observation* – a realization of random variable with unknown (at least partially) probability distribution and want to decide, based on this observation, on two or more hypotheses on the actual distribution of the observed variable. A convenient for us formal setting is as follows:

HYPOTHESIS TESTING

Given are:

- Observation space Ω , where the observed random variable (r.v.) takes its values;
- L families \mathcal{P}_{ℓ} of probability distributions on Ω . We associate with these families L hypotheses $H_1, ..., H_L$, with H_{ℓ} stating that the probability distribution P of the observed r.v. belongs to the family \mathcal{P}_{ℓ} (shorthand: $H_{\ell}: P \in \mathcal{P}_{\ell}$). We shall say that the distributions from \mathcal{P}_{ℓ} obey hypothesis H_{ℓ} .

Hypothesis H_{ℓ} is called *simple*, if \mathcal{P}_{ℓ} is a singleton, and is called *composite* otherwise.

Our goal is, given an observation – a realization ω of the r.v. in question – to decide which one of the hypotheses is true.

2.1.2 Tests

Informally, a *test* is an inference procedure one can use in the above testing problem. Formally, a test for this testing problem is a function $\mathcal{T}(\omega)$ of $\omega \in \Omega$; the value $\mathcal{T}(\omega)$ of this function at a point ω is some subset of the set $\{1, ..., L\}$:

 $\mathcal{T}(\Omega) \subset \{1, \dots, L\}.$

Given observation ω , the test accepts all hypotheses H_{ℓ} with $\ell \in \mathcal{T}(\omega)$ and rejects all hypotheses H_{ℓ} with $\ell \notin \mathcal{T}(\omega)$. We call a test *simple*, if $\mathcal{T}(\omega)$ is a singleton for every ω , that is, whatever be the observation, the test accepts exactly one of the hypotheses $H_1, ..., H_L$ and rejects all other hypotheses.

Note: what we have defined is a *deterministic* test. Sometimes we shall consider also randomized tests, where the set of accepted hypotheses is a (deterministic) function of observation ω and of a realization θ of independent of ω random parameter (which w.l.o.g. can be assumed to be uniformly distributed on [0, 1]). Thus, in a randomized test, the inference depends both on the observation ω and the outcome θ of "flipping a coin," while in a deterministic test the inference depends on observation only. In fact, randomized testing can be reduced to deterministic one. To this end it suffices to pass from our "actual" observation ω to new observation $\omega_{+} = (\omega, \theta)$, where $\theta \sim Uniform[0, 1]$ is independent of ω ; the ω -component of our new observation ω_{+} is, as before, generated by "the nature," and the θ -component is generated by ourselves. Now, given families \mathcal{P}_{ℓ} , $1 \leq \ell \leq L$, of probability distributions on the original observation space Ω , we can associate with them families $\mathcal{P}_{\ell,+} = \{P \times Uniform[0,1] : P \in \mathcal{P}_{\ell}\}$ of probability distributions on our new observation space $\Omega_+ = \Omega \times [0,1]$; clearly, to decide on the hypotheses associated with the families \mathcal{P}_{ℓ} via observation ω is the same as to decide on the hypotheses associated with the families $\mathcal{P}_{\ell,+}$ of our new observation ω_+ , and deterministic tests for the latter testing problem are exactly the randomized tests for the former one.

2.1.3 Testing from repeated observations

There are situations where an inference can be based on several observations $\omega_1, ..., \omega_K$ rather than on a single one. Our related setup is as follows:

We are given L families \mathcal{P}_{ℓ} , $\ell = 1, ..., L$, of probability distributions on

observation space Ω and a collection

$$\omega^{K} = (\omega_{1}, ..., \omega_{K}) \in \Omega^{K} = \underbrace{\Omega \times ... \times \Omega}_{K}$$

and want to make conclusions on how the distribution of ω^K "is positioned" w.r.t. the families \mathcal{P}_{ℓ} , $1 \leq \ell \leq L$.

Specifically, we are interested in three situations of this type, specifically, as follows.

2.1.3.1 Stationary K-repeated observations

In the case of stationary K-repeated observations $\omega_1, ..., \omega_K$ are independently of each other drawn from a distribution P. Our goal is to decide, given ω^K , on the hypotheses $P \in \mathcal{P}_{\ell}, \ell = 1, ..., L$.

Equivalently: Families \mathcal{P}_{ℓ} of probability distributions of $\omega \in \Omega$, $1 \leq \ell \leq L$, give rise to the families

$$\mathcal{P}_{\ell}^{\odot,K} = \{ P^{K} = \underbrace{P \times \dots \times P}_{K} : P \in \mathcal{P}_{\ell} \}$$

of probability distributions on Ω^K ; we refer to the families $\mathcal{P}_{\ell}^{\odot,K}$ as to *K*-th diagonal powers of the family \mathcal{P}_{ℓ} . Given observation $\omega^K \in \Omega^K$, we want to decide on the hypotheses

$$H^{\odot,K}_{\ell}: \omega^K \sim P^K \in \mathcal{P}^{\odot,K}_{\ell}, \ 1 \leq \ell \leq L.$$

2.1.3.2 Semi-stationary K-repeated observations

In the case of semi-stationary K-repeated observations, "the nature" selects somehow a sequence $P_1, ..., P_K$ of distributions on Ω , and then draws, *independently across* k, observations ω_k , k = 1, ..., K, from these distributions:

 $\omega_k \sim P_k$ are independent across $k \leq K$

Our goal is to decide, given $\omega^K = (\omega_1, ..., \omega_K)$, on the hypotheses $\{P_k \in \mathcal{P}_\ell, 1 \leq k \leq K\}, \ell = 1, ..., L$.

Equivalently: Families \mathcal{P}_{ℓ} of probability distributions of $\omega \in \Omega$, $1 \leq \ell \leq L$, give rise to the families

$$\mathcal{P}_{\ell}^{\oplus,K} = \{ P^K = P_1 \times \dots \times P_K : P_k \in \mathcal{P}_{\ell}, 1 \le k \le K \}$$

of probability distributions on Ω^K . Given observation $\omega^K \in \Omega^K$, we want to decide on the hypotheses

$$H_{\ell}^{\oplus,K}: \omega^K \sim P^K \in \mathcal{P}_{\ell}^{\oplus,K}, \ 1 \le \ell \le L.$$

In the sequel, we refer to families $\mathcal{P}_{\ell}^{\oplus,K}$ as to *K*-th direct powers of the families \mathcal{P}_{ℓ} .

HYPOTHESIS TESTING

A closely related notion is the one of *direct product*

$$\mathcal{P}_{\ell}^{\oplus,K} = \bigoplus_{k=1}^{K} \mathcal{P}_{\ell,k}$$

of K families $\mathcal{P}_{\ell,k}$, of probability distributions on Ω_k , over k = 1, ..., K. By definition,

$$\mathcal{P}_{\ell}^{\oplus, \kappa} = \{ P^{\kappa} = P_1 \times \dots \times P_K : P_k \in \mathcal{P}_{\ell, k}, 1 \le k \le K \}.$$

2.1.3.3 Quasi-stationary K-repeated observations

Quasi-stationary K-repeated observations $\omega_1 \in \Omega, ..., \omega_K \in \Omega$ stemming from a family \mathcal{P} of probability distributions on an observation space Ω are generated as follows:

"In the nature" there exists random sequence $\zeta^{K} = (\zeta_{1}, ..., \zeta_{K})$ of "driving factors" such that for every k, ω_{k} is a deterministic function of $\zeta_{1}, ..., \zeta_{k}$:

$$\omega_k = \theta_k(\zeta_1, ..., \zeta_k)$$

and the conditional, $\zeta_1, ..., \zeta_{k-1}$ given, distribution $P_{\omega_k|\zeta_1,...,\zeta_{k-1}}$ of ω_k always (i.e., for all $\zeta_1, ..., \zeta_{k-1}$) belongs to \mathcal{P} .

With the above mechanism, the collection $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ has some distribution P^{K} which depends on the distribution of driving factors and on functions $\theta_{k}(\cdot)$. We denote by $\mathcal{P}^{\otimes,K}$ the family of all distributions P^{K} which can be obtained in this fashion and we refer to random observations ω^{K} with distribution P^{K} of the just define type as to quasi-stationary K-repeated observations stemming from \mathcal{P} . The quasi-stationary version of our hypothesis testing problem reads: Given L families \mathcal{P}_{ℓ} of probability distributions \mathcal{P}_{ℓ} , $\ell = 1, ..., L$, on Ω and an observation $\omega^{K} \in \Omega^{K}$, decide on the hypotheses

$$H_{\ell}^{\otimes,K} = \{ P^K \in \mathcal{P}_{\ell}^{\otimes,K} \}, \, 1 \le \ell \le K$$

on the distribution P^K of the observation ω^K .

A closely related notion is the one of quasi-direct product

$$\mathcal{P}_{\ell}^{\otimes,K} = \bigotimes_{k=1}^{K} \mathcal{P}_{\ell,k}$$

of K families $\mathcal{P}_{\ell,k}$, of probability distributions on Ω_k , over k = 1, ..., K. By definition, $\mathcal{P}_{\ell}^{\otimes,K}$ is comprised of all probability distributions of random sequences $\omega^K = (\omega_1, ..., \omega_K), \ \omega_k \in \Omega_k$, which can be generated as follows: "in the nature" there exists a random sequence $\zeta^K = (\zeta_1, ..., \zeta_K)$ of "driving factors" such that for every $k \leq K, \ \omega_k$ is a deterministic function of $\zeta^k = (\zeta_1, ..., \zeta_k)$, and conditional, ζ^{k-1} being given, distribution of ω_k always belongs to $\mathcal{P}_{\ell,k}$.

The above description of quasi-stationary K-repeated observations seems to be too complicated; well, this is what happens in some important applications, e.g., in hidden Markov chain. Here $\Omega = \{1, ..., d\}$ is a finite set, and $\omega_k \in \Omega$, k = 1, 2, ..., are generated as follows: "in the nature there" exists a Markov chain with D-element state space S split into d non-overlapping bins, and ω_k is the serial number $\beta(\eta)$

of the bin to which the state η_k of the chain belongs. Now, every column Q^j of the transition matrix Q of the chain (this column is a probability distribution on $\{1, ..., D\}$) generates a probability distribution P_j on Ω , specifically, the distribution of $\beta(\eta), \eta \sim Q^j$. Now, a family \mathcal{P} of distributions on Ω induces a family $\mathcal{Q}[\mathcal{P}]$ of all $D \times D$ stochastic matrices Q for which all D distributions $P^j, j = 1, ..., D$, belong to \mathcal{P} . When $Q \in \mathcal{Q}[\mathcal{P}]$, observations $\omega_k, k = 1, 2, ...$ clearly are given by the above "quasi-stationary mechanism" with η_k in the role of driving factors and \mathcal{P} in the role of \mathcal{P}_{ℓ} . Thus, in the situation in question, given L families $\mathcal{P}_{\ell}, \ell = 1, ..., L$ of probability distributions on \mathcal{S} , deciding on hypotheses $Q \in \mathcal{Q}[\mathcal{P}_{\ell}], \ell = 1, ..., L$, on the transition matrix Q of the Markov chain underlying our observations reduces to hypothesis testing via quasi-stationary K-repeated observations.

2.1.4 Risk of a simple test

Let \mathcal{P}_{ℓ} , $\ell = 1, ..., L$, be families of probability distributions on observation space Ω ; these families give rise to hypotheses

$$H_{\ell}: P \in \mathcal{P}_{\ell}, \ \ell = 1, \dots, L$$

on the distribution P of a random observation $\omega \sim P$. We are about to define the risks of a simple test \mathcal{T} deciding on the hypotheses H_{ℓ} , $\ell = 1, ..., L$, via observation ω ; recall that simplicity means that as applied to an observation, our test accepts exactly one hypothesis and rejects all other hypotheses.

Partial risks Risk_{ℓ}($\mathcal{T}|H_1, ..., H_L$) are the worst-case, over $P \in \mathcal{P}_\ell$, *P*-probabilities for \mathcal{T} to reject ℓ -th hypothesis when it is true, that is, when $\omega \sim P$:

$$\operatorname{Risk}_{\ell}(\mathcal{T}|H_1, ..., H_L) = \sup_{P \in \mathcal{P}_{\ell}} \operatorname{Prob}_{\omega \sim P} \left\{ \omega : \mathcal{T}(\omega) \neq \{\ell\} \right\}, \ \ell = 1, ..., L.$$

Note that for ℓ fixed, ℓ -th partial risk depends on how we order the hypotheses; when reordering them, we should reorder risks as well. In particular, for a test \mathcal{T} deciding on two hypotheses H, H' we have

$$\operatorname{Risk}_1(\mathcal{T}|H, H') = \operatorname{Risk}_2(\mathcal{T}|H', H).$$

Total risk $\operatorname{Risk}_{\operatorname{tot}}(\mathcal{T}|H_1, ..., H_L)$ is the sum of all L partial risks:

$$\operatorname{Risk}_{\operatorname{tot}}(\mathcal{T}|H_1,...,H_L) = \sum_{\ell=1}^L \operatorname{Risk}_{\ell}(\mathcal{T}|H_1,...,H_L).$$

Risk Risk $(\mathcal{T}|H_1, ..., H_L)$ is the maximum of all L partial risks:

$$\operatorname{Risk}(\mathcal{T}|H_1,...,H_L) = \max_{1 \le \ell \le L} \operatorname{Risk}_{\ell}(\mathcal{T}|H_1,...,H_L).$$

Note that at the first glance, we have defined risks for single-observation tests only; in fact, we have defined them for tests based on stationary, semi-stationary, and

HYPOTHESIS TESTING

quasi-stationary K-repeated observations as well, since, as we remember from Section 2.1.3, the corresponding testing problems, after redefining observations and families of probability distributions (ω^{K} in the role of ω and, say, $\mathcal{P}_{\ell}^{\oplus,K} = \bigoplus_{k=1}^{K} \mathcal{P}_{\ell}$ in the role of \mathcal{P}_{ℓ}), become single-observation testing problems.

Pay attention to the following two important observations:

- Partial risks of a simple test are defined in the worst-case-oriented fashion: as the worst, over the true distributions P of observations compatible with the hypothesis in question, probability to reject this hypothesis
- Risks of a simple test say what happens, statistically speaking, when the true distribution P of observation obeys one of the hypotheses in question, and say nothing on what happens when P does not obey neither one of the L hypotheses.

Remark 2.1. "The smaller are hypotheses, the less are risks." Specifically, given families of probability distributions $\mathcal{P}_{\ell} \subset \mathcal{P}'_{\ell}, \ell = 1, ..., L$, on observation space Ω , along with hypotheses $H_{\ell} : P \in \mathcal{P}_{\ell}, H'_{\ell} : P \in \mathcal{P}'_{\ell}$ on the distribution P of an observation $\omega \in \Omega$, every test \mathcal{T} deciding on the "larger" hypotheses $H'_1, ..., H'_L$ can be considered as a test deciding on smaller hypotheses $H_1, ..., H_L$ as well, and the risks of the test when passing from larger hypotheses to smaller ones can only drop down:

$$\mathcal{P}_{\ell} \subset \mathcal{P}'_{\ell}, 1 \leq \ell \leq L \Rightarrow \operatorname{Risk}(\mathcal{T}|H_1, ..., H_L) \leq \operatorname{Risk}(\mathcal{T}|H'_1, ..., H'_L).$$

For example, families of probability distributions \mathcal{P}_{ℓ} , $1 \leq \ell \leq L$, on Ω and a positive integer K induce three families of hypotheses on a distribution P^{K} of K-repeated observations:

$$H_{\ell}^{\odot,K}K: P^{K} \in \mathcal{P}_{\ell}^{\odot,K}, H_{\ell}^{\oplus,K}: P^{K} \in \mathcal{P}_{\ell}^{\oplus,K} = \bigoplus_{k=1}^{K} \mathcal{P}_{\ell},$$
$$H_{\ell}^{\otimes,K}: P^{K} \in \mathcal{P}_{\ell}^{\otimes,K} = \bigotimes_{k=1}^{K} \mathcal{P}_{\ell}, 1 \le \ell \le L,$$

(see Section 2.1.3), and clearly

$$\mathcal{P}_{\ell}^{K} \subset \mathcal{P}_{\ell}^{\oplus,K} \subset \mathcal{P}_{\ell}^{\otimes,K};$$

it follows that when passing from quasi-stationary K-repeated observations to semistationary K-repeated, and then to stationary K-repeated observations, the risks of a test can only go down.

2.1.5 Two-point lower risk bound

The following observation is nearly evident:

Proposition 2.2. Consider two simple hypotheses $H_1 : P = P_1$ and $H_2 : P = P_2$ on the distribution P of observation $\omega \in \Omega$, and assume that P_1 , P_2 have densities p_1 , p_2 w.r.t. some reference measure Π on Ω ¹². Then for any simple test \mathcal{T}

¹² This assumption is w.l.o.g. – we can take, as Π , the sum of the measures P_1 and P_2 .

deciding on H_1, H_2 it holds

$$\operatorname{Risk}_{\operatorname{tot}}(\mathcal{T}|H_1, H_2) \ge \int_{\Omega} \min[p_1(\omega), p_2(\omega)] \Pi(d\omega).$$
(2.1)

Note that the right hand side in this relation is independent of how Π is selected.

Proof. Consider a simple test \mathcal{T} , perhaps a randomized one, and let $\pi(\omega)$ be the probability for this test to accept H_1 and reject H_2 when the observation is ω ; since the test is simple, the probability for \mathcal{T} to accept H_2 and to reject H_1 , observation being ω , is $1 - \pi(\omega)$. Consequently,

$$\begin{aligned} \operatorname{Risk}_{1}(\mathcal{T}|H_{1}, H_{2}) &= \int_{\Omega} (1 - \pi(\omega)) p_{1}(\omega) \Pi(d\omega), \\ \operatorname{Risk}_{2}(\mathcal{T}|H_{1}, H_{2}) &= \int_{\Omega} \pi(\omega) p_{2}(\omega) \Pi(d\omega), \end{aligned}$$

whence

$$\operatorname{Risk_{tot}}(\mathcal{T}|H_1, H_2) = \int_{\Omega} [(1 - \pi(\omega))p_1(\omega) + \pi(\omega)p_2(\omega)]\Pi(d\omega) \\ \geq \int_{\Omega} \min[p_1(\omega), p_2(\omega)]\Pi(d\omega). \square$$

Remark 2.3. Note that the lower risk bound (2.1) is achievable; given an observation ω , the corresponding test \mathcal{T} accepts H_1 with probability 1 (i.e., $\pi(\omega) = 1$ when $p_1(\omega) > p_2(\omega)$), accepts H_2 when $p_1(\omega) < p_2(\omega)$ (i.e., $\pi(\omega) = 0$ when $p_1(\omega) < p_2(\omega)$) and accepts H_1 and H_2 with probabilities 1/2 in the case of tie (i.e., $\pi(\omega) = 1/2$ when $p_1(\omega) = p_2(\omega)$); this is nothing but maximum likelihood test naturally adjusted to account for ties.

Example 2.4. Let $\Omega = \mathbf{R}^d$, let the reference measure Π be the Lebesgue measure on \mathbf{R}^d , and let $p_{\chi}(\cdot) = \mathcal{N}(\mu_{\chi}, I_d)$, be the Gaussian densities on \mathbf{R}^d with unit covariance and means $\mu_{\chi}, \chi = 1, 2$. In this case, assuming $\mu_1 \neq \mu_2$, the recipe from Remark 2.3 reduces to the following:

Let

$$\phi_{1,2}(\omega) = \frac{1}{2} [\mu_1 - \mu_2]^T [\omega - w], \ w = \frac{1}{2} [\mu_1 + \mu_2].$$
(2.2)

Consider the simple test \mathcal{T} which, given an observation ω , accepts $H_1: p = p_1$ and rejects $H_2: p = p_2$ when $\phi_{1,2}(\omega) \geq 0$, otherwise accepts H_2 and rejects H_1 . For this test,

$$\operatorname{Risk}_{1}(\mathcal{T}|H_{1}, H_{2}) = \operatorname{Risk}_{2}(\mathcal{T}|H_{1}, H_{2}) = \operatorname{Risk}(\mathcal{T}|H_{1}, H_{2})$$

= $\frac{1}{2}\operatorname{Risk}_{\operatorname{tot}}(\mathcal{T}|H_{1}, H_{2}) = \operatorname{Erf}(\frac{1}{2}\|\mu_{1} - \mu_{2}\|_{2}),$ (2.3)

where

$$\operatorname{Erf}(\delta) = \frac{1}{\sqrt{2\pi}} \int_{\delta}^{\infty} e^{-s^2/2} ds \qquad (2.4)$$

is the error function, and the test is optimal in terms of its risk and its total risk.

Note that optimality of \mathcal{T} in terms of total risk is given by Proposition 2.2 and Remark 2.3; optimality in terms of risk is ensured by optimality in terms of total risk combined with the first equality in (2.3).

Example 2.4 admits an immediate and useful extension:

HYPOTHESIS TESTING



Figure 2.1: "Gaussian Separation" (Example 2.5): Optimal test deciding on whether the mean of Gaussian r.v. belongs to the dark red (H_1) or to the dark blue (H_2) domains. Dark and light red: acceptance domain for H_1 . Dark and light blue: acceptance domain for H_2 .

Example 2.5. Let $\Omega = \mathbf{R}^d$, let the reference measure Π be the Lebesgue measure on \mathbf{R}^d , and let M_1 , M_2 be two nonempty closed convex sets in \mathbf{R}^d with empty intersection and such that the convex optimization program

$$\min_{\mu_1,\mu_2} \left\{ \|\mu_1 - \mu_2\|_2 : \mu_\chi \in M_\chi, \ \chi = 1,2 \right\}$$
(*)

has an optimal solution μ_1^*, μ_2^* (this definitely is the case when at least one of the sets M_1, M_2 is bounded). Let

$$\phi_{1,2}(\omega) = \frac{1}{2} [\mu_1^* - \mu_2^*]^T [\omega - w], \ w = \frac{1}{2} [\mu_1^* + \mu_2^*], \tag{2.5}$$

and let the simple test \mathcal{T} deciding on the hypotheses

$$H_1: p = \mathcal{N}(\mu, I_d)$$
 with $\mu \in M_1$, $H_2: p = \mathcal{N}(\mu, I_d)$ with $\mu \in M_2$

be as follows (see Figure 2.1): given an observation ω , \mathcal{T} accepts H_1 and rejects H_2 when $\phi_{1,2}(\omega) \geq 0$, otherwise accepts H_2 and rejects H_1 . Then

$$\operatorname{Risk}_{1}(\mathcal{T}|H_{1}, H_{2}) = \operatorname{Risk}_{2}(\mathcal{T}|H_{1}, H_{2}) = \operatorname{Risk}(\mathcal{T}|H_{1}, H_{2}) = \frac{1}{2}\operatorname{Risk}_{\operatorname{tot}}(\mathcal{T}|H_{1}, H_{2}) = \operatorname{Erf}(\frac{1}{2}\|\mu_{1}^{*} - \mu_{2}^{*}\|_{2}),$$
(2.6)

and the test is optimal in terms of its risk and its total risk.

Justification of Example 2.5 is immediate. Let e be the $\|\cdot\|_2$ -unit vector with

the same direction as the one of $\mu_1^* - \mu_2^*$, and let $\xi[\omega] = e^T(\omega - w)$. From optimality conditions for (*) it follows that

$$e^T \mu \ge e^T \mu_1^* \ \forall \mu \in M_1 \ \& \ e^T \mu \le e^T \mu_2^* \ \forall \mu \in M_2.$$

As a result, if $\mu \in M_1$ and the density of ω is $p_{\mu} = \mathcal{N}(\mu, I_d)$, the random variable $\xi[\omega]$ is scalar Gaussian random variable with unit variance and expectation $\geq \delta := \frac{1}{2} \|\mu_1^* - \mu_2^*\|_2$, implying that p_{μ} -probability for $\xi[\omega]$ to be negative (which is exactly the same as the p_{μ} -probability for \mathcal{T} to reject H_1 and accept H_2) is at most $\mathrm{Erf}(\delta)$. Similarly, when $\mu \in M_2$ and the density of ω is $p_{\mu} = \mathcal{N}(\mu, I_d)$, $\xi[\omega]$ is scalar Gaussian random variable with unit variance and expectation $\leq -\delta$, implying that the p_{μ} -probability for $\xi[\omega]$ to be nonnegative (which is exactly the same as the probability for \mathcal{T} to reject H_2 and accept H_1) is at most $\mathrm{Erf}(\delta)$. These observations imply the validity of (2.6); optimality in terms of risks follows from the fact that risks of a simple test deciding on our now – composite – hypotheses H_1 , H_1 on the density p of observation ω can be only larger than the risks of a simple test deciding on our now – composite – hypotheses H_1 , H_1 on the density p of observation ω can be only larger than the risk of a simple test deciding on H_1, H_2 ; with this in mind, the announced optimalities of \mathcal{T} in terms of risks are immediate consequences of (2.6).

We remark that the (nearly self-evident) result stated in Example 2.5 seems first been noticed in [27].

Example 2.5 allows for substantial extensions in two directions: first. it turns out that the "Euclidean separation" underlying the test built in this example can be used to decide on hypotheses on location of a "center" of *d*-dimensional distribution far beyond the Gaussian observation model considered in this example; this extension will be our goal in the next Section, based on recent paper [70]. A less straightforward and, we believe, more instructive extensions, originating from [64], will be considered in Section 2.3.

2.2 HYPOTHESIS TESTING VIA EUCLIDEAN SEPARATION

2.2.1 Situation

In this section, we will be interested in testing hypotheses

$$H_{\ell}: P \in \mathcal{P}_{\ell}, \ell = 1, ..., L$$
 (2.7)

on the probability distribution of a random observation ω in the situation where the families of distributions \mathcal{P}_{ℓ} are obtained from the probability distributions from a given family \mathcal{P} by shifts. Specifically, we are given

• A family \mathcal{P} of probability distributions on $\Omega = \mathbf{R}^d$ such that all distributions from \mathcal{P} possess densities with respect to the Lebesgue measure on \mathbf{R}^n , and these densities are even functions on \mathbf{R}^{d-13} ;

¹³Allowing for a slight abuse of notation, we write $P \in \mathcal{P}$, where P is a probability distribution, to express the fact that P belongs to \mathcal{P} (no abuse of notation so far), and write $p(\cdot) \in \mathcal{P}$ (this is the abuse of notation), where $p(\cdot)$ is the density of a probability distribution P, to express

HYPOTHESIS TESTING

• A collection $X_1, ..., X_L$ of nonempty closed and convex subsets of \mathbf{R}^d , with at most one of the sets unbounded.

These data specify L families \mathcal{P}_{ℓ} of distributions on \mathbf{R}^d ; \mathcal{P}_{ℓ} is comprised of distributions of random vectors of the form $x + \xi$, where $x \in X_{\ell}$ is deterministic, and ξ is random with distribution from \mathcal{P} . Note that with this setup, deciding upon hypotheses (2.7) via observation $\omega \sim P$ is exactly the same as to decide, given observation

$$\omega = x + \xi, \tag{2.8}$$

where x is a deterministic "signal" and ξ is random noise with distribution P known to belong to \mathcal{P} , on the "position" of x w.r.t. $X_1, ..., X_L$; ℓ -th hypothesis H_{ℓ} merely says that $x \in H_{\ell}$. The latter allows us to write down ℓ -th hypothesis as $H_{\ell} : x \in X_{\ell}$ (of course, this shorthand makes sense only within the scope of our current "signal plus noise" setup).

2.2.2 Pairwise Hypothesis Testing via Euclidean Separation

2.2.2.1 The simplest case

Consider nearly the simplest case of the situation from Section 2.2.1, one with L = 2, $X_1 = \{x^1\}$ and $X_2 = \{x^2\}$, $x^1 \neq x^2$, are singletons, and \mathcal{P} also is a singleton; moreover, the probability density of the only distribution from \mathcal{P} is of the form

 $p(u) = f(||u||_2), \ f(\cdot)$ is a strictly monotonically increasing function on the nonnegative ray. (2.9)

This situation is a generalization of the one considered in Example 2.4, where we dealt with the special case of f, namely, with

$$p(u) = (2\pi)^{-d/2} e^{-u^T u/2}.$$

In the case in question our goal is to decide on two simple hypotheses $H_{\chi} : p(u) = f(||u - x^{\chi}||_2), \ \chi = 1, 2$, on the density of observation (2.8). Let us set

$$\delta = \frac{1}{2} \|x^1 - x^2\|_2, \ e = \frac{x^1 - x^2}{\|x^1 - x^2\|_2}, \ \phi(\omega) = e^T \omega - \underbrace{\frac{1}{2} e^T [x^1 + x^2]}_{\underbrace{2}}, \tag{2.10}$$

and consider the test \mathcal{T} which, given observation $\omega = x + \xi$, accepts the hypothesis $H_1: x = x^1$ when $\phi(\omega) \ge 0$, and accepts the hypothesis $H_2: x = x^2$ otherwise.

the fact that $P \in \mathcal{P}$.



We have (cf. Example 2.4)

$$\operatorname{Risk}_{1}(\mathcal{T}|H_{1}, H_{2}) = \int_{\substack{\omega:\phi(\omega)<0\\ \\ = \int_{\substack{\omega:\phi(\omega)\geq 0\\ \\ \omega:\phi(\omega)\geq 0}}} p_{1}(\omega)d\omega = \int_{\substack{u:e^{T}u\geq\delta\\ \\ u:e^{T}u\geq\delta}} f(\|u\|_{2})du$$

Since p(u) is strictly decreasing function of $||u||_2$, we have also

$$\min[p_1(u), p_2(u)] = \begin{cases} p_1(u), & \phi(u) \ge 0\\ p_2(u), & \phi(u) \le 0 \end{cases},$$

whence

$$\operatorname{Risk}_{1}(\mathcal{T}|H_{1}, H_{2}) + \operatorname{Risk}_{2}(\mathcal{T}|H_{1}, H_{2}) = \int_{\substack{\omega:\phi(\omega)<0\\ \prod\\ \mathbf{R}^{d}}} p_{1}(\omega)d\omega + \int_{\substack{\omega:\phi(\omega)\geq0\\ \omega:\phi(\omega)\geq0}} p_{2}(\omega)d\omega$$

Invoking Proposition 2.2, we conclude that the test \mathcal{T} is the minimum risk simple test deciding on H_1 , H_2 , and the risk of this test is

$$\operatorname{Risk}(\mathcal{T}|H_1, H_2) = \int_{u:e^T u \ge \delta} f(\|u\|_2) du.$$
(2.11)

2.2.2.2 Extension

Now consider a slightly more complicated case of the situation from Section 2.2.1, the one with L = 2 and nonempty and nonintersecting closed convex sets X_1, X_2 , one of the sets being bounded; as about \mathcal{P} , we still assume that it is a singleton, and the density of the only distribution from \mathcal{P} is of the form 2.9. Our now situation is an extension of the one from Example 2.5. By the same reasons as in the case of the latter Example, with X_1, X_2 as above, the convex minimization problem

$$Opt = \min_{x^1 \in X_1, x^2 \in X_2} \frac{1}{2} \|x^1 - x^2\|_2$$
(2.12)
is solvable, and denoting by (x_*^1, x_*^2) an optimal solution and setting

$$\phi(\omega) = e^T \omega - c, \ e = \frac{x_*^1 - x_*^2}{\|x_*^1 - x_*^2\|_2}, \ c = \frac{1}{2} e^T [x_*^1 + x_*^2]$$
(2.13)

the stripe $\{\omega : -\text{Opt} \le \phi(x) \le \text{Opt}\}$ separates X_1 and X_2 :



Proposition 2.6. Let X_1, X_2 be nonempty and nonintersecting closed convex sets in \mathbb{R}^d , one of the sets being bounded. With Opt and $\phi(\cdot)$ given by (2.12) – (2.13), let us split the width 2Opt of the stripe { $\omega : -Opt \le \phi(\omega) \le Opt$ } separating X_1 and X_2 into two nonnegative parts:

$$\delta_1 \ge 0, \delta_2 \ge 0, \, \delta_1 + \delta_2 = 2\text{Opt} \tag{2.15}$$

and consider simple test \mathcal{T} deciding on the hypotheses $H_1: x \in X_1, H_2: x \in X_2$ via observation (2.8) by accepting H_1 when

$$\phi(\omega) \ge \frac{1}{2} [\delta_2 - \delta_1]$$

and accepting H_2 otherwise. Then

$$\operatorname{Risk}_{\chi}(\mathcal{T}|H_1, H_2) \le \int_{\delta_{\chi}}^{\infty} \gamma(s) ds, \ \chi = 1, 2,$$
(2.16)

where $\gamma(\cdot)$ is the univariate marginal density of ξ , that is, probability density of the scalar random variable $h^T\xi$, where $\|h\|_2 = 1$ (note that due to (2.9), $\gamma(\cdot)$ is independent of how we select h with $\|h\|_2 = 1$).

In addition, when $\delta_1 = \delta_2 = \text{Opt}$, \mathcal{T} is the minimum risk test deciding on H_1 ,

 H_2 . The risk of this test is

$$\operatorname{Risk}(\mathcal{T}|H_1, H_2) = \int_{O_{pt}}^{\infty} \gamma(s) ds.$$
(2.17)

Proof. By (2.9) and (2.14), for $x \in X_1$ we have (see the picture above):

$$\operatorname{Prob}_{\xi \sim p(\cdot)} \left\{ \phi(x+\xi) < \frac{1}{2} [\delta_2 - \delta_1] \right\} \leq \operatorname{Prob}_{\xi \sim p(\cdot)} \left\{ [-e]^T \xi \ge \delta_1 \right\} = \int_{\delta_1}^{\infty} \gamma(s) ds;$$

by "symmetric" reasoning, for $x \in X_2$ we have

$$\operatorname{Prob}_{\xi \sim p(\cdot)} \left\{ \phi(x+\xi) \ge \frac{1}{2} [\delta_2 - \delta_1] \right\}] \le \operatorname{Prob}_{\xi \sim p(\cdot)} \left\{ e^T \xi \ge \delta_2 \right\} = \int_{\delta_2}^{\infty} \gamma(s) ds,$$

and we arrive at (2.16). The fact that in the case of $\delta_1 = \delta_2 = \text{Opt}$ our test \mathcal{T} becomes the minimum risk test deciding on composite hypotheses H_1, H_2 , same as (2.16), are readily given by the fact that due to the analysis in Section 2.2.2.1, the minimal, over all possible tests, risk of deciding on two simple hypotheses $H'_1 : x = x^1_*, H'_2 : x = x^2_*$ is given by (2.11), that is, is equal to $\int_{\text{Opt}}^{\infty} \gamma(s) ds$ (note that e in (2.11) by construction is a $\|\cdot\|_2$ -unit vector), that is, it is equal to the already justified upper bound (2.17) on the risk of the test \mathcal{T} deciding on the larger than H'_{χ} composite hypotheses $H_{\chi}, \chi = 1, 2$.

2.2.2.3 Further extensions: spherical families of distributions

Now let us assume that we are in the situation of Section 2.2.1 with L = 2and nonempty closed, convex and non-intersecting X_1 , X_2 , one of the sets being bounded, exactly what we have assumed in Section 2.2.2.2. What we intend to do now, is to relax the restrictions on the family \mathcal{P} of noise distributions, which in Section 2.2.2.2 was just a singleton with density which is a strictly decreasing function of the $\|\cdot\|_2$ -norm. Observe that as far as the density $p(\cdot)$ of noise is concerned, justification of the upper risk bound (2.16) in Proposition 2.6 used the only fact that whenever $h \in \mathbf{R}^d$ is a $\|\cdot\|_2$ -unit vector and $\delta \geq 0$, we have $\int_{h^T u \geq \delta} p(u) du \leq \int_{\delta}^{\infty} \gamma(s) ds$, with the even univariate probability density $\gamma(\cdot)$ specified in Proposition. We use this observation to extend our construction to *spherical families of probability densities*.

2.2.2.3.A. Spherical families of probability densities. Let $\gamma(\cdot)$ be an even probability density on the axis such that there is no neighbourhood of the origin where $\gamma = 0$ almost surely. We associate with γ the spherical family of densities \mathcal{P}^d_{γ} comprised of all probability densities $p(\cdot)$ on \mathbf{R}^d such that

A. $p(\cdot)$ is even **B.** Whenever $e \in \mathbf{R}^d$, $||e||_2 = 1$, and $\delta \ge 0$, we have

$$\operatorname{Prob}_{\xi \sim P} \{\xi : e^T \xi \ge \delta\} \le P_{\gamma}(\delta) := \int_{\delta}^{\infty} \gamma(s) ds.$$
(2.18)

Geometrically: $p(\cdot)$ -probability for $\xi \sim p(\cdot)$ to belong to a half-space not containing origin does not exceed $P_{\gamma}(\delta)$, where δ is the $\|\cdot\|_2$ -distance from the origin to the half-space.

Note that density (2.9) belongs to the family \mathcal{P}^d_{γ} with $\gamma(\cdot)$ defined in Proposition 2.6; the resulting γ , in addition to being an even density, is strictly monotonically decreasing on the nonnegative ray. When speaking about general-type spherical families \mathcal{P}^d_{γ} , we do not impose monotonicity requirements on $\gamma(\cdot)$. If a spherical family \mathcal{P}^d_{γ} includes a density $p(\cdot)$ of the form (2.9) such that $\gamma(\cdot)$ is the induced by $p(\cdot)$ univariate marginal density, as in Proposition 2.6, we say that \mathcal{P}^d_{γ} has a cap, and this cap is $p(\cdot)$.

2.2.2.3.B. Example: Gaussian mixtures. Let $\eta \sim \mathcal{N}(0,\Theta)$, where the $d \times d$ covariance matrix Θ satisfies $\Theta \leq I_d$, and let Z be an independent of η positive scalar random variable. Gaussian mixture of Z and η (or, better to say, of the distribution P_Z of Z and the distribution $\mathcal{N}(0,\Theta)$) is the probability distribution of the random vector $\xi = \sqrt{Z\eta}$. Examples of Gaussian mixtures include

- Gaussian distribution $\mathcal{N}(0,\Theta)$ (take Z identically equal to 1),
- multidimensional Student's t-distribution with $\nu \in \{1, 2, ...\}$ degrees of freedom and "covariance structure" Θ ; here Z is given by the requirement that ν/Z has χ^2 -distribution with ν degrees of freedom.

An immediate observation (see Exercise 2.55) is that with γ given by the distribution P_Z of Z according to

$$\gamma_Z(s) = \int_{z>0} \frac{1}{\sqrt{2\pi z}} e^{-\frac{s^2}{2z}} P_Z(dz), \qquad (2.19)$$

the distribution of random variable $\sqrt{Z\eta}$, with $\eta \sim \mathcal{N}(0,\Theta)$, $\Theta \leq I_d$, independent of Z, belongs to the family $\mathcal{P}_{\gamma_Z}^d$, and the family $\mathcal{P}_{\gamma_Z}^d$ has a cap, specifically, the Gaussian mixture of P_Z and $\mathcal{N}(0, I_d)$.

Another example of this type: Gaussian mixture of a distribution P_Z of random variable Z taking values in (0, 1] and a distribution $\mathcal{N}(0, \Theta)$ with $\Theta \leq I_d$ belongs to the spherical family $\mathcal{P}^d_{\gamma_G}$ associated with the standard univariate Gaussian density

$$\gamma_{\mathcal{G}}(s) = \frac{1}{\sqrt{2\pi}} \mathrm{e}^{-s^2/2}$$

This family has a cap, specifically, the standard Gaussian *d*-dimensional distribution $\mathcal{N}(0, I_d)$.

2.2.3.C. Main result. Looking at the proof of Proposition 2.6, we arrive at the following

Proposition 2.7. Let X_1, X_2 be nonempty and nonintersecting closed convex sets in \mathbb{R}^d , one of the sets being bounded, and let \mathcal{P}^d_{γ} be a spherical family of probability distributions. With Opt and $\phi(\cdot)$ given by (2.12) – (2.13), let us split the width 2Opt of the stripe { $\omega : -\text{Opt} \leq \phi(\omega) \leq \text{Opt}$ } separating X_1 and X_2 into two nonnegative parts:

$$\delta_1 \ge 0, \delta_2 \ge 0, \ \delta_1 + \delta_2 = 2\text{Opt} \tag{2.20}$$

and consider simple test \mathcal{T} deciding on the hypotheses $H_1 : x \in X_1, H_2 : x \in X_2$ via observation (2.8) by accepting H_1 when

$$\phi(\omega) \ge \frac{1}{2} [\delta_2 - \delta_1]$$

and accepting H_2 otherwise. Then

$$\operatorname{Risk}_{\chi}(\mathcal{T}|H_1, H_2) \leq \int_{\delta_{\chi}}^{\infty} \gamma(s) ds, \ \chi = 1, 2$$
(2.21)

In addition, when $\delta_1 = \delta_2 = \text{Opt}$ and \mathcal{P}^d_{γ} has a cap, \mathcal{T} is the minimum risk test deciding on H_1 , H_2 . The risk of this test is

$$\operatorname{Risk}(\mathcal{T}|H_1, H_2) = P_{\gamma}(\operatorname{Opt}) := \int_{O_{pt}}^{\infty} \gamma(s) ds.$$
(2.22)

To illustrate the power of Proposition 2.7, consider the case when γ is the function (2.19) stemming from Student's *t*-distribution on \mathbf{R}^d with ν degrees of freedom. It is known that in this case γ is the density of univariate Student's *t*-distribution with ν degrees of freedom:

$$\gamma(s) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)\sqrt{\pi\nu}} (1+s^2/\nu)^{-\frac{\nu+1}{2}},$$

where $\Gamma(\cdot)$ is Euler's Gamma function. When $\nu = 1$, $\gamma(\cdot)$ is just the heavy tailed (no expectation!) standard Cauchy density $\frac{1}{\pi}(1+s^2)^{-1}$. Same as in this "extreme case," multidimensional Student's distributions have relatively heavy tails (the heavier the less is ν) and as such are of interest in Finance.

2.2.3 Euclidean Separation, Repeated Observations, and Majority Tests

Assume that $X_1, X_2, \mathcal{P}^d_{\gamma}$ are as in the premise of Proposition 2.7 and K-repeated observations are allowed, K > 1. An immediate attempt to reduce the situation to the single-observation case by calling K-repeated observation $\omega^K = (\omega_1, ..., \omega_K)$ our new observation and thus reducing testing via repeated observations to the single-observation case seemingly fails: already in the simplest case of stationary K-repeated observations this reduction would require replacing the family \mathcal{P}^d_{γ} with the family of product distributions $\underbrace{P \times ... \times P}_{K}$ stemming from $P \in \mathcal{P}^d_{\gamma}$, and it

is unclear how to apply to the resulting single-observation testing problem our machinery based on Euclidean separation. Instead, let us use K-step majority test.

We are in the situation when our inference should be based on observations

$$\omega^K = (\omega_1, \omega_2, ..., \omega_K), \tag{2.23}$$

and decide on hypotheses \mathcal{H}_1 , \mathcal{H}_2 on the distribution Q^K of ω^K , and we are interested in the following 3 cases:

- **S** [stationary K-repeated observations, cf. Section 2.1.3.1]: $\omega_1, ..., \omega_K$ are drawn independently of each other from the same distribution Q, that is, Q^K is the product distribution $Q \times ... \times Q$. Further, under hypothesis $\mathcal{H}_{\chi}, \chi = 1, 2, Q$ is the distribution of random variable $\omega = x + \xi$, where $x \in X_{\chi}$ is deterministic, and the distribution P of ξ belongs to the family \mathcal{P}_{γ}^d ;
- **SS** [semi-stationary K-repeated observations, cf. Section 2.1.3.2]: there are two deterministic sequences, one of signals $\{x_k\}_{k=1}^K$, another of distributions $\{P_k \in \mathcal{P}_q^{\mathcal{A}}\}_{k=1}^K$, and $\omega_k = x_k + \xi_k$, $1 \leq k \leq K$, with $\xi_k \sim P_k$ independent across k. Under hypothesis \mathcal{H}_{χ} , all signals $x_k, k \leq K$, belong to X_{χ} .
- **QS** [quasi-stationary K-repeated observations, cf. Section 2.1.3.3]: "in the nature" there exists a random sequence of driving factors $\zeta^k = (\zeta_1, ..., \zeta_K)$ such that observation ω_k , for every k, is a deterministic function of $\zeta^k = (\zeta_1, ..., \zeta_k)$: $\omega_k = \theta_k(\zeta^k)$. On the top of it, under ℓ -th hypothesis \mathcal{H}_ℓ , for all $k \leq K$ and all ζ^{k-1} the conditional, ζ^{k-1} being given, distribution of ω_k belong to the family \mathcal{P}_ℓ of distributions of all random vectors of the form $x + \xi$, where $x \in X_\ell$ is deterministic, and ξ is random noise with distribution from \mathcal{P}_{γ}^d .

2.2.3.2 Majority Test

2.2.3.2.A. The construction of *K*-observation majority test is very natural. We use Euclidean separation to build simple single-observation test \mathcal{T} deciding on hypotheses H_{χ} : $x \in X_{\chi}$, $\chi = 1, 2$, via observation $\omega = x + \xi$, where x is deterministic, and the distribution of noise ξ belongs to \mathcal{P}_{γ}^d . \mathcal{T} is given by the construction from Proposition 2.7 applied with $\delta_1 = \delta_2 = \text{Opt.}$ The summary of our actions is as follows:

$$X_{1}, X_{2} \Rightarrow \begin{cases} \text{Opt} = \min_{x^{1} \in X_{1}, x^{2} \in X_{2}} \frac{1}{2} \|x^{1} - x^{2}\|_{2} \\ (x_{*}^{1}, x_{*}^{2}) \in \operatorname{Argmin}_{x^{1} \in X_{1}, x^{2} \in X_{2}} \frac{1}{2} \|x^{1} - x^{2}\|_{2} \\ e = \frac{x_{*}^{1} - x_{*}^{2}}{\|x_{*}^{1} - x_{*}^{2}\|_{2}}, c = \frac{1}{2} e^{T} [x_{*}^{1} + x_{*}^{2}] \\ \Rightarrow \qquad \phi(\omega) = e^{T} \omega - c \end{cases}$$

$$(2.24)$$

Majority test $\mathcal{T}_{K}^{\text{maj}}$, as applied to K-repeated observation $\omega^{K} = (\omega_{1}, ..., \omega_{K})$, builds the K reals $v_{k} = \phi(\omega_{k})$. If at least K/2 of these reals are nonnegative, the test accepts \mathcal{H}_{1} and rejects \mathcal{H}_{2} ; otherwise the test accepts \mathcal{H}_{2} and rejects \mathcal{H}_{1} .

2.2.3.2.B. Risk analysis. We intend to carry out the risk analysis for the case QS of quasi-stationary K-repeated observations; this analysis automatically applies to the cases SS of stationary and S of K-repeated stationary/semi-stationary observations, which are special cases of QS.

Proposition 2.8. With $X_1, X_2, \mathcal{P}^d_{\gamma}$ obeying the premise of Proposition 2.7, in the case **QS** of quasi-stationary observations the risk of K-observation Majority test \mathcal{T}_K^{maj} can be bounded as

$$\operatorname{Risk}(\mathcal{T}_{K}^{maj}|\mathcal{H}_{1},\mathcal{H}_{2}) \leq \epsilon_{K} \equiv \sum_{K/2 \leq k \leq K} \binom{K}{k} \epsilon_{\star}^{k} (1-\epsilon_{\star})^{K-k}, \ \epsilon_{\star} = \int_{Opt}^{\infty} \gamma(s) ds.$$
(2.25)

Proof. Here we restrict ourselves to the case SS of semi-stationary K-repeated observations. In "full generality," that is, in the case QS of quasi-stationary K-repeated observations, Proposition will be proved in Section 2.11.2.

Assume that \mathcal{H}_1 takes place, so that (recall that we are in the **SS** case!) $\omega_k = x_k + \xi_k$ with some deterministic $x_k \in X_1$ and independent across k noises $\xi_k \sim P_k$, for some deterministic sequence $P_k \in \mathcal{P}_{\gamma}^d$. Let us fix $\{x_k \in X_1\}_{k=1}^K$ and $\{P_k \in \mathcal{P}_{\gamma}^d\}_{k=1}^K$. Then the random reals $v_k = \phi(\omega_k = x_k + \xi_k)$ are independent across k, and so are the Boolean random variables

$$\chi_k = \begin{cases} 1, & v_i < 0\\ 0, & v_i \ge 0 \end{cases}$$

 $\chi_k = 1$ if and only if test \mathcal{T} , as applied to observation ω_k , rejects hypothesis $H_1: x_k \in X_1$. By Proposition 2.7, P_k -probability p_k of the event $\chi_k = 1$ is at most ϵ_{\star} . Further, by construction of the Majority test, if $\mathcal{T}_K^{\text{maj}}$ rejects the true hypothesis \mathcal{H}_1 , then the number of k's with $\chi_k = 1$ is $\geq K/2$. Thus, with $\{x_k \in X_1\}$ and $P_k \in \mathcal{P}_{\gamma}^d$ the probability to reject \mathcal{H}_1 is not greater than the probability of the event

In K independent coin tosses, with probability $p_k \leq \epsilon_*$ to get head in k-th toss, the total number of heads is $\geq K/2$.

The probability of this event clearly does not exceed the right hand side in (2.25), implying that $\operatorname{Risk}_1(\mathcal{T}_K^{\operatorname{maj}}|\mathcal{H}_1,\mathcal{H}_2) \leq \epsilon_k$. "Symmetric" reasoning yields

$$\operatorname{Risk}_2(\mathcal{T}_K^{\operatorname{maj}}|\mathcal{H}_1,\mathcal{H}_2) \leq \epsilon_k$$

completing the proof of (2.25).

Corollary 2.9. Under the premise of Proposition 2.8, the upper bounds ϵ_K on the risk of the K-observation Majority test goes to 0 exponentially fast as $K \to \infty$.

Indeed, we are in the situation of Opt > 0, so that $\epsilon_{\star} < \frac{1}{2}$ ¹⁴.

Remark 2.10. When proving (**SS**-version of) Proposition 2.8, we have used "evident" observation as follows:

(#) Let $\chi_1,..., \chi_K$ be independent random variables taking values 0 and 1, and let the probabilities p_k for χ_k to be 1 be upper-bounded by some $\epsilon \in [0,1]$ for all k. Then for every fixed M the probability of the event "at least M of $\chi_1,...,\chi_K$ are equal to 1" is upper-bounded by the probability $\sum_{M \leq k \leq K} {K \choose k} \epsilon^k (1-\epsilon)^{K-k}$ of the same event in the case when $p_k = \epsilon$ for all k.

62

¹⁴Recall that we have assumed from the very beginning that γ is an even probability density on the axis, and there is no neighbourhood of the origin where $\gamma = 0$ a.s.

If there are evident facts in Math, (#) definitely is one of them. Nevertheless, why (#) is true?

Reader is kindly asked to prove (#). For your information: design of proof took about 10-minute effort of the authors (a bit too much for an evident statement); the results of their effort can be found in Section 2.11.2.

2.2.3.2.C. Near-optimality. We are about to show that under appropriate assumptions, the majority test $\mathcal{T}_{K}^{\text{maj}}$ is near-optimal. The precise statement is as follows:

Proposition 2.11. Let $X_1, X_2, \mathcal{P}^d_{\gamma}$ obey the premise of Proposition 2.7. Assume that the spherical family \mathcal{P}_{γ} and positive reals D, α, β are such that

$$\beta D \le \frac{1}{4},\tag{2.26}$$

$$\int_{0}^{\delta} \gamma(s) ds \ge \beta \delta, \quad 0 \le \delta \le D, \tag{2.27}$$

and \mathcal{P}_{γ} contains a density $q(\cdot)$ such that

$$\int_{\mathbf{R}^n} \sqrt{q(\xi - e)q(\xi + e)} d\xi \ge \exp\{-\alpha e^T e\} \ \forall (e : \|e\|_2 \le D).$$
(2.28)

Let, further, the sets X_1 , X_2 be such that Opt as given by (2.12) satisfies the relation

$$Opt \le D. \tag{2.29}$$

Given tolerance $\epsilon \in (0, 1/5)$, the risk of K-observation majority test \mathcal{T}_{K}^{maj} utilizing QS observations ensures the relation

$$K \ge K^* := \left\lfloor \frac{\ln(1/\epsilon)}{2\beta^2 \operatorname{Opt}^2} \right\rfloor \implies \operatorname{Risk}(\mathcal{T}_K^{maj} | \mathcal{H}_1, \mathcal{H}_2) \le \epsilon$$
(2.30)

(here |x| stands for the smallest integer $\geq x \in \mathbf{R}$). In addition, for every K-observation test \mathcal{T}_K utilizing stationary repeated observations and satisfying

$$\operatorname{Risk}(\mathcal{T}_K | \mathcal{H}_1, \mathcal{H}_2) \leq \epsilon$$

it holds

$$K \ge K_* := \frac{\ln\left(\frac{1}{4\epsilon}\right)}{2\alpha\delta^2}.\tag{2.31}$$

As a result, the majority test $\mathcal{T}_{K^*}^{maj}$ for (\mathcal{S}_{K^*}) has risk at most ϵ and is near-optimal, in terms of the required number of observations, among all tests with risk $\leq \epsilon$: the number K of observations in such test satisfies the relation

$$K^*/K \le \theta := K^*/K_* = O(1)\frac{\alpha}{\beta^2}.$$

Proof of Proposition is the subject of Exercise 2.56.

Illustration. Given $\nu \geq 1$, consider the case when $\mathcal{P} = \mathcal{P}_{\gamma}$ is the spherical family

with n-variate (spherical) Student's distribution in the role of the cap, so that

$$\gamma(s) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)(\pi\nu)^{1/2}} \left[1 + s^2/\nu\right]^{-(\nu+1)/2}$$
(2.32)

It is easily seen (see Exercise 2.56) that \mathcal{P} contains the $\mathcal{N}(0, \frac{1}{2}I_n)$ density $q(\cdot)$, implying that setting

$$D = 1, \ \alpha = 1, \ \beta = \frac{1}{7},$$

one ensures relations (2.26), (2.27) and (2.29). As a result, when Opt as yielded by (2.12) is ≤ 1 , the non-optimality factor θ of the majority test $\mathcal{T}_{K^*}^{\text{maj}}$ as defined in Proposition 2.11 does not exceed O(1).

2.2.4 From Pairwise to Multiple Hypotheses Testing

2.2.4.1 Situation

Assume we are given L families of probability distributions \mathcal{P}_{ℓ} , $1 \leq \ell \leq L$, on observation space Ω , and observe a realization of random variable $\omega \sim P$ taking values in Ω . Given ω , we want to decide on the L hypotheses

$$H_{\ell}: P \in \mathcal{P}_{\ell}, \ 1 \le \ell \le L. \tag{2.33}$$

Our *ideal goal* would be to find a low-risk simple test deciding on the hypotheses. However, it may happen that the "ideal goal" is not achievable, for example, when some pairs of families \mathcal{P}_{ℓ} have nonempty intersections. When $\mathcal{P}_{\ell} \cap \mathcal{P}_{\ell'} \neq \emptyset$ for some $\ell \neq \ell'$, there is no way to decide on the hypotheses with risk < 1/2.

But: Impossibility to decide reliably on all L hypotheses "individually" does not mean that no meaningful inferences can be done.

For example, consider the 3 colored rectangles on the plane:



and 3 hypotheses, with H_{ℓ} , $1 \leq \ell \leq 3$, stating that our observation is $\omega = x + \xi$ with deterministic "signal" x belonging to ℓ -th rectangle and $\xi \sim \mathcal{N}(0, \sigma^2 I_2)$. Whatever small σ be, no test can decide on the 3 hypotheses with risk < 1/2; e.g., there is no way to decide reliably on H_1 vs. H_2 . However, we may hope that when σ is small (or when repeated observations are allowed), observations allow us to discard reliably some of the hypotheses; for example, when the signal "is brown" (i.e., H_1 holds true), we hardly can discard reliably the hypothesis H_2 stating that the signal "is green," but hopefully can discard reliably H_3 (that is, infer that the signal is not blue).

When handling multiple hypotheses which cannot be reliably decided upon "as they are," it makes sense to speak about *testing the hypotheses "up to closeness.*"

2.2.4.2Closeness relation and "up to closeness" risks

Closeness relation, or simply *closeness* C on a collection of L hypotheses $H_1, ..., H_L$ is defined as some set of pairs (ℓ, ℓ') with $1 \leq \ell, \ell' \leq L$. We interpret the relation $(\ell, \ell') \in \mathcal{C}$ as the fact that the hypotheses H_{ℓ} and H'_{ℓ} are close to each other. Sometimes we shall use the words " ℓ and ℓ' are/are not C-close to each other" as an equivalent form of "hypotheses H_{ℓ} , $H_{\ell'}$ are/are not \mathcal{C} -close to each other." We always assume that

- \mathcal{C} contains all "diagonal pairs" $(\ell, \ell), 1 \leq \ell \leq L$ ("every hypothesis is close to itself");
- $(\ell, \ell') \in \mathcal{C}$ is and only if $(\ell', \ell) \in \mathcal{C}$ ("closeness is a symmetric relation").

Note that by symmetry of \mathcal{C} , the relation $(\ell, \ell') \in \mathcal{T}$ is in fact a property of unordered pair $\{\ell, \ell'\}$.

"Up to closeness" risks. Let \mathcal{T} be a test deciding on L hypotheses $H_1, ..., H_L$, see (2.33); given observation ω , \mathcal{T} accepts all hypotheses H_{ℓ} with indexes $\ell \in \mathcal{T}(\omega)$ and rejects all other hypotheses. We say that ℓ -th partial C-risk of test T is $\leq \epsilon$, if whenever H_{ℓ} is true: $\omega \sim P \in \mathcal{P}_{\ell}$, the P-probability of the event

$$\mathcal{T} \ accepts \ H_{\ell}: \ \ell \in \mathcal{T}(\omega)$$

and
all hypotheses $H_{\ell'}$ accepted by \mathcal{T} are \mathcal{C} -close to $H_{\ell}: \ (\ell, \ell') \in \mathcal{C}, \forall \ell' \in \mathcal{T}(\omega)$

is at least 1 -

 ℓ -th partial \mathcal{C} -risk Risk $^{\mathcal{C}}_{\ell}(\mathcal{T}|H_1,...,H_L)$ of \mathcal{T} is the smallest ϵ with the outlined property, or, equivalently,

$$\operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}|H_1, ..., H_L) = \sup_{P \in \mathcal{P}_{\ell}} \operatorname{Prob}_{\omega \sim P} \left\{ [\ell \notin \mathcal{T}(\omega)] \text{ or } [\exists \ell' \in \mathcal{T}(\omega) : (\ell, \ell') \notin \mathcal{C}] \right\}$$

 \mathcal{C} -risk Risk^{\mathcal{C}}($\mathcal{T}|H_1, ..., H_L$) of \mathcal{T} is the largest of the partial \mathcal{C} -risks of the test:

$$\operatorname{Risk}^{\mathcal{C}}(\mathcal{T}|H_1,...,H_L) = \max_{1 \le \ell \le L} \operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}|H_1,...,H_L).$$

Observe that when \mathcal{C} is the "strictest possible" closeness, that is, $(\ell, \ell') \in \mathcal{C}$ if and only if $\ell = \ell'$, then a test \mathcal{T} deciding on $H_1, ..., H_L$ up to closeness \mathcal{C} with risk ϵ is, *basically*, the same as a simple test deciding on $H_1, ..., H_L$ with risk $\leq \epsilon$. Indeed, a test with the latter property clearly decides on $H_1, ..., H_L$ with C-risk $\leq \epsilon$. The inverse statement, taken literally, is not true, since even with our "as strict as possible" closeness, a test \mathcal{T} with \mathcal{C} -risk $\leq \epsilon$ not necessarily is simple. However, we can enforce \mathcal{T} to be simple, specifically, to accept a once for ever fixed hypothesis, say, H_1 , and only it, when the set of hypotheses accepted by \mathcal{T} "as is" is not a singleton, otherwise accept exactly the same hypothesis as \mathcal{T} . The modified test already is simple, and clearly its C-risk does not exceed the one of T.

2.2.4.3 Multiple Hypothesis Testing via pairwise tests

Assume that for every unordered pair $\{\ell, \ell'\}$ with $(\ell, \ell') \notin \mathcal{C}$ we are given a simple test $\mathcal{T}_{\{\ell,\ell'\}}$ deciding on H_{ℓ} vs. $H_{\ell'}$ via observation ω .

Our goal is to "assemble" the tests $\mathcal{T}_{\{\ell,\ell'\}}$, $(\ell,\ell') \notin \mathcal{C}$, into a test \mathcal{T} deciding on $H_1..., H_L$ up to closeness \mathcal{C} .

The construction we intend to use is as follows:

- For $1 \leq \ell, \ell' \leq L$, we define functions $T_{\ell\ell'}(\omega)$ as follows:
 - when $(\ell, \ell') \in \mathcal{C}$, we set $T_{\ell\ell'}(\cdot) \equiv 0$.
 - when $(\ell, \ell') \notin \mathcal{C}$, so that $\ell \neq \ell'$, we set

$$T_{\ell\ell'}(\omega) = \begin{cases} 1, & \mathcal{T}_{\{\ell,\ell'\}}(\omega) = \{\ell\} \\ -1, & \mathcal{T}_{\{\ell,\ell'\}}(\omega) = \{\ell'\} \end{cases}$$
(2.34)

Note that $\mathcal{T}_{\{\ell,\ell'\}}$ is a simple test, so that $T_{\ell\ell'}(\cdot)$ is well defined and takes values ± 1 when $(\ell,\ell') \notin \mathcal{C}$ and 0 when $(\ell,\ell') \in \mathcal{C}$.

Note that by construction and since C is symmetric, we have

$$T_{\ell\ell'}(\omega) \equiv -T_{\ell'\ell}(\omega), \ 1 \le \ell, \ell' \le L.$$
(2.35)

• The test \mathcal{T} is as follows: given observation ω , we build the $L \times L$ matrix $T(\omega) = [T_{\ell\ell'}(\omega)]$ and accept exactly those of the hypotheses H_{ℓ} for which ℓ -th row in $T(\omega)$ is nonnegative.

Observation 2.12. When \mathcal{T} accepts some hypothesis H_{ℓ} , all hypotheses accepted by \mathcal{T} are C-close to H_{ℓ} .

Indeed, if ω is such that $\ell \in \mathcal{T}(\omega)$, then the ℓ -th row in $T(\omega)$ is nonnegative. If now ℓ' is not \mathcal{C} -close to ℓ , we have $T_{\ell\ell'}(\omega) \ge 0$ and $T_{\ell\ell'}(\omega) \in \{-1, 1\}$, whence $T_{\ell\ell'}(\omega) = 1$. Consequently, by (2.35) it holds $T_{\ell'\ell}(\omega) = -1$, implying that ℓ' -th row in $T(\omega)$ is not nonnegative, and thus $\ell' \notin \mathcal{T}(\omega)$.

Risk analysis. For $(\ell, \ell') \notin C$, let

$$\epsilon_{\ell\ell'} = \operatorname{Risk}_{1}(\mathcal{T}_{\{\ell,\ell'\}}|H_{\ell},H_{\ell'}) = \sup_{P\in\mathcal{P}_{\ell}}\operatorname{Prob}_{\omega\sim P}\{\ell \notin \mathcal{T}_{\{\ell,\ell'\}}(\omega)\}$$

$$= \sup_{P\in\mathcal{P}_{\ell}}\operatorname{Prob}_{\omega\sim P}\{T_{\ell\ell'}(\omega) = -1\} = \sup_{P\in\mathcal{P}_{\ell}}\operatorname{Prob}_{\omega\sim P}\{T_{\ell'\ell}(\omega) = 1\}$$

$$= \sup_{P\in\mathcal{P}_{\ell}}\operatorname{Prob}_{\omega\sim P}\{\ell' \in \mathcal{T}_{\{\ell,\ell'\}}(\omega)\}$$

$$= \operatorname{Risk}_{2}(\mathcal{T}_{\{\ell,\ell'\}}|H_{\ell'},H_{\ell}).$$
(2.36)

Proposition 2.13. For the just defined test \mathcal{T} it holds

$$\forall \ell \leq L : \operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}|H_1, ..., H_L) \leq \epsilon_{\ell} := \sum_{\ell' : (\ell, \ell') \notin \mathcal{C}} \epsilon_{\ell \ell'}.$$
(2.37)

Proof. Let us fix ℓ , let H_{ℓ} be true, and let $P \in \mathcal{P}_{\ell}$ be the distribution of observation ω . Set $I = \{\ell' \leq L : (\ell, \ell') \notin \mathcal{C}\}$. For $\ell' \in I$, let $E_{\ell'}$ be the event

$$\{\omega: T_{\ell\ell'}(\omega) = -1\}.$$

We have $\operatorname{Prob}_{\omega \sim P}(E_{\ell'}) \leq \epsilon_{\ell \ell'}$ (by definition of $\epsilon_{\ell \ell'}$), whence

$$\operatorname{Prob}_{\omega \sim P}\left(\underbrace{\cup_{\ell' \in I} E_{\ell'}}_{E}\right) \leq \epsilon_{\ell}.$$

When the event E does not take place, we have $T_{\ell\ell'}(\omega) = 1$ for all $\ell' \in I$, so that $T_{\ell\ell'}(\omega) \ge 0$ for all $\ell', 1 \le \ell' \le L$, whence $\ell \in \mathcal{T}(\omega)$. By Observation 2.12, the latter inclusion implies that

$$\{\ell \in \mathcal{T}(\omega) \& \{(\ell, \ell') \in \mathcal{C} \forall \ell' \in \mathcal{T}(\omega)\}.$$

Invoking the definition of partial C-risk, we get

$$\operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}|H_1, ..., H_L) \leq \operatorname{Prob}_{\omega \sim P}(E) \leq \epsilon_{\ell}.$$

2.2.4.4 Testing Multiple Hypotheses via Euclidean separation

Situation. We are given L nonempty and closed convex sets $X_{\ell} \subset \Omega = \mathbf{R}^d$, $1 \leq \ell \leq L$, with at least L-1 of the sets being bounded, and a spherical family of probability distributions \mathcal{P}^d_{γ} . These data define L families \mathcal{P}_{ℓ} of probability distributions on \mathbf{R}^d ; the family \mathcal{P}_{ℓ} , $1 \leq \ell \leq L$, is comprised of probability distributions of all random vectors of the form $x + \xi$, where deterministic x ("signal") belongs to X_{ℓ} , and ξ is random noise with distribution from \mathcal{P}^d_{γ} . Given positive integer K, we can speak about L hypotheses on the distribution P^K of K-repeated observation $\omega^K = (\omega_1, ..., \omega_K)$, with \mathcal{H}_{ℓ} stating that ω^K is a quasi-stationary K-repeated observation associated with \mathcal{P}_{ℓ} . In other words $\mathcal{H}_{\ell} = H_{\ell}^{\otimes, K}$, see Section 2.1.3.3. Finally, we are given a closeness \mathcal{C} .

Our goal is to decide on the hypotheses $\mathcal{H}_1, ..., \mathcal{H}_L$ up to closeness \mathcal{C} via K-repeated observation ω^K . Note that this is a natural extension of the case **QS** of pairwise testing from repeated observations considered in Section 2.2.3 (there L = 2 and \mathcal{C} is the only meaningful closeness on a two-hypotheses set: $(\ell, \ell') \in \mathcal{C}$ is and only if $\ell = \ell'$).

Standing Assumption which is by default in force everywhere in this Section is:

Whenever ℓ, ℓ' are not C-close: $(\ell, \ell') \notin C$, the sets $X_{\ell}, X_{\ell'}$ do not intersect.

Strategy: We intend to attack the above testing problem by assembling pairwise Euclidean separation Majority tests via the construction from Section 2.2.4.3.

Building blocks to be assembled are Euclidean separation K-observation pairwise Majority tests built for the pairs \mathcal{H}_{ℓ} , $\mathcal{H}_{\ell'}$ of hypotheses with *not* close to each other ℓ and ℓ' , that is, with $(\ell, \ell') \notin \mathcal{C}$. These tests are built as explained in Section 2.2.3.2; for reader's convenience, here is the construction. For a pair $(\ell, \ell') \notin \mathcal{C}$, we

1. Find the optimal value $Opt_{\ell\ell'}$ and an optimal solution $(u_{\ell\ell'}, v_{\ell\ell'})$ to the convex optimization problem

$$Opt_{\ell\ell'} = \min_{u \in X_{\ell'}, v \in X_{\ell'}} \frac{1}{2} \|u - v\|_2,$$
(2.38)

The latter problem is solvable, since we have assumed from the very beginning that X_{ℓ} , X'_{ℓ} are nonempty, closed and convex, and that at least one of these sets is bounded;

2. Set

$$e_{\ell\ell'} = \frac{u_{\ell\ell'} - v_{\ell\ell'}}{\|u_{\ell\ell'} - v_{\ell\ell'}\|_2}, \, c_{\ell\ell'} = \frac{1}{2} e_{\ell\ell'}^T [u_{\ell\ell'} + v_{\ell\ell'}], \, \phi_{\ell\ell'}(\omega) = e_{\ell\ell'}^T \omega - c_{\ell\ell'}.$$

Note that the construction makes sense, since by our Standing Assumption for ℓ , ℓ' in question X_{ℓ} and $X_{\ell'}$ do not intersect. Further, $e_{\ell\ell'}$ and $c_{\ell\ell'}$ clearly depend solely on (ℓ, ℓ') , but not on how we select an optimal solution $(u_{\ell\ell'}, v_{\ell\ell'})$ to (2.38). Finally, we have

$$e_{\ell\ell'} = -e_{\ell'\ell}, c_{\ell\ell'} = -c_{\ell'\ell}, \phi_{\ell\ell'}(\cdot) \equiv -\phi_{\ell'\ell}(\cdot)$$

- 3. We consider separately the case of K = 1 and the case of K > 1. Specifically,
 - a) when K = 1, we select somehow nonnegative reals $\delta_{\ell\ell'}$, $\delta_{\ell'\ell}$ such that

$$\delta_{\ell\ell'} + \delta_{\ell'\ell} = 2\mathrm{Opt}_{\ell\ell'} \tag{2.39}$$

and specify the single-observation simple test $\mathcal{T}_{\ell\ell'}$ deciding on the hypotheses $\mathcal{H}_{\ell}, \mathcal{H}_{\ell'}$ according to

$$\mathcal{T}_{\ell\ell'}(\omega) = \begin{cases} \{\ell\}, & \phi_{\ell\ell'}(\omega) \ge \frac{1}{2} [\delta_{\ell'\ell} - \delta_{\ell\ell'}] \\ \{\ell'\}, & \text{otherwise} \end{cases};$$

Note that by Proposition 2.7, setting

$$P_{\gamma}(\delta) = \int_{\delta}^{\infty} \gamma(s) ds, \qquad (2.40)$$

we have

$$\begin{aligned}
\operatorname{Risk}_{1}(\mathcal{T}_{\ell\ell'}|\mathcal{H}_{\ell},\mathcal{H}_{\ell'}) &\leq P_{\gamma}(\delta_{\ell\ell'}) \\
\operatorname{Risk}_{2}(\mathcal{T}_{\ell\ell'}|\mathcal{H}_{\ell},\mathcal{H}_{\ell'}) &\leq P_{\gamma}(\delta_{\ell'\ell}) \\
\operatorname{Risk}_{1}(\mathcal{T}_{\ell'\ell}|\mathcal{H}_{\ell'},\mathcal{H}_{\ell}) &\leq P_{\gamma}(\delta_{\ell'\ell}) \\
\operatorname{Risk}_{2}(\mathcal{T}_{\ell'\ell}|\mathcal{H}_{\ell'},\mathcal{H}_{\ell}) &\leq P_{\gamma}(\delta_{\ell\ell'})
\end{aligned} (2.41)$$

b) when K > 1, we specify K-observation simple test $\mathcal{T}_{\ell\ell' K}$ deciding on \mathcal{H}_{ℓ} , $\mathcal{H}_{\ell'}$ according to

$$\mathcal{T}_{\ell\ell'}(\omega^k = (\omega_1, ..., \omega_k)) = \begin{cases} \{\ell\}, & \operatorname{Card}\{k \le K : \phi_{\ell\ell'} \ge 0\} \ge K/2, \\ \{\ell'\}, & \operatorname{otherwise} \end{cases}$$

Note that by Proposition 2.8 we have

$$\operatorname{Risk}(\mathcal{T}_{\ell\ell'K}|\mathcal{H}_{\ell},\mathcal{H}_{\ell}') \leq \epsilon_{\ell\ell'K} := \sum_{K/2 \leq k \leq K} {K \choose k} \epsilon_{\star\ell\ell'}^k (1 - \epsilon_{\star\ell\ell'})^{K-k},$$

$$\epsilon_{\star\ell\ell'} = P_{\gamma}(\operatorname{Opt}_{\ell\ell'}) = \epsilon_{\star\ell'\ell}.$$
(2.42)

Assembling building blocks, case of K = 1. In the case of K = 1, we specify the simple pairwise tests $\mathcal{T}_{\{\ell,\ell'\}}$, $(\ell,\ell') \notin \mathcal{C}$, participating in the construction of the multi-hypothesis test presented in Section 2.2.4.3, as follows. Given unordered pair $\{\ell,\ell'\}$ with $(\ell,\ell') \notin \mathcal{C}$ (which is exactly the same as $(\ell',\ell) \notin \mathcal{C}$), we arrange ℓ,ℓ' in

ascending order, thus arriving at ordered pair $(\bar{\ell}, \bar{\ell}')$, and set

$$\mathcal{T}_{\{\ell,\ell'\}}(\cdot) = \mathcal{T}_{\bar{\ell}\bar{\ell}'}(\cdot),$$

with the right hand side tests defined as explained above. We then assemble, as explained in Section 2.2.4.3, the tests $\mathcal{T}_{\{\ell,\ell'\}}$ into a single-observation test \mathcal{T}_1 deciding on hypotheses $\mathcal{H}_1, ..., \mathcal{H}_L$. Looking at (2.36) and (2.41), we conclude that for the just defined tests $\mathcal{T}_{\{\ell,\ell'\}}$ and the associated with the tests $\mathcal{T}_{\{\ell,\ell'\}}$, via (2.36), quantities $\epsilon_{\ell\ell'}$ it holds

$$(\ell, \ell') \notin \mathcal{C} \Rightarrow \epsilon_{\ell\ell'} \le P_{\gamma}(\delta_{\ell\ell'}).$$
 (2.43)

Invoking Proposition 2.13, we get

Proposition 2.14. In the situation described in the beginning of Section 2.2.4.4 and under Standing Assumption, the C-risks of the just defined test \mathcal{T}_1 , whatever be the choice of nonnegative $\delta_{\ell\ell'}$, $(\ell, \ell') \notin C$, satisfying (2.39), can be upper-bounded as

$$\operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}_{1}|\mathcal{H}_{1},...,\mathcal{H}_{L}) \leq \sum_{\ell':(\ell,\ell')\notin\mathcal{C}} P_{\gamma}(\delta_{\ell\ell'}).$$
(2.44)

with $P_{\gamma}(\cdot)$ given by (2.40).

Case of K = 1 (continued): Optimizing the construction. We can try to optimize the risk bounds (2.44) over the parameters $\delta_{\ell\ell'}$ of the construction. The first question to be addressed here is what to minimize – we have several risks! A natural model here is as follows. Let us fix a nonnegative $M \times L$ weight matrix W and M-dimensional positive profile vector w, and solve the optimization problem

$$\min_{\substack{t,\{\delta_{\ell\ell'}:(\ell,\ell')\notin\mathcal{C}\\t\,:\,\,}} \left\{ t: \begin{array}{l} W \cdot \left[\sum_{\ell':(\ell,\ell')\notin\mathcal{C}} P_{\gamma}(\delta_{\ell\ell'})\right]_{\ell=1}^{L} \leq tw \\ \delta_{\ell\ell'} \geq 0, \delta_{\ell\ell'} + \delta_{\ell'\ell} = 2\mathrm{Opt}_{\ell\ell'}, \ (\ell,\ell')\notin\mathcal{C} \end{array} \right\}.$$
(2.45)

For example, when M = 1 and w = 1, we are minimizing weighted sum of (upper bounds on) partial *C*-risks of our test, and when *W* is a diagonal matrix with positive diagonal entries and *w* is the all-ones vector, we are minimizing the largest of scaled partial risks. Note that when $P_{\gamma}(\cdot)$ is convex on \mathbf{R}_+ , or, which is the same, $\gamma(\cdot)$ is nonincreasing in \mathbf{R}_+ , (2.45) is a convex, and thus efficiently solvable, problem.

Assembling building blocks, case of K > 1. We again pass from our building blocks – K-observation simple pairwise tests $\mathcal{T}_{\ell\ell'K}$, $(\ell, \ell') \notin C$, we have already specified, to tests $\mathcal{T}_{\{\ell,\ell'\}} = \mathcal{T}_{\ell\ell'K}$, with $\ell = \min[\ell, \ell']$ and $\ell' = \max[\ell, \ell']$, and then apply to the resulting tests the construction from Section 2.2.4.3, arriving at K-observation multi-hypothesis test \mathcal{T}_K . By Proposition 2.8, the quantities $\epsilon_{\ell\ell'}$ associated with the tests $\mathcal{T}_{\{\ell,\ell'\}}$ via (2.36) satisfy the relation

$$(\ell, \ell') \notin \mathcal{C} \Rightarrow \epsilon_{\ell\ell'} \leq \sum_{K/2 \leq k \leq K} \binom{K}{k} [P_{\gamma}(\operatorname{Opt}_{\ell\ell'})]^{k} [1 - P_{\gamma}(\operatorname{Opt}_{\ell\ell'})]^{K-k}, \quad (2.46)$$

which combines with Proposition 2.13 to imply

Proposition 2.15. Let the situation described in the beginning of Section 2.2.4.4 take place, and let K > 1. Under Standing Assumption, the C-risks of the just defined test \mathcal{T}_K can be upper-bounded as

$$\operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}_{1}|\mathcal{H}_{1},...,\mathcal{H}_{L}) \leq \sum_{\ell':(\ell,\ell')\notin\mathcal{C}} P_{\gamma}(\delta_{\ell\ell'}) \sum_{K/2 \leq k \leq K} \binom{K}{k} [P_{\gamma}(\operatorname{Opt}_{\ell\ell'})]^{k} [1 - P_{\gamma}(\operatorname{Opt}_{\ell\ell'})]^{K-k},$$
(2.47)

with $P_{\gamma}(\cdot)$ given by (2.40) and $\operatorname{Opt}_{\ell\ell'}$ given by (2.38).

Note that by Standing Assumption the quantities $P_{\gamma}(\text{Opt}_{\ell\ell'})$ are < 1/2, so that the risks $\text{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}_K|H_1,...,H_L)$ go to 0 exponentially fast as $K \to \infty$.

2.3 DETECTORS AND DETECTOR-BASED TESTS

2.3.1 Detectors and their risks

Let Ω be an observation space, and \mathcal{P}_{χ} , $\chi = 1, 2$, be two families of probability distributions on Ω . By definition a *detector* associated with Ω is a real-valued function $\phi(\omega)$ of Ω . We associate with a detector ϕ and families \mathcal{P}_{χ} , $\chi = 1, 2$, risks defined as follows:

$$\begin{array}{rcl} \operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}] &=& \sup_{P \in \mathcal{P}_{1}} \int_{\Omega} \exp\{-\phi(\omega)\} P(d\omega) & (a) \\ \operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}] &=& \sup_{P \in \mathcal{P}_{2}} \int_{\Omega} \exp\{\phi(\omega)\} P(d\omega) & (b) \\ \operatorname{Risk}[\phi|\mathcal{P}_{1},\mathcal{P}_{2}] &=& \max[\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}], \operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}]] & (c) \end{array}$$

Given a detector ϕ , we can associate with it simple test \mathcal{T}_{ϕ} deciding, via observation $\omega \sim P$, on the hypotheses

$$H_1: P \in \mathcal{P}_1, \ H_2: P \in \mathcal{P}_2; \tag{2.49}$$

specifically, given observation $\omega \in \Omega$, the test \mathcal{T}_{ϕ} accepts H_1 and rejects H_2 whenever $\phi(\omega) \geq 0$, otherwise the test accepts H_2 and rejects H_1 .

Let us make the following immediate observation:

Proposition 2.16. Let Ω be an observation space, \mathcal{P}_{χ} , $\chi = 1, 2$, be two families of probability distributions on Ω , and ϕ be a detector. The risks of the test \mathcal{T}_{ϕ} associated with this detector satisfy

$$\operatorname{Risk}_{1}(\mathcal{T}_{\phi}|H_{1}, H_{2}) \leq \operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}];$$

$$\operatorname{Risk}_{2}(\mathcal{T}_{\phi}|H_{1}, H_{2}) \leq \operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}].$$
(2.50)

Proof. Let $\omega \sim P \in \mathcal{P}_1$. Then the *P*-probability of the event $\{\omega : \phi(\omega) < 0\}$ does not exceed Risk_ $[\phi|\mathcal{P}_1]$, since on the set $\{\omega : \phi(\omega) < 0\}$ the integrand in (2.48.*a*) is > 1, and this integrand is nonnegative everywhere, so that the integral in (2.48.*a*) is $\geq P\{\omega : \phi(\omega) < 0\}$. Recalling what \mathcal{T}_{ϕ} is, we see that the *P*-probability to reject H_1 is at most Risk_ $[\phi|\mathcal{P}_1]$, implying the first relation in (2.50). By similar argument, with (2.48.*b*) in the role of (2.48.*a*), when $\omega \sim P \in \mathcal{P}_2$, the *P*-probability of the event $\{\omega : \phi(\omega) \geq 0\}$ is upper-bounded by Risk_ $[\phi|\mathcal{P}_2]$, implying the second relation in (2.50).

2.3.2 Detector-based tests

Our current goal is to establish some basic properties of detector-based tests.

2.3.2.1 Structural properties of risks

Observe that the fact that ϵ_1 and ϵ_2 are upper bounds on the risks of a detector are expressed by system of *convex* constraints

$$\sup_{P \in \mathcal{P}_1} \int_{\Omega} \exp\{-\phi(\omega)\} P(d\omega) \le \epsilon_1 \quad (a)$$

$$\sup_{P \in \mathcal{P}_2} \int_{\Omega} \exp\{\phi(\omega)\} P(d\omega) \le \epsilon_2 \quad (b)$$
(2.51)

on ϵ_1 , ϵ_2 and $\phi(\cdot)$; this observation is useful, but not too useful, since the convex constraints in question usually are infinite-dimensional when $\phi(\cdot)$ is so, and are semiinfinite (suprema, over parameter ranging in infinite set, of parametric families of convex constraints), provided \mathcal{P}_1 or \mathcal{P}_2 are of infinite cardinalities; constraints of this type can be intractable computationally.

Another important observation is that the distributions P enter the constraints linearly; as a result, when passing from families of probability distributions \mathcal{P}_1 , \mathcal{P}_2 to their convex hulls, the risks of a detector remain intact.

2.3.2.2 Renormalization

Let Ω , \mathcal{P}_1 , \mathcal{P}_2 be the same as in Section 2.3.1, and let ϕ be a detector. When shifting this detector by a real a – passing from ϕ to the detector

$$\phi_a(\omega) = \phi(\omega) - a$$

- the risks clearly are updated as follows:

$$\begin{aligned} \operatorname{Risk}_{-}[\phi_{a}|\mathcal{P}_{1}] &= e^{a}\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}],\\ \operatorname{Risk}_{+}[\phi_{a}|\mathcal{P}_{2}] &= e^{-a}\operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}]. \end{aligned} (2.52)$$

We see that

When speaking about risks of a detector, what matters is the product

$$\operatorname{Risk}_{\odot}[\phi|\mathcal{P}_1, \mathcal{P}_2] := \operatorname{Risk}_{-}[\phi|\mathcal{P}_1]\operatorname{Risk}_{+}[\phi|\mathcal{P}_2]$$

of the risks, not these risks individually: by shifting the detector, we can redistribute this product between the factors in any way we want. In particular, we can always shift a detector to make it balanced, i.e., satisfying

$$\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}] = \operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}] = \operatorname{Risk}[\phi|\mathcal{P}_{1},\mathcal{P}_{2}].$$

When deciding on the hypotheses

$$H_1: P \in \mathcal{P}_1, \ H_2: P \in \mathcal{P}_2$$

on the distribution P of observation, the risk of the test \mathcal{T}_{ϕ} associated with a balanced detector ϕ is bounded by the risk Risk $[\phi|\mathcal{P}_1, \mathcal{P}_2]$ of the detector:

 $\operatorname{Risk}(\mathcal{T}_{\phi}|H_1, H_2) := \max\left[\operatorname{Risk}_1(\mathcal{T}_{\phi}|H_1, H_2), \operatorname{Risk}_2(\mathcal{T}_{\phi}|H_1, H_2)\right] \le \operatorname{Risk}[\phi|\mathcal{P}_1, \mathcal{P}_2].$

2.3.2.3 Detector-based testing from repeated observations

We are about to show that detector-based tests are perfectly well suited for passing from inferences based on *single* observation to those based on *repeated* observations.

Given K observation spaces Ω_k , $1 \leq k \leq K$, each equipped with pair $\mathcal{P}_{k,1}$, $\mathcal{P}_{k,2}$ of families of probability distributions, we can build a new observation space

$$\boldsymbol{\Omega}^{K} = \boldsymbol{\Omega}_{1} \times ... \times \boldsymbol{\Omega}_{K} = \{\boldsymbol{\omega}^{K} = (\boldsymbol{\omega}_{1},...,\boldsymbol{\omega}_{K}): \boldsymbol{\omega}_{k} \in \boldsymbol{\Omega}_{k}, k \leq K\}$$

and equip it with two families \mathcal{P}_{χ}^{K} , $\chi = 1, 2$, of probability distributions; distributions from \mathcal{P}_{χ}^{K} are exactly the product-type distributions $P = P_1 \times ... \times P_K$ with all factors P_k taken from $\mathcal{P}_{k,\chi}$. Observations $\omega^K = (\omega_1, ..., \omega_K)$ from Ω^K drawn from a distribution $P = P_1 \times ... \times P_K \in \mathcal{P}_{\chi}^{K}$ are nothing but collections of observations $\omega_k, k = 1, ..., K$, drawn, independently of each other, from distributions P_k . Now, given detectors $\phi_k(\cdot)$ on observation spaces Ω_k and setting

$$\phi^{(K)}(\omega^K) = \sum_{k=1}^K \phi_k(\omega_k) : \Omega^K \to \mathbf{R},$$

we clearly have

$$\operatorname{Risk}_{-}[\phi^{(K)}|\mathcal{P}_{1}^{K}] = \prod_{\substack{k=1\\K}}^{K} \operatorname{Risk}_{-}[\phi_{k}|\mathcal{P}_{k,1}],$$

$$\operatorname{Risk}_{+}[\phi^{(K)}|\mathcal{P}_{2}^{K}] = \prod_{\substack{k=1\\k=1}}^{K} \operatorname{Risk}_{+}[\phi_{k}|\mathcal{P}_{k,2}].$$
(2.53)

Let us look at some useful consequences of (2.53).

Stationary K-repeated observations. Consider the case of Section 2.1.3.1: we are given an observation space Ω and a positive integer K, and what we observe, is a sample $\omega^K = (\omega_1, ..., \omega_K)$ with $\omega_1, ..., \omega_K$ drawn, independently of each other, from some distribution P on Ω . Let now \mathcal{P}_1 , \mathcal{P}_2 , be two families of probability distributions on Ω ; we can associate with these families two hypotheses, $H_1^{\odot,K}$, $H_2^{\odot,K}$, on the distribution of K-repeated observation $\omega^K = (\omega_1, ..., \omega_K)$, with $H_{\chi}^{\odot,K}$ stating that $\omega_1, ..., \omega_K$ are drawn, independently of each other, from a distribution $P \in \mathcal{P}_{\chi}$. Given a detector ϕ on Ω , we can associate with it the detector

$$\phi^{(K)}(\omega^K) = \sum_{k=1}^K \phi(\omega_k)$$

on

$$\Omega^K:\underbrace{\Omega\times\ldots\times\Omega}_K.$$

Combining (2.53) and Proposition 2.16, we arrive at the following nice result:

Proposition 2.17. The risks of the simple test $\mathcal{T}_{\phi^{(K)}}$ deciding, given K-repeated observation $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ on the hypotheses

 $H_1^{\odot,K}: \omega_k, k \leq K$, are independently of each other drawn from a distribution $P \in \mathcal{P}_1$ $H_2^{\odot,K}: \omega_k, k \leq K$, are independently of each other drawn from a distribution $P \in \mathcal{P}_2$

according to the rule

$$\phi^{(K)}(\omega^K) := \sum_{k=1}^K \phi(\omega_k) \left\{ \begin{array}{ll} \geq 0 & \Rightarrow & accept \ H_1^{\odot,K} \\ < 0 & \Rightarrow & accept \ H_2^{\odot,K} \end{array} \right.$$

admit the upper bounds

$$\operatorname{Risk}_{1}(\mathcal{T}_{\phi^{(K)}}|H_{1}^{\odot,K}, H_{2}^{\odot,K}) \leq (\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}])^{K} \operatorname{Risk}_{2}(\mathcal{T}_{\phi^{(K)}}|H_{1}^{\odot,K}, H_{2}^{\odot,K}) \leq (\operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}])^{K}$$

$$(2.54)$$

Semi- and Quasi-Stationary K-repeated observations. Recall that Semi-Stationary and Quasi-Stationary K-repeated observations associated with a family \mathcal{P} of distributions on observation space Ω were defined in Sections 2.1.3.2 and 2.1.3.3, respectively. It turns out that Proposition 2.17 extends to quasi-stationary K-repeated observations:

Proposition 2.18. Let Ω be an observation space, \mathcal{P}_{χ} , $\chi = 1, 2$ be families of probability distributions on Ω , $\phi : \Omega \to \mathbf{R}$ be a detector, and K be a positive integer.

Families \mathcal{P}_{χ} , $\chi = 1, 2$, give rise to two hypotheses on the distribution P^{K} of quasi-stationary K-repeated observation ω^{K} :

$$H_{\chi}^{\otimes,K}: P^K \in \mathcal{P}_{\chi}^{\otimes,K} = \bigotimes_{k=1}^K \mathcal{P}_{\chi}, \ \chi = 1,2$$

(see Section 2.1.3.3), and ϕ gives rise to the detector

$$\phi^{(K)}(\omega^K) := \sum_{k=1}^K \phi(\omega_k).$$

The risks of the detector $\phi^{(K)}$ on the families $\mathcal{P}_{\chi}^{\otimes,K}$, $\chi = 1, 2$, can be upper-bounded as follows:

$$\begin{aligned}
\operatorname{Risk}_{-}[\phi^{(K)}|\mathcal{P}_{1}^{\otimes,K}] &\leq \left(\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}]\right)^{K}, \\
\operatorname{Risk}_{+}[\phi^{(K)}|\mathcal{P}_{2}^{\otimes,K}] &\leq \left(\operatorname{Risk}_{-}[\phi|\mathcal{P}_{2}]\right)^{K}.
\end{aligned}$$
(2.55)

Further, the detector $\phi^{(K)}$ induces simple test $\mathcal{T}_{\phi^{(K)}}$ deciding on $H_{\chi}^{\otimes,K}$, $\chi = 1, 2$ as follows: given ω^{K} , the test accepts $H_{1}^{\otimes,K}$ when $\phi^{(K)}(\omega^{K}) \geq 0$, and accepts $H_{2}^{\otimes,K}$ otherwise. The risks of this test can be upper-bounded as

$$\operatorname{Risk}_{1}(\mathcal{T}_{\phi^{(K)}}|H_{1}^{\otimes,K}, H_{2}^{\otimes,K}) \leq (\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}])^{K},$$

$$\operatorname{Risk}_{2}(\mathcal{T}_{\phi^{(K)}}|H_{1}^{\otimes,K}, H_{2}^{\otimes,K}) \leq (\operatorname{Risk}_{+}[\phi|\mathcal{P}_{2}])^{K}.$$
(2.56)

Finally, the above results remain intact when passing from quasi-stationary to semistationary K-repeated observations (that is, when replacing $\mathcal{P}_{\chi}^{\otimes,K}$ with $\mathcal{P}_{\chi}^{\oplus,K} = \bigoplus_{k=1}^{K} \mathcal{P}_{\chi}$ and $H_{\chi}^{\otimes,K}$ with the hypotheses $H_{\chi}^{\oplus,K}$ stating that the distribution of ω^{K} belongs to $\mathcal{P}_{\chi}^{\oplus,K}$, $\chi = 1, 2$).

Proof. All we need is to verify (2.55) – in view of Proposition 2.16, all other claims in Proposition 2.18 are immediate consequences of (2.55) and the inclusions $\mathcal{P}_{\chi}^{\oplus,K} \subset \mathcal{P}_{\chi}^{\otimes,K}, \ \chi = 1, 2$. Verification of (2.55) is as follows. Let $P^{K} \in \mathcal{P}_{1}^{\otimes,K}$, and let P^{K} be the distribution of random sequence $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ generated as follows: there exists a random sequence of driving factors $\zeta_{1}, ..., \zeta_{K}$ such that ω_{k} is a deterministic function of $\zeta^{k} = (\zeta_{1}, ..., \zeta_{k})$:

$$\omega_k = \theta_k(\zeta_1, \dots, \zeta_k),$$

and the conditional, $\zeta_1, ..., \zeta_{k-1}$ being given, distribution $P_{\omega_k|\zeta^{k-1}}$ belongs to \mathcal{P}_1 . Let P_{ζ^k} be the distribution of the first k driving factors, and $P_{\zeta_k|\zeta^{k-1}}$ be the conditional, $\zeta_1, ..., \zeta_{k-1}$ being given, distribution of ζ_k . Let us set

$$\psi^{(k)}(\zeta_1, ..., \zeta_k) = \sum_{t=1}^k \phi(\theta_t(\zeta_1, ..., \zeta_t)),$$

so that

$$\int_{\Omega^{K}} \exp\{-\phi^{(K)}(\omega^{k})\} P^{K}(d\omega^{K}) = \int \exp\{-\psi^{(K)}(\zeta^{K})\} P_{\zeta^{K}}(d\zeta^{K}).$$
(2.57)

On the other hand, denoting $C_0 = 1$, we have

$$C_{k} := \int \exp\{-\psi^{(k)}(\zeta^{k})\} P_{\zeta^{k}}(d\zeta^{k}) = \int \exp\{-\psi^{(k-1)}(\zeta^{k-1}) - \phi(\theta_{k}(\zeta^{k}))\} P_{\zeta^{k}}(d\zeta^{k})$$

$$= \int \exp\{-\psi^{(k-1)}(\zeta^{k-1})\} \underbrace{\left[\int \exp\{-\phi(\theta_{k}(\zeta^{k}))\} P_{\zeta_{k}|\zeta^{k-1}}(d\zeta_{k})\right]}_{= \int_{\Omega} \exp\{-\phi(\omega_{k})\} P_{\omega_{k}|\zeta^{k-1}}(d\omega_{k})} P_{\zeta^{k-1}}(d\zeta^{k-1})$$

$$\underset{(*)}{\leq} \operatorname{Risk}_{[\phi|\mathcal{P}_{1}]} \int \exp\{-\psi^{(k-1)}(\zeta^{k-1})\} P_{\zeta^{k-1}}(d\zeta^{k-1}) = \operatorname{Risk}_{[\phi|\mathcal{P}_{1}]} C_{k-1},$$

where (*) is due to the fact that the distribution $P_{\omega_k|\zeta^{k-1}}$ belongs to \mathcal{P}_1 . From the resulting recurrence we get

$$C_K \leq (\operatorname{Risk}_{-}[\phi|\mathcal{P}_1])^K$$
,

which combines with (2.57) to imply that

$$\int_{\Omega^{K}} \exp\{-\phi^{(K)}(\omega^{k})\} P^{K}(d\omega^{K}) \le (\operatorname{Risk}_{-}[\phi|\mathcal{P}_{1}])^{K}.$$

The latter inequality holds true for every distribution $P^K \in \mathcal{P}_{\chi}^{\otimes,K}$, and the first inequality in (2.55) follows. The second inequality in (2.55) is given by completely similar reasoning, with \mathcal{P}_2 in the role of \mathcal{P}_1 , and $-\phi$, $-\phi^{(K)}$ in the roles of ϕ , $\phi^{(K)}$, respectively.

The fact that observations ω_k under hypotheses $H_{\ell}^{\otimes,K}$, $\ell = 1, 2$ are related to "constant in time" families \mathcal{P}_{ℓ} has no importance here, and in fact the proof of Proposition 2.18 after absolutely evident modifications of wording allows to justify the following "non-stationary" version of Proposition:

Proposition 2.19. For k = 1, ..., K, let Ω_k be observation spaces, $\mathcal{P}_{\chi,k}$, $\chi = 1, 2$ be families of probability distributions on Ω_k , and $\phi_k : \Omega_k \to \mathbf{R}$ be detectors.

Families $\mathcal{P}_{\chi,k}$, $\chi = 1, 2$, give rise to quasi-direct products (see Section 2.1.3.3) $\mathcal{P}_{\chi}^{\otimes,K} = \bigotimes_{k=1}^{K} \mathcal{P}_{\chi,k}$ of the families $\mathcal{P}_{\chi,k}$ over $1 \leq k \leq K$, and thus to two hypotheses on the distribution P^{K} of observation $\omega^{K} = (\omega_{1}, ..., \omega_{K}) \in \Omega^{K} = \Omega_{1} \times ... \times \Omega_{K}$:

$$H_{\chi}^{\otimes,K}: P^K \in \mathcal{P}_{\chi}^{\otimes,K}, \, \chi = 1, 2.$$

and detectors ϕ_k , $1 \leq k \leq K$, give rise to the detector

$$\phi^K(\omega^K) := \sum_{k=1}^K \phi_k(\omega_k).$$

The risks of the detector ϕ^K on the families $\mathcal{P}_{\chi}^{\otimes,K}$, $\chi = 1, 2$, can be upper-bounded as follows:

$$\operatorname{Risk}_{-}[\phi^{K}|\mathcal{P}_{1}^{\otimes,K}] \leq \prod_{k=1}^{K} \operatorname{Risk}_{-}[\phi|\mathcal{P}_{1,k}], \\ \operatorname{Risk}_{+}[\phi^{K}|\mathcal{P}_{2}^{\otimes,K}] \leq \prod_{k=1}^{K} \operatorname{Risk}_{+}[\phi|\mathcal{P}_{2,K}].$$

$$(2.58)$$

Further, the detector ϕ^{K} induces simple test $\mathcal{T}_{\phi^{(K)}}$ deciding on $H_{\chi}^{\otimes,K}$, $\chi = 1, 2$ as follows: given ω^{K} , the test accepts $H_{1}^{\otimes,K}$ when $\phi^{K}(\omega^{K}) \geq 0$, and accepts $H_{2}^{\otimes,K}$ otherwise. The risks of this test can be upper-bounded as

$$\operatorname{Risk}_{1}(\mathcal{T}_{\phi^{K}}|H_{1}^{\otimes,K}, H_{2}^{\otimes,K}) \leq \prod_{k=1}^{K} \operatorname{Risk}_{-}[\phi|\mathcal{P}_{1,k}], \\ \operatorname{Risk}_{2}(\mathcal{T}_{\phi^{(K)}}|H_{1}^{\otimes,K}, H_{2}^{\otimes,K}) \leq \prod_{k=1}^{K} \operatorname{Risk}_{+}[\phi|\mathcal{P}_{2,k}].$$

$$(2.59)$$

Finally, the above results remain intact when passing from quasi-direct products to direct products of the families of distributions in question (that is, when replacing $\mathcal{P}_{\chi}^{\otimes,K}$ with $\mathcal{P}_{\chi}^{\oplus,K} = \bigoplus_{k=1}^{K} \mathcal{P}_{\chi,k}$ and $H_{\chi}^{\otimes,K}$ with the hypotheses $H_{\chi}^{\oplus,K}$ stating that the distribution of ω^{K} belongs to $\mathcal{P}_{\chi}^{\oplus,K}$, $\chi = 1, 2$).

2.3.2.4 Limits of performance of detector-based tests

We are about to demonstrate that as far as limits of performance of pairwise simple detector-based tests are concerned, these tests are nearly as good as simple tests can be.

Proposition 2.20. Let Ω be an observation space, and \mathcal{P}_{χ} , $\chi = 1, 2$, be families of probability distributions on Ω . Assume that for some $\epsilon \in (0, 1/2)$ "in the nature" there exists a simple test (deterministic or randomized) deciding on the hypotheses

$$H_1: P \in \mathcal{P}_1, \ H_2: P \in \mathcal{P}_2$$

on the distribution P of observation ω with risks $\leq \epsilon$:

$$\operatorname{Risk}_1(\mathcal{T}|H_1, H_2) \leq \epsilon \& \operatorname{Risk}_2(\mathcal{T}|H_1, H_2) \leq \epsilon$$

Then there exists a detector-based test \mathcal{T}_{ϕ} deciding on the same pair of hypotheses

with risk "comparable" with ϵ :

$$\operatorname{Risk}_{1}(\mathcal{T}_{\phi}|H_{1},H_{2}) \leq \epsilon^{+} \& \operatorname{Risk}_{2}(\mathcal{T}_{\phi}|H_{1},H_{2}) \leq \epsilon^{+}, \ \epsilon^{+} = 2\sqrt{\epsilon(1-\epsilon)}.$$
(2.60)

Proof. Let us prove the claim in the case when the test \mathcal{T} is deterministic; the case when this test is randomized is the subject of Exercise 2.64.

Let $\Omega_{\chi}, \chi = 1, 2$, be the sets of $\omega \in \Omega$ such that \mathcal{T} as "feeded" by observation ω accepts H_{χ} . Since \mathcal{T} is simple, Ω_1, Ω_2 split Ω into two non-overlapping parts, and since the risks of \mathcal{T} are $\leq \epsilon$, we have

(a)
$$\epsilon_2(P) := P\{\Omega_2\} \le \epsilon \,\forall P \in \mathcal{P}_1$$

(a) $\epsilon_1(P) := P\{\Omega_1\} \le \epsilon \,\forall P \in \mathcal{P}_2$

Let $\delta = \sqrt{(1-\epsilon)/\epsilon}$, so that $\delta \ge 1$ due to $0 < \epsilon \le 1/2$, and let

$$\psi(\omega) = \begin{cases} \delta, & \omega \in \Omega_1 \\ 1/\delta, & \omega \in \Omega_2 \end{cases}, \ \phi(\omega) = \ln(\psi(\omega)).$$

When $P \in \mathcal{P}_1$, we have

$$\int_{\Omega} \exp\{-\phi(\omega)\}P(d\omega) = \frac{1}{\delta}P\{\Omega_1\} + \delta P\{\Omega_2\} = \frac{1}{\delta} + \underbrace{\left[\delta - \frac{1}{\delta}\right]}_{\geq 0} \epsilon_2(P) \le \frac{1}{\delta} + \left[\delta - \frac{1}{\delta}\right]\epsilon = \epsilon^+,$$

whence $\operatorname{Risk}_{-}[\phi|\mathcal{P}_1] \leq \epsilon^+$. Similarly, when $P \in \mathcal{P}_2$, we have

$$\int_{\Omega} \exp\{\phi(\omega)\} P(d\omega) = \delta P\{\Omega_1\} + \frac{1}{\delta} P\{\Omega_2\} = \underbrace{\left[\delta - \frac{1}{\delta}\right]}_{\geq 0} \epsilon_1(P) + \frac{1}{\delta} \leq \left[\delta - \frac{1}{\delta}\right] \epsilon + \frac{1}{\delta} = \epsilon^+$$

whence $\operatorname{Risk}_{+}[\phi|\mathcal{P}_2] \leq \epsilon^+$.

Discussion. Proposition 2.20 states that we can restrict ourselves with detectorbased tests at the price of passing from risk ϵ exhibited by "the best test existing in the nature" to "comparable" risk $\epsilon^+ = 2\sqrt{\epsilon(1-\epsilon)}$. What we buy when sticking to detector-based tests are nice properties listed in Sections 2.3.2.1 - 2.3.2.3 and possibility to compute under favorable circumstances, see below, the best, in terms of their risk, among the detector-based tests; optimizing risk of a detector-based test turns out to be an essentially more realistic task than optimizing risk of a generaltype test. This being said, one can argue that treating ϵ and ϵ^+ "comparable" is a too optimistic attitude; for example, risk level $\epsilon = 0.01$ seems to be much more attractive than $[0.01]^+ \approx 0.2$. While passing from a test \mathcal{T} with risk 0.01 to a detector-based test \mathcal{T}_{ϕ} with risk 0.2 could indeed be a "heavy toll," there is some comfort in the fact that passing from a single observation to three of them (i.e., to 3-repeated, stationary or non-stationary alike, version of the original observation scheme), we can straightforwardly convert \mathcal{T}_{ϕ} into a test with risk $(0.2)^3 = 0.008 <$ 0.01, and passing to 6 observations, to make the risk less than 0.0001. On the other hand, seemingly the only way to convert a general-type single-observation test \mathcal{T} with risk 0.01 into a multi-observation test with essentially smaller risk is to pass to a Majority version of \mathcal{T} , see Section 2.2.3.2¹⁵. Computation shows that

 $^{^{15}}$ In Section 2.2.3.2, we dealt with "signal plus noise" observations and with specific test ${\cal T}$

2.4 SIMPLE OBSERVATION SCHEMES

2.4.1 Simple observation schemes – Motivation

A natural conclusion one can extract from the previous Section is that it makes sense, to say the least, to learn how to build detector-based tests with minimal risk. Thus, we arrive at the following design problem:

Given an observation space Ω and two families, \mathcal{P}_1 and \mathcal{P}_2 , of probability distributions on Ω , solve the optimization problem

$$Opt = \min_{\phi:\Omega \to \mathbf{R}} \max\left[\underbrace{\sup_{P \in \mathcal{P}_1} \int_{\Omega} e^{-\phi(\omega)} P(d\omega)}_{F[\phi]}, \underbrace{\sup_{P \in \mathcal{P}_2} \int_{\Omega} e^{\phi(\omega)} P(d\omega)}_{G[\phi]}\right] \quad (2.61)$$

While being convex, problem (2.61) typically is computationally intractable. First, it is infinite-dimensional – candidate solutions are multivariate functions; how to represent them in a computer, not speaking of how to optimize over them? Besides, the objective to be optimized is expressed in terms of suprema of infinitely many (provided \mathcal{P}_1 and/or \mathcal{P}_2 are infinite) expectations, and computing just a single expectation can be a difficult task... We are about to consider "favorable" cases – simple observation schemes – where (2.61) is efficiently solvable.

To arrive at the notion of a simple observation scheme, consider the case when all distributions from \mathcal{P}_1 , \mathcal{P}_2 admit densities taken w.r.t. some reference measure Π on Ω , and these densities are parameterized by "parameter" μ running through some parameter space \mathcal{M} , so that \mathcal{P}_1 is comprised of all distributions with densities $p_{\mu}(\cdot)$ and μ belonging to some subset M_1 of \mathcal{M} , while \mathcal{P}_2 is comprised of distributions with densities $p_{\mu}(\cdot)$ and μ belonging to another subset, M_2 , of \mathcal{M} . To save words, we shall identify distributions with their densities taken w.r.t. Π , so that

$$\mathcal{P}_{\chi} = \{ p_{\mu} : \mu \in M_{\chi} \}, \, \chi = 1, 2,$$

where $\{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}$ is a given "parametric" family of probability densities. Quotation marks in "parametric" reflect the fact that at this point in time, the "parameter" μ can be infinite-dimensional (e.g, we can parameterise a density by itself), so that assuming "parametric" representation of the distributions from \mathcal{P}_1 , \mathcal{P}_2 in fact does not restrict generality.

Our first observation is that in our "parametric" setup, we can rewrite problem

given by Euclidean separation. Straightforward inspection of the construction and the proof of Proposition 2.8 makes it clear that the construction is applicable to a whatever simple test \mathcal{T} , and that the risk of the resulting multi-observation test obeys the upper bound in (2.25), with the risk of \mathcal{T} in the role of ϵ_{\star} .

78

(2.61) equivalently as

$$\ln(\text{Opt}) = \min_{\phi:\Omega \to \mathbf{R}} \sup_{\mu \in M_1, \nu \in M_2} \underbrace{\frac{1}{2} \left[\ln\left(\int_{\Omega} e^{-\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega)\right) + \ln\left(\int_{\Omega} e^{\phi(\omega)} p_{\nu}(\omega) \Pi(d\omega)\right) \right]}_{\Phi(\phi;\mu,\nu)}.$$
(2.62)

Indeed, when shifting ϕ by a constant: $\phi(\cdot) \mapsto \phi(\cdot) - a$, the positive quantities $F[\phi]$ and $G[\phi]$ participating in (2.61) are multiplied by e^a and e^{-a} , respectively, and their product remains intact. It follows that to minimize over ϕ the maximum of $F[\phi]$ and $G[\phi]$ (this is what (2.61) wants of us) is exactly the same as to minimize over ϕ the quantity $H[\phi] := \sqrt{F[\phi]G[\phi]}$. Indeed, a candidate solution ϕ to the problem $\min_{\phi} H[\phi]$ can be balanced – shifted by a constant to ensure $F[\phi] = G[\phi]$, and this balancing does not change $H[\cdot]$; as a result, minimizing H over all ϕ is the same as minimizing H over balanced ϕ , and the latter problem clearly is equivalent to (2.61). It remains to note that (2.62) is nothing but the problem of minimizing $\ln(H[\phi])$.

Now, (2.62) is a min-max problem – a problem of the generic form

$$\min_{u \in U} \max_{v \in V} \Psi(u, v).$$

Problems of this type (at least, finite-dimensional ones) are computationally tractable when the domain of the minimization argument is convex and the cost function Ψ is convex in the minimization argument (this indeed is the case for (2.62)), and the domain of the maximization argument is convex, and the cost function is concave in this argument (this not necessarily is the case for (2.62)). Simple observation schemes we are about to define are, essentially, the schemes where the just outlined requirements of finite dimensionality and convexity-concavity indeed are met.

2.4.2 Simple observation schemes – Definition

Consider the situation where we are given

- 1. A Polish (complete separable metric) observation space Ω equipped with σ -finite σ -additive Borel reference measure Π such that the support of Π is the entire Ω . Those not fully comfortable with some of the notions from the previous sentence can be assured that the only observation spaces we indeed shall deal with are pretty simple:
 - $\Omega = \mathbf{R}^d$ equipped with the Lebesgue measure Π , and
 - a finite or countable set Ω which is discrete (distances between distinct points are equal to 1) and is equipped with the counting measure Π .
- 2. A parametric family $\{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}$ of probability densities, taken w.r.t. Π , such that
 - the space \mathcal{M} of parameters is a convex set in some \mathbb{R}^n which coincides with its relative interior,
 - the function $p_{\mu}(\omega) : \mathcal{M} \times \Omega \to \mathbf{R}$ is continuous in (μ, ω) and positive everywhere.
- 3. A finite-dimensional linear subspace \mathcal{F} of the space of continuous functions on

LECTURE 2

 Ω such that

- \mathcal{F} contains constants,
- all functions of the form $\ln(p_{\mu}(\omega)/p_{\nu}(\omega))$ with $\mu, \nu \in \mathcal{M}$ are contained in \mathcal{F} ;
- for every $\phi(\cdot) \in \mathcal{F}$, the function

$$\ln\left(\int_{\Omega} e^{\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega)\right)$$

is real-valued and *concave* on \mathcal{M} .

In this situation we call the collection

$$(\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F})$$

a simple observation scheme (s.o.s. for short).

Nondegenerate simple o.s. We call a simple observation scheme *nondegenerate*, if the mapping $\mu \mapsto p_{\mu}$ is an embedding: whenever $\mu, \mu' \in \mathcal{M}$ and $\mu \neq \mu'$, we have $p_{\mu} \neq p_{\mu'}$.

2.4.3 Simple observation schemes – Examples

We are about to list basic examples of s.o.s.'s.

2.4.3.1 Gaussian observation scheme

In Gaussian o.s.,

- the observation space (Ω, Π) is the space \mathbf{R}^d with Lebesgue measure,
- the family $\{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}$ is the family of Gaussian densities $\mathcal{N}(\mu, \Theta)$, with fixed positive definite covariance matrix Θ , distributions from the family are parameterized by their expectations μ . Thus,

$$\mathcal{M} = \mathbf{R}^d, \ p_\mu(\omega) = \frac{1}{(2\pi)^{d/2} \sqrt{\operatorname{Det}(\Theta)}} \exp\{-\frac{(\omega-\mu)^T \Theta^{-1}(\omega-\mu)}{2}\};$$

• the family \mathcal{F} is the family of all affine functions on \mathbf{R}^d .

It is immediately seen that Gaussian o.s. meets all requirements imposed on a simple o.s. For example,

$$\ln(p_{\mu}(\omega)/p_{\nu}(\omega)) = (\nu - \mu)^{T} \Theta^{-1} \omega + \frac{1}{2} \left[\nu^{T} \Theta^{-1} \nu - \mu^{T} \Theta^{-1} \mu \right]$$

is an affine function of ω and thus belongs to \mathcal{F} . Besides this, a function $\phi(\cdot) \in \mathcal{F}$ is affine: $\phi(\omega) = a^T \omega + b$, implying that

$$\begin{aligned} f(\mu) &:= & \ln\left(\int_{\mathbf{R}^d} e^{\phi(\omega)} p_{\mu}(\omega) d\omega\right) = \ln\left(\mathbf{E}_{\xi \sim \mathcal{N}(0, I_d)} \left\{\exp\{a^T(\Theta^{1/2}\xi + \mu) + b\}\right\}\right) \\ &= & a^T \mu + b + const, \\ const &= & \ln\left(\mathbf{E}_{\xi \sim \mathcal{N}(0, I_d)} \left\{\exp\{a^T \Theta^{1/2}\xi\}\right\}\right) = \frac{1}{2}a^T \Theta a \end{aligned}$$

is affine (and thus concave) function of μ .

As we remember from Lecture 1, Gaussian o.s. is responsible for the standard *signal processing* model where one is given a noisy observation

$$\omega = Ax + \xi \qquad [\xi \sim \mathcal{N}(0, \Theta)]$$

of the image Ax of unknown signal $x \in \mathbf{R}^n$ under linear transformation with known $d \times n$ sensing matrix, and the goal is to infer from this observation some knowledge about x. In this situation, a hypothesis that x belongs to some set X translates into the hypothesis that the observation ω is drawn from Gaussian distribution with known covariance matrix Θ and expectation known to belong to the set $M = \{\mu = Ax : x \in X\}$, so that deciding on various hypotheses on where x is located reduces to deciding on hypotheses on the distribution in Gaussian o.s.

2.4.3.2 Poisson observation scheme

In Poisson observation scheme,

- the observation space Ω is the set \mathbf{Z}^d_+ of *d*-dimensional vectors with nonnegative integer entries, and this set is equipped with the counting measure,
- the family $\{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}$ is the family of product-type Poisson distributions with positive parameters. In other words,

$$\mathcal{M} = \{\mu \in \mathbf{R}^d : \mu > 0\}, p_{\mu}(\omega) = \frac{\mu_1^{\omega_1} \mu_2^{\omega_2} \dots \mu_d^{\omega_d}}{\omega_1! \omega_2! \dots \omega_d!} e^{-\mu_1 - \mu_2 - \dots - \mu_d}, \, \omega \in \mathbf{Z}_+^d,$$

that is, random variable $\omega \sim p_{\mu}$, $\mu \in \mathcal{M}$, is *d*-dimensional vector with independent random entries, and *i*-th of the entries is $\omega_i \sim \text{Poisson}(\mu_i)$;

• the space \mathcal{F} is comprised of affine functions on \mathbf{Z}_d^+ .

It is immediately seen that Poisson o.s. is simple. For example,

$$\ln(p_{\mu}(\omega)/p_{\nu}(\omega)) = \sum_{i=1}^{d} \ln(\mu_{i}/\nu_{i})\omega_{i} - \sum_{i=1}^{d} [\mu_{i} - \nu_{i}]$$

is affine function of ω and thus belongs to \mathcal{F} . Besides this, a function $\phi \in \mathcal{F}$ is affine: $\phi(\omega) = a^T \omega + b$, implying that the function

$$\begin{aligned} f(\mu) &:= \ln\left(\int_{\Omega} e^{\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega)\right) = \ln\left(\sum_{\omega \in \mathbf{Z}_{+}^{d}} e^{a^{T}\omega + b} \prod_{i=1}^{d} \frac{\mu_{i}^{\omega_{i}} e^{-\mu_{i}}}{\omega_{i}!}\right) \\ &= b + \ln\left(\prod_{i=1}^{d} \left[e^{-\mu_{i}} \sum_{s=0}^{\infty} \frac{[e^{a_{i}} \mu_{i}]^{s}}{s!}\right]\right) = b + \sum_{i=1}^{d} \ln(\exp\{e^{a_{i}} \mu_{i} - \mu_{i}\}) \\ &= \sum_{i} [e^{a_{i}} - 1] \mu_{i} + b \end{aligned}$$

is affine (and thus concave) function of μ .

Poisson observation scheme is responsible for *Poisson Imaging*. This is the situation where there are n "sources of customers;" arrivals of customers at source i are independent of what happens at other sources, and inter-arrival times at source j are independent random variables with exponential, with parameter λ_j , random variables, so that the number of customers arriving at source j in a unit time interval is Poisson random variable with parameter λ_j . Now, there are d "servers", and a customer arrived at source j is dispatched to server i with some given probability A_{ij} , $\sum_i A_{ij} \leq 1$; with probability $1 - \sum_i A_{ij}$, such a customer

leaves the system. Needless to say, the dispatches are independent of each other and of the arrival processes. What we observe is the vector $\omega = (\omega_1, ..., \omega_d)$, where ω_i is the number of customers dispatched to server *i* on the time horizon [0, 1]. It is easy to verify that in the just described situation, the entries ω_i in ω indeed are independent of each other Poisson random variables with Poisson parameters

$$\mu_i = \sum_{j=1}^n A_{ij} \lambda_j$$

In what is called *Poisson Imaging*, one is given a random observation ω of the above type along with *sensing matrix* $A = [A_{ij}]$, and the goal is to use the observation to infer conclusions on the parameter $\mu = A\lambda$ and underlying this parameter "signal" λ .

Poisson imaging is has several important applications¹⁶, for example, in Positron Emission Tomography (PET).



In PET, a patient is injected radioactive tracer and is placed in PET tomograph, which can be thought of as a cylinder with surface split into small detector cells. The tracer disintegrates, and every disintegration act produces a positron which immediately annihilates with a nearby electron, producing two γ -quants flying at the speed of light in two opposite directions along a line ("line of response" – LOR) with completely random orientation. Eventually, each of the γ -quants hits its own detector cell. When two detector cells are "simultaneously" hit (in fact - hit within a short time interval, like 10^{-8} sec), this event – *coincidence* – and the serial number of the *bin* (pair of detectors) where the hits were observed are registered; observing a coincidence in some bin, we know that somewhere on the line linking the detector cells from the bin a disintegration act took place. The data collected in a PET study are the numbers of coincidences registered in every one of the bins; discretizing the field of view (patient's body) into small 3D cubes (voxels), an accurate enough model of the data is a realization ω of random vector with independent Poisson entries $\omega_i \sim \text{Poisson}(\mu_i)$, with μ_i given by

$$\mu_i = \sum_{j=1}^n p_{ij} \lambda_j$$

where λ_i is proportional to the amount of tracer in voxel j, and p_{ij} is the probability

 $^{^{16}\}text{in}$ all these applications, the signal λ we ultimately are interested in is an image, this is where "Imaging" comes from.

for LOR emanating from voxel j to be registered in bin i (these probabilities can be computed given the geometry of PET device). The tracer is selected in such a way that in the body it concentrates in the areas of interest (say, the areas of high metabolic activity when tumor is sought), and the goal of the study is to infer from the observation ω conclusions on the density of the tracer. The characteristic feature of PET as compared to other types of tomography is that with properly selected tracer, this technique allows to visualize metabolic activity, and not only the anatomy of tissues in the body. Now, PET fits perfectly well the above "dispatching customers" story, with disintegration acts taking place in voxel j in the role of customers arriving in location j and bins in the role of servers; the arrival intensities are (proportional to) the amounts λ_j of tracer in voxels, and the random dispatch of customers to servers corresponds to random orientation of LOR's (in reality, the nature draws their directions from the uniform distribution on the unit sphere in 3D).

It is worthy of noting that there are two other real life applications of Poisson Imaging: Large Binocular Telescope and Nanoscale Fluorescent Microscopy ¹⁷.

2.4.3.3 Discrete observation scheme

In Discrete observation scheme,

- the observation space is a finite set $\Omega = \{1, ..., d\}$ equipped with counting measure,
- the family $\{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}$ is comprised of all non-vanishing distributions on Ω , that is,

$$\mathcal{M} = \{ \mu \in \mathbf{R}^d : \mu > 0, \sum_{\omega \in \Omega} \mu_\omega = 1 \}, \ p_\mu(\omega) = \mu_\omega, \omega \in \Omega;$$

• \mathcal{F} is the space of all real-valued functions on the finite set Ω .

Clearly, Discrete o.s. is simple; for example, the function

$$f(\mu) := \ln\left(\int_{\Omega} e^{\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega)\right) = \ln\left(\sum_{\omega \in \Omega} e^{\phi(\omega)} \mu_{\omega}\right)$$

indeed is concave in $\mu \in \mathcal{M}$.

2.4.3.4 Direct products of simple observation schemes

Given K simple observation schemes

$$\mathcal{O}_k = (\Omega_k, \Pi_k; \{p_{\mu,k}(\cdot) : \mu \in \mathcal{M}_k\}; \mathcal{F}_k), \ 1 \le k \le K,$$

¹⁷Large Binocular Telescope is a cutting edge instrument for high-resolution optical/infrared astronomical imaging; it is the subject of huge ongoing international project, see http://www.lbto.org. Nanoscale Fluorescent Microscopy (a.k.a. Poisson Biophotonics) is a revolutionary tool for cell imaging trigged by the advent of techniques [15, 73, 76, 132] (2014 Nobel Prize in Chemistry) allowing to break the diffraction barrier and to view biological molecules "at work" at a resolution 10-20 nm, yielding entirely new insights into the signalling and transport processes within cells.

we can define their *direct product*

$$\mathcal{O}^{K} = \prod_{k=1}^{K} \mathcal{O}_{k} = (\Omega^{K}, \Pi^{K}; \{p_{\mu} : \mu \in \mathcal{M}^{K}\}; \mathcal{F}^{K})$$

by modeling the situation where our observation is a tuple $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ with components ω_{k} yielded, independently of each other, by observation schemes \mathcal{O}_{k} , namely, as follows:

- The observation space Ω^{K} is the direct product of observations spaces $\Omega_{1}, ..., \Omega_{K}$, and the reference measure Π^{K} is the product of the measures $\Pi_{1}, ..., \Pi_{K}$;
- The parameter space \mathcal{M}^K is the direct product of partial parameter spaces $\mathcal{M}_1, ..., \mathcal{M}_K$, and the distribution $p_{\mu}(\omega^K)$ associated with parameter

$$\mu = (\mu_1, \mu_2, ..., \mu_K) \in \mathcal{M}^K = \mathcal{M}_1 \times ... \times \mathcal{M}_K$$

is the probability distribution on Ω^{K} with the density

$$p_{\mu}(\omega^{K}) = \prod_{k=1}^{K} p_{\mu,k}(\omega_{k})$$

w.r.t. Π^{K} . In other words, random observation $\omega^{K} \sim p_{\mu}$ is a sample of observations $\omega_{1}, ..., \omega_{K}$, drawn, independently of each other, from the distributions $p_{\mu_{1},1}, p_{\mu_{2},2}, ..., p_{\mu_{K},K}$;

• The space \mathcal{F}^{K} is comprised of all *separable* functions

$$\phi(\omega^K) = \sum_{k=1}^K \phi_k(\omega_k)$$

with $\phi_k(\cdot) \in \mathcal{F}_k, 1 \leq k \leq K$.

It is immediately seen that the direct product of simple observation o.s.'s is simple. When all factors \mathcal{O}_k , $1 \leq k \leq K$, are identical to simple o.s.

$$\mathcal{O} = (\Omega, \Pi; \{ p_{\mu} : \mu \in \mathcal{M} \}; \mathcal{F}),$$

the direct product of the factors can be "truncated" to yield the K-th power (called also the stationary K-repeated version) of \mathcal{O} , denoted

$$[\mathcal{O}]^K = (\Omega^K, \Pi^K; \{p_\mu^{(K)} : \mu \in \mathcal{M}\}; \mathcal{F}^{(K)})$$

and defined as follows:

• Ω^K and Π^K are exactly the same as in the direct product:

$$\Omega^K = \underbrace{\Omega \times \ldots \times \Omega}_K, \ \Pi^K = \underbrace{\Pi \times \ldots \times \Pi}_K;$$

• the parameter space is \mathcal{M} rather than the direct product of K copies of \mathcal{M} , and

the densities are

$$p_{\mu}^{(K)}(\omega^{K} = (\omega_{1}, ..., \omega_{K})) = \prod_{k=1}^{K} p_{\mu}(\omega_{k});$$

in other words, random observations $\omega^K \sim p_{\mu}^{(K)}$ are K-element samples with components drawn, independently of each other, from p_{μ} ;

• the space $\mathcal{F}^{(K)}$ is comprised of separable functions

$$\phi^{(K)}(\omega^K) = \sum_{k=1}^K \phi(\omega_k)$$

with identical components belonging to \mathcal{F} (i.e., $\phi \in \mathcal{F}$).

It is immediately seen that a power of simple o.s. is simple.

Remark 2.21. Gaussian, Poisson and Discrete o.s.'s clearly are nondegenerate. It is also clear that the direct product of nondegenerate o.s.'s is nondegenerate.

2.4.4 Simple observation schemes – Main result

We are about to demonstrate that when deciding on *convex*, in some precise sense to be specified below, hypotheses in *simple* observation schemes, optimal detectors can be found efficiently by solving *convex-concave saddle point problems*.

We start with "executive summary" on convex-concave saddle point problems.

2.4.4.1 Executive summary of convex-concave saddle point problems

The results to follow are absolutely standard, and their proofs can be found in all textbooks on the subject, see, e.g., [11, Section D.4].

Let U and V be nonempty sets, and $\Phi : U \times V \to \mathbf{R}$ be a function. These data define an *antagonistic game* of two players, I and II, where player I selects a point $u \in U$, and player II selects a point $v \in V$; as an outcome of these selections, player I pays to player II the sum $\Phi(u, v)$. Clearly, the player I is interested to minimize this payment, and player II – to maximize the payment. The data U, V, Φ are known to the players in advance, and the question is, what should be their selections.

When the player I makes his selection u first, and player II makes his selection v with u already known, player I should be ready to pay for a selection $u \in U$ the toll as large as

$$\overline{\Phi}(u) = \sup_{v \in V} \Phi(u, v).$$

In this situation, a risk-averse player I would select u by minimizing the above worst-case payment, by solving the *primal* problem

$$Opt(P) = \inf_{u \in U} \overline{\Phi}(u) = \inf_{u \in U} \sup_{v \in V} \Phi(u, v)$$
(P)

associated with the data U, V, Φ .

Similarly, if player II makes his selection v first, and player I selects u after v becomes known, player II should be ready to get, as a result of selecting $v \in V$, the

amount as small as

$$\underline{\Phi}(v) = \inf_{u \in U} \phi(u, v).$$

In this situation, a risk-averse player II would select v by maximizing the above worst-case payment, by solving the *dual* problem

$$Opt(D) = \sup_{v \in V} \underline{\Phi}(v) = \sup_{v \in V} \inf_{u \in U} \Phi(u, v) \tag{D}$$

Intuitively, the first situation is less preferable for player I than the second one, so that his guaranteed payment in the first situation, that is, Opt(P), should be \geq his guaranteed payment, Opt(D), in the second situation:

$$\operatorname{Opt}(P) := \inf_{u \in U} \sup_{v \in V} \Phi(u, v) \ge \sup_{v \in V} \inf_{u \in U} \Phi(u, v) =: \operatorname{Opt}(D);$$
(2.63)

this fact, called *Weak Duality*, indeed is true.

The central question related to the game is what should the players do when making their selections simultaneously, with no knowledge of what is selected by the adversary. There is a case when this question has a completely satisfactory answer – this is the case where Φ has a *saddle point* on $U \times V$.

Definition 2.22. A point $(u_*, v_*) \in U \times V$ is called a saddle point ¹⁸ of function $\Phi(u, v) : U \times V \to \mathbf{R}$, if Φ as a function of $u \in U$ attains at this point its minimum, and as a function of $v \in V$ – its maximum, that is, if

$$\Phi(u, v_*) \ge \Phi(u_*, v_*) \ge \Phi(u_*, v) \ \forall (u \in U, v \in V).$$

From the viewpoint of our game, a saddle point (u_*, v_*) is an equilibrium: when one of the players sticks to the selection stemming from this point, the other one has no incentive to deviate from his selection stemming from the point: if player II selects v_* , there is no reason for player I to deviate from selecting u_* , since with another selection, his loss (the payment) can only increase; similarly, when player I selects u_* , there is no reason for player II to deviate from v_* , since with any other selection, his gain (the payment) can only decrease. As a result, if the cost function Φ has a saddle point on $U \times V$, this saddle point (u_*, v_*) can be considered as a solution to the game, as the pair of preferred selections of rational players. It can be easily seen that while Φ can have many saddle points, the values of Φ at all these points are equal to each other, let us denote their common value by SadVal. If (u_*, v_*) is a saddle point and player I selects $u = u_*$, his worst, over selections $v \in V$ of player II, loss is SadVal, and if player I selects a whatever $u \in U$, his worst-case, over the selections of player II, loss can be only \geq SadVal. Similarly, when player II selects $v = v_*$, his worst-case, over the selections of player I, gain is SadVal, and if player II selects a whatever $v \in V$, his worst-case, over the selections of player I, gain can be only < SadVal.

Existence of saddle points of Φ (min in $u \in U$, max in $v \in V$) can be expressed in terms of the primal problem (P) and the dual problem (P):

¹⁸ more precisely, "saddle point (min in $u \in U$, max in $v \in V$);" we will usually skip the clarification in parentheses, since it always will be clear from the context what are the minimization variables and what are the maximization ones.

Proposition 2.23. Φ has saddle point (min in $u \in U$, max in $v \in V$) if and only if problems (P) and (D) are solvable with equal optimal values:

$$\operatorname{Opt}(P) := \inf_{u \in U} \sup_{v \in V} \Phi(u, v) = \sup_{v \in V} \inf_{u \in U} \Phi(u, v) =: \operatorname{Opt}(D).$$
(2.64)

Whenever this is the case, the saddle points of Φ are exactly the pairs (u_*, v_*) comprised of optimal solutions to problems (P) and (D), and the value of Φ at every one of these points is the common value SadVal of Opt(P) and Opt(D).

Existence of a saddle point of a function is "rare commodity;" the standard sufficient condition for it is convexity-concavity of Φ coupled with convexity of U and V; the precise statement is as follows:

Theorem 2.24. [Sion-Kakutani, see, e.g., [11, Theorems D.4.3, D.4.4]] Let $U \subset \mathbf{R}^m, V \subset \mathbf{R}^n$ be nonempty closed convex sets, with W bounded, and let $\Phi : U \times V \to \mathbf{R}$ be continuous function which is convex in $u \in U$ for every fixed $v \in V$, and is concave in $v \in V$ for every fixed $u \in U$. Then the equality (2.64) holds true (although it may happen that $Opt(P) = Opt(D) = -\infty$).

If, in addition, Φ is coercive in u, meaning that the level sets

$$\{u \in U : \Phi(u, v) \le a\}$$

are bounded for every $a \in \mathbf{R}$ and $v \in V$ (equivalently: for every $v \in V$, $\Phi(u_i, v) \to +\infty$ along every sequence $u_i \in U$ going to ∞ : $||u_i|| \to \infty$ as $i \to \infty$), then Φ admits saddle points (min in $u \in U$, max in $v \in V$).

Note that the "true" Sion-Kakutani Theorem is a bit stronger than Theorem 2.24; the latter, however, covers all our related needs.

2.4.4.2 Main Result

Theorem 2.25. Let

$$\mathcal{O} = (\Omega, \Pi; \{ p_{\mu} : \mu \in \mathcal{M} \}; \mathcal{F})$$

be a simple observation scheme, and let M_1, M_2 be nonempty compact convex subsets of \mathcal{M} . Then

(i) The function

$$\Phi(\phi, [\mu; \nu]) = \frac{1}{2} \left[\ln \left(\int_{\Omega} e^{-\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega) \right) + \ln \left(\int_{\Omega} e^{\phi(\omega)} p_{\nu}(\omega) \Pi(d\omega) \right) \right] :$$

$$\mathcal{F} \times (M_1 \times M_2) \to \mathbf{R}$$
(2.65)

is continuous on its domain, is convex in $\phi(\cdot) \in \mathcal{F}$, concave in $[\mu; \nu] \in M_1 \times M_2$, and possesses a saddle point (min in $\phi \in \mathcal{F}$, max in $[\mu; \nu] \in M_1 \times M_2$) ($\phi_*(\cdot), [\mu_*; \nu_*]$) on $\mathcal{F} \times (M_1 \times M_2)$. ϕ_* w.l.o.g. can be assumed to satisfy the relation¹⁹

$$\int_{\Omega} \exp\{-\phi_*(\omega)\} p_{\mu_*}(\omega) \Pi(d\omega) = \int_{\Omega} \exp\{\phi_*(\omega)\} p_{\mu_*}(\omega) \Pi(d\omega).$$
(2.66)

¹⁹Note that \mathcal{F} contains constants, and shifting by a constant the ϕ -component of a saddle point of Φ and keeping its $[\mu; \nu]$ -component intact, we clearly get another saddle point of Φ .

Denoting the common value of the two quantities in (2.66) by ε_{\star} , the saddle point value

$$\min_{\phi \in \mathcal{F}} \max_{[\mu;\nu] \in M_1 \times M_2} \Phi(\phi, [\mu;\nu])$$

is $\ln(\varepsilon_{\star})$. Besides this, setting $\phi_{\star}^{a}(\cdot) = \phi_{\star}(\cdot) - a$, one has

implying, in view of Proposition 2.16, that when deciding via an observation $\omega \in \Omega$ on the hypotheses

$$H_{\chi}: \omega \sim p_{\mu} \text{ with } \mu \in M_{\chi}, \quad \chi = 1, 2,$$

the risks of the simple test $\mathcal{T}_{\phi^a_*}$ based on the detector ϕ^a_* can be upper-bounded as follows:

$$\operatorname{Risk}_1(\mathcal{T}_{\phi^a_*}|H_1, H_2) \le \exp\{a\}\varepsilon_\star, \ \operatorname{Risk}_2(\mathcal{T}_{\phi^a_*}|H_1, H_2) \le \exp\{-a\}\varepsilon_\star.$$
(2.68)

Besides this, ϕ_*, ε_* form an optimal solution to the optimization problem

$$\min_{\phi,\epsilon} \left\{ \epsilon : \int_{\Omega} e^{-\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega) \le \epsilon \,\forall \mu \in M_1 \\ \int_{\Omega} e^{\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega) \le \epsilon \,\forall \mu \in M_2 \right\}$$
(2.69)

(the minimum in (2.69) is taken over all $\epsilon > 0$ and all Π -measurable functions $\phi(\cdot)$, not just over $\phi \in \mathcal{F}$).

(ii) The dual problem associated with the saddle point data Φ , \mathcal{F} , $M_1 \times M_2$ is

$$\max_{\mu \in M_1, \nu \in M_2} \underline{\Phi}(\mu, \nu) := \inf_{\phi \in \mathcal{F}} \Phi(\phi; [\mu; \nu]).$$
(D)

The objective in this problem is in fact the logarithm of Hellinger affinity of p_{μ} and p_{ν} :

$$\underline{\Phi}(\mu,\nu) = \ln\left(\int_{\Omega} \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}\Pi(d\omega)\right),\tag{2.70}$$

and this objective is concave and continuous on $M_1 \times M_2$.

The (μ, ν) -components of saddle points of Φ are exactly the maximizers (μ_*, ν_*) of the concave function $\underline{\Phi}$ on $M_1 \times M_2$. Given such a maximizer $[\mu_*; \nu_*]$ and setting

$$\phi_*(\omega) = \frac{1}{2} \ln(p_{\mu_*}(\omega)/p_{\nu_*}(\omega))$$
(2.71)

we get a saddle point $(\phi_*, [\mu_*; \nu_*])$ of Φ satisfying (2.66).

(iii) Let $[\mu_*; \nu_*]$ be a maximizer of $\underline{\Phi}$ over $M_1 \times M_2$. Let, further, $\epsilon \in [0, 1/2]$ be such that there exists a (whatever, perhaps randomized) test for deciding via observation $\omega \in \Omega$ on two simple hypotheses

$$(A): \omega \sim p(\cdot) := p_{\mu_*}(\cdot), \quad (B): \omega \sim q(\cdot) := p_{\nu_*}(\cdot)$$
(2.72)

with total risk $\leq 2\epsilon$. Then

$$\varepsilon_{\star} \leq 2\sqrt{\epsilon(1-\epsilon)}.$$

In other words, if the simple hypotheses (A), (B) can be decided, by a whatever test, with total risk 2ϵ , then the risks of the simple test with detector ϕ_* given by (2.71) 88

LECTURE 2

on the composite hypotheses H_1 , H_2 do not exceed $2\sqrt{\epsilon(1-\epsilon)}$.

Proof. 1⁰. Since \mathcal{O} is a simple o.s., the function $\Phi(\phi, [\mu; \nu])$ given by (2.65) is a well defined real-valued function on $\mathcal{F} \times (\mathcal{M} \times \mathcal{M})$ which is concave in $[\mu; \nu]$; convexity of the function in $\phi \in \mathcal{F}$ is evident. Since both \mathcal{F} and \mathcal{M} are convex sets coinciding with their relative interiors, convexity-concavity and real valuedness of Φ on $\mathcal{F} \times (\mathcal{M} \times \mathcal{M})$ imply the continuity of Φ on the indicated domain. As a consequence, Φ is convex-concave continuous real-valued function on $\mathcal{F} \times (\mathcal{M}_1 \times \mathcal{M}_2)$.

Now let

$$\underline{\Phi}(\mu,\nu) = \inf_{\phi \in \mathcal{F}} \Phi(\phi, [\mu;\nu]).$$
(2.73)

Note that $\underline{\Phi}$, being the infimum of a family of concave functions of $[\mu; \nu] \in \mathcal{M} \times \mathcal{M}$, is concave on $\mathcal{M} \times \mathcal{M}$. We claim that for $\mu, \nu \in \mathcal{M}$ the function

$$\phi_{\mu,\nu}(\omega) = \frac{1}{2} \ln(p_{\mu}(\omega)/p_{\nu}(\omega))$$

(which, by definition of a simple o.s., belongs to \mathcal{F}) is an optimal solution to the right hand side minimization problem in (2.73), so that

$$\forall (\mu \in M_1, \nu \in M_2) :$$

$$\underline{\Phi}([x; y]) := \inf_{\phi \in \mathcal{F}} \Phi(\phi, [\mu; \nu]) = \Phi(\phi_{\mu, \nu}, [\mu; \nu]) = \ln\left(\int_{\Omega} \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}\Pi(d\omega)\right).$$
(2.74)

Indeed, we have

$$\exp\{-\phi_{\mu,\nu}(\omega)\}p_{\mu}(\omega) = \exp\{\phi_{\mu,\nu}(\omega)\}p_{\nu}(\omega) = g(\omega) := \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)},$$

whence $\Phi(\phi_{\mu,\nu}, [\mu; \nu]) = \ln \left(\int_{\Omega} g(\omega) \Pi(d\omega) \right)$. On the other hand, for $\phi(\cdot) = \phi_{\mu,\nu}(\cdot) + \delta(\cdot) \in \mathcal{F}$ we have

$$\begin{split} &\int_{\Omega} g(\omega) \Pi(d\omega) = \int_{\Omega} \left[\sqrt{g(\omega)} \exp\{-\delta(\omega)/2\} \right] \left[\sqrt{g(\omega)} \exp\{\delta(\omega)/2\} \right] \Pi(d\omega) \\ (a) &\leq \left(\int_{\Omega} g(\omega) \exp\{-\delta(\omega)\} \Pi(d\omega) \right)^{1/2} \left(\int_{\Omega} g(\omega) \exp\{\delta(\omega)\} \Pi(d\omega) \right)^{1/2} \\ &= \left(\int_{\Omega} \exp\{-\phi(\omega)\} p_{\mu}(\omega) \Pi(d\omega) \right)^{1/2} \left(\int_{\Omega} \exp\{\phi(\omega)\} p_{\nu}(\omega) \Pi(d\omega) \right)^{1/2} \\ (b) &\Rightarrow \ln \left(\int_{\Omega} g(\omega) \Pi(d\omega) \right) \leq \Phi(\phi, [\mu; \nu]), \end{split}$$

and thus $\Phi(\phi_{\mu,\nu}, [\mu; \nu]) \leq \Phi(\phi, [\mu; \nu])$ for every $\phi \in \mathcal{F}$.

Remark 2.26. Note that the above reasoning did not use the fact that the minimization in the right hand side of (2.73) is over $\phi \in \mathcal{F}$; in fact, this reasoning shows that $\phi_{\mu,\nu}(\cdot)$ minimizes $\Phi(\phi, [\mu; \nu])$ over all functions ϕ for which the integrals $\int_{\Omega} \exp\{-\phi(\omega)\}p_{\mu}(\omega)\Pi(d\omega)$ and $\int_{\Omega} \exp\{\phi(\omega)\}p_{\nu}(\omega)\Pi(d\omega)$ exist.

Remark 2.27. Note that the inequality in (b) can be equality only when the inequality in (a) is so. In other words, if $\bar{\phi}$ is a minimizer of $\Phi(\phi, [\mu; \nu])$ over $\phi \in \mathcal{F}$, setting $\delta(\cdot) = \bar{\phi}(\cdot) - \phi_{\mu,\nu}(\cdot)$, the functions $\sqrt{g(\omega)} \exp\{-\delta(\omega)/2\}$ and $\sqrt{g(\omega)} \exp\{\delta(\omega)/2\}$, considered as elements of $L_2[\Omega, \Pi]$, are proportional to each other. Since g is positive and g, δ are continuous, while the support of Π is the entire Ω , this " L_2 -proportionality" means that the functions in question differ by a constant factor, or, which is the same, that $\delta(\cdot)$ is constant. Thus, the minimizers of $\Phi(\phi, [\mu; \nu])$ over $\phi \in \mathcal{F}$ are exactly the functions of the form $\phi(\omega) = \phi_{\mu,\nu}(\omega) + const$.

2⁰. We are about to verify that $\Phi(\phi, [\mu; \nu])$ has a saddle point (min in $\phi \in \mathcal{F}$, max in $[\mu; \nu] \in M_1 \times M_2$). Indeed, observe, first, that on the domain of Φ it holds

$$\Phi(\phi(\cdot) + a, [\mu; \nu]) = \Phi(\phi(\cdot), [\mu; \nu]) \ \forall (a \in \mathbf{R}, \phi \in \mathcal{F}).$$

$$(2.75)$$

Let us select somehow $\bar{\mu} \in \mathcal{M}$, and let $\bar{\Pi}$ be the measure on Ω with density $p_{\bar{\mu}}$ w.r.t. Π . For $\phi \in \mathcal{F}$, the integrals $\int_{\Omega} e^{\pm \phi(\omega)} \bar{\Pi}(d\omega)$ are finite (since \mathcal{O} is simple), implying that

 $\phi \in L_1[\Omega, \overline{\Pi}]$; note also that $\overline{\Pi}$ is a probabilistic measure. Let now $\mathcal{F}_0 = \{\phi \in \mathcal{F} : \int_\Omega \phi(\omega)\overline{\Pi}(d\omega) = 0\}$, so that \mathcal{F}_0 is a linear subspace in \mathcal{F} , and all functions $\phi \in \mathcal{F}$ can be obtained by shifts of functions from \mathcal{F}_0 by constants. Invoking (2.75), to prove the existence of a saddle point of Φ on $\mathcal{F} \times (M_1 \times M_2)$ is exactly the same as to prove the existence of a saddle point of Φ on $\mathcal{F}_0 \times (M_1 \times M_2)$. Let us verify that $\Phi(\phi, [\mu; \nu])$ indeed has a saddle point on $\mathcal{F}_0 \times (M_1 \times M_2)$. $M_1 \times M_2$ is a convex compact set, and Φ is continuous on $\mathcal{F}_0 \times (M_1 \times M_2)$ and convex-concave; invoking Sion-Kakutani Theorem we see that all we need in order to verify the existence of a saddle point is to show that Φ is coercive in the first argument, that is, for every fixed $[\mu; \nu] \in M_1 \times M_2$ one has $\Phi(\phi, [\mu; \nu]) \to +\infty$ as $\phi \in \mathcal{F}_0$ and $\|\phi\| \to \infty$ (whatever be the norm $\|\cdot\|$ on \mathcal{F}_0 ; recall that \mathcal{F}_0 is a finite-dimensional linear space). Setting

$$\Theta(\phi) = \Phi(\phi, [\mu; \nu]) = \frac{1}{2} \left[\ln \left(\int_{\omega} e^{-\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega) \right) + \ln \left(\int_{\omega} e^{\phi(\omega)} p_{\nu}(\omega) \Pi(d\omega) \right) \right]$$

and taking into account that Θ is convex and finite on \mathcal{F}_0 , in order to prove that Θ is coercive, it suffices to verify that $\Theta(t\phi) \to \infty$, $t \to \infty$, for every nonzero $\phi \in \mathcal{F}_0$, which is evident: since $\int_{\Omega} \phi(\omega) \overline{\Pi}(d\omega) = 0$ and ϕ is nonzero, we have $\int_{\Omega} \max[\phi(\omega), 0] \overline{\Pi}(d\omega) = \int_{\Omega} \max[-\phi(\omega), 0] \overline{\Pi}(d\omega) > 0$, whence $\phi > 0$ and $\phi < 0$ on sets of Π -positive measure, so that $\Theta(t\phi) \to \infty$ as $t \to \infty$ due to the fact that both $p_{\mu}(\cdot)$ and $p_{\nu}(\cdot)$ are positive everywhere.

3⁰. Now let $(\phi_*(\cdot); [\mu_*; \nu_*])$ be a saddle point of Φ on $\mathcal{F} \times (M_1 \times M_2)$. Shifting, if necessary, $\phi_*(\cdot)$ by a constant (by (2.75), this does not affect the fact that $(\phi_*, [\mu_*; \nu_*])$ is a saddle point of Φ), we can assume that

$$\varepsilon_* := \int_{\Omega} \exp\{-\phi_*(\omega)\} p_{\mu_*}(\omega) \Pi(d\omega) = \int_{\Omega} \exp\{\phi_*(\omega)\} p_{\nu_*}(\omega) \Pi(d\omega), \qquad (2.76)$$

so that the saddle point value of Φ is

$$\Phi_* := \max_{[\mu;\nu] \in M_1 \times M_2} \min_{\phi \in \mathcal{F}} \Phi(\phi, [\mu;\nu]) = \Phi(\phi_*, [\mu_*;\nu_*]) = \ln(\varepsilon_\star).$$
(2.77)

as claimed in item (i) of Theorem.

Now let us prove (2.67). For $\mu \in M_1$, we have

$$\begin{aligned} \ln(\varepsilon_{\star}) &= \Phi_{\star} \ge \Phi(\phi_{\star}, [\mu; \nu_{\star}]) \\ &= \frac{1}{2} \ln\left(\int_{\Omega} \exp\{-\phi_{\star}(\omega)\} p_{\mu}(\omega) \Pi(d\omega)\right) + \frac{1}{2} \ln\left(\int_{\Omega} \exp\{\phi_{\star}(\omega)\} p_{\nu_{\star}}(\omega) \Pi(d\omega)\right) \\ &= \frac{1}{2} \ln\left(\int_{\Omega} \exp\{-\phi_{\star}(\omega)\} p_{\mu}(\omega) P(d\omega)\right) + \frac{1}{2} \ln(\varepsilon_{\star}), \end{aligned}$$

whence $\ln\left(\int_{\Omega} \exp\{-\phi_*^a(\omega)\}p_\mu(\omega)\Pi(d\omega)\right) = \ln\left(\int_{\Omega} \exp\{-\phi_*(\omega)\}p_\mu(\omega)P(d\omega)\right) + a \le \ln(\varepsilon_\star) + a$, and (2.67.*a*) follows. Similarly, when $\nu \in M_2$, we have

$$\begin{aligned} \ln(\varepsilon_{\star}) &= \Phi_{\star} \geq \Phi(\phi_{\star}, [\mu_{\star}; \nu]) \\ &= \frac{1}{2} \ln\left(\int_{\Omega} \exp\{-\phi_{\star}(\omega)\} p_{\mu_{\star}}(\omega) \Pi(d\omega)\right) + \frac{1}{2} \ln\left(\int_{\Omega} \exp\{\phi_{\star}(\omega)\} p_{\nu}(\omega) \Pi(d\omega)\right) \\ &= \frac{1}{2} \ln(\varepsilon_{\star}) + \frac{1}{2} \ln\left(\int_{\Omega} \exp\{\phi_{\star}(\omega)\} p_{\nu}(\omega) \Pi(d\omega)\right), \end{aligned}$$

so that $\ln\left(\int_{\Omega} \exp\{\phi_*^a(\omega)\}p_{\nu}(\omega)\Pi(d\omega)\right) = \ln\left(\int_{\Omega} \exp\{\phi_*(\omega)\}p_{\nu}(\omega)\Pi(d\omega)\right) - a \le \ln(\varepsilon_{\star}) - a$, and (2.67.b) follows.

We have proved all claims in item (i), except for the claim that the just defined ϕ_*, ε_* form an optimal solution to (2.69). Note that by (2.67) as applied with a = 0, the pair in question is feasible for (2.69). Assuming that the problem admits a feasible solution $(\bar{\phi}, \epsilon)$ with $\epsilon < \varepsilon_*$, let us lead this assumption to a contradiction. Note that $\bar{\phi}$ should be such that

$$\int_{\Omega} e^{-\bar{\phi}(\omega)} p_{\mu_{*}}(\omega) \Pi(d\omega) < \varepsilon_{\star} \& \int_{\Omega} e^{\bar{\phi}(\omega)} p_{\nu_{*}}(\omega) \Pi(d\omega) < \varepsilon_{\star},$$

and consequently $\Phi(\bar{\phi}, [\mu_*; \nu_*]) < \ln(\varepsilon_*)$. On the other hand, Remark 2.26 says that $\Phi(\bar{\phi}, [\mu_*; \nu_*])$ cannot be less than $\min_{\phi \in \mathcal{F}} \Phi(\phi, [\mu_*; \nu_*])$, and the latter quantity is $\Phi(\phi_*, [\mu_*; \nu_*])$ due to the fact that $(\phi_*, [\mu_*; \nu_*])$ is a saddle point of Φ on $\mathcal{F} \times (M_1 \times M_2)$. Thus, assuming that the optimal value in (2.69) is $< \varepsilon_*$, we conclude that $\Phi(\phi_*, [\mu_*; \nu_*]) \leq \Phi(\bar{\phi}, [\mu_*; \nu_*]) < \ln(\varepsilon_*)$, contradicting (2.77). Item (i) of Theorem 2.25 is proved.

4⁰. Let us prove item (ii) of Theorem 2.25. Relation (2.70) and concavity of the right hand side of this relation in $[\mu; \nu]$ were already proved; moreover, these relations were proved in the range $\mathcal{M} \times \mathcal{M}$ of $[\mu; \nu]$. Since this range coincides with its relative interior, the real-valued concave function $\underline{\Phi}$ is continuous in $\mathcal{M} \times \mathcal{M}$ and thus is continuous in $M_1 \times M_2$. Next, let ϕ_* be the ϕ -component of a saddle point of Φ on $\mathcal{F} \times (M_1 \times M_2)$ (we already know that a saddle point exists). Invoking Proposition 2.23, the $[\mu; \nu]$ -components of saddle points of Φ on $\mathcal{F} \times (M_1 \times M_2)$ are exactly the maximizers of $\underline{\Phi}$ on $M_1 \times M_2$. Let $[\mu_*; \nu_*]$ be such a maximizer; by the same Proposition 2.23, $(\phi_*, [\mu_*; \nu_*])$ is a saddle point of Φ , whence $\Phi(\phi, [\mu_*; \nu_*])$ attains its minimum over $\phi \in \mathcal{F}$ at $\phi = \phi_*$. We have also seen that $\Phi(\phi, [\mu_*; \nu_*])$ attains its minimum over $\phi \in \mathcal{F}$ at $\phi = \phi_{\mu_*,\nu_*}$. These observations combine with Remark 2.27 to imply that ϕ_* and ϕ_{μ_*,ν_*} differ by a constant, which, in view of (2.75), means that $(\phi_{\mu_*,\nu_*}, [\mu_*; \nu_*])$ is a saddle point of Φ along with $(\phi_*, [\mu_*; \nu_*])$. (ii) is proved.

5⁰. It remains to prove item (iii) of Theorem 2.25. In the notation from (iii), simple hypotheses (A) and (B) can be decided with the total risk $\leq 2\epsilon$, and therefore, by Proposition 2.2,

$$2\bar{\epsilon} := \int_{\Omega} \min[p(\omega), q(\omega)] \Pi(d\omega) \le 2\epsilon.$$

On the other hand, we have seen that the saddle point value of Φ is $\ln(\varepsilon_{\star})$; since $[\mu_{*};\nu_{\star}]$ is a component of a saddle point of Φ , it follows that $\min_{\phi \in \mathcal{F}} \Phi(\phi, [\mu_{*};\nu_{\star}]) = \ln(\varepsilon_{\star})$. The left hand side in this equality, as we know from item 1⁰, is $\Phi(\phi_{x_{\star},y_{\star}}, [x_{*};y_{\star}])$, and we arrive at $\ln(\varepsilon_{\star}) = \Phi(\frac{1}{2}\ln(p_{\mu_{\star}}(\cdot)/p_{\nu_{\star}}(\cdot)), [\mu_{*};\nu_{\star}]) = \ln\left(\int_{\Omega} \sqrt{p_{\mu_{\star}}(\omega)p_{\nu_{\star}}(\omega)}\Pi(d\omega)\right)$, so that $\varepsilon_{\star} = \int_{\Omega} \sqrt{p_{\mu_{\star}}(\omega)p_{\nu_{\star}}(\omega)}\Pi(d\omega) = \int_{\Omega} \sqrt{p(\omega)q(\omega)}\Pi(d\omega)$. We now have

$$\begin{split} \varepsilon_{\star} &= \int_{\Omega} \sqrt{p(\omega)q(\omega)} \Pi(d\omega) = \int_{\Omega} \sqrt{\min[p(\omega), q(\omega)]} \sqrt{\max[p(\omega), q(\omega)]} \Pi(d\omega) \\ &\leq \left(\int_{\Omega} \min[p(\omega), q(\omega)] \Pi(d\omega)\right)^{1/2} \left(\int_{\Omega} \max[p(\omega), q(\omega)] \Pi(d\omega)\right)^{1/2} \\ &= \left(\int_{\Omega} \min[p(\omega), q(\omega)] \Pi(d\omega)\right)^{1/2} \left(\int_{\Omega} (p(\omega) + q(\omega) - \min[p(\omega), q(\omega)]) \Pi(d\omega)\right)^{1/2} \\ &= \sqrt{2\overline{\epsilon}(2 - 2\overline{\epsilon})} \leq 2\sqrt{(1 - \epsilon)\epsilon}, \end{split}$$

where the concluding inequality is due to $\bar{\epsilon} \leq \epsilon \leq 1/2$. (iii) is proved, and the proof of Theorem 2.25 is complete.

Remark 2.28. Assume that we are under the premise of Theorem 2.25 and that the simple o.s. in question is nondegenerate (see Section 2.4.2). Then $\varepsilon_{\star} < 1$ if and only if the sets M_1 and M_2 do not intersect.

Indeed, by Theorem 2.25.i, $\ln(\varepsilon_{\star})$ is the saddle point value of $\Phi(\phi, [\mu; \nu])$ on $\mathcal{F} \times (M_1 \times M_2)$, or, which is the same by Theorem 2.25.ii, the maximum of the function (2.70) on $M_1 \times M_2$; since saddle points exist, this maximum is achieved at some pair $[\mu; \nu] \in M_1 \times M_2$. Since (2.70) clearly is ≤ 0 , we conclude that $\varepsilon_{\star} \leq 1$ and the equality takes place if and only if $\int_{\Omega} \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}\Pi(d\omega) = 1$ for some $\mu \in M_1$ and $\nu \in M_2$, or, which is the same, $\int_{\Omega} (\sqrt{p_{\mu}(\omega)} - \sqrt{p_{\nu}(\omega)})^2 \Pi(d\omega) = 0$ for these μ and ν . Since $p_{\mu}(\cdot)$ are continuous and the support of Π is the entire Ω , the latter can happen if and only if $p_{\mu} = p_{\nu}$ for our μ, ν , or, by nondegeneracy of \mathcal{O} , if and only if $M_1 \cap M_2 \neq \emptyset$.

2.4.5 Simple observation schemes – Examples of optimal detectors

Theorem 2.25.i states that when the observation scheme

$$\mathcal{O} = (\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F})$$

is simple and we are interested to decide on a pair of hypotheses on the distribution of observation $\omega \in \Omega$,

$$H_{\chi}: \omega \sim p_{\mu}$$
 with $\mu \in M_{\chi}, \chi = 1, 2$

and the hypotheses are convex, meaning that the underlying parameter sets M_{χ} are convex and compact, building optimal, in terms of its risk, detector ϕ_* – that is, solving (in general, semi-infinite and infinite-dimensional) optimization problem (2.69) reduces to solving the usual finite-dimensional convex problem. Specifically, an optimal solution (ϕ_*, ε_*) can be built as follows:

1. We solve optimization problem

$$Opt = \max_{\mu \in M_1, \nu \in M_2} \left[\underline{\Phi}(\mu, \nu) := \ln \left(\int_{\Omega} \sqrt{p_{\mu}(\omega) p_{\nu}(\omega)} \Pi(d\omega) \right) \right];$$
(2.78)

of maximizing Hellinger affinity (the quantity under the logarithm) of a pair of distributions obeying H_1 and H_2 , respectively; for a simple o.s., the objective in this problem is concave and continuous, and optimal solutions do exist;

2. (Any) optimal solution $[\mu_*; \nu_*]$ to (2.78) gives rise to an optimal detector ϕ_* and its risk ε_* , according to

$$\phi_*(\omega) = \frac{1}{2} \ln\left(\frac{p_{\mu_*}(\omega)}{p_{\nu_*}(\omega)}\right), \ \varepsilon_\star = \exp\{\text{Opt}\}.$$
(2.79)

The risks of the simple test \mathcal{T}_{ϕ_*} associated with the above detector and deciding on H_1 , H_2 , satisfy the bounds

$$\max\left[\operatorname{Risk}_{1}(\mathcal{T}_{\phi_{*}}|H_{1},H_{2}),\operatorname{Risk}_{2}(\mathcal{T}_{\phi_{*}}|H_{1},H_{2})\right] \leq \varepsilon_{\star},\tag{2.80}$$

and the test is *near-optimal*, meaning that whenever the hypotheses H_1 , H_2 (and in fact – even two simple hypotheses stating that $\omega \sim p_{\mu_*}$ and $\omega \sim p_{\nu_*}$, respectively) can be decided upon by a test with total risk $\leq 2\epsilon$, \mathcal{T}_{ϕ_*} exhibits "comparable" risk:

$$\varepsilon_{\star} \le 2\sqrt{\epsilon(1-\epsilon)}.\tag{2.81}$$

Note that the test \mathcal{T}_{ϕ_*} is just the maximum likelihood test induced by the probability densities p_{μ_*} and p_{ν_*} .

Note that after we know that (ϕ_*, ε_*) form an optimal solution to (2.69), some kind of near-optimality of the test \mathcal{T}_{ϕ_*} is guaranteed already by Proposition 2.20; specifically, by this Proposition, whenever in the nature there exists a test \mathcal{T} which decides on H_1, H_2 with risks Risk₁, Risk₂ bounded by some $\epsilon \leq 1/2$, the upper bound ε_* on the risks of \mathcal{T}_{ϕ_*} can be bounded according to (2.81). Our now nearoptimality statement is a bit stronger: first, we allow \mathcal{T} to have the total risk $\leq 2\epsilon$, which is weaker than to have both risks $\leq \epsilon$; second, and more important, now 2ϵ should upper-bound the total risk of \mathcal{T} on a pair of *simple* hypotheses "embedded" into the hypotheses H_1, H_2 ; both these modifications extend the family of tests \mathcal{T} to which we compare the test \mathcal{T}_{ϕ_*} , and thus enrich the comparison.

Let us look how the above recipe works for our basic simple o.s.'s.

2.4.5.1 Gaussian o.s.

When \mathcal{O} is a Gaussian o.s., that is, $\{p_{\mu} : \mu \in \mathcal{M}\}\$ are Gaussian densities with expectations $\mu \in \mathcal{M} = \mathbf{R}^d$ and common positive definite covariance matrix Θ , and \mathcal{F} is the family of affine functions on $\Omega = \mathbf{R}^d$,

- M_1, M_2 can be arbitrary nonempty convex compact subsets of \mathbf{R}^d ,
- problem (2.78) becomes the convex optimization problem

$$Opt = -\min_{\mu \in M_1, \nu \in M_2} \frac{(\mu - \nu)^T \Theta^{-1} (\mu - \nu)}{8}$$
(2.82)

the optimal detector φ_{*} and the upper bound ε_{*} on its risks given by an optimal solution (μ_{*}, ν_{*}) to (2.82) are

$$\phi_*(\omega) = \frac{1}{2} [\mu_* - \nu_*]^T \Theta^{-1} [\omega - w], \ w = \frac{1}{2} [\mu_* + \nu_*]$$

$$\varepsilon_* = \exp\{-\frac{[\mu_* - \nu_*]\Theta^{-1} [\mu_* - \nu_*]}{8}\}$$
(2.83)

Note that when $\Theta = I_d$, the test \mathcal{T}_{ϕ_*} becomes exactly the optimal test from Example 2.4. The upper bound on the risks of this test established in Example 2.4 (in our present notation, this bound is $\operatorname{Erf}(\frac{1}{2} \| \mu_* - \nu_* \|_2)$) is slightly better than the bound $\varepsilon_* = \exp\{-\|\mu_* - \nu_*\|_2^2/8\}$ given by (2.83) when $\Theta = I_d$. Note, however, that when speaking about the distance $\delta = \|\mu_* - \nu_*\|_2$ between M_1 and M_2 allowing for a test with risks $\leq \epsilon \ll 1$, the results of Example 2.4) and (2.83) say nearly the same: Example 2.4 says that δ should be $\geq 2\operatorname{ErfInv}(\epsilon)$, where $\operatorname{ErfInv}(\epsilon)$ is the Inverse Error function:

$$\operatorname{Erf}(\operatorname{ErfInv}(\epsilon)) \equiv \epsilon, \ 0 < \epsilon < 1,$$

and (2.83) says that δ should be $\geq 2\sqrt{2\ln(1/\epsilon)}$. When $\epsilon \to +0$, the ratio of these two lower bounds on δ tends to 1.

It should be noted that our general construction of optimal detectors as applied to Gaussian o.s. and a pair of convex hypotheses results in *exactly* optimal test and can be analyzed directly, without any "science" (see Example 2.4).

2.4.5.2 Poisson o.s.

When \mathcal{O} is a Poisson o.s., that is, $\mathcal{M} = \mathbf{R}_{++}^d$ is the interior of nonnegative orthant in \mathbf{R}^d , and $p_{\mu}, \mu \in \mathcal{M}$, is the density

$$p_{\mu}(\omega) = \prod_{i} \left(\frac{\mu_{i}^{\omega_{i}}}{\omega_{i}!} e^{-\mu_{i}} \right), \ \omega = (\omega_{!}, ..., \omega_{d}) \in \mathbf{Z}_{+}^{d}$$

taken w.r.t. the counting measure Π on $\Omega = \mathbf{Z}_{+}^{d}$, and \mathcal{F} is the family of affine functions on Ω , the recipe from the beginning of Section 2.4.5 reads as follows:

• M_1, M_2 can be arbitrary nonempty convex compact subsets of $\mathbf{R}_{++}^d = \{x \in \mathbf{R}^d : x > 0\};$
• problem (2.78) becomes the convex optimization problem

Opt =
$$-\min_{\mu \in M_1, \nu \in M_2} \frac{1}{2} \sum_{i=1}^{a} \left(\sqrt{\mu_i} - \sqrt{\nu_i} \right)^2;$$
 (2.84)

the optimal detector φ_{*} and the upper bound ε_{*} on its risks given by an optimal solution (μ^{*}, ν^{*}) to (2.84) are

$$\phi_{*}(\omega) = \frac{1}{2} \sum_{i=1}^{d} \ln\left(\frac{\mu_{i}^{*}}{\nu_{i}^{*}}\right) \omega_{i} + \frac{1}{2} \sum_{i=1}^{d} [\nu_{i}^{*} - \mu_{i}^{*}], \qquad (2.85)$$

$$\varepsilon_{\star} = e^{\text{Opt}}$$

2.4.5.3 Discrete o.s.

When \mathcal{O} is a Discrete o.s., that is, $\Omega = \{1, ..., d\}$, Π is a counting measure on Ω , $\mathcal{M} = \{\mu \in \mathbf{R}^d : \mu > 0, \sum_i \mu_i = 1\}$ and

$$p_{\mu}(\omega) = \mu_{\omega}, \, \omega = 1, ..., d, \, \mu \in \mathcal{M},$$

the recipe from the beginning of Section 2.4.5 reads as follows:

- M_1, M_2 can be arbitrary nonempty convex compact subsets of the relative interior \mathcal{M} of the probabilistic simplex,
- problem (2.78) is equivalent to the convex program

$$\varepsilon_{\star} = \max_{\mu \in M_1, \nu \in M_2} \sum_{i=1}^d \sqrt{\mu_i \nu_i}; \qquad (2.86)$$

• the optimal detector ϕ_* given by an optimal solution (μ^*, ν^*) to (2.84) is

$$\phi_*(\omega) = \frac{1}{2} \ln \left(\frac{\mu_{\omega}^*}{\nu_{\omega}^*} \right), \qquad (2.87)$$

and the upper bound ε_{\star} on the risks of this detector is given by (2.86).

2.4.5.4 K-th power of simple o.s.

Recall that K-th power of a simple o.s. $\mathcal{O} = (\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F})$ (see Section 2.4.3.4) is the o.s.

$$[\mathcal{O}]^K = (\Omega^K, \Pi^K; \{p_\mu^{(K)}: \mu \in \mathcal{M}\}; \mathcal{F}^{(K)})$$

where Ω^{K} is the direct product of K copies of Ω , Π^{K} is the product of K copies of Π , the densities $p_{\mu}^{(K)}$ are product densities induced by K copies of density p_{μ} , $\mu \in \mathcal{M}$:

$$p_{\mu}^{(K)}(\omega^{K} = (\omega_{1}, ..., \omega_{K})) = \prod_{k=1}^{K} p_{\mu}(\omega_{k}),$$

and $\mathcal{F}^{(K)}$ is comprised of functions

$$\phi^{(K)}(\omega^K = (\omega_1, ..., \omega_K)) = \sum_{k=1}^K \phi(\omega_k)$$

stemming from functions $\phi \in \mathcal{F}$. Clearly, $[\mathcal{O}]^K$ is the observation scheme describing the stationary K-repeated observations $\omega^K = (\omega_1, ..., \omega_K)$ with ω_k stemming from the o.s. \mathcal{O} , see Section 2.3.2.3. As we remember, $[\mathcal{O}]^K$ is simple provided that \mathcal{O} is so.

Assuming \mathcal{O} simple, it is immediately seen that as applied to the o.s. $[\mathcal{O}]^K$, the recipe from the beginning of Section 2.4.5 reads as follows:

- M_1, M_2 can be arbitrary nonempty convex compact subsets of \mathcal{M} , and the corresponding hypotheses, $H_{\chi}^K, \chi = 1, 2$, state that the components ω_k of observation $\omega^K = (\omega_1, ..., \omega_K)$ are independently of each other drawn from distribution p_{μ} with $\mu \in M_1$ (hypothesis H_1^K) or $\mu \in M_2$ (hypothesis H_2^K).
- problem (2.78) is the convex program

$$Opt(K) = \max_{\mu \in M_1, \nu \in M_2} \underbrace{\ln\left(\int_{\Omega^K} \sqrt{p_{\mu}^{(K)}(\omega^K)p_{\nu}^{(K)}(\omega^K)}\Pi^K(d\Omega)\right)}_{\equiv K \ln\left(\int_{\Omega} \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}\Pi(d\omega)\right)}$$
(D_K)

implying that any optimal solution to the "single-observation" problem (D_1) associated with M_1 , M_2 is optimal for the "K-observation" problem (D_K) associated with M_1 , M_2 , and Opt(K) = KOpt(1);

• the optimal detector $\phi_*^{(K)}$ given by an optimal solution (μ_*, ν_*) to (D_1) (this solution is optimal for (D_K) as well) is

$$\phi_*^{(K)}(\omega^K) = \sum_{k=1}^K \phi_*(\omega_k),
\phi_*(\omega) = \frac{1}{2} \ln\left(\frac{p_{\mu_*}(\omega)}{p_{\nu_*}(\omega)}\right),$$
(2.88)

and the upper bound $\varepsilon_{\star}(K)$ on the risks of the detector $\phi_{\star}^{(K)}$ on the pair of families of distributions obeying hypotheses H_1^K , resp., H_2^K , is

$$\varepsilon_{\star}(K) = e^{\operatorname{Opt}(K)} = e^{K\operatorname{Opt}(1)} = [\epsilon_{\star}(1)]^{K}.$$
(2.89)

The just outlined results on powers of simple observation schemes allow to express near-optimality of detector-based tests in simple o.s.'s in a nicer form, specifically, as follows.

Proposition 2.29. Let $\mathcal{O} = (\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F})$ be a simple observation scheme, M_1 , M_2 be two nonempty convex compact subsets of \mathcal{M} , and let (μ_*, ν_*) be an optimal solution to the convex optimization problem (cf. Theorem 2.25)

Opt =
$$\max_{\mu \in M_1, \nu \in M_2} \ln \left(\int_{\Omega} \sqrt{p_{\mu}(\omega) p_{\nu}(\omega)} \Pi(d\omega) \right).$$

Let ϕ_* and ϕ_*^K be single- and K-observation detectors induced by (μ_*, ν_*) via (2.88). Let $\epsilon \in (0, 1/2)$, and assume that for some positive integer K in the nature

exists a simple test \mathcal{T}^K deciding via K i.i.d. observations $\omega^K = (\omega_1, ..., \omega_K)$ with $\omega_k \sim p_{\mu}$, for some unknown $\mu \in \mathcal{M}$, on the hypotheses

$$H_{\chi}^{(K)}: \mu \in M_{\chi}, \ \chi = 1, 2,$$

with risks $Risk_1$, $Risk_2$ not exceeding ϵ . Then setting

$$K_{+} = \int \frac{2}{1 - \ln(4(1 - \epsilon)) / \ln(1/\epsilon)} K \lfloor,$$

the simple test $\mathcal{T}_{\phi_*^{(K_+)}}$ utilizing K_+ i.i.d. observations decides on $H_1^{(K_+)}$, $H_2^{(K_+)}$ with risks $\leq \epsilon$. Note that K_+ "is of order of K:" $K_+/K \rightarrow 2$ as $\epsilon \rightarrow +0$.

Proof. Applying item (iii) of Theorem 2.25 to the simple o.s. $[\mathcal{O}]^K$, we see that what above was called $\varepsilon_{\star}(K)$ satisfies

$$\varepsilon_{\star}(K) \le 2\sqrt{\epsilon(1-\epsilon)}.$$

By (2.89), we conclude that $\varepsilon_{\star}(1) \leq \left(2\sqrt{\epsilon(1-\epsilon)}\right)^{1/K}$, whence, by the same (2.89), $\varepsilon_{\star}(T) \leq \left(2\sqrt{\epsilon(1-\epsilon)}\right)^{T/K}$, T = 1, 2, ...; plugging in this bound $T = K_+$, we get (check it!) the inequality $\varepsilon_{\star}(K_+) \leq \epsilon$. It remains to recall that $\varepsilon_{\star}(K_+)$ upperbounds the risks of the test $\mathcal{T}_{\phi^{(K_+)}}$ when deciding on $H_1^{(K_+)}$ vs. $H_2^{(K_+)}$.

2.5 TESTING MULTIPLE HYPOTHESES

So far, we focused on detector-based tests deciding on pairs of hypotheses, and our "constructive" results were restricted to pairs of *convex* hypotheses dealing with a simple o.s.

$$\mathcal{O} = (\Omega, \Pi; \{ p_{\mu} : \mu \in \mathcal{M} \}; \mathcal{F}), \tag{2.90}$$

convexity of a hypothesis meaning that the family of probability distributions obeying the hypothesis is $\{p_{\mu} : \mu \in X\}$ associated with a convex (in fact, convex compact) set $X \subset \mathcal{M}$.

In this Section, we will be interested in pairwise testing *unions* of convex hypotheses and testing *multiple* (more than two) hypotheses.

2.5.1 Testing unions

2.5.1.1 Situation and goal

Let Ω be an observation space, and assume we are given two finite collections of families of probability distributions on Ω : families of *red* distributions \mathcal{R}_i , $1 \leq i \leq r$, and families of *blue* distributions \mathcal{B}_j , $1 \leq j \leq b$. These families give rise to r red and b blue hypotheses on the distribution P of an observation $\omega \in \Omega$, specifically,

$$R_i: P \in \mathcal{R}_i \text{ (red hypotheses)} \text{ and } B_j: P \in \mathcal{B}_j \text{ (blue hypotheses)}$$

Assume that for every $i \leq r, j \leq B$ we have at our disposal a simple detector-based test \mathcal{T}_{ij} capable to decide on R_i vs B_j ; what we want is to assemble these tests into a test \mathcal{T} deciding on the union R of red hypotheses vs. the union B of blue ones:

$$R: P \in \mathcal{R} := \bigcup_{i=1}^{r} \mathcal{R}_{i}, \ B: P \in \mathcal{B} := \bigcup_{j=1}^{b} \mathcal{B}_{j},$$

where P, as always, stands for the probability distribution of observation $\omega \in \Omega$.

Our motivation primarily stems from the case where R_i and B_j are convex hypotheses in a simple o.s. (2.90):

$$\mathcal{R}_i = \{p_\mu : \mu \in M_i\}, \, \mathcal{B}_j = \{p_\mu : \mu \in N_j\},$$

where M_i and N_j are convex compact subsets of \mathcal{M} . In this case we indeed know how to build near-optimal tests deciding on R_i vs. B_j , and the question we have posed becomes, how to assemble these tests into a test deciding on R vs. B, with

$$R: P \in \mathcal{R} = \{p_{\mu} : \mu \in X\}, X = \bigcup_{i} M_{i}, B: P \in \mathcal{B} = \{p_{\mu} : \mu \in Y\}, Y = \bigcup_{i} N_{j};$$

while structure of R, B is similar to the one of R_i , B_j , there is a significant difference: the sets X, Y are, in general, non-convex, and therefore the techniques we have developed fail to address testing R vs. B directly.

2.5.1.2 The construction

In the just described situation, let ϕ_{ij} be the detectors underlying the tests \mathcal{T}_{ij} ; w.l.o.g., we can assume these detectors balanced (see Section 2.3.2.2) with some risks ϵ_{ij} :

$$\int_{\Omega} e^{-\phi_{ij}(\omega)} P(d\omega) \leq \epsilon_{ij} \,\forall P \in \mathcal{R}_i \\ \int_{\Omega} e^{\phi_{ij}(\omega)} P(d\omega) \leq \epsilon_{ij} \,\forall P \in \mathcal{B}_j \ \right\}, \ 1 \leq i \leq r, 1 \leq j \leq b.$$
 (2.91)

Let us assemble the detectors ϕ_{ij} into a detector for R, B as follows:

$$\phi(\omega) = \max_{1 \le i \le r} \min_{1 \le j \le b} [\phi_{ij} - \alpha_{ij}], \qquad (2.92)$$

where the *shifts* α_{ij} are construction's parameters.

Proposition 2.30. The risks of ϕ on R, B can be bounded as

$$\forall P \in \mathcal{R} : \quad \int_{\Omega} e^{-\phi(\omega)} P(d\omega) \le \max_{i \le r} \left[\sum_{j=1}^{b} \epsilon_{ij} e^{\alpha_{ij}} \right]$$

$$\forall P \in \mathcal{B} : \quad \int_{\Omega} e^{\phi(\omega)} P(d\omega) \le \max_{j \le b} \left[\sum_{i=1}^{r} \epsilon_{ij} e^{-\alpha_{ij}} \right]$$

$$(2.93)$$

Thus, the risks of ϕ on R, B are upper-bounded by the quantity

$$\varepsilon_{\star} = \max\left[\max_{i \le r} \left[\sum_{j=1}^{b} \epsilon_{ij} e^{\alpha_{ij}}\right], \max_{j \le b} \left[\sum_{i=1}^{r} \epsilon_{ij} e^{-\alpha_{ij}}\right]\right], \quad (2.94)$$

whence the risks of the based on the detector ϕ simple test \mathcal{T}_{ϕ} deciding on R, B are upper-bounded by ε_{\star} .

Proof. Let $P \in \mathcal{R}$, so that $P \in \mathcal{R}_{i_*}$ for some $i_* \leq r$. Then

$$\begin{split} &\int_{\Omega} e^{-\phi(\omega)} P(d\omega) = \int_{\Omega} e^{\min_{i \leq r} \max_{j \leq b} [-\phi_{ij}(\omega) + \alpha_{ij}]} P(d\omega) \\ &\leq \int_{\Omega} e^{\max_{j \leq b} [-\phi_{i*j}(\omega) + \alpha_{i*j}]} P(d\omega) \leq \sum_{j=1}^{b} \int_{\Omega} e^{-\phi_{i*j}(\omega) + \alpha_{i*j}} P(d\omega) \\ &= \sum_{j=1}^{b} \exp^{\alpha_{i*j}} \int_{\Omega} e^{-\phi_{i*j}(\omega)} P(d\omega) \\ &\leq \sum_{j=1}^{b} \epsilon_{i*j} e^{\alpha_{i*j}} \text{ [by (2.91) due to } P \in \mathcal{R}_{i*}] \\ &\leq \max_{i \leq r} \left[\sum_{j=1}^{b} \epsilon_{ij} e^{\alpha_{ij}} \right] \end{split}$$

Now let $P \in \mathcal{B}$, so that $P \in \mathcal{B}_{j_*}$ for some j_* . We have

$$\begin{split} &\int_{\Omega} e^{\phi(\omega)} P(d\omega) = \int_{\Omega} e^{\max_{i \leq r} \min_{j \leq b} [\phi_{ij}(\omega) - \alpha_{ij}]} P(d\omega) \\ &\leq \int_{\Omega} e^{\max_{i \leq r} [\phi_{ij_*}(\omega) - \alpha_{ij_*}]} P(d\omega) \leq \sum_{i=1}^{r} \int_{\Omega} e^{\phi_{ij_*}(\omega) - \alpha_{ij_*}} P(d\omega) \\ &= \sum_{i=1}^{r} \exp^{-\alpha_{ij_*}} \int_{\Omega} e^{\phi_{ij_*}(\omega)} P(d\omega) \\ &\leq \sum_{i=1}^{r} \epsilon_{ij_*} e^{\alpha_{ij_*}} \text{ [by (2.91) due to } P \in \mathcal{B}_{j_*}] \\ &\leq \max_{j \leq b} \left[\sum_{i=1}^{r} \epsilon_{ij} e^{-\alpha_{ij}} \right] \end{split}$$

(2.93) is proved. The remaining claims in Proposition are readily given by (2.93) combined with Proposition 2.16. $\hfill \Box$

Optimal choice of shift parameters. The detector and the test considered in Proposition 2.30, same as the resulting risk bound ε_{\star} , depend on the shifts α_{ij} . We are about to optimize the risk bound w.r.t. these shifts. To this end, consider the $r \times b$ matrix

$$E = [\epsilon_{ij}]_{\substack{i \leq j \\ i \leq j}}$$

and the symmetric $(r+b) \times (r+b)$ matrix

$$\mathcal{E} = \left[\begin{array}{c|c} & E \\ \hline & E^T & \\ \hline \end{array} \right]$$

As it is well known, the eigenvalues of the symmetric matrix \mathcal{E} are comprised of the pairs $(\sigma_s, -\sigma_s)$, where σ_s are the singular values of E, and several zeros; in particular, the leading eigenvalue of \mathcal{E} is the spectral norm $||E||_{2,2}$ (the largest singular value) of matrix E. Further, E is a matrix with positive entries, so that \mathcal{E} is a symmetric entrywise nonnegative matrix. By Perron-Frobenius Theorem, the leading eigenvector of this matrix can be selected to be nonnegative. Denoting this nonnegative eigenvector [g; h] with r-dimensional g and b-dimensional h, and setting $\rho = ||E||_{2,2}$, we have

$$\begin{array}{lll}
\rho g &=& Eh\\
\rho h &=& E^T g
\end{array} \tag{2.95}$$

Observe that $\rho > 0$ (evident), whence both g and h are nonzero (since otherwise (2.95) would imply g = h = 0, which is impossible – the eigenvector [g; h] is nonzero). Since h and g are nonzero nonnegative vectors, $\rho > 0$ and E is entrywise positive, (2.95) says that g and h are strictly positive vectors. The latter allows to define shifts α_{ij} according to

$$\alpha_{ij} = \ln(h_j/g_i). \tag{2.96}$$

With these shifts, we get

$$\max_{i \leq r} \left[\sum_{j=1}^{b} \epsilon_{ij} e^{\alpha_{ij}} \right] = \max_{i \leq r} \sum_{j=1}^{b} \epsilon_{ij} h_j / g_i = \max_{i \leq r} (Eh)_i / g_i = \max_{i \leq r} \rho = \rho$$

(we have used the first relation in (2.95)) and

$$\max_{j \le b} \left[\sum_{i=1}^r \epsilon_{ij} e^{-\alpha_{ij}} \right] = \max_{j \le b} \sum_{i=1}^r \epsilon_{ij} g_i / h_j = \max_{j \le b} [E^T g]_j / h_j = \max_{j \le b} \rho = \rho$$

(we have used the second relation in (2.95)). The bottom line is as follows:

Proposition 2.31. In the situation and the notation from Section 2.5.1.1, the risks of the detector (2.92) with shifts (2.95), (2.96) on the families \mathcal{R} , \mathcal{B} do not exceed the quantity

$$||E:=[\epsilon_{ij}]_{i\leq r,j\leq b}||_{2,2}.$$

As a result, the risks of the simple test \mathcal{T}_{ϕ} deciding on the hypotheses R, B, does not exceed $||E||_{2,2}$ as well.

In fact, the shifts in the above proposition are the best possible; this is an immediate consequence of the following simple fact:

Proposition 2.32. Let $\mathcal{E} = [e_{ij}]$ be nonzero entrywise nonnegative $n \times n$ symmetric matrix. Then the optimal value in the optimization problem

$$Opt = \min_{\alpha_{ij}} \left\{ \max_{i \le n} \sum_{j=1}^{n} e_{ij} e^{\alpha_{ij}} : \alpha_{ij} = -\alpha_{ji} \right\}$$
(*)

is equal to $\|\mathcal{E}\|_{2,2}$. When the Perron-Frobenius eigenvector f of \mathcal{E} can be selected positive, the problem is solvable, and an optimal solution is given by

$$\alpha_{ij} = \ln(f_j/f_i), \ 1 \le i, j \le n.$$
 (2.97)

Proof. Let us prove, first, that $Opt \leq \rho := ||\mathcal{E}||_{2,2}$. Given $\epsilon > 0$, we clearly can find an entrywise nonnegative symmetric matrix \mathcal{E}' with entries e'_{ij} in-between e_{ij} and $e_{ij} + \epsilon$ such that the Perron-Frobenius eigenvector f of \mathcal{E}' can be selected positive (it suffices, e.g., to set $e'_{ij} = e_{ij} + \epsilon$). Selecting α_{ij} according to (2.97), we get a feasible solution to (*) such that

$$\forall i: \sum_{j} e_{ij} e^{\alpha_{ij}} \leq \sum_{j} e'_{ij} f_j / f_i = \|\mathcal{E}'\|_{2,2},$$

implying that $\text{Opt} \leq \|\mathcal{E}'\|_{2,2}$. Passing to limit as $\epsilon \to +0$, we get $\text{Opt} \leq \|\mathcal{E}\|_{2,2}$. As a byproduct of our reasoning, we see that if \mathcal{E} admits a positive Perron-Frobenius eigenvector f, then (2.97) yields a feasible solution to (*) with the value of the objective equal to $\|\mathcal{E}\|_{2,2}$.

It remain to prove that $\text{Opt} \geq \|\mathcal{E}\|_{2,2}$. Assume that this is not the case, so that (*) admits a feasible solution $\hat{\alpha}_{ij}$ such that

$$\widehat{\rho} := \max_{i} \sum_{j} e_{ij} e^{\widehat{\alpha}_{ij}} < \rho := \|\mathcal{E}\|_{2,2}.$$

Perturbing \mathcal{E} a little bit, we can make this matrix symmetric and entrywise positive, and still satisfying the above strict inequality; to save notation, assume that already the original \mathcal{E} is entrywise positive. Let f be a positive Perron-Frobenius eigenvector of \mathcal{E} , and let, as above, $\alpha_{ij} = \ln(f_j/f_i)$, so that

$$\sum_{j} e_{ij} e^{\alpha_{ij}} = \sum_{j} e_{ij} f_j / f_i = \rho \ \forall i.$$

Setting $\delta_{ij} = \hat{\alpha}_{ij} - \alpha_{ij}$, we conclude that the convex functions

$$\theta_i(t) = \sum_j e_{ij} e^{\alpha_{ij} + t\delta_{ij}}$$

all are equal to ρ as t = 0, and all are $\leq \hat{\rho} < \rho$ as t = 1, implying that $\theta_i(1) < \theta_i(0)$ for every *i*. The latter, in view of convexity of $\theta_i(\cdot)$, implies that

$$\theta_i'(0) = \sum_j e_{ij} e^{\alpha_{ij}} \delta_{ij} = \sum_j e_{ij} (f_j/f_i) \delta_{ij} < 0 \ \forall i.$$

Multiplying the resulting inequalities by f_i^2 and summing up over *i*, we get

$$\sum_{i,j} e_{ij} f_i f_j \delta_{ij} < 0,$$

which is impossible: we have $e_{ij} = e_{ji}$ and $\delta_{ij} = -\delta_{ji}$, implying that the left hand side in the latter inequality is 0.

2.5.2 Testing multiple hypotheses "up to closeness"

So far, we have considered detector-based simple tests deciding on pairs of hypotheses, specifically, convex hypotheses in simple o.s.'s (Section 2.4.4) and unions of convex hypotheses (Section 2.5.1)²⁰. Now we intend to consider testing of multiple (perhaps more than 2) hypotheses "up to closeness;" the latter notion was introduced in Section 2.2.4.2.

2.5.2.1 Situation and goal

Let Ω be an observation space, and let a collection $\mathcal{P}_1, ..., \mathcal{P}_L$ of families of probability distributions on Ω be given. As always, families \mathcal{P}_ℓ give rise to hypotheses

$$H_{\ell}: P \in \mathcal{P}_{\ell}$$

on the distribution P of observation $\omega \in \Omega$. Assume also that we are given a closeness relation C on $\{1, ..., L\}$; recall that a closeness relation, formally, is some set of pairs of indexes $(\ell, \ell') \in \{1, ..., L\}$; we interpret the inclusion $(\ell, \ell') \in C$ as the

 $^{^{20}}$ strictly speaking, in Section 2.5.1 it was not explicitly stated that the unions under consideration involve convex hypotheses in simple o.s.'s; our emphasis was on how to decide on a pair of union-type hypotheses given pairwise detectors for "red" and "blue" components of the unions from the pair. Note, however, that as of now, the only situation where we indeed have at our disposal good pairwise detectors for red and blue components is the one where these components are convex hypotheses in a good o.s.

100

LECTURE 2



Figure 2.2: 11 hypotheses on the location of the mean μ of observation $\omega \sim \mathcal{N}(\mu, I_2)$, each stating that μ belongs to the polygon of specific color.

fact that hypothesis H_{ℓ} "is close" to hypothesis $H_{\ell'}$. When $(\ell, \ell') \in \mathcal{C}$, we say that ℓ' is close (or \mathcal{C} -close) to ℓ . We always assume that

- \mathcal{C} contains the diagonal: $(\ell, \ell) \in \mathcal{C}$ for every $\ell \leq L$ ("each hypothesis is close to itself"), and
- C is symmetric: whenever $(\ell, \ell') \in C$, we have also $(\ell', \ell) \in C$ ("if ℓ -th hypothesis is close to ℓ' -th one, then ℓ' -th hypothesis is close to ℓ -th one").

Recall that a test \mathcal{T} deciding on the hypotheses $H_1, ..., H_L$ via observation $\omega \in \Omega$ is a procedure which, given on input $\omega \in \Omega$, builds some set $\mathcal{T}(\omega) \subset \{1, ..., L\}$, accepts all hypotheses H_ℓ with $\ell \in \mathcal{T}(\omega)$, and rejects all other hypotheses.

Risks of an "up to closeness" test. The notion of C-risk of a test was introduced in Section 2.2.4.2; we reproduce it here for reader's convenience. Given closeness Cand a test T, we define the C-risk

$$\operatorname{Risk}^{\mathcal{C}}(\mathcal{T}|H_1,...,H_L)$$

of \mathcal{T} as the smallest $\epsilon \geq 0$ such that

Whenever an observation ω is drawn from a distribution $P \in \bigcup_{\ell} \mathcal{P}_{\ell}$, and ℓ_* is such that $P \in \mathcal{P}_{\ell_*}$ (i.e., hypothesis H_{ℓ_*} is true), the P-probability of the event " $\ell_* \notin \mathcal{T}(\omega)$ ("true hypothesis H_{ℓ_*} is not accepted") <u>or</u> there exists ℓ' <u>not close to ℓ </u> such that $H_{\ell'}$ is accepted" is <u>at most</u> ϵ .

Equivalently:

 $\operatorname{Risk}^{\mathcal{C}}(\mathcal{T}|H_1,...,H_L) \leq \epsilon$ if and only if the following takes place:

Whenever an observation ω is drawn from a distribution $P \in \bigcup_{\ell} \mathcal{P}_{\ell}$, and ℓ_* is such that $P \in \mathcal{P}_{\ell_*}$ (i.e., hypothesis H_{ℓ_*} is true), the P-probability of the event

 $\ell_* \in \mathcal{T}(\omega)$ ("the true hypothesis H_{ℓ_*} is accepted") and $\ell' \in \mathcal{T}(\omega)$ implies that $(\ell, \ell') \in \mathcal{C}$ ("all accepted hypotheses are \mathcal{C} -close to the true hypothesis H_{ℓ_*} ") is at least $1 - \epsilon$.

For example, consider 11 colored polygons presented on Figure 2.2 and associate with them 11 hypotheses on 2D "signal plus noise" observation $\omega = x + \xi$, $\xi \sim \mathcal{N}(0, I_2)$, with ℓ -th hypothesis stating that x belongs to ℓ -th polygon. When defining closeness \mathcal{C} on the collection of 11 hypotheses presented on Figure 2.2 as

"two hypotheses are close if and only if the corresponding color polygons intersect"

the fact that a test \mathcal{T} has \mathcal{C} -risk ≤ 0.01 implies, in particular, that if the probability

distribution P underlying the observed ω "is black," (i.e., the mean of ω belongs to the black polygon), then with P-probability at least 0.99 the list of accepted hypotheses will include the black one, and the only other hypotheses in this list will be among the red, yellow and light-blue ones.

2.5.2.2 "Building blocks" and construction

The construction we are about to present is, essentially, the one used in Section 2.2.4.3 as applied to detector-generated tests; this being said, the presentation to follow is self-contained.

Building blocks for our construction are pairwise detectors $\phi_{\ell\ell'}(\omega)$, $1 \leq \ell \leq \ell' \leq L$, for pairs \mathcal{P}_{ℓ} , $\mathcal{P}_{\ell'}$ along with (upper bounds on) the risks $\epsilon_{\ell\ell'}$ of these detectors:

$$\begin{array}{ll} \forall (P \in \mathcal{P}_{\ell}) : & \int_{\Omega} e^{-\phi_{\ell\ell'}(\omega)} P(d\omega) \leq \epsilon_{\ell\ell'} \\ \forall (P \in \mathcal{P}_{\ell'}) : & \int_{\Omega} e^{\phi_{\ell\ell'}(\omega)} P(d\omega) \leq \epsilon_{\ell\ell'} \end{array} \right\}, \ 1 \leq \ell < \ell' \leq L.$$

Setting

$$\phi_{\ell'\ell}(\omega) = -\phi_{\ell\ell'}(\omega), \ \epsilon_{\ell'\ell} = \epsilon_{\ell\ell'}, \ 1 \le \ell < \ell' \le L, \ \phi_{\ell\ell}(\omega) \equiv 0, \\ \epsilon_{\ell\ell} = 1, \ 1 \le \ell \le L,$$

we get what we shall call balanced system of detectors $\phi_{\ell\ell'}$ and risks $\epsilon_{\ell\ell'}$, $1 \leq \ell, \ell' \leq L$, for the collection $\mathcal{P}_1, ..., \mathcal{P}_L$, meaning that

$$(a): \quad \phi_{\ell\ell'}(\omega) + \phi_{\ell'\ell}(\omega) \equiv 0, \ \epsilon_{\ell\ell'} = \epsilon_{\ell'\ell}, \ 1 \leq \ell, \ell' \leq L (b): \quad \forall P \in \mathcal{P}_{\ell}: \int_{\Omega} e^{-\phi_{\ell\ell'}(\omega)} P(d\omega) \leq \epsilon_{\ell\ell'}, \ 1 \leq \ell, \ell' \leq L.$$
(2.98)

Given closeness \mathcal{C} , we associate with it the symmetric $L \times L$ matrix **C** given by

$$\mathbf{C}_{\ell\ell'} = \begin{cases} 0, & (\ell, \ell') \in \mathcal{C} \\ 1, & (\ell, \ell') \notin \mathcal{C} \end{cases}$$
(2.99)

Test $\mathcal{T}_{\mathcal{C}}$. Let a collection of shifts $\alpha_{\ell\ell'} \in \mathbf{R}$ satisfying the relation

$$\alpha_{\ell\ell'} = -\alpha_{\ell'\ell}, \ 1 \le \ell, \ell' \le L \tag{2.100}$$

be given. The detectors $\phi_{\ell\ell'}$ and the shifts $\alpha_{\ell\ell'}$ specify a test $\mathcal{T}_{\mathcal{C}}$ deciding on hypotheses $H_1, ..., H_L$; specifically, given an observation ω , the test $\mathcal{T}_{\mathcal{C}}$ accepts exactly those hypotheses H_ℓ for which $\phi_{\ell\ell'}(\omega) - \alpha_{\ell\ell'} > 0$ whenever ℓ' is not \mathcal{C} -close to ℓ :

$$\mathcal{T}_{\mathcal{C}}(\omega) = \{\ell : \phi_{\ell\ell'}(\omega) - \alpha_{\ell\ell'} > 0 \ \forall (\ell' : (\ell, \ell') \notin \mathcal{C})\}.$$
(2.101)

Proposition 2.33. (i) The C-risk of the just defined test \mathcal{T}_{C} is upper-bounded by the quantity

$$\varepsilon[\alpha] = \max_{\ell \le L} \sum_{\ell'=1}^{L} \epsilon_{\ell\ell'} \mathbf{C}_{\ell\ell'} e^{\alpha_{\ell\ell'}}$$

with \mathbf{C} given by (2.99).

(ii) The infimum, over shifts α satisfying (2.100), of the risk bound $\varepsilon[\alpha]$ is the

102

LECTURE 2

quantity

$$\varepsilon_{\star} = \|\mathcal{E}\|_{2,2},$$

where the $L \times L$ symmetric entrywise nonnegative matrix \mathcal{E} is given by

 $\mathcal{E} = [e_{\ell\ell'} := \epsilon_{\ell\ell'} \mathbf{C}_{\ell\ell'}]_{\ell,\ell' < L} \,.$

Assuming \mathcal{E} admits a strictly positive Perron-Frobenius vector f, an optimal choice of the shifts is

$$\alpha_{\ell\ell'} = \ln(f_{\ell'}/f_\ell), 1 \le \ell, \ell' \le L,$$

resulting in $\varepsilon[\alpha] = \varepsilon_{\star} = \|\mathcal{E}\|_{2,2}$.

Proof. (i): Setting

$$\bar{\phi}_{\ell\ell'}(\omega) = \phi_{\ell\ell'}(\omega) - \alpha_{\ell\ell'}, \ \bar{\epsilon}_{\ell\ell'} = \epsilon_{\ell\ell'} e^{\alpha_{\ell\ell'}},$$

(2.98), (2.100) imply that

$$(a): \quad \bar{\phi}_{\ell\ell'}(\omega) + \bar{\phi}_{\ell'\ell}(\omega) \equiv 0, 1 \leq \ell, \ell' \leq L (b): \quad \forall P \in \mathcal{P}_{\ell}: \int_{\Omega} e^{-\bar{\phi}_{\ell\ell'}(\omega)} P(d\omega) \leq \bar{\epsilon}_{\ell\ell'}, 1 \leq \ell, \ell' \leq L.$$
(2.102)

Now let ℓ_* be such that the distribution P of observation ω belongs to \mathcal{P}_{ℓ_*} . Then the P-probability of the event $\bar{\phi}_{\ell_*\ell'}(\omega) \leq 0$ is, for every $\ell', \leq \bar{\epsilon}_{\ell_*\ell'}$ by (2.102.*b*), whence the P-probability of the event

$$E_* = \{ \omega : \exists \ell' : (\ell_*, \ell') \notin \mathcal{C} \& \bar{\phi}_{\ell_* \ell'}(\omega) \le 0 \}$$

is upper-bounded by

$$\sum_{\ell':(\ell_*,\ell')\notin\mathcal{C}}\bar{\epsilon}_{\ell_*\ell'} = \sum_{\ell'=1}^L \mathbf{C}_{\ell_*\ell'}\epsilon_{\ell_*\ell'}e^{\alpha_{\ell_*\ell'}} \le \varepsilon[\alpha].$$

Assume that E_* does not take place (as we have seen, this indeed is so with P-probability $\geq 1 - \varepsilon[\alpha]$). Then $\bar{\phi}_{\ell_*\ell'}(\omega) > 0$ for all ℓ' such that $(\ell_*, \ell') \notin C$, implying, first, that H_{ℓ_*} is accepted by our test. Second, $\bar{\phi}_{\ell'\ell_*}(\omega) = -\bar{\phi}_{\ell_*\ell'}(\omega) < 0$ whenever $(\ell_*, \ell') \notin C$, or, which is the same due to the symmetry of closeness, whenever $(\ell', \ell_*) \notin C$, implying that the test \mathcal{T}_C rejects the hypothesis $H_{\ell'}$ when ℓ' is not C-close to ℓ_* . Thus, the P-probability of the event " H_{ℓ_*} is accepted, and all accepted hypotheses are C-close to H_{ℓ_*} " is at least $1 - \varepsilon[\alpha]$. We conclude that the C-risk Risk^C($\mathcal{T}_C|H_1, ..., H_L$) of the test \mathcal{T}_C is at most $\varepsilon[\alpha]$. (i) is proved. (ii) is readily given by Proposition 2.32.

2.5.2.3 Testing multiple hypotheses via repeated observations

In the situation of Section 2.5.2.1, given a balanced system of detectors $\phi_{\ell\ell'}$ and risks $\epsilon_{\ell\ell'}$, $1 \leq \ell, \ell' \leq L$ for the collection $\mathcal{P}_1, ..., \mathcal{P}_L$ (see (2.98)) and a positive integer K, we can

• pass from detectors $\phi_{\ell\ell'}$ and risks $\epsilon_{\ell\ell'}$ to the entities

$$\phi_{\ell\ell'}^{(K)}(\omega^K = (\omega_1, ..., \omega_K)) = \sum_{k=1}^K \phi_{\ell\ell'}(\omega_k), \ \epsilon_{\ell\ell'}^{(K)} = \epsilon_{\ell\ell'}^K, \ 1 \le \ell, \ell' \le L$$

• associate with the families \mathcal{P}_{ℓ} families $\mathcal{P}_{\ell}^{(K)}$ of probability distributions underlying quasi-stationary K-repeated versions of observations $\omega \sim P \in \mathcal{P}_{\ell}$, see Section 2.3.2.3, and thus arrive at hypotheses $H_{\ell}^{K} = \mathcal{H}_{\ell}^{\otimes,K}$ stating that the distribution P^{K} of K-repeated observation $\omega^{K} = (\omega_{1}, ..., \omega_{K}), \omega_{k} \in \Omega$, belongs to the family $\mathcal{P}_{\ell}^{\otimes,K} = \bigotimes_{k=1}^{K} \mathcal{P}_{\ell}$, see Section 2.1.3.3, associated with \mathcal{P}_{ℓ} .

Invoking Proposition 2.18 and (2.98), we arrive at the following analogy of (2.98):

$$(a): \quad \phi_{\ell\ell\ell'}^{(K)}(\omega^K) + \phi_{\ell'\ell}^{(K)}(\omega^K) \equiv 0, \ \epsilon_{\ell\ell'}^{(K)} = \epsilon_{\ell'\ell}^{(K)} = \epsilon_{\ell\ell'}^K, \ 1 \le \ell, \ell' \le L (b): \quad \forall P^K \in \mathcal{P}_{\ell}^{(K)}: \int_{\Omega^K} e^{-\phi_{\ell\ell'}^{(K)}(\omega^K)} P^K(d\omega^K) \le \epsilon_{\ell\ell'}^{(K)}, \ 1 \le \ell, \ell' \le L.$$

$$(2.103)$$

Given shifts $\alpha_{\ell\ell'}$ satisfying (2.100) and applying the construction from Section 2.5.2.2 to these shifts and our new detectors and risks, we arrive at the test test $\mathcal{T}_{\mathcal{C}}^{K}$ deciding on hypotheses $H_{1}^{K}, ..., H_{L}^{K}$ via K-repeated observation ω^{K} ; specifically, given an observation ω^{K} , the test $\mathcal{T}_{\mathcal{C}}^{K}$ accepts exactly those hypotheses H_{ℓ}^{K} for which $\phi_{\ell\ell'}^{(K)}(\omega^{K}) - \alpha_{\ell\ell'} > 0$ whenever ℓ' is not \mathcal{C} -close to ℓ :

$$\mathcal{T}_{\mathcal{C}}^{K}(\omega^{K}) = \{\ell : \phi_{\ell\ell'}^{(K)}(\omega^{K}) - \alpha_{\ell\ell'} > 0 \ \forall (\ell' : (\ell, \ell') \notin \mathcal{C})\},$$
(2.104)

Invoking Proposition 2.33, we arrive at

Proposition 2.34. (i) The C-risk of the just defined test $\mathcal{T}_{\mathcal{C}}^{K}$ is upper-bounded by the quantity

$$\varepsilon[\alpha, K] = \max_{\ell \leq L} \sum_{\ell'=1}^{L} \epsilon_{\ell\ell'}^{K} \mathbf{C}_{\ell\ell'} e^{\alpha_{\ell\ell'}}.$$

(ii) The infimum, over shifts α satisfying (2.100), of the risk bound $\varepsilon[\alpha, K]$ is the quantity

$$\varepsilon_{\star}(K) = \|\mathcal{E}^{(K)}\|_{2,2},$$

where the $L \times L$ symmetric entrywise nonnegative matrix $\mathcal{E}^{(K)}$ is given by

$$\mathcal{E}^{(K)} = \left[e_{\ell\ell'}^{(K)} := \epsilon_{\ell\ell'}^K \mathbf{C}_{\ell\ell'} \right]_{\ell,\ell' \le L}$$

Assuming $\mathcal{E}^{(K)}$ admits a strictly positive Perron-Frobenius vector f, an optimal choice of the shifts is

$$\alpha_{\ell\ell'} = \ln(f_\ell/f_{\ell'}), 1 \le \ell, \ell' \le L,$$

resulting in $\varepsilon[\alpha, K] = \varepsilon_{\star}(K) = \|\mathcal{E}^{(K)}\|_{2,2}$.

2.5.2.4 Consistency and near-optimality

Observe that when the closeness C is such that $\epsilon_{\ell\ell'} < 1$ whenever ℓ , ℓ' are not C-close to each other, the entries on the matrix $\mathcal{E}^{(K)}$ exponentially fast go to 0 as

 $K \to \infty$, whence the *C*-risk of test $\mathcal{T}_{\mathcal{C}}^{K}$ also goes to 0 as $K \to \infty$; this is called *consistency*. When, in addition, \mathcal{P}_{ℓ} correspond to convex hypotheses in a simple o.s., the test $\mathcal{T}_{\mathcal{C}}^{K}$ possesses certain near-optimality properties similar to those stated in Proposition 2.29

Proposition 2.35. Consider the special case of the situation from Section 2.5.2.1 where, given a simple o.s. $\mathcal{O} = (\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F})$, the families \mathcal{P}_{ℓ} of probability distributions are of the form $\mathcal{P}_{\ell} = \{p_{\mu} : \mu \in N_{\ell}\}$, where $N_{\ell}, 1 \leq \ell \leq L$, are nonempty convex compact subsets of \mathcal{M} . Let also the pairwise detectors $\phi_{\ell\ell'}$ and their risks $\epsilon_{\ell\ell'}$ underlying the construction from Section 2.5.2.2 be obtained by applying Theorem 2.25 to the pairs $N_{\ell}, N_{\ell'}$, so that for $1 \leq \ell < \ell' \leq L$ one has

$$\phi_{\ell\ell'}(\omega) = \frac{1}{2} \ln(p_{\mu_{\ell,\ell'}}(\omega)/p_{\nu_{\ell,\ell'}}(\omega)), \ \epsilon_{\ell\ell'} = \exp\{\operatorname{Opt}_{\ell\ell'}\},$$

where

$$\operatorname{Opt}_{\ell\ell'} = \min_{\mu \in N_{\ell}, \nu \in N_{\ell'}} \ln\left(\int_{\Omega} \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)} \Pi(d\omega)\right)$$

and $(\mu_{\ell\ell'}, \nu_{\ell\ell'})$ form an optimal solution to the right hand side optimization problem.

Assume that for some positive integer K_* in the nature there exists a test \mathcal{T}^{K_*} capable to decide with C-risk $\epsilon \in (0, 1/2)$, via stationary K_* -repeated observation ω^{K_*} , on the hypotheses $H_{\ell}^{(K_*)}$, stating that the components in ω^{K_*} are drawn, independently of each other, from a distribution $P \in \mathcal{P}_{\ell}$, $\ell = 1, ..., L$, and let

$$K = \int 2 \frac{1 + \ln(L - 1)/\ln(1/\epsilon)}{1 - \ln(4(1 - \epsilon))/\ln(1/\epsilon)} K_* \lfloor.$$
(2.105)

Then the test $\mathcal{T}_{\mathcal{C}}^{K}$ yielded by the construction from Section 2.5.2.2 as applied to the above $\phi_{\ell\ell'}$, $\epsilon_{\ell\ell'}$ and trivial shifts $\alpha_{\ell\ell'} \equiv 0$ decides on the hypotheses H_{ℓ}^{K} , see Section 2.5.2.3, via quasi-stationary K-repeated observations ω^{K} , with C-risk $\leq \epsilon$.

Note that $K/K_* \to 2$ as $\epsilon \to +0$.

Proof. Let

$$\bar{\epsilon} = \max_{\ell,\ell'} \left\{ \epsilon_{\ell\ell'} : \ell < \ell' \text{ and } \ell, \ell' \text{ are not } \mathcal{C}\text{-close to each other} \right\}.$$

Denoting by (ℓ_*, ℓ'_*) the maximizer in the right hand side maximization, note that \mathcal{T}^{K_*} induces a simple test \mathcal{T} capable to decide via stationary K_* -repeated observations ω^K on the pair of hypotheses $H_{\ell_*}^{(K_*)}$, $H_{\ell'_*}^{(K_*)}$ with risks $\leq \epsilon$ (it suffices to make \mathcal{T} to accept the first of the hypotheses in the pair and reject the second one whenever \mathcal{T}^{K_*} on the same observation accepts $H_{\ell_*}^{(K_*)}$, otherwise \mathcal{T} rejects the first hypothesis in the pair and accepts the second one). This observation, by the same argument as in the proof of Proposition 2.29, implies that $\bar{\epsilon}^{K_*} \leq 2\sqrt{\epsilon(1-\epsilon)} < 1$, whence all entries in the matrix $\mathcal{E}^{(K)}$ do not exceed $\bar{\epsilon}^{(K/K_*)}$, implying by Proposition 2.33 that the \mathcal{C} -risk of the test $\mathcal{T}_{\mathcal{C}}^K$ does not exceed

$$\epsilon(K) := (L-1)[2\sqrt{\epsilon(1-\epsilon)}]^{K/K_*}.$$

It remains to note that for K given by (2.105) one has $\epsilon(K) \leq \epsilon$.

Remark 2.36. Note that the tests $\mathcal{T}_{\mathcal{C}}$ and $\mathcal{T}_{\mathcal{C}}^{K}$ we have built, may, depending on

observations, accept no hypotheses at all, which sometimes is undesirable. Clearly, every test deciding on multiple hypotheses up to C-closeness always can be modified to ensure that a hypothesis always is accepted; to this end, it suffices to accept exactly those hypotheses, if any, which are accepted by our original test, and accept, say, hypothesis # 1 when the original test accepts no hypotheses. It is immediate to see that the C-risk of the modified test cannot be larger than the one of the original test.

2.5.3 Illustration: Selecting the best among a family of estimates

Let us illustrate our machinery for multiple hypothesis testing by applying it to the situation as follows:

We are given:

- a simple nondegenerate observation scheme $\mathcal{O} = (\Omega, \Pi; \{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}; \mathcal{F}),$
- a seminorm $\|\cdot\|$ on \mathbf{R}^n , ²¹
- a convex compact set $X \subset \mathbf{R}^n$ along with a collection of M points $x_i \in \mathbf{R}^n$, $1 \leq i \leq M$ and a positive D such that the $\|\cdot\|$ -diameter of the set $X^+ = X \cup \{x_i : 1 \leq i \leq M\}$ is at most D:

$$||x - x'|| \le D \ \forall (x, x' \in X^+),$$

- an affine mapping $x \mapsto A(x)$ from \mathbf{R}^n into the embedding space of \mathcal{M} such that $A(x) \in \mathcal{M}$ for all $x \in \mathcal{M}$,
- a tolerance $\epsilon \in (0, 1)$.

We observe K-element sample $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ of independent across k observations

$$\omega_k \sim p_{A(x_*)}, \ 1 \le k \le K,$$
 (2.106)

where $x_* \in \mathbf{R}^n$ is unknown signal known to belong to X. Our "ideal goal" is to use ω^K in order to identify, with probability $\geq 1 - \epsilon$, the $\|\cdot\|$ -closest to x_* point among the points $x_1, ..., x_M$.

This just outlined goal often is too ambitious, and in the sequel we focus on the relaxed goal as follows:

Given a positive integer N and a "resolution" $\theta > 1$, consider the grid

$$\Gamma = \{r_j = D\theta^{-j}, 0 \le j \le N\}$$

and let

$$\rho(x) = \min\left\{\rho_j \in \Gamma : \rho_j \ge \min_{1 \le i \le M} \|x - x_i\|\right\}.$$

Given design parameters $\alpha \geq 1, \beta \geq 0$, we want to specify volume of obser-

²¹A seminorm on \mathbb{R}^n is defined by exactly the same requirements as a norm, except that now we allow zero seminorms for some nonzero vectors. Thus, a seminorm on \mathbb{R}^n is a nonnegative function $\|\cdot\|$ which is even and homogeneous: $\|\lambda x\| = |\lambda| \|x\|$ and satisfies the triangle inequality $\|x + y\| \leq \|x\| + \|y\|$. A universal example is $\|x\| = \|Bx\|_o$, where $\|\cdot\|_o$ is a norm on some \mathbb{R}^m and B is an $m \times n$ matrix; whenever this matrix has a nontrivial kernel, $\|\cdot\|$ is a seminorm rather than a norm.

vations K and an inference routine $\omega^K \mapsto i_{\alpha,\beta}(\omega^K) \in \{1, ..., M\}$ such that

$$\forall (x_* \in X) : \operatorname{Prob}\{\|x_* - x_{i_{\alpha,\beta}(\omega^K)}\| > \alpha \rho(x_*) + \beta\} \ge 1 - \epsilon.$$
(2.107)

Note that when passing from the "ideal" to the relaxed goal, the simplification is twofold: first, we do not care about the precise distance $\min_i ||x_* - x_i||$ from x_* to $\{x_1, ..., x_M\}$, all we look at is the best upper bound $\rho(x_*)$ on this distance from the grid Γ ; second, we allow factor α and additive term β in mimicking the (discretized) distance $\rho(x_*)$ by $||x_* - x_{i_{\alpha,\beta}(\omega^{\kappa})}||$.

The problem we have posed is rather popular in Statistics; its origin usually looks as follows: x_i are candidate estimates of x_* yielded by a number of a priori "models" of x_* and perhaps some preliminary noisy observations of x_* . Given x_i and a matrix B, we want to select among the vectors Bx_i the (nearly) best, w.r.t. a given norm $\|\cdot\|_o$, approximation of Bx_* , utilizing additional observations ω^K of the signal. To bring this problem into our framework, it suffices to specify the seminorm as $\|x\| = \|Bx\|_o$. We shall see in the mean time that in the context of this problem, the above "discretization of distances" is, for all practical purposes, irrelevant: the dependence of the volume of observations on N is just logarithmic, so that we can easily handle fine grid, like the one with $\theta = 1.001$ and $\theta^{-N} = 10^{-10}$. As about factor α and additive term β , they indeed could be "expensive in terms of applications," but the "nearly ideal" goal of making α close to 1 and β close to 0 is in many cases too ambitious to be achievable.

2.5.3.1 The construction

Let us associate with $i \leq M$ and $j, 0 \leq j \leq N$, hypothesis H_{ij} stating that the independent across k observations ω_k , see (2.106), stem from $x_* \in X_{ij} = \{x \in X :$ $\|x - x_i\| \leq r_j\}$. Note that the sets X_{ij} are convex and compact. We denote by \mathcal{J} the set of all pairs (i, j), for which $i \in \{1, ..., M\}$, $j \in \{0, 1, ..., N\}$, and $X_{ij} \neq \emptyset$. Further, we define closeness $\mathcal{C}_{\alpha,\beta}$ on the set of hypotheses $H_{ij}, (i, j) \in \mathcal{J}$, as follows:

 $(ij, i'j') \in \mathcal{C}_{\alpha\beta}$ if and only if

$$\|x_i - x_{i'}\| \le \bar{\alpha}(r_j + r_{j'}) + \beta, \ \bar{\alpha} = \frac{\alpha - 1}{2}.$$
 (2.108)

(here and in what follows, $k\ell$ denotes the ordered pair (k, ℓ)).

Applying Theorem 2.25, we can build, in a computation-friendly fashion, the system $\phi_{ij,i'j'}(\omega), ij, i'j' \in \mathcal{J}$, of optimal balanced detectors for the hypotheses H_{ij} along with the risks of these detectors, so that

$$\begin{array}{ll} (a) & \phi_{ij,i'j'}(\omega) \equiv -\phi_{i'j',ij}(\omega) \,\forall (ij,i'j' \in \mathcal{J}) \\ (b) & \int_{\Omega} e^{-\phi_{ij,i'j'}(\omega)} p_{A(x)}(\omega) \Pi(d\omega) \leq \epsilon_{ij,i'j'} \,\forall (ij \in \mathcal{J}, i'j' \in \mathcal{J}, x \in X_{ij}) \end{array}$$

$$(2.109)$$

Let us say that a pair (α, β) is *admissible*, if $\alpha \ge 1, \beta \ge 0$ and

$$\forall ((i,j) \in \mathcal{J}, (i',j') \in \mathcal{J}, (ij,i'j') \notin \mathcal{C}_{\alpha,\beta}) : A(X_{ij}) \cap A(X_{i'j'}) = \emptyset.$$
(2.110)

Note that checking admissibility of a given pair (α, β) is a computationally tractable task.

Given an admissible par (α, β) , we associate with it positive integer $K = K(\alpha, \beta)$

107

and inference $\omega^K \mapsto i_{\alpha,\beta}(\omega^K)$ as follows:

- 1. $K = K(\alpha, \beta)$ is the smallest integer such that the detector-based test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}}^{K}$ yielded by the machinery of Section 2.5.2.3 decides on the hypotheses $H_{ij}, ij \in \mathcal{J}$, with $\mathcal{C}_{\alpha,\beta}$ -risk not exceeding ϵ . Note that by admissibility, $\epsilon_{ij,i'j'} < 1$ whenever $(ij, i'j') \notin \mathcal{C}_{\alpha,\beta}$, so that $K(\alpha, \beta)$ is well defined.
- 2. Given observation ω^{K} , $K = K(\alpha, \beta)$, we define $i_{\alpha,\beta}(\omega^{K})$ as follows:
 - a) We apply to ω^{K} the test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}}^{K}$. If the test accepts no hypothesis (case A), $i_{\alpha\beta}(\omega^{K})$ is undefined. The observations ω^{K} resulting in case A comprise some set, which we denote by \mathcal{B} ; given ω^{K} , we can recognize whether or not $\omega^{K} \in \mathcal{B}$.
 - b) When $\omega^K \notin \mathcal{B}$, the test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}}^K$ accepts some of the hypotheses H_{ij} , let the set of their indexes ij be $\mathcal{J}(\omega^K)$; we select from the pairs $ij \in \mathcal{J}(\omega^K)$ one with the largest j, and set $i_{\alpha,\beta}(\omega^K)$ to be equal to the first component, and $j_{\alpha,\beta}(\omega^K)$ to be equal to the selected pair.

We are about to prove the following

Proposition 2.37. Assuming (α, β) admissible, for the just defined inference $\omega^K \mapsto i_{\alpha,\beta}(\omega^K)$ and for every $x_* \in X$, denoting by $P_{x_*}^K$ the distribution of stationary K-repeated observation ω^K stemming from x_* one has

$$\|x_* - x_{i_{\alpha,\beta}(\omega^K)}\| \le \alpha \rho(x_*) + \beta. \tag{2.111}$$

with $P_{x_*}^K$ -probability at least $1 - \epsilon$.

Proof. Let us fix $x_* \in X$, let $j_* = j_*(x_*)$ be the largest $j \leq N$ such that $r_j \geq \min_{i \leq M} ||x_* - x_i||$; note that j_* is well defined due to $r_0 = D \geq ||x_* - x_1||$. We set

$$r_{j_*} = \min_j \{r_j : r_j \ge \min_i \|x_* - x_i\|\} = \rho(x_*)$$

and specify $i_* = i_*(x_*) \leq M$ in such a way that

$$\|x_* - x_{i_*}\| \le r_{j_*}.\tag{2.112}$$

Note that i_* is well defined and that observations (2.106) stemming from x_* obey the hypothesis $H_{i_*i_*}$.

Let \mathcal{E} be the set of those ω^K for which the predicate

 \mathcal{P} : As applied to observation ω^{K} , the test $\mathcal{T}_{\mathcal{C}\alpha,\beta}^{K}$ accepts $H_{i_{*}j_{*}}$, and all hypotheses accepted by the test are $\mathcal{C}_{\alpha,\beta}$ -close to $H_{i_{*}j_{*}}$

holds true. Taking into account that the $\mathcal{C}_{\alpha,\beta}$ -risk of $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}}^{K}$ does not exceed ϵ and that the hypothesis $H_{i_*j_*}$ is true, the $P_{x_*}^{K}$ -probability of the event \mathcal{E} is at least $1-\epsilon$.

Let observation ω^K satisfy

$$\omega^K \in \mathcal{E}.\tag{2.113}$$

Then

1. The test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}}^{K}$ accepts the hypothesis $H_{i_{*}j_{*}}$, that is, $\omega^{K} \notin \mathcal{B}$. By construction of $i_{\alpha,\beta}(\omega^{K})j_{\alpha,\beta}(\omega^{K})$ (see the rule 2b above) and due to the fact that $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}}^{K}$ accepts $H_{i_{*}j_{*}}$, we have $j_{\alpha,\beta}(\omega^{K}) \geq j_{*}$.

2. The hypothesis $H_{i_{\alpha,\beta}(\omega^{K})j_{\alpha,\beta}(\omega^{K})}$ is $\mathcal{C}_{\alpha,\beta}$ -close to $H_{i_{*}j_{*}}$, so that

$$\|x_{i_*} - x_{i_{\alpha,\beta}(\omega^K)}\| \le \bar{\alpha}(r_{j_*} + r_{j_{\alpha,\beta}(\omega^K)}) + \beta \le 2\bar{\alpha}r_{j_*} + \beta = 2\bar{\alpha}\rho(x_*) + \beta, \quad (2.114)$$

where the concluding inequality is due to the fact that, as we have already seen, $j_{\alpha,\beta}(\omega^K) \ge j_*$ when (2.113) takes place.

Invoking (2.112), we conclude that with $P_{x_*}^K$ -probability at least $1 - \epsilon$ it holds

$$\|x_* - x_{i_{\alpha,\beta}(\omega^K)}\| \le (2\bar{\alpha} + 1)\rho(x_*) + \beta = \alpha\rho(x_*) + \beta,$$
(2.115)

where the concluding equality is due to the definition of $\bar{\alpha}$.

2.5.3.2 A modification

From the computational viewpoint, a shortcoming of the construction presented in the previous Section is the necessity to operate with M(N + 1) hypotheses, which could require computing as many as $O(M^2N^2)$ detectors. We are about to present a modified construction, where we deal at most N + 1 times with just Mhypotheses at a time (i.e., with the total of at most $O(M^2N)$ detectors). The idea is to replace simultaneous processing of all hypotheses H_{ij} , $ij \in \mathcal{J}$, with processing them in stages j = 0, 1, ..., with stage j operating only with the hypotheses H_{ij} , i = 1, ..., M.

The implementation of this idea is as follows. In the situation of Section 2.5.3, given the same entities Γ , (α, β) , H_{ij} , X_{ij} , $ij \in \mathcal{J}$, as in the beginning of Section 2.5.3.1 and specifying closeness $\mathcal{C}_{\alpha,\beta}$ according to (2.108), we now act as follows.

Preprocessing. We look, one by one, at j = 0, 1, ..., N, and for such a j,

- 1. identify the set $\mathcal{I}_j = \{i \leq M : X_{ij} \neq \emptyset\}$ and stop if this set is empty. If this set is nonempty, we
- 2. specify closeness $C_{\alpha\beta}^{j}$ on the set of hypotheses H_{ij} , $i \in \mathcal{I}_{j}$ as a "slice" of the closeness $C_{\alpha,\beta}$:

 H_{ij} and $H_{i'j}$ (equivalently, *i* and *i'*) are $\mathcal{C}^{j}_{\alpha,\beta}$ -close to each other if (ij, i'j) are $\mathcal{C}_{\alpha,\beta}$ -close, that is,

$$\|x_i - x_{i'}\| \le 2\bar{\alpha}r_j + \beta, \ \bar{\alpha} = \frac{\alpha - 1}{2}.$$

3. build the optimal detectors $\phi_{ij,i'j}$, along with their risks $\epsilon_{ij,i'j}$, for all $i, i' \in \mathcal{I}_j$ such that $(i, i') \notin \mathcal{C}^j_{\alpha,\beta}$.

If for a pair i, i' of this type it happens that $\epsilon_{ij,i'j} = 1$, that is, $A(X_{ij}) \cap A(X_{i'j}) \neq \emptyset$, we claim that (α, β) is inadmissible and stop. Otherwise we find the smallest $K = K_j$ such that the spectral norm of the symmetric $M \times M$ matrix E^{jK} with the entries

$$E_{ii'}^{jK} = \begin{cases} \epsilon_{ij,i'j}^{K}, & i \in \mathcal{I}_j, i' \in \mathcal{I}_j, (i,i') \notin \mathcal{C}_{\alpha,\beta}^{j} \\ 0, & \text{otherwise} \end{cases}$$

does not exceed $\bar{\epsilon} = \epsilon/(N+1)$. We then use the machinery of Section 2.5.2.3 to build detector-based test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}^{j}}^{K_{j}}$ which decides on the hypotheses H_{ij} , $i \in \mathcal{I}_{j}$, with $\mathcal{C}_{\alpha\beta}^{j}$ -risk not exceeding $\bar{\epsilon}$.

It may happen that the outlined process stops when processing some value \overline{j} of j; if this does not happen, we set $\overline{j} = N + 1$. Now, if the process does stop, and stops with the claim that (α, β) is inadmissible, we call (α, β) inadmissible and terminate – in this case we fail to produce a desired inference; note that if this is the case, (α, β) is inadmissible in the sense of Section 2.5.3.1 as well. When we do not stop with inadmissibility claim, we call (α, β) admissible, and in this case we do produce an inference, specifically, as follows.

Processing observations.

- 1. We set $\overline{\mathcal{J}} = \{0, 1, ..., \hat{j} = \bar{j} 1\}, K = K(\alpha, \beta) = \max_{0 \le j \le \hat{j}} K^j$. Note that $\overline{\mathcal{J}}$ is nonempty due to $\bar{j} > 0$.²²
- 2. Given observation ω^K with independent across k components stemming from unknown signal $x_* \in X$ according to (2.106), we act as follows.
 - a) We set $\widehat{\mathcal{I}}_{-1}(\omega^K) = \{1, ..., M\} = \mathcal{I}_0.$
 - b) We look, one by one, at the values $j = 0, 1, ..., \hat{j}$. When processing j, we already have at our disposal subsets $\widehat{\mathcal{I}}_k(\omega^K) \subset \{1, ..., M\}, -1 \leq k < j$, and act as follows:
 - i. we apply the test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}^{j}}^{K_{j}}$ to the initial K_{j} components of the observation ω^{K} . Let $\mathcal{I}_{j}^{+}(\omega^{K})$ be the set of hypotheses H_{ij} , $i \in \mathcal{I}_{j}$, accepted by the test.
 - ii. it may happen that $\mathcal{I}_i^+(\omega^K) = \emptyset$; if it is so, we terminate.
 - iii. if $\mathcal{I}_{j}^{+}(\omega^{K})$ is nonempty, we look, one after one, at indexes $i \in \mathcal{I}_{j}^{+}(\omega^{K})$ and for such an *i*, check, for every $\ell \in \{-1, 0, ..., j - 1\}$, whether $i \in \widehat{\mathcal{I}}_{\ell}(\omega^{K})$. If it is the case for every $\ell \in \{-1, 0, ..., j - 1\}$, we call index *i* good.
 - iv. if good indexes in $\mathcal{I}_{j}^{+}(\omega^{K})$ are discovered, we define $\widehat{\mathcal{I}}_{j}(\omega^{K})$ as the set of these good indexes and process to the next value of j (if $j < \hat{j}$), or terminate (if $j = \hat{j}$). If there are no good indexes in $\mathcal{I}_{j}^{+}(\omega^{K})$, we terminate.
 - c) Upon termination, we have at our disposal a collection $\widehat{\mathcal{I}}_{j}(\omega^{K})$, $0 \leq j \leq \widetilde{j}(\omega^{K})$, of all sets $\widehat{\mathcal{I}}_{j}(\omega^{K})$ we have built (this collection can be empty, which we encode by setting $\widetilde{j}(\omega^{K}) = -1$). When $\widetilde{j}(\omega^{K}) = -1$, our inference remains undefined. Otherwise we select from the set $\widehat{\mathcal{I}}_{\widetilde{j}(\omega^{K})}(\omega^{K})$ an index $i_{\alpha,\beta}(\omega^{K})$, say, the smallest one, and claim that the point $x_{i_{\alpha,\beta}(\omega^{K})}$ is the "nearly closest" to x_{*} point among $x_{1}, ..., x_{M}$.

We have the following analogy of Proposition 2.37:

Proposition 2.38. Assuming (α, β) admissible, for the just defined inference $\omega^K \mapsto i_{\alpha,\beta}(\omega^K)$ and for every $x_* \in X$, denoting by $P_{x_*}^K$ the distribution of stationary K-repeated observation ω^K stemming from x_* one has

 $P_{x_*}^K\left\{\omega^K: i_{\alpha,\beta}(\omega^K) \text{ is well defined and } \|x_* - x_{i_{\alpha,\beta}(\omega^K)}\| \le \alpha\rho(x_*) + \beta\right\} \ge 1 - \epsilon.$ (2.116)

²²All the sets X_{i0} contain X and thus are nonempty, so that $\mathcal{I}_0 = \{1, ..., M\} \neq \emptyset$, and thus we cannot stop at step j = 0 due to $\mathcal{I}_0 = \emptyset$; and another possibility to stop at step j = 0 is ruled out by the fact that we are in the case when (α, β) is admissible.

Proof. Let us fix the signal $x_* \in X$ underlying observations ω^K . Same as in the proof of Proposition 2.37, let j_* be such that $\rho(x_*) = r_{j_*}$, and let $i_* \leq M$ be such that $x_* \in X_{i_*j_*}$; clearly, i_* and j_* are well defined, and the hypotheses H_{i_*j} , $0 \leq j \leq j_*$, are true. In particular, $X_{i_*j} \neq \emptyset$ when $j \leq j_*$, implying that $i_* \in \mathcal{I}_j$, $0 \leq j \leq j_*$, whence also $\hat{j} \geq j_*$.

For $0 \leq j \leq j_*$, let \mathcal{E}_j be the set of all realizations of ω^K such that

$$i_* \in \mathcal{I}_j^+(\omega^K) \& \{(i_*, i) \in \mathcal{C}_{\alpha, \beta}^j \,\forall i \in \mathcal{I}_j^+(\omega^K)\}.$$

Since $\mathcal{C}_{\alpha,\beta}^{j}$ -risk of the test $\mathcal{T}_{\mathcal{C}_{\alpha,\beta}^{j}}^{K_{j}}$ is $\leq \bar{\epsilon}$, we conclude that the $P_{x_{*}}^{K}$ -probability of \mathcal{E}_{j} is at least $1 - \bar{\epsilon}$, whence the $P_{x_{*}}^{K}$ -probability of the event

$$\mathcal{E} = \bigcap_{j=0}^{j_*} \mathcal{E}_j$$

is at least $1 - (N+1)\vec{\epsilon} = 1 - \epsilon$.

Now let

$$\omega^K \in \mathcal{E}.$$

Then, by the definition of \mathcal{E}_j , $j \leq j_*$,

- When $j \leq j_*$, we have $i_* \in \mathcal{I}_j^+(\omega^K)$, whence, by evident induction in $j, i_* \in \widehat{\mathcal{I}}_j(\omega^K)$ for all $j \leq j_*$.
- From the above item, $\tilde{j}(\omega^K) \geq j_*$; in particular, $i := i_{\alpha,\beta}(\omega^K)$ is well defined and turned out to be good at step $\tilde{j} \geq j_*$, implying that $i \in \hat{\mathcal{I}}_{j_*}(\omega^K) \subset \mathcal{I}_{j_*}^+(\omega^K)$.

Thus, $i \in \mathcal{I}_{j_*}^+(\omega^K)$, which combines with the definition of \mathcal{E}_{j_*} to imply that *i* and i_* are $\mathcal{C}_{\alpha,\beta}^{j_*}$ -close to each other, whence

$$\|x_{i(\alpha,\beta)(\omega^{K})} - x_{i_*}\| \le 2\bar{\alpha}r_{j_*} + \beta = 2\bar{\alpha}\rho(x_*) + \beta,$$

resulting in the desired relation

$$\|x_{i(\alpha,\beta)(\omega^{K})} - x_{*}\| \leq 2\bar{\alpha}\rho(x_{*}) + \beta + \|x_{i_{*}} - x_{*}\| \leq [2\bar{\alpha}+1]\rho(x_{*}) + \beta = \alpha\rho(x_{*}) + \beta. \square$$

2.5.3.3 "Near-optimality"

We augment the above simple constructions with the following

Proposition 2.39. Let in the nature for some positive integer \bar{K} , $\epsilon \in (0, 1/2)$ and a pair $(a,b) \geq 0$ there exists an inference $\omega^{\bar{K}} \mapsto i(\omega^{\bar{K}}) \in \{1, ..., M\}$ such that whenever $x_* \in X$, we have

$$\operatorname{Prob}_{\omega^{\bar{K}} \sim P_{\pi}^{\bar{K}}} \left\{ \left\| x_* - x_{i(\omega^{\bar{K}})} \right\| \le a\rho(x_*) + b \right\} \ge 1 - \epsilon.$$

Then the pair ($\alpha = 2a + 3, \beta = 2b$) is admissible in the sense of Section 2.5.3.1 (and thus – in the sense of Section 2.5.3.2), and for both our constructions – the

one from Section 2.5.3.1 and the one from Section 2.5.3.2) - one has

$$K(\alpha,\beta) \leq Ceil\left(2\frac{1+\ln(M(N+1))/\ln(1/\epsilon)}{1-\frac{\ln(4(1-\epsilon))}{\ln(1/\epsilon)}}\bar{K}\right);$$
(2.117)

Proof. Consider the situation of Section 2.5.3.1 (the situation of Section 2.5.3.2 can be processed in a completely similar fashion). Observe that with α, β as above, there exists a simple test deciding on a pair of hypotheses $H_{ij}, H_{i'j'}$ which are not $\mathcal{C}_{\alpha,\beta}$ -close to each other via stationary \bar{K} -repeated observation $\omega^{\bar{K}}$ with risk $\leq \epsilon$. Indeed, the desired test \mathcal{T} is as follows: given $ij \in \mathcal{J}, i'j' \in \mathcal{J}$, and observation $\omega^{\bar{K}}$, we compute $i(\omega^{\bar{K}})$ and accept H_{ij} if and only if $||x_{i(\omega^{\bar{K}})} - x_i|| \leq (a+1)r_j + b$, and accept $H_{i'j'}$ otherwise. Let us check that the risk of this test indeed is at most ϵ . Assume, first, that H_{ij} takes place. The $P_{x_*}^{\bar{K}}$ -probability of the event $\mathcal{E} : ||x_{i(\omega^{\bar{K}})} - x_*|| \leq a\rho(x_*) + b$ is at lest $1 - \epsilon$ due to the origin of $i(\cdot)$, and $||x_i - x_*|| \leq r_j$ since H_{ij} takes place, implying that $\rho(x_*) \leq r_j$ by the definition of $\rho(\cdot)$. Thus, in the case of \mathcal{E} it holds

$$\|x_{i(\omega^{\bar{K}})} - x_i\| \le \|x_{i(\omega^{\bar{K}})} - x_*\| + \|x_i - x_*\| \le a\rho(x_*) + b + r_j \le (a+1)r_j + b.$$

We conclude that if H_{ij} is true and $\omega^{\bar{K}} \in \mathcal{E}$, then the test \mathcal{T} accepts H_{ij} , and thus the $P_{x_*}^{\bar{K}}$ -probability for the simple test \mathcal{T} not to accept H_{ij} when the hypothesis takes place is $\leq \epsilon$.

Now let $H_{i'j'}$ take place, and let \mathcal{E} be the same event as above. When $\omega^{\bar{K}} \in \mathcal{E}$, which happens with the $P_{x_*}^{\bar{K}}$ -probability at least $1 - \epsilon$, we by exactly the same reasons as above have $\|x_{i(\omega\bar{K})} - x_{i'}\| \leq (a+1)r_{j'} + b$. It follows that when $H_{i'j'}$ takes place and $\omega^{\bar{K}} \in \mathcal{E}$, we have $\|x_{i(\omega\bar{K})} - x_i\| > (a+1)r_j + b$, since otherwise we would have

$$\begin{aligned} \|x_i - x_{i'}\| &\leq \|x_{i(\omega^{\bar{K}})} - x_i\| + \|x_{i(\omega^{\bar{K}})} - x_{i'}\| \leq (a+1)r_j + b + (a+1)r_{j'} + b \\ &\leq (a+1)(r_j + r_{j'}) + 2b = \frac{\alpha - 1}{2}(r_j + r_{j'}) + \beta, \end{aligned}$$

which contradicts the fact that ij and i'j' are not $\mathcal{C}_{\alpha,\beta}$ -close. Thus, whenever $H_{i'j'}$ holds true and \mathcal{E} takes place, we have $||x_{i(\omega\bar{K})} - x_i|| > (a+1)r_j + b$, implying that \mathcal{T} accepts $H_{i'j'}$. Thus, the $P_{x_*}^{\bar{K}}$ -probability not to accept $H_{i'j'}$ when the hypotheses if true is at most ϵ . From the just established fact that whenever $(ij, i'j') \notin \mathcal{C}_{\alpha,\beta}$, the hypotheses H_{ij} , $H_{i'j'}$ can be decided upon, via \bar{K} observations, with risk $\leq \epsilon < 0.5$ it follows that for ij, i'j' in question, the sets $A(X_{ij})$ and $A(X_{i'j'})$ do not intersect, so that (α, β) is an admissible pair.

Same as in the proof of Proposition 2.35, by basic properties of simple observation schemes, the fact that the hypotheses H_{ij} , $H_{i'j'}$ with $(ij, i'j') \notin C_{\alpha,\beta}$ can be decided upon via \bar{K} -repeated observations (2.106) with risk $\leq \epsilon < 1/2$ implies that $\epsilon_{ij,i'j'} \leq [2\sqrt{\epsilon(1-\epsilon)}]^{1/\bar{K}}$, whence, again by basic results on simple observation scheme (look once again at the proof of Proposition 2.35), the $C_{\alpha,\beta}$ -risk of K-observation detector-based test \mathcal{T}_K deciding on the hypotheses H_{ij} , $ij \in \mathcal{J}$, up to closeness $\mathcal{C}_{\alpha,\beta}$ does not exceed $\operatorname{Card}(\mathcal{J})[2\sqrt{\epsilon(1-\epsilon)}]^{K/\bar{K}} \leq M(N + 1)[2\sqrt{\epsilon(1-\epsilon)}]^{K/\bar{K}}$, and (2.117) follows. \Box

Comment. Proposition 2.39 says that in our problem, the "statistical toll" for quite large values of N and M is quite moderate: with $\epsilon = 0.01$, resolution $\theta = 1.001$

(which for all practical purposes is the same as no discretization of distances at all), D/r_N as large as 10^{10} , and M as large as 10,000, (2.117) reads $K = \text{Ceil}(10.7\bar{K})$ – not a disaster! The actual statistical toll in our construction is in replacing the "existing in the nature" a and b with $\alpha = 2\alpha + 3$ and $\beta = 2b$. And of course there is a huge computational toll for large M and N: we need to operate with large (albeit polynomial in M, N) number of hypotheses and detectors.

2.5.3.4 Numerical illustration

The toy problem we use to illustrate the approach presented in this Section is as follows:

A signal $x_* \in \mathbf{R}^n$ (it makes sense to think of x_* as of the restriction on the equidistant *n*-point grid in [0, 1] of a function of continuous argument $t \in [0, 1]$) is observed according to

$$\omega = Ax_* + \xi, \ \xi \sim \mathcal{N}(0, \sigma^2 I_n), \tag{2.118}$$

where A is "discretized integration:"

$$(Ax)_s = \frac{1}{n} \sum_{j=1}^s x_s, \ s = 1, ..., n.$$

We want to approximate x in the discrete version of L_1 -norm

$$\|y\| = \frac{1}{n} \sum_{s=1}^{n} |y_s|, y \in \mathbf{R}^n$$

by a low order polynomial.

In order to build the approximation, we use a single observation ω , stemming from x_* according to (2.118), to build 5 candidate estimates x_i , i = 1, ..., 5 of x_* . Specifically, x_i is the Least Squares polynomial, of degree $\leq i - 1$, approximation of x:

$$x_i = \operatorname*{argmin}_{y \in \mathcal{P}_{i-1}} \|Ay - \omega\|_2^2,$$

where \mathcal{P}_{κ} is the linear space of algebraic polynomials, of degree $\leq \kappa$, of discrete argument s varying in $\{1, 2, ..., n\}$. After the candidate estimates are built, we use additional K observations (2.118) "to select the model" – to select among our estimates the $\|\cdot\|$ -closest to x_* .

In the experiment to be reported we used n = 128 and $\sigma = 0.01$. The true signal x_* is plotted in magenta on the top of Figure 2.3; it is discretization of function of continuous argument $t \in [0, 1]$ which is linear, with slope 1, to the left of t = 0.5, and is linear, with slope -1, to the right of t = 0.5; at t = 0.5, the function has a jump. A priori information on the true signal is that it belongs to the box $\{x \in \mathbf{R}^n : ||x||_{\infty} \leq 1\}$. Sample polynomial approximations x_i of x_* , $1 \leq i \leq 5$, are plotted in blue on the top of Figure 2.3; their actual $|| \cdot ||$ -distances to x_* are as follows:

i	1	2	3	4	5	
$\ x_i - x_*\ $	0.534	0.354	0.233	0.161	0.172	



Figure 2.3: Signal (top, magenta) and its candidate estimates (top,blue). Bottom: the primitive of the signal.

As usual, the reliability tolerance ϵ was set to 0.01. We used N = 22 and $\theta = 2^{1/4}$, $\alpha = 3$, $\beta = 0.05$, resulting in K = 3. In a series of 1000 simulations of the resulting inference, all 1000 results correctly identified the $\|\cdot\|$ -closest to x_* candidate estimate, specifically, x_4 , in spite of the factor $\alpha = 3$ in (2.111). Surprisingly, the same holds true when we use the resulting inference with the reduced values of K, namely, K = 1 and K = 2, although the theoretical reliability guarantees deteriorate: with K = 1 and K = 2, theory guarantees the validity of (2.111) with probabilities 0.77 and 0.97, respectively.

2.6 SEQUENTIAL HYPOTHESIS TESTING

2.6.1 Motivation: Election Polls

Consider the question as follows:

One of L candidates for an office is about to be selected by population-wide majority vote. Every member of the population votes for exactly one of the candidates. How to predict the winner via an opinion poll?

A (naive) model of situation could be as follows. Let us represent the preference of a particular voter by his *preference vector* – basic orth e in \mathbf{R}^L with unit entry in a position ℓ meaning that the voter is about to vote for the ℓ -th candidate. The entries μ_{ℓ} in the average μ , over the population, of these vectors are the fractions of votes in favor of ℓ -th candidate, and the elected candidate is the one "indexing" the largest of μ_{ℓ} 's. Now assume that we select at random, from the uniform distribution, a member of the population and observe his preference vector. Our observation ω is a realization of discrete random variable taking values in the set $\Omega = \{e_1, ..., e_L\}$

of basic orths in \mathbf{R}^{L} , and μ is the distribution of ω (technically, the density of this distribution w.r.t. the counting measure Π on Ω). Selecting a small threshold δ and assuming that the true – unknown to us – μ is such that the largest entry in μ is at least by δ larger than every other entry and that $\mu_{\ell} \geq \frac{1}{N}$ for all ℓ , N being the population size²³, the fact that ℓ -th candidate wins the elections means that

$$\mu \in M_{\ell} = \{ \mu \in \mathbf{R}^d : \mu_i \ge \frac{1}{N}, \sum_i \mu_i = 1, \mu_\ell \ge \mu_i + \delta \forall (i \neq \ell) \}$$

$$\subset \mathcal{M} = \{ \mu \in \mathbf{R}^d : \mu > 0, \sum_i \mu_i = 1 \}.$$

In an (idealized) poll, we select at random a number K of voters and observe their preferences, thus arriving at a sample $\omega^K = (\omega_1, ..., \omega_K)$ of observations drawn, independently of each other, from unknown distribution μ on Ω , with μ known to belong to $\bigcup_{\ell=1}^{L} M_{\ell}$, and to predict the winner is the same as to decide on Lconvex hypotheses, $H_1, ..., H_L$, in the Discrete o.s., with H_{ℓ} stating that $\omega_1, ..., \omega_K$ are drawn, independently of each other, from a distribution $\mu \in M_{\ell}$. What we end up with, is the problem of deciding on L convex hypotheses in the Discrete o.s. with L-element Ω via stationary K-repeated observations.

Illustration. Consider two-candidate elections; now the goal of a poll is, given K independent of each other realizations $\omega_1, ..., \omega_K$ of random variable ω taking value $\chi = 1, 2$ with probability $\mu_{\chi}, \mu_1 + \mu_2 = 1$, to decide what is larger, μ_1 or μ_2 . As explained above, we select somehow a threshold δ and impose on the unknown μ a priori assumption that the gap between the largest and the next largest (in our case – just the smallest) entry of μ is at least δ , thus arriving at two hypotheses:

$$H_1: \mu_1 \ge \mu_2 + \delta, \quad H_2: \mu_2 \ge \mu_1 + \delta,$$

which is the same as

$$H_1: \mu \in M_1 = \{\mu: \mu_1 \ge \frac{1+\delta}{2}, \mu_2 \ge 0, \mu_1 + \mu_2 = 1\}, H_2: \mu \in M_2 = \{\mu: \mu_2 \ge \frac{1+\delta}{2}, \mu_1 \ge 0, \mu_1 + \mu_2 = 1\}.$$

We now want to decide on these two hypotheses from stationary K-repeated observations. We are in the case of simple (specifically, Discrete) o.s.; the optimal detector as given by Theorem 2.25 stems from the optimal solution (μ^*, ν^*) to the convex optimization problem

$$\varepsilon_{\star} = \max_{\mu \in M_1, \nu \in M_2} \left[\sqrt{\mu_1 \nu_1} + \sqrt{\mu_2 \nu_2} \right], \tag{2.119}$$

the optimal balanced single-observation detector is

$$\phi_*(\omega) = f_*^T \omega, \ f_* = \frac{1}{2} [\ln(\mu_1^*/\nu_1^*); \ln(\mu_2^*/\nu_2^*)]$$

(recall that we encoded observations ω_k by basic orths from \mathbf{R}^2), the risk of this detector being ε_{\star} . In other words,

$$\begin{split} \mu^* &= [\frac{1+\delta}{2}; \frac{1-\delta}{2}], \, \nu^* = [\frac{1-\delta}{2}; \frac{1+\delta}{2}], \, \varepsilon_\star = \sqrt{1-\delta^2}, \\ f_* &= \frac{1}{2} \left[\ln((1+\delta)/(1-\delta)); \ln((1-\delta)/(1+\delta)) \right]. \end{split}$$

²³ with the size N of population in the range of tens of thousands and δ like 1/N, both these assumptions seem to be quite realistic.

The optimal balanced K-observation detector and its risk are

$$\phi_*^{(K)}(\underbrace{\omega_1, ..., \omega_K}_{\omega^K}) = f_*^T(\omega_1 + ... + \omega_K), \ \varepsilon_*^{(K)} = (1 - \delta^2)^{K/2}.$$

The near-optimal K-observation test $\mathcal{T}_{\phi_*}^K$ accepts H_1 and rejects H_2 if $\phi_*^{(K)}(\omega^K) \geq 0$, otherwise it accepts H_2 and rejects H_1 . Both risks of this test do not exceed $\varepsilon_*^{(K)}$.

Given risk level ϵ , we can identify the minimal "poll size" K for which the risks Risk₁, Risk₂ of the test $\mathcal{T}_{\phi_*}^K$ do not exceed ϵ . This poll size depends on ϵ and on our a priory "hypotheses separation" parameter δ : $K = K_{\epsilon}(\delta)$. Some impression on this size can be obtained from Table 2.1, where, as in all subsequent "election illustrations," ϵ is set to 0.01. We see that while poll sizes for "landslide" elections are surprisingly low, reliable prediction of the results of "close run" elections requires surprisingly high sizes of the polls. Note that this phenomenon reflects reality (to the extent at which the reality is captured by our model²⁴); indeed, from Proposition 2.29 we know that our poll size is within an explicit factor, depending solely on ϵ , from the "ideal" poll sizes – the smallest ones which allow to decide upon H_1 , H_2 with risk $\leq \epsilon$. For $\epsilon = 0.01$, this factor is about 2.85, meaning that when $\delta = 0.01$, the ideal poll size is larger than 32,000. In fact, we can build more accurate lower bounds on the sizes of ideal polls, specifically, as follows. When computing the optimal detector ϕ_* , we get, as a byproduct, two distributions, μ^* , ν^* obeying H_1, H_2 , respectively. Denoting by μ_K^*, ν_K^* the distributions of K-element i.i.d. samples drawn from μ^* and ν^* , the risk of deciding on two simple hypotheses on the distribution of ω^{K} , stating that this distribution is μ_{K}^{*} , respectively, ν_{K}^{*} can be only smaller than the risk of deciding on H_1 , H_2 via K-repeated stationary observations. On the other hand, the former risk can be lower-bounded by one half of the total risk of deciding on our two simple hypotheses, and the latter risk admits a sharp lower bound given by Proposition 2.2, namely,

$$\sum_{i_1,\dots,i_K\in\{1,2\}} \min\left[\prod_{\ell} \mu_{i_\ell}^*, \prod_{\ell} \nu_{i_\ell}^*\right] = \mathbf{E}_{(i_1,\dots,i_K)} \left\{\min\left[\prod_{\ell} (2\mu_{i_\ell}^*), \prod_{\ell} (2\nu_{i_\ell}^*)\right]\right\},$$

with the expectation taken w.r.t independent tuples of K integers taking values 1 and 2 with probabilities 1/2. Of course, when K is in the range of few tens and more, we cannot compute the above 2^{K} -term sum exactly; however, we can use Monte Carlo simulation in order to estimate the sum reliably within moderate accuracy, like 0.005, and use this estimate to lower-bound the value of K for which "ideal" K-observation test decides on H_1 , H_2 with risks ≤ 0.01 . Here are the resulting lower bounds (along with upper bounds stemming from the data in Table

²⁴in actual opinion polls, additional information is used; for example, in reality voters can be split into groups according to their age, sex, education, income, etc., etc., with variability of preferences within a group essentially lower than across the entire population; when planning a poll, respondents are selected at random within these groups, with a prearranged number of selections in every group, and their preferences are properly weighted, yielding more accurate predictions as compared to the case when the respondents are selected from the uniform distribution. In other words, in actual polls a non-trivial a priori information on the "true" distribution of preferences is used – something we do not have in our naive model.

δ	0.5623	0.3162	0.1778	0.1000	0.0562	0.0316	0.0177	0.0100	-
$K_{0.01}(\delta), L = 2$	25	88	287	917	2908	9206	29118	92098	-
$K_{0.01}(\delta), L = 5$	32	114	373	1193	3784	11977	37885	119745	

Table 2.1: Sample of values of poll size $K_{0.01}(\delta)$ as a function of δ for 2-candidate (L = 2) and 5-candidate (L = 5) elections. Values of δ form a decreasing geometric progression with ratio $10^{-1/4}$.

2.1):

δ	0.5623	0.3162	0.1778	0.1000	0.0562	0.0316	0.0177	0.0100
$\underline{K}, \overline{K}$	14, 25	51,88	166, 287	534,917	1699,2908	5379,9206	17023, 29122	53820,92064
Lower (<u>K</u>) and upper (\overline{K}) bounds on the "ideal" poll sizes								

We see that the poll sizes as yielded by our machinery are within factor 2 of the "ideal" poll sizes.

Clearly, the outlined approach can be extended to L-candidate elections with $L \geq 2$. We model the corresponding problem as the one where we need to decide, via stationary K-repeated observations drawn from unknown probability distribution μ on L-element set, on L hypotheses

$$H_{\ell}: \mu \in M_{\ell} = \{ \mu \in \mathbf{R}^{d} : \mu_{i} \ge \frac{1}{N}, i \le L, \sum_{i} \mu_{i} = 1, \mu_{\ell} \ge \mu_{\ell'} + \delta \,\forall (\ell' \neq \ell) \}, \, \ell \le L;$$
(2.120)

here $\delta > 0$ is a selected in advance threshold small enough to believe that the actual preferences of the voters correspond to $\mu \in \bigcup_{\ell} M_{\ell}$. Defining closeness C in the strongest possible way $-H_{\ell}$ is close to $H_{\ell'}$ if and only if $\ell = \ell'$, predicting the outcome of elections with risk ϵ becomes the problem of deciding upon our multiple hypotheses with C-risk $\leq \epsilon$, and we can use the pairwise detectors yielded by Theorem 2.25 to identify the smallest possible $K = K_{\epsilon}$ such that the test \mathcal{T}_{C}^{K} from Section 2.5.2.3 is capable to decide upon our L hypotheses with C-risk $\leq \epsilon$. Numerical illustration of the performance of this approach in 5-candidate elections is presented in Table 2.1 (where ϵ is set to 0.01).

2.6.2 Sequential hypothesis testing

In view of the above analysis, when predicting outcomes of "close run" elections, huge poll sizes are a must. It, however, does not mean that nothing can be done in order to build more reasonable opinion polls. The classical related statistical idea, going back to Wald [144] is to pass to *sequential tests* where the observations are processed one by one, and at every time we either accept some of our hypotheses and terminate, or conclude that the observations obtained so far are insufficient to make a reliable inference and pass to the next observation. The idea is that a properly built sequential test, while still ensuring a desired risk, will be able to make "early decisions" in the case when the distribution underlying observations is "well inside" the true hypothesis and thus is far from the alternatives. Let us show how to utilize our machinery in building a sequential test for the problem of predicting the outcome of *L*-candidate elections; thus, our goal is, given a small threshold δ , to decide upon *L* hypotheses (2.120). Let us act as follows.



Figure 2.4: 3-candidate hypotheses in probabilistic simplex Δ_3 :

[green]	M_1	dark green + light green: candidate A wins with margin $\geq \delta_S$
[green]	M_1^s	dark green: candidate A wins with margin $\geq \delta_s > \delta_S$
[red]	M_2	dark red + pink: candidate B wins with margin $\geq \delta_S$
[red]	M_2^s	dark red: candidate B wins with margin $\geq \delta_s > \delta_S$
[blue]	M_3	dark blue + light blue: candidate C wins with margin $\geq \delta_S$
[blue]	M_3^s	dark blue: candidate C wins with margin $\geq \delta_s > \delta_S$
closen	ess: hy	vpotheses in the tuple $\{G_{2\ell-1}^s : \mu \in M_\ell, G_{2\ell}^s : \mu \in M_\ell^s, 1 < \ell <$

 C_s closeness: hypotheses in the tuple $\{G_{2\ell-1}^s : \mu \in M_\ell, G_{2\ell}^s : \mu \in M_\ell^s, 1 \leq \ell \leq 3\}$ are *not* C_s -close to each other, if the corresponding *M*-sets are of different colors and at least one the sets is dark-painted, like M_1^s and M_2 , but not M_1 and M_2 .

- 1. We select a factor $\theta \in (0, 1)$, say, $\theta = 10^{-1/4}$, and consider thresholds $\delta_1 = \theta$, $\delta_2 = \theta \delta_1$, $\delta_3 = \theta \delta_2$, and so on, until for the first time we get a threshold $\leq \delta$; to save notation, we assume that this threshold is exactly δ , and let the number of the thresholds be S.
- 2. We split somehow (e.g., equally) the risk ϵ which we want to guarantee into S portions ϵ_s , $1 \le s \le S$, so that ϵ_s are positive and

$$\sum_{s=1}^{S} \epsilon_s = \epsilon$$

3. For $s \in \{1, 2, ..., S\}$, we define, along with the hypotheses H_{ℓ} , the hypotheses

 $H^{s}_{\ell}: \mu \in M^{s}_{\ell} = \{ \mu \in M_{\ell}: \mu_{\ell} \ge \mu_{\ell'} + \delta_{s}, \, \forall (\ell' \neq \ell) \}, \, \ell = 1, ..., L,$

see Figure 2.4, and introduce 2L hypotheses $G_{2\ell-1}^s = H_\ell$, and $G_{2\ell}^s = H_\ell^s$, $1 \le \ell \le L$. It is convenient to color these hypotheses in L colors, with $G_{2\ell-1}^s = H_\ell$ and $G_{2\ell}^s = H_\ell^s$ assigned color ℓ . We define also *s*-th closeness C_s as follows:

When s < S, hypotheses G_i^s and G_j^s are \mathcal{C}_s -close to each other if either they are of the same color, or they are of different colors and both of them have odd indexes (that is, one of them is H_ℓ , and another one is $H_{\ell'}$ with $\ell \neq \ell'$).

When s = S (in this case $G_{2\ell-1}^S = H_\ell = G_{2\ell}^S$), hypotheses G_ℓ^S and $G_{\ell'}^S$ are \mathcal{C}_S -close to each other if and only if they are of the same color, i.e., both coincide with the same hypothesis H_ℓ .

Observe that G_i^s is a convex hypothesis:

$$G_i^s : \mu \in Y_i^s \qquad [Y_{2\ell-1}^s = M_\ell, Y_{2\ell}^s = M_\ell^s]$$

The key observation is that when G_i^s and G_j^s are not \mathcal{C}_s -close, the sets Y_i^s and Y_j^s are "separated" by at least δ_s , meaning that for some vector $e \in \mathbf{R}^L$ with just two nonnegative entries, equal to 1 and -1, we have

$$\min_{\mu \in Y_i^s} e^T \mu \ge \delta_s + \max_{\mu \in Y_j^s} e^T \mu.$$
(2.121)

Indeed, let G_i^s and G_j^s be not \mathcal{C}_s -close to each other. That means that the hypotheses are of different colors, say, ℓ and $\ell' \neq \ell$, and at least one of them has even index; w.l.o.g. we can assume that the even-indexed hypothesis is G_i^s , so that

$$Y_i^s \subset \{\mu : \mu_\ell - \mu_{\ell'} \ge \delta_s\},\$$

while Y_j^s is contained in the set $\{\mu : \mu_{\ell'} \ge \mu_{\ell}\}$. Specifying *e* as the vector with just two nonzero entries, ℓ -th equal to 1 and ℓ' -th equal to -1, we ensure (2.121).

4. For $1 \leq s \leq S$, we apply the construction from Section 2.5.2.3 to identify the smallest K = K(s) for which the test \mathcal{T}_s yielded by this construction as applied to stationary K-repeated observation allows to decide on the hypotheses $G_1^s, ..., G_{2L}^s$ with \mathcal{C}_s -risk $\leq \epsilon_s$; the required K exists due to the already mentioned separation of members in a pair of not \mathcal{C}_s -close hypotheses G_i^s, G_j^s . It is easily seen that $K(1) \leq K(2) \leq ... \leq K(S-1)$; however, it may happen that K(S-1) > K(S), the reason being that \mathcal{C}_S is defined differently than \mathcal{C}_s with s < S. We set

$$\mathcal{S} = \{ s \le S : K(s) \le K(S) \}.$$

For example, this is what we get in *L*-candidate Opinion Poll problem when S = 8, $\delta = \delta_S = 0.01$, and for properly selected ϵ_s with $\sum_{s=1}^{8} \epsilon_s = 0.01$:

$\mid L \mid$	K(1)	K(2)	K(3)	K(4)	K(5)	K(6)	K(7)	K(8)	
2	177	617	1829	5099	15704	49699	153299	160118	
5	208	723	2175	6204	19205	60781	188203	187718	
$S = 8, \delta_s = 10^{-s/4}.$									

 $S = \{1, 2, ..., 8\}$ when L = 2 and $S = \{1, 2, ..., 6\} \cup \{8\}$ when L = 5.

5. Our sequential test \mathcal{T}_{seq} works in attempts $s \in S$ – it tries to make conclusions after observing $K(s), s \in S$, realizations ω_k of ω . At s-th attempt, we apply the test \mathcal{T}_s to the collection $\omega^{K(s)}$ of observations obtained so far to decide on hypotheses $G_1^s, ..., G_{2L}^s$. If \mathcal{T}_s accepts some of these hypotheses and all accepted hypotheses are of the same color, let it be ℓ , the sequential test accepts the hypothesis H_ℓ and terminates, otherwise we continue to observe the realizations of ω (when s < S) or terminate with no hypotheses accepted/rejected (when s = S).

It is easily seen that the risk of the outlined sequential test \mathcal{T}_{seq} does not exceed ϵ , meaning that whatever be a distribution $\mu \in \bigcup_{\ell=1}^{L} M_{\ell}$ underlying observations $\omega_1, \omega_2, ..., \omega_{K(S)}$ and ℓ_* such that $\mu \in M_{\ell_*}$, the μ -probability of the event \mathcal{T}_{seq} accepts exactly one hypothesis, namely, H_{ℓ_*}

is at least $1 - \epsilon$.

Indeed, observe, first, that the sequential test always accepts at most one of the hypotheses $H_1, ..., H_L$. Second, let $\omega_k \sim \mu$ with μ obeying H_{ℓ_*} . Consider events

 $E_s, s \in \mathcal{S}$, defined as follows:

- when s < S, E_s is the event "the test \mathcal{T}_s as applied to observation $\omega^{K(s)}$ does not accept the true hypothesis $G^s_{2\ell_*-1} = H_{\ell_*}$ ";
- E_S is the event "as applied to observation $\omega^{K(S)}$, the test \mathcal{T}_S does not accept the true hypothesis $G_{2\ell_*-1}^S = H_{\ell_*}$ or accepts a hypothesis not \mathcal{C}_S -close to $G_{2\ell_*-1}^S$."

Note that by our selection of K(s)'s, the μ -probability of E_s does not exceed ϵ_s , so that the μ -probability of no one of the events E_s , $s \in S$, taking place is at least $1 - \epsilon$. To justify the above claim on the risk of our sequential test, all we need is to verify that when no one of the events E_s , $s \in S$, takes place, then the sequential test accepts the true hypothesis H_{ℓ_s} . Verification is immediate: let the observations be such that no one of the events E_s takes place. We claim that in this case

(a) The sequential test does accept a hypothesis – if this does not happen at s-th attempt with some s < S, it definitely happens at S-th attempt.

- Indeed, since E_S does not take place, T_S accepts G^S_{2ℓ*-1} and all other hypotheses, if any, accepted by T_S are C_S-close to G^S_{2ℓ*-1}, implying by construction of C_S that T_S does accept hypotheses, and all these hypotheses are of the same color, that is, the sequential test at S-th attempt does accept a hypothesis.
 (b) The sequential test does not accept a wrong hypothesis.
- Indeed, assume that the sequential test accepts a wrong hypothesis, $H_{\ell'}$, $\ell' \neq \ell_*$, and it happens at s-th attempt, and let us lead this assumption to a contradiction. Observe that under our assumption the test \mathcal{T}_s as applied to observation $\omega^{K(s)}$ does accept some hypothesis G_i^s , but does not accept the true hypothesis $G_{2\ell_*-1}^s = H_{\ell_*}$ (indeed, assuming the latter hypothesis to be accepted, its color, which is ℓ_* , should be the same as the color ℓ' of G_i^s (we are in the case when the sequential test accepts $H_{\ell'}$ at s-th attempt!); since in fact $\ell' \neq \ell_*$, the above assumption leads to a contradiction). On the other hand, we are in the case when E_s does not take place, that is, \mathcal{T}_s does accept the true hypothesis $G_{2\ell_*-1}^s$, and we arrive at the desired contradiction.
- (a) and (b) provide us with a verification we were looking for.

Discussion and illustration. It can be easily seen that when $\epsilon_s = \epsilon/S$ for all s, the worst-case duration K(S) of our sequential test is within a logarithmic in SL factor of the duration of any other test capable to decide on our L hypotheses with risk ϵ . At the same time it is easily seen that when the distribution μ of our observation is "deeply inside" some set M_{ℓ} , specifically, $\mu \in M_{\ell}^s$ for some $s \in S$, s < S, then the μ -probability to terminate not later than after just K(s) realizations ω_k of $\omega \sim \mu$ are observed and to infer correctly what is the true hypothesis is at least $1 - \epsilon$. Informally speaking, in the case of "landslide" elections, a reliable prediction of elections' outcome will be made after a relatively small number of respondents are interviewed.

Indeed, let $s \in S$ and $\omega_k \sim \mu \in M^s_{\ell}$, so that μ obeys the hypothesis $G^s_{2\ell}$. Consider the s events E_t , $1 \leq t \leq s$, defined as follows:

- For t < s, E_t occurs when the sequential test terminates at attempt t with accepting, instead of H_ℓ , wrong hypothesis $H_{\ell'}$, $\ell' \neq \ell$. Note that E_t can take place only when \mathcal{T}_t does not accept the true hypothesis $G_{2\ell}^s = H_\ell^s$ (why?), and μ -probability of this outcome is $\leq \epsilon_t$.
- E_s occurs when \mathcal{T}_s does not accept the true hypothesis $G_{2\ell}^s$ or accepts it along with some hypothesis G_j^s , $1 \leq j \leq 2L$, of color different from ℓ . Note that we are in the situation where the hypothesis $G_{2\ell}^s$ is true, and, by construction of \mathcal{C}_s , all hypotheses \mathcal{C}_s -close to $G_{2\ell}^s$ are of the same color ℓ as $G_{2\ell}^s$. Recalling what \mathcal{C}_s -risk is and that the \mathcal{C}_s -risk of \mathcal{T}_s is $\leq \epsilon_s$, we conclude that the μ -probability

120

of E_s is at most ϵ_s .

The bottom line is that μ -probability of the event $\bigcup_{t\leq s} E_t$ is at most $\sum_{t=1}^s \epsilon_t \leq \epsilon$; by construction of the sequential test, if the event $\bigcup_{t\leq s} E_t$ does not take place, the test terminates in course of the first s attempts with accepting the correct hypothesis H_{ℓ} . Our claim is justified.

Numerical illustration. To get an impression of the "power" of sequential hypothesis testing, here are the data on the durations of non-sequential and sequential tests with risk $\epsilon = 0.01$ for various values of δ ; in the sequential tests, $\theta = 10^{-1/4}$ is used. The worst-case data for 2-candidate and 5-candidate elections are as follows (below "volume" stands for the number of observations used by test)

δ	0.5623	0.3162	0.1778	0.1000	0.0562	0.0316	0.0177	0.0100		
K, L = 2	25	88	287	917	2908	9206	29118	92098		
S & K(S), L = 2	1&25	2&152	3&499	4&1594	5&5056	6&16005	7&50624	8&160118		
K, L = 5	32	114	373	1193	3784	11977	37885	119745		
S & K(S), L = 5	1&32	2&179	3&585	4&1870	5&5931	6&18776	7&59391	8&187720		
Volum	Volume of non-sequential test (K) number of stages (S) and worst-case volume									

(K(S)) of sequential test as functions of threshold $\delta = \delta_S$. Risk ϵ is set to 0.01.

As it should be, the worst-case volume of sequential test is essentially worse than the volume of the non-sequential test²⁵. This being said, let us look what happens in the "average," rather than the worst, case, specifically, let us look at the empirical distribution of the volume when the distribution μ of observations is selected in the *L*-dimensional probabilistic simplex $\Delta_L = \{\mu \in \mathbf{R}^L : \mu \ge 0, \sum_{\ell} \mu_{\ell} = 1\}$ at random. Here is the empirical statistics of test volume obtained when drawing μ from the uniform distribution on $\bigcup_{\ell \le L} M_{\ell}$ and running the sequential test²⁶ on observations drawn from the selected μ :

L	risk	median	mean	50%	55%	60%	65%
2	0.0010	177	9182	177	177	177	397
5	0.0040	1449	18564	1449	2175	2175	4189
L	70%	75%	80%	85%	90%	95%	100%
2	617	617	1223	1829	8766	87911	160118
5	6204	12704	19205	39993	60781	124249	187718

The data on empirical risk (column "risk") and volume (columns "median...100%") of Sequential test. Column "X%": empirical X%-quantile of test volume.

The data in the table are obtained from 1,000 experiments. We see that with the Sequential test, "typical" numbers of observations before termination are much less than the worst-case values of these numbers. For example, in as much as 80% of experiments these numbers were below quite reasonable levels, at least in the case L = 2. Of course, what is "typical," and what is not, depends on how we generate μ 's (scientifically speaking, this is called "prior Bayesian distribution"); were our generation more likely to produce "close run" distributions, the advantages of sequential decision making would be reduced. This ambiguity is, however, unavoidable when attempting to go beyond worst-case-oriented analysis.

²⁵the reason is twofold: first, for s < S we pass from deciding on L hypotheses to deciding on 2L of them; second, the desired risk ϵ is now distributed among several tests, so that each of them should be more reliable than the non-sequential test with risk ϵ .

 $^{^{26} {\}rm corresponding \ to} \ \delta = 0.01, \ \theta = 10^{-1/4} \ {\rm and} \ \epsilon = 0.01$

2.6.3 Concluding remarks

Application of our machinery to sequential hypothesis testing is in no sense restricted to the simple election model considered so far. A natural general setup we can handle is as follows:

We are given a simple observation scheme \mathcal{O} and a number L of related convex hypotheses, colored in d colors, on the distribution of an observation, with distributions obeying hypotheses of different colors being distinct from each other. Given risk level ϵ , we want to infer $(1-\epsilon)$ -reliably the color of the distribution underlying observations (i.e., the color of the hypothesis obeyed by this distribution) from stationary K-repeated observations, utilizing as small number of observations as possible.

For detailed description of our related constructions and results, an interested reader is referred to [87].

2.7 MEASUREMENT DESIGN IN SIMPLE OBSERVATION SCHEMES

2.7.1 Motivation: Opinion Polls revisited

Consider the same situation as in Section 2.6.1 – we want to use opinion poll to predict the winner in a population-wide elections with L candidates. When addressing this situation earlier, no essential a priori information on the distribution of voters' preferences was available. Now consider the case when the population is split into I groups (according to age, sex, income, etc., etc.), with *i*-th group forming fraction θ_i of the entire population, and we have at our disposal, at least for some i, a nontrivial a priori information about the distribution p^i of the preferences across group # i (ℓ -th entry p_{ℓ}^i in p^i is the fraction of voters of group i voting for candidate ℓ). For example, we could know in advance that at least 90% of members of group #1 vote for candidate #1, and at least 85% of members of group #2 vote for candidate #2; no information of this type for group #3 is available. In this situation it would be wise to select respondents in the poll via two-stage procedure, first – selecting at random, with probabilities q_1, \ldots, q_I , the group from which the next respondent will be picked, and second – selecting the respondent from this group at random according to the uniform distribution on the group. When q_i are proportional to the sizes of the groups (i.e., $q_i = \theta_i$ for all i), we come back to selecting respondents at random from the uniform distribution over the entire population; the point, however, is that in the presence of a priori information, it makes sense to use q_i different from θ_i , specifically, to make the ratios q_i/θ_i "large" or "small" depending on whether a priori information on group #i is poor or rich.

The story we just have told is an example of situation when we can "design measurements" – draw observations from a distribution which partly is under our control. Indeed, what in fact happens in the story, is the following. "In the nature" there exist I probabilistic vectors $p^1, ..., p^I$ of dimension L representing distributions of voting preferences within the corresponding groups; the distribution of preferences across the entire population is $p = \sum_i \theta_i p^i$. With two-stage selection of respondents, the outcome of a particular interview becomes a pair (i, ℓ) , with i

identifying the group to which the respondent belongs, and ℓ identifying the candidate preferred by this respondent. In subsequent interviews, the pairs (i, ℓ) – these are our observations – are drawn, independently of each other, from the probability distribution on the pairs (i, ℓ) , $i \leq I$, $\ell \leq L$, with the probability of an outcome (i, ℓ) equal to

 $p(i,\ell) = q_i p_\ell^i.$

Thus, we find ourselves in the situation of stationary repeated observations stemming from the Discrete o.s. with observation space Ω of cardinality *IL*; the distribution from which the observations are drawn is a probabilistic vector μ of the form

$$\mu = Ax$$

where

• $x = [p^1; ...; p^I]$ is the "signal" underlying our observations and representing the preferences of the population; this signal is selected by the nature in the known to us set \mathcal{X} defined in terms of our a priori information on $p^1, ..., p^I$:

$$\mathcal{X} = \{ x = [x^1; ...; x^I] : x^i \in \Pi_i, 1 \le i \le I \},$$
(2.122)

where Π_i are the sets, given by our a priori information, of possible values of the preference vectors p^i of the voters from *i*-th group. In the sequel, we assume that Π_i are convex compact subsets in the positive part $\Delta_L^o = \{p \in \mathbf{R}^L : p > 0, \sum_{\ell} p_{\ell} = 1\}$ of the *L*-dimensional probabilistic simplex;

• A is "sensing matrix" which, to some extent, is under our control; specifically,

$$A[x^1; ...; x^I] = [q_1 x^1; q_2 x^2; ...; q_I x^I],$$
(2.123)

with $q = [q_1; ...; q_I]$ fully controlled by us (up to the fact that q must be a probabilistic vector).

Note that in the situation under consideration the hypotheses we want to decide upon can be represented by convex sets in the space of signals, with particular hypothesis stating that the observations stem from a distribution μ on Ω , with μ belonging to the image of some convex compact set $X_{\ell} \subset \mathcal{X}$ under the mapping $x \mapsto \mu = Ax$. For example, the hypotheses

$$H_{\ell}: \mu \in M_{\ell} = \{ \mu \in \mathbf{R}^{L}: \sum_{i} \mu_{i} = 1, \mu_{i} \ge \frac{1}{N}, \mu_{\ell} \ge \mu_{\ell}' + \delta, \ \ell' \neq \ell \}, \ 1 \le \ell \le L$$

considered in Section 2.6.1 can be expressed in terms of the signal $x = [x^1; ...; x^I]$:

$$H_{\ell}: \mu = Ax, x \in X_{\ell} = \left\{ x = [x^{1}; ...; x^{I}]: \begin{array}{l} x^{i} \ge 0, \sum_{\ell} x^{i}_{\ell} = 1 \forall i \le I \\ \sum_{i} \theta_{i} x^{i}_{\ell} \ge \sum_{i} \theta_{i} x^{i}_{\ell'} + \delta \forall (\ell' \neq \ell) \\ \sum_{i} \theta_{i} x^{i}_{j} \ge \frac{1}{N}, \forall j \end{array} \right\}.$$

$$(2.124)$$

The challenge we intend to address is as follows: so far, we were interested in inferences from observations drawn from distributions selected "by nature." Now our goal is to make inferences from observations drawn from a distribution selected partly by the nature and partly by us: the nature selects the signal x, we select from

some set matrix A, and the observations are drawn from the distribution Ax. As a result, we arrive at a completely new for us question: how to utilize the freedom in selecting A in order to improve our inferences (this is somehow similar to what in statistics is called "design of experiments.")

2.7.2 Measurement Design: SetUp

In what follows we address measurement design in simple observation schemes, and our setup is as follows (to make our intensions transparent, we illustrate our general setup by explaining how it should be specified to cover the outlined two-stage Design of Opinion Polls – DOP for short).

Given are

• simple observation scheme $\mathcal{O} = (\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F})$, specifically, Gaussian, Poisson or Discrete one, with $\mathcal{M} \subset \mathbf{R}^d$.

In DOP, \mathcal{O} is the Discrete o.s. with $\Omega = \{(i, \ell) : 1 \leq i \leq I, 1 \leq \ell \leq L\}$, that is, points of Ω are the potential outcomes "reference group, preferred candidate" of individual interviews.

• a nonempty closed convex signal space $\mathcal{X} \subset \mathbf{R}^n$, along with L nonempty convex compact subsets X_{ℓ} of \mathcal{X} , $\ell = 1, ..., L$.

In DOP, \mathcal{X} is the set (2.122) comprised by tuples of allowed distributions of voters' preferences from various groups, and X_{ℓ} are the sets (2.124) of signals associated with the hypotheses H_{ℓ} we intend to decide upon.

• a nonempty convex compact set \mathcal{Q} in some \mathbb{R}^N along with a continuous mapping $q \mapsto A_q$ acting from \mathcal{Q} into the space of $d \times n$ matrices such that

$$\forall (x \in \mathcal{X}, q \in \mathcal{Q}) : A_q x \in \mathcal{M}.$$
(2.125)

In DOP, Q is the set of probabilistic vectors $q = [q_1; ...; q_I]$ specifying our measurement design, and A_q is the matrix of the mapping (2.123).

• a closeness C on the set $\{1, ..., L\}$ (that is, a set C of pairs (i, j) with $1 \le i, j \le L$ such that $(i, i) \in C$ for all $i \le L$ and $(j, i) \in C$ whenever $(i, j) \in C$), and a positive integer K.

In DOP, the closeness S is as strict as it could be -i is close to j if and only if $i = j^{27}$, and K is the total number of interviews in the poll.

We can associate with $q \in \mathcal{Q}$ and every one of X_{ℓ} , $\ell \leq L$, nonempty convex compact sets M_{ℓ}^{q} in the space \mathcal{M} :

$$M^q_\ell = \{A_q x : x \in X_\ell\}$$

and hypotheses H_{ℓ}^q on K-repeated stationary observations $\omega^K = (\omega_1, ..., \omega_K)$, with H_{ℓ}^q stating that ω_k , k = 1, ..., K, are drawn, independently of each other, from a distribution $\mu \in M_{\ell}^q$, $\ell = 1, ..., L$. Closeness \mathcal{C} can be thought of as closeness on the collection of hypotheses $H_1^q, H_2^q, ..., H_L^q$. Given $q \in \mathcal{Q}$, we can use the construction from Section 2.5.2 in order to build the test $\mathcal{T}_{\phi_*}^K$ deciding on the hypotheses H_{ℓ}^q up to closeness \mathcal{C} , the \mathcal{C} -risk of the test being the smallest allowed by the construction. Note that this \mathcal{C} -risk depends on q; the "Measurement Design" (MD for short)

 $^{^{27}}$ this closeness makes sense when the goal of the poll is to predict the winner; less ambitious goal, like to decide whether the winner will or will not belong to a particular set of candidates, would require weaker closeness.

problem we are about to consider is to select $q \in \mathcal{Q}$ which minimizes the C-risk of the associated test $\mathcal{T}_{\phi_{\omega}}^{K}$.

2.7.3 Formulating the MD problem

By Proposition 2.34, the C-risk of the test $\mathcal{T}_{\phi_*}^K$ is upper-bounded by the spectral norm of the symmetric entrywise nonnegative $L \times L$ matrix

$$E^{(K)}(q) = \left[\epsilon_{\ell\ell'}(q)\right]_{\ell,\ell'},$$

and this is what we intend to minimize in our MD problem. In the above formula, $\epsilon_{\ell\ell'}(q) = \epsilon_{\ell'\ell}(q)$ are zeros when $(\ell, \ell') \in \mathcal{C}$; when $(\ell, \ell') \notin \mathcal{C}$ and $1 \leq \ell < \ell' \leq L$, the quantities $\epsilon_{\ell\ell'}(q) = \epsilon_{\ell'\ell}(q)$ are defined depending on what is the simple o.s. \mathcal{O} . Specifically,

• In the case of *Gaussian* observation scheme (see Section 2.4.5.1), restriction (2.125) does not restrict the dependence A_q on q at all (modulo the default restriction that A_q is a continuous in $q \in \mathcal{Q} \ d \times n$ matrix), and

$$\epsilon_{\ell\ell'}(q) = \exp\{KOpt_{\ell\ell'}(q)\}\$$

where

$$Opt_{\ell\ell'}(q) = \max_{x \in X_{\ell'}, y \in X_{\ell'}} - [A_q(x-y)]^T \Theta^{-1}[A_q(x-y)]$$
(G_q)

and Θ is the common covariance matrix of the Gaussian densities forming the family $\{p_{\mu} : \mu \in \mathcal{M}\};\$

• In the case of Poisson o.s. (see Section 2.4.5.2), restriction (2.125) requires from $A_q x$ to be positive vector whenever $q \in Q$ and $x \in \mathcal{X}$, and

$$\epsilon_{\ell\ell'}(q) = \exp\{KOpt_{\ell\ell'}(q)\},\$$

where

$$Opt_{\ell\ell'}(q) = \max_{x \in X_{\ell}, y \in X_{\ell'}} \left[\sum_{i} \sqrt{[A_q x]_i [A_q y]_i} - \frac{1}{2} \sum_{i} [A_q x]_i - \frac{1}{2} \sum_{i} [A_q y]_i \right]; \ (P_q)$$

• In the case of Discrete o.s. (see Section 2.4.5.3), restriction (2.125) requires from $A_q x$ to be a positive probabilistic vector whenever $q \in Q$ and $x \in \mathcal{X}$, and

$$\epsilon_{\ell\ell'}(q) = \left[\operatorname{Opt}_{\ell\ell'}(q)\right]^K,$$

where

$$\operatorname{Opt}_{\ell\ell'}(q) = \max_{x \in X_{\ell}, y \in X_{\ell'}} \sum_{i} \sqrt{[A_q x]_i [A_q y]_i}.$$
 (D_q)

The summary of above observations is as follows. The norm $||E^{(K)}||_{2,2}$ – the quantity we are interested to minimize in $q \in Q$ – as a function of $q \in Q$ is of the form

$$\Psi(q) = \psi(\underbrace{\{\operatorname{Opt}_{\ell\ell'}(q) : (\ell, \ell') \notin \mathcal{C}\}}_{\overline{\operatorname{Opt}}(q)})$$
(2.126)

where the outer function ψ is real-valued convex and nondecreasing in every one

of its arguments explicitly given function on \mathbb{R}^N (N is the cardinality of the set of pairs $(\ell, \ell'), 1 \leq \ell, \ell' \leq L$, with $(\ell, \ell') \notin C$). Indeed, denoting by $\Gamma(S)$ the spectral norm of $d \times d$ matrix S, note that Γ is convex function of S, and this function is nondecreasing in every one of the entries of S, provided that S is restricted to be entrywise nonnegative²⁸. $\psi(\cdot)$ is obtained from $\Gamma(S)$ by substitution, instead of entries $S_{\ell\ell'}$ of S, everywhere convex, nonnegative and nondecreasing functions of new variables $\vec{z} = \{z_{\ell\ell'} : (\ell, \ell') \notin C, 1 \leq \ell, \ell' \leq L\}$, specifically

- when $(\ell, \ell') \in \mathcal{C}$, we set $S_{\ell\ell'}$ to zero;
- when $(\ell, \ell') \notin C$, we set $S_{\ell\ell'} = \exp\{Kz_{\ell\ell'}\}$ in the case of Gaussian and Poisson o.s.'s, and set $S_{\ell\ell'} = \max[0, z_{\ell\ell'}]^K$, in the case of Discrete o.s.

As a result, we indeed get a convex and nondecreasing in every one of its arguments function ψ of $\vec{z} \in \mathbf{R}^N$.

Now, the Measurement Design problem we want to solve reads

$$Opt = \min_{q \in \mathcal{Q}} \psi(\overline{Opt}(q)); \qquad (2.127)$$

As we remember, the entries in the inner function $\overline{\text{Opt}}(q)$ are optimal values of solvable *convex* optimization problems and as such are efficiently computable. When these entries are also *convex* functions of $q \in \mathcal{Q}$, the objective in (2.127), due to the already established convexity and monotonicity properties of ψ , is a convex function of q, meaning that (2.127) is a convex and thus efficiently solvable problem. On the other hand, when some of the entries in $\overline{Opt}(q)$ are nonconvex in q, we hardly could expect (2.127) to be an easy-to-solve problem. Unfortunately, convexity of the entries in $\overline{\operatorname{Opt}}(q)$ in q turns out to be a "rare commodity." For example, we can verify by inspection that the objectives in (G_q) , (P_q) , (D_q) as a functions of A_q (not of q!) are *concave* rather than convex, so that the optimal values in the problems, as a functions of q, are maxima, over the parameters, of parametric families of concave functions of A_q (the parameter in these parametric families are the optimization variables in $(G_q) - (D_q)$ and as such as a functions of A_q hardly are convex. And indeed, as a matter of fact, the MD problem usually is nonconvex and difficult to solve. We intend to consider "Simple case" where this difficulty does not arise, specifically, the case where the objectives of the optimization problems specifying $\operatorname{Opt}_{\ell\ell'}(q)$ are affine in q; in this case, $\operatorname{Opt}_{\ell\ell'}(q)$ as a function of q is the maximum, over the parameters (optimization variables in the corresponding problems), of parametric families of affine functions of q and as such is convex.

Our current goal is to understand what our sufficient condition for tractability of the MD problem – affinity in q of the objectives in the respective problems $(G_q), (P_q), (D_q)$ – actually means, and to show that this, by itself quite restrictive, assumption indeed takes place in some important applications.

 $^{^{28}}$ monotonicity follows from the fact that for an entrywise nonnegative S, we have

 $^{\|}S\|_{2,2} = \max_{x,y} \{x^T Sy : \|x\|_2 \le 1, \|y\|_2 \le 1\} = \max_{x,y} \{x^T Sy : \|x\|_2 \le 1, \|y\|_2 \le 1, x \ge 0, y \ge 0\}$

2.7.3.1 Simple case, Discrete o.s.

Looking at the optimization problem (D_q) , we see that the simplest way to ensure that its objective is affine in q is to assume that

$$A_q = \text{Diag}\{Bq\}A,\tag{2.128}$$

where A is some fixed $d \times n$ matrix, and B is some fixed $d \times (\dim q)$ matrix such that Bq is positive whenever $q \in Q$. On the top of this, we should ensure that when $q \in Q$ and $x \in \mathcal{X}$, $A_q x$ is a positive probabilistic vector; this amounts to some restrictions linking Q, \mathcal{X}, A , and B.

Illustration. An instructive example of the Simple case of Measurement Design in Discrete o.s. is the "Opinion Poll" problem with a priori information presented in Section 2.7.1: the voting population is split into I groups, with *i*-th group constituting fraction θ_i of the entire population. In *i*-th group, the distribution of voters' preferences is represented by unknown L-dimensional probabilistic vector $x^i = [x_1^i; ...; x_L^i]$ (L is the number of candidates, x_ℓ^i is the fraction of voters in *i*-th group intending to vote for ℓ -th candidate), known to belong to a given convex compact subset Π_i of the "positive part" $\Delta_L^o = \{x \in \mathbf{R}^L : x > 0, \sum_\ell x_\ell = 1\}$ of the L-dimensional probabilistic simplex. We are given threshold $\delta > 0$ and want to decide on L hypotheses H_1, \ldots, H_L , with H_ℓ stating that the population-wide vector $y = \sum_{i=1}^{I} \theta_i x^i$ of voters' preferences belongs to the closed convex set

$$Y_{\ell} = \{ y = \sum_{i=1}^{I} \theta_{i} x^{i} : x^{i} \in \Pi_{i}, 1 \le i \le I, y_{\ell} \ge y_{\ell'} + \delta, \forall (\ell' \ne \ell) \};$$

note that Y_{ℓ} is the image, under the linear mapping

$$[x^1;...;x^I]\mapsto y(x)=\sum_i\theta_ix^i$$

of the compact convex set

$$X_{\ell} = \{x = [x^1; ...; x^I] : x^i \in \Pi_i, \ 1 \le i \le I, y_{\ell}(x) \ge y_{\ell'}(x) + \delta, \ \forall (\ell' \ne \ell) \}$$

which is a subset of the convex compact set

$$\mathcal{X} = \{ x = [x^1; ...; x^I] : x^i \in \Pi_i, \, 1 \le i \le I \}.$$

k-th poll interview is organized as follows:

We draw at random a group among the I groups of voters, with probability q_i to draw *i*-th group, and then draw at random, from the uniform distribution on the group, the respondent to be interviewed. The outcome of the interview – our observation ω_k – is the pair (i, ℓ) , where *i* is the group to which the respondent belongs, and ℓ is the candidate preferred by the respondent.

This results in a sensing matrix A_q , see (2.123), which is in the form of (2.128), specifically,

$$A_q = \text{Diag}\{q_1 I_L, q_2 I_L, \dots, q_I I_L\} \qquad [q \in \mathbf{\Delta}_I]$$

the outcome of k-th interview is drawn at random from the discrete probability distribution $A_q x$, where $x \in \mathcal{X}$ is the "signal" summarizing voters' preferences in the groups.

Given total number of observations K, our goal is to decide with a given risk ϵ on our L hypotheses; whether this goal is or is not achievable, it depends on K and on A_q . What we want, is to find q for which the above goal is achievable with as small K as possible; in the case in question, this reduces to solving, for various trial values of K, problem (2.127), which under the circumstances is an explicit *convex* optimization problem.

To get an impression of the potential of Measurement Design, we present a sample of numerical results. In all reported experiments, we used $\delta = 0.05$ and $\epsilon = 0.01$. The sets Π_i , $1 \leq i \leq I$, were generated as follows: we pick at random a probabilistic vector \bar{p}^i of dimension L, and Π_i was the intersection of the box $\{p : \bar{p}_{\ell} - u_i \leq p_{\ell} \leq \bar{p}_{\ell} + u_i\}$ centered at \bar{p} with the probabilistic simplex Δ_L , where u_i , i = 1, ..., I, are prescribed "uncertainty levels;" note that uncertainty level $u_i \geq 1$ is the same as absence of any a priori information on the preferences of voters from *i*-th group.

The results of our numerical experiments are as follows:

	I	Uncertainty levels u	Group sizes θ	K _{ini}	q_{opt}	K _{opt}
2	2	[0.03; 1.00]	[0.500; 0.500]	1212	[0.437; 0.563]	1194
2	2	[0.02; 1.00]	[0.500; 0.500]	2699	[0.000; 1.000]	1948
3	3	[0.02; 0.03; 1.00]	[0.333; 0.333; 0.333]	3177	[0.000; 0.455; 0.545]	2726
5	4	[0.02; 0.02; 0.03; 1.00]	[0.250; 0.250; 0.250; 0.250]	2556	[0.000; 0.131; 0.322; 0.547]	2086
5	4	[1.00; 1.00; 1.00; 1.00]	[0.250; 0.250; 0.250; 0.250]	4788	[0.250; 0.250; 0.250; 0.250]	4788

Effect of measurement design. $K_{\rm ini}$ and $K_{\rm opt}$ are the poll sizes required for 0.99-reliable prediction of the winner when $q = \theta$ and $q = q_{\rm opt}$, respectively.

We see that measurement design allows to reduce (for some data – quite significantly) the volume of observations as compared to the straightforward selecting the respondents uniformly across the entire population. To compare our current model and results with those from Section 2.6.1, note that now we have more a priori information on the true distribution of voting preferences due to some a priori knowledge of preferences within groups, which allows to reduce the poll sizes with both straightforward and optimal measurement design²⁹. The differences between $K_{\rm ini}$ and $K_{\rm opt}$ is fully due to measurement design.

Comparative drug study. A related to DOP and perhaps more interesting Simple case of the Measurement Design in Discrete o.s. is as follows. Let us speak about L competing drugs rather than L competing candidates running for an office, and population of patients the drugs are aimed to help rather than population of voters. For the sake of simplicity, assume that when a particular drug is administered to a particular patient, the outcome is binary: (positive) "effect" or "no effect" (what follows can be easily extended to the case of non-binary categorial outcomes, like "strong positive effect," "weak positive effect," "negative effect," and alike). Our goal is to organize a clinical study in order to make inferences on comparative drug efficiency, measured by the percentage of patients on which a particular drug has effect. The difference with organizing opinion poll is that now we cannot just ask a respondent what are his/her preferences; we are supposed to administer to a participant of the study a single drug on our choice and to look at the result.

 $^{^{29}}$ To illustrate this point, look at the last two lines in the table: utilizing a priori information allows to reduce the poll size from 4788 to 2556 even with the straightforward measurement design.

As in the DOP problem, we assume that the population of patients is split into I groups, with *i*-th group comprising fraction θ_i of the entire population.

We model the situation as follows. We associate with a patient Boolean vector of dimension 2L, with ℓ -th entry in the vector equal to 1 or 0 depending on whether drug $\# \ell$ has effect on the patient, and the $(L + \ell)$ -th entry complementing the ℓ -th one to 1 (that is, if ℓ -th entry is χ , then $(L + \ell)$ -th entry is $1 - \chi$). Let x^i be the average of these vectors over patients from group i. We define "signal" x underlying our measurements as the vector $[x^1; ...; x^I]$ and assume that our a priori information allows to localize x in a closed convex subset \mathcal{X} of the set

$$\mathcal{Y} = \{x = [x^1; ...; x^I] : x^i \ge 0, x^i_{\ell} + x^i_{L+\ell} = 1, 1 \le i \le I, 1 \le \ell \le L\}$$

to which all our signals belong by construction. Note that the vector

$$y = Bx = \sum_{i} \theta_i x^i$$

can be treated as "population-wise distribution of drug effects:" y_{ℓ} , $\ell \leq L$, is the fraction, in the entire population of patients, of those patients on whom drug ℓ has effect, and $y_{L+\ell} = 1 - y_{\ell}$. As a result, typical hypotheses related to comparison of the drugs, like "drug ℓ has effect on a larger, at least by margin δ , percentage of patients than drug ℓ' ," become convex hypotheses on the signal x. In order to test hypotheses of this type, we can use two-stage procedure for observing drug effects, namely, as follows.

To get a particular observation, we select at random, with probability $q_{i\ell}$, pair (i,ℓ) from the set $\{(i,\ell): 1 \leq i \leq I, 1 \leq \ell \leq L\}$, select a patient from group i according to the uniform distribution on the group, administer the patient drug ℓ and check whether the drug has effect on the patient. Thus, a single observation is a triple (i,ℓ,χ) , where $\chi = 0$ when the administered drug has no effect on the patient, and $\chi = 1$ otherwise. The probability to get observation $(i,\ell,1)$ is $q_{i\ell}x_{\ell}^i$, and the probability to get observation $(i,\ell,0)$ is $q_{i\ell}x_{L+\ell}^i$. Thus, we arrive at the Discrete o.s. where the distribution μ of observations is of the form $\mu = A_q x$, with the rows in A_q indexed by triples $\omega = (i,\ell,\chi) \in \Omega := \{1,2,...,I\} \times \{1,2,...,L\} \times \{0,1\}$ and given by

$$(A_q[x^1;...;x^I])_{i,\ell,\chi} = \begin{cases} q_{i\ell}x^i_{\ell}, & \chi = 1\\ q_{i\ell}x^i_{L+\ell}, & \chi = 0 \end{cases}$$

Specifying the set \mathcal{Q} of allowed measurement designs q as a closed convex subset of the set of all non-vanishing discrete probability distributions on the set $\{1, 2, ..., I\} \times \{1, 2, ..., L\}$, we find ourselves in the Simple case, as defined by (2.128), of Discrete o.s., and $A_q x$ is a probabilistic vector whenever $q \in \mathcal{Q}$ and $x \in \mathcal{Y}$.

2.7.3.2 Simple case, Poisson o.s.

Looking at the optimization problem (P_q) , we see that the simplest way to ensure that its objective is, same as in the case of Discrete o.s., to assume that

$$A_q = \text{Diag}\{Bq\}A,$$

where A is some fixed $d \times n$ matrix, and B is some fixed $d \times (\dim q)$ matrix such that Bq is positive whenever $q \in Q$. On the top of this, we should ensure that when $q \in Q$ and $x \in \mathcal{X}$, $A_q x$ is a positive vector; this amounts to some restrictions
linking $\mathcal{Q}, \mathcal{X}, A$, and B.

Application Example: PET with time control. Positron Emission Tomography was already mentioned, as an example of Poisson o.s., in Section 2.4.3.2. As explained in the latter Section, in PET we observe a random vector $\omega \in \mathbf{R}^d$ with independent entries $[\omega]_i \sim \text{Poisson}(\mu_i), 1 \leq i \leq d$, where the vector of parameters $\mu = [\mu_1; ..., \mu_d]$ of the Poisson distributions is the linear image $\mu = A\lambda$ of unknown "signal" (tracer's density in patient's body) λ belonging to some known subset Λ of \mathbf{R}^D_+ , with entrywise nonnegative matrix A; our goal is to make inferences about λ . Now, in actual PET scan, patient's position w.r.t. the scanner is not the same during the entire study; the position is kept fixed within *i*-th time period, $1 \leq i \leq I$, and changes from period to period in order to expose to the scanner the entire "area of interest"



For example, with the scanner shown on the picture, during PET study the imaging table with the patient will be shifted several times along the axis of the scanning ring. As a result, observed vector ω can be split into blocks ω^i , i = 1, ..., I, of data acquired during *i*-th period, $1 \leq i \leq I$; on the closest inspection, the corresponding block μ^i in μ is

$$\mu^i = q_i A_i \lambda,$$

where A_i is a known in advance entrywise nonnegative matrix, and q_i is the duration of *i*-th period. In principle, q_i could be treated as nonnegative design variables subject to the "budget constraint" $\sum_{i=1}^{I} q_i = T$, where *T* is the total duration of the study³⁰, and perhaps some other convex constraints, say, positive lower bounds on q_i . It is immediately seen that the outlined situation is exactly as is required in the Simple case of Poisson o.s.

2.7.3.3 Simple case, Gaussian o.s.

Looking at the optimization problem (G_q) , we see that the simplest way to ensure that its objective is affine in q is to assume that the covariance matrix Θ is diagonal, and

$$A_q = \operatorname{Diag}\{\sqrt{q_1}, \dots, \sqrt{q_d}\}A\tag{2.129}$$

where A is a fixed $d \times n$ matrix, and q runs through a convex compact subset of \mathbf{R}^{d}_{+} .

It turns out that there are situations where assumption (2.129) makes perfect

 $^{^{30}}T$ cannot be too large; aside of other considerations, the tracer disintegrates, and its density can be considered as nearly constant only on a properly restricted time horizon.

sense. Let us start with preamble. In Gaussian o.s.

$$\omega = Ax + \xi \left[A \in \mathbf{R}^{d \times n}, \xi \sim \mathcal{N}(0, \Sigma), \Sigma = \text{Diag}\{\sigma_1^2, ..., \sigma_d^2\}\right]$$
(2.130)

the "physics" behind the observations in many cases is as follows. There are d sensors (receivers), *i*-th registering continuous time analogous input depending linearly on the underlying observations signal x; on the time horizon on which the measurements are taken, this input is constant in time and is registered by *i*-th sensor on time interval Δ_i . The deterministic component of the measurement registered by sensor *i* is the integral of the corresponding input taken over Δ_i , and the stochastic component of the measurement is obtained by integrating over the same interval white Gaussian noise. As far as this noise is concerned, the only thing which matters is that when the white noise affecting *i*-th sensor is integrated over a time interval Δ , the result is random Gaussian variable with zero mean and variance $\sigma_i^2 |\Delta|$ ($|\Delta|$ is the length of Δ), and the random variables obtained by integrating white noise over non-overlapping segments are independent. Besides this, we assume that the noisy components of measurements are independent across the sensors.

Now, there could be two basic versions of the just outlined situation, both leading to the same observation model (2.130). In the first, "parallel," version, all d sensors work in parallel on the same time horizon of duration 1. In the second, "sequential," version, the sensors are activated and scanned one by one, each working unit time; thus, here the full time horizon is d, and the sensors are registering their respective inputs on consecutive time intervals of duration 1 each. In this second "physical" version of Gaussian o.s., we can, in principle, allow for sensors to register their inputs on consecutive time segments of varying durations $q_1 \ge 0, q_2 \ge 0, ..., q_d \ge 0$, with the additional to nonnegativity restriction that our total time budget is respected: $\sum_{i} q_i = d$ (and perhaps with some other convex constraints on q_i). Let us look what is the observation scheme we end up with. Assuming that (2.130) represents correctly our observations in the reference case where all $|\Delta_i|$ are equal to 1, the deterministic component of the measurement registered by sensor i in time interval of duration q_i will be $q_i \sum_j a_{ij} x_j$, and the standard deviation of the noisy component will be $\sigma_i \sqrt{q_i}$, so that the measurements become

$$z_i = \sigma_i \sqrt{q_i} \zeta_i + q_i \sum_j a_{ij} x_j, \ i = 1, ..., d,$$

with independent of each other standard (zero mean, unit variance) Gaussian noises ζ_i . Now, since we know q_i , we can scale the latter observations by making the standard deviation of the noisy component the same σ_i as in the reference case; specifically, we lose nothing when assuming that our observations are

$$\omega_i = z_i / \sqrt{q_i} = \underbrace{\sigma_i \zeta_i}_{\xi_i} + \sqrt{q_i} \sum_j a_{ij} x_j,$$

or, equivalently,

$$\omega = \xi + \underbrace{\text{Diag}\{\sqrt{q_1}, \dots, \sqrt{q_d}\}A}_{A_q} x, \ \xi \sim \mathcal{N}(0, \text{Diag}\{\sigma_1^2, \dots, \sigma_d^2\}) \qquad [A = [a_{ij}]]$$

where q is allowed to run through a convex compact subset \mathcal{Q} of the simplex $\{q \in \mathbf{R}^d_+ : \sum_i q_i = d\}$. Thus, if the "physical nature" of a Gaussian o.s. is sequential, then, making, as is natural under the circumstances, the activity times of the sensors our design variables, we arrive at (2.129), and, as a result, end up with easy-to-solve Measurements Design problem.

2.8 AFFINE DETECTORS BEYOND SIMPLE OBSERVATION SCHEMES

On a closer inspection, the "common denominator" of our basic simple o.s.'s – Gaussian, Poisson and Discrete ones, is that in all these cases the minimal risk detector for a pair of convex hypotheses is *affine*. At the first glance, this indeed is so for the Gaussian and the Poisson o's"s, where \mathcal{F} is comprised of affine functions on the corresponding observation space Ω (\mathbf{R}^d for Gaussian o.s., and $\mathbf{Z}^d_+ \subset \mathbf{R}^d$ for Poisson o.s.), but is *not* so for the Discrete o.s. – in the latter case, $\Omega = \{1, ..., d\}$, and \mathcal{F} is comprised of all functions on Ω , while "affine functions on $\Omega = \{1, ..., d\}$ " merely make no sense. Note, however, that we can encode (and from now on indeed encode) the points i = 1, ..., d of *d*-element set by basic orths $e_i = [0; ...; 0; 1; 0; ...; 0] \in \mathbf{R}^d$ in \mathbf{R}^d , thus making our observation space $\Omega = \{1, ..., d\}$ a subset of \mathbf{R}^d . With this encoding, *every* real valued function on $\{1, ..., d\}$ becomes restriction on Ω of an affine function. Note that when passing from our basic simple o.s.'s to their direct products, the minimum risk detectors for pairs of convex hypotheses remain affine.

Now, good in our context news about simple o.s.'s state that

- A) the best with the smallest possible risk *affine* detector, same as its risk, can be efficiently computed;
- B) the smallest risk *affine* detector from A) is the best, in terms of risk, detector available under the circumstances, so that the associated test is near-optimal.

Note that as far as practical applications of the detector-based hypothesis testing are concerned, one "can survive" without B) (near-optimality of the constructed detectors), while A) is a must.

In this Section we focus on families of probability distributions obeying A). This class turns out to be incomparably larger than what was defined as simple o.s.'s in Section 2.4; in particular, it includes nonparametric families of distributions. Staying within this much broader class, we still are able to construct in a computationally efficient way the best affine detectors for a pair of "convex", in certain precise sense, hypotheses, along with valid upper bounds on the risks of the detectors. What we, in general, cannot claim anymore, is that the tests associated with the above detectors are near-optimal. This being said, we believe that investigating possibilities for building tests and quantifying their performance in a computationally friendly manner is of value even when we cannot provably guarantee near-optimality of these tests. The results to follow originate from [89, 88].

2.8.1 Situation

In what follows, we fix observation space $\Omega = \mathbf{R}^d$, and let \mathcal{P}_j , $1 \leq j \leq J$, be given families of probability distributions on Ω . Put broadly, our goal still is, given a random observation $\omega \sim P$, where $P \in \bigcup_{j \leq J} \mathcal{P}_j$, to decide upon the hypotheses $H_j : P \in \mathcal{P}_j$, j = 1, ..., J. We intend to address this goal in the case when the families \mathcal{P}_j are simple – they are comprised of distributions for which moment-generating functions admit an explicit upper bound.

2.8.1.1 Preliminaries: Regular data and associated families of distributions

Regular data is defined as a triple $\mathcal{H}, \mathcal{M}, \Phi(\cdot, \cdot)$, where

- \mathcal{H} is a nonempty closed convex set in $\Omega = \mathbf{R}^d$ symmetric w.r.t. the origin,
- \mathcal{M} is a closed convex set in some \mathbf{R}^n ,
- $\Phi(h;\mu): \mathcal{H} \times \mathcal{M} \to \mathbf{R}$ is a continuous function convex in $h \in \mathcal{H}$ and concave in $\mu \in \mathcal{M}$.

Regular data $\mathcal{H}, \mathcal{M}, \Phi(\cdot, \cdot)$ define two families of probability distributions on Ω :

• the family of *regular* distributions

$$\mathcal{R} = \mathcal{R}[\mathcal{H}, \mathcal{M}, \Phi]$$

comprised of all probability distributions P on Ω such that

$$\forall h \in \mathcal{H} \exists \mu \in \mathcal{M} : \ln\left(\int_{\Omega} \exp\{h^T \omega\} P(d\omega)\right) \le \Phi(h;\mu).$$
(2.131)

• the family of *simple* distributions

$$S = S[\mathcal{H}, \mathcal{M}, \Phi]$$

comprised of probability distributions P on Ω such that

$$\exists \mu \in \mathcal{M} : \forall h \in \mathcal{H} : \ln\left(\int_{\Omega} \exp\{h^T \omega\} P(d\omega)\right) \le \Phi(h;\mu).$$
(2.132)

Recall that beginning with Section 2.3, the starting point in all our constructions is a "plausibly good" detector-based test which, given two families \mathcal{P}_1 and \mathcal{P}_2 of distributions with common observation space, and repeated observations $\omega_1, ..., \omega_t$ drawn from a distribution $P \in \mathcal{P}_1 \cup \mathcal{P}_2$, decides whether $P \in \mathcal{P}_1$ or $P \in \mathcal{P}_2$. Our interest in the families of regular/simple distributions stems from the fact that when the families \mathcal{P}_1 and \mathcal{P}_2 are of this type, building such a test reduces to solving a convex-concave saddle point problem and thus can be carried out in a computationally efficient manner. We postpone the related construction and analysis to Section 2.8.2, and continue with presenting some basic examples of families of simple and regular distributions along with a simple "calculus" of these families.

2.8.1.2 Basic examples of simple families of probability distributions

2.8.1.2.A. Sub-Gaussian distributions: Let $\mathcal{H} = \Omega = \mathbf{R}^d$, \mathcal{M} be a closed convex subset of the set $\mathcal{G}_d = \{\mu = (\theta, \Theta) : \theta \in \mathbf{R}^d, \Theta \in \mathbf{S}^d_+\}$, where \mathbf{S}^d_+ is cone of positive semidefinite matrices in the space \mathbf{S}^d of symmetric $d \times d$ matrices, and let

$$\Phi(h;\theta,\Theta) = \theta^T h + \frac{1}{2}h^T \Theta h.$$

In this case, $S[\mathcal{H}, \mathcal{M}, \Phi]$ contains all *sub-Gaussian* distributions P on \mathbf{R}^d with sub-Gaussianity parameters from \mathcal{M} .

Recall that a distributions P on $\Omega = \mathbf{R}^d$ is called sub-Gaussian with sub-Gaussianity parameters $\theta \in \mathbf{R}^d$ and $\Theta \in \mathbf{S}^d_+$, if

$$\mathbf{E}_{\omega \sim P}\{\exp\{h^T\omega\}\} \le \exp\{\theta^T h + \frac{1}{2}h^T\Theta h\} \ \forall h \in \mathbf{R}^d.$$
(2.133)

Whenever this is the case, θ is the expected value of P. We shall use the notation $\xi \sim SG(\theta, \Theta)$ as a shortcut for the sentence "random vector ξ is sub-Gaussian with parameters θ , Θ ." It is immediately seen that when $\xi \sim \mathcal{N}(\theta, \Theta)$, we have also $\xi \sim SG(\theta, \Theta)$, and (2.133) in this case is an identity rather than inequality.

In particular, $S[H, M, \Phi]$ contains all Gaussian distributions $\mathcal{N}(\theta, \Theta)$ with $(\theta, \Theta) \in \mathcal{M}$.

2.8.1.2.B. Poisson distributions: Let $\mathcal{H} = \Omega = \mathbf{R}^d$, let \mathcal{M} be a closed convex subset of *d*-dimensional nonnegative orthant \mathbf{R}^d_+ , and let

$$\Phi(h = [h_1; ...; h_d]; \mu = [\mu_1; ...; \mu_d]) = \sum_{i=1}^d \mu_i [\exp\{h_i\} - 1] : \mathcal{H} \times \mathcal{M} \to \mathbf{R}.$$

The family $S[\mathcal{H}, \mathcal{M}, \Phi]$ contains all product-type Poisson distributions $Poisson[\mu]$ with vectors μ of parameters belonging to \mathcal{M} ; here $Poisson[\mu]$ is the distribution of random *d*-dimensional vector with independent of each other entries, *i*-th entry being Poisson random variable with parameter μ_i .

2.8.1.2.C. Discrete distributions. Consider a discrete random variable taking values in *d*-element set $\{1, 2, ..., d\}$, and let us think of such a variable as of random variable taking values $e_i \in \mathbf{R}^d$, i = 1, ..., d, where $e_i = [0; ...; 0; 1; 0; ...; 0]$ (1 in position *i*) are standard basic orths in \mathbf{R}^d . Probability distribution of such a variable can be identified with a point $\mu = [\mu_1; ...; \mu_d]$ from the *d*-dimensional probabilistic simplex

$$\mathbf{\Delta}_d = \{ \nu \in \mathbf{R}^d_+ : \sum_{i=1}^d \nu_i = 1 \},\$$

where μ_i is the probability for the variable to take value e_i . With these identifica-

tions, setting $\mathcal{H} = \mathbf{R}^d$, specifying \mathcal{M} as a closed convex subset of $\mathbf{\Delta}_d$ and setting

$$\Phi(h = [h_1; ...; h_d]; \mu = [\mu_1; ...; \mu_d]) = \ln\left(\sum_{i=1}^d \mu_i \exp\{h_i\}\right),$$

the family $S[\mathcal{H}, \mathcal{M}, \Phi]$ contains distributions of all discrete random variables taking values in $\{1, ..., d\}$ with probabilities $\mu_1, ..., \mu_d$ comprising a vector from \mathcal{M} .

2.8.1.2.D. Distributions with bounded support. Consider the family $\mathcal{P}[X]$ of probability distributions supported on a closed and bounded convex set $X \subset \Omega = \mathbf{R}^d$, and let

$$\phi_X(h) = \max_{x \in X} h^T x$$

be the support function of X. We have the following result (to be refined in Section 2.8.1.3):

Proposition 2.40. For every $P \in \mathcal{P}[X]$ it holds

$$\forall h \in \mathbf{R}^d : \ln\left(\int_{\mathbf{R}^d} \exp\{h^T \omega\} P(d\omega)\right) \le h^T e[P] + \frac{1}{8} \left[\phi_X(h) + \phi_X(-h)\right]^2, \quad (2.134)$$

where $e[P] = \int_{\mathbf{R}^d} \omega P(d\omega)$ is the expectation of P, and the right hand side function in (2.134) is convex. As a result, setting

$$\mathcal{H} = \mathbf{R}^{d}, \ \mathcal{M} = X, \ \Phi(h;\mu) = h^{T}\mu + \frac{1}{8} \left[\phi_{X}(h) + \phi_{X}(-h)\right]^{2},$$

we get regular data such that $\mathcal{P}[X] \subset \mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$.

For proof, see Section 2.11.3

2.8.1.3 Calculus of regular and simple families of probability distributions

Families of regular and simple distributions admit "fully algorithmic" calculus, with the main calculus rules as follows.

2.8.1.3.A. Direct summation. For $1 \leq \ell \leq L$, let regular data $\mathcal{H}_{\ell} \subset \Omega_{\ell} = \mathbf{R}^{d_{\ell}}$, $\mathcal{M}_{\ell} \subset \mathbf{R}^{n_{\ell}}, \Phi_{\ell}(h_{\ell}; \mu_{\ell}) : \mathcal{H}_{\ell} \times \mathcal{M}_{\ell} \to \mathbf{R}$ be given. Let us set

$$\Omega = \Omega_1 \times \ldots \times \Omega_L = \mathbf{R}^d, \ d = d_1 + \ldots + d_L,
\mathcal{H} = \mathcal{H}_1 \times \ldots \times \mathcal{H}_L = \{h = [h^1; \ldots; h^L] : h^\ell \in \mathcal{H}_\ell, \ell \leq L\},
\mathcal{M} = \mathcal{M}_1 \times \ldots \times \mathcal{M}_L = \{\mu = [\mu^1; \ldots; \mu^L] : \mu^\ell \in \mathcal{M}^\ell, \ell \leq L\} \subset \mathbf{R}^n,
n = n_1 + \ldots + n_L,
\Phi(h = [h^1; \ldots; h^L]; \mu = [\mu^1; \ldots; \mu^L]) = \sum_{\ell=1}^L \Phi_\ell(h^\ell; \mu^\ell) : \mathcal{H} \times \mathcal{M} \to \mathbf{R}.$$

Then \mathcal{H} is a symmetric w.r.t. the origin closed convex set in $\Omega = \mathbf{R}^d$, \mathcal{M} is a nonempty closed convex set in \mathbf{R}^n , $\Phi : \mathcal{H} \times \mathcal{M} \to \mathbf{R}$ is a continuous convexconcave function, and clearly

• the family $\mathcal{R}[\mathcal{H}, \mathcal{M}, \Phi]$ contains all product-type distributions $P = P_1 \times ... \times P_L$ on $\Omega = \Omega_1 \times ... \times \Omega_L$ with $P_\ell \in \mathcal{R}[\mathcal{H}_\ell, \mathcal{M}_\ell, \Phi_\ell], 1 \le \ell \le L$;

• the family $S[\mathcal{H}, \mathcal{M}, \Phi]$ contains all product-type distributions $P = P_1 \times ... \times P_L$ on $\Omega = \Omega_1 \times ... \times \Omega_L$ with $P_\ell \in S[\mathcal{H}_\ell, \mathcal{M}_\ell, \Phi_\ell], 1 \le \ell \le L$.

2.8.1.3.B. Mixing. For $1 \leq \ell \leq L$, let regular data $\mathcal{H}_{\ell} \subset \Omega = \mathbf{R}^{d}$, $\mathcal{M}_{\ell} \subset \mathbf{R}^{n_{\ell}}$, $\Phi_{\ell}(h_{\ell};\mu_{\ell}): \mathcal{H}_{\ell} \times \mathcal{M}_{\ell} \to \mathbf{R}$ be given, with compact \mathcal{M}_{ℓ} . Let also $\nu = [\nu_{1};...;\nu_{L}]$ be a probabilistic vector. For a tuple $P^{L} = \{P_{\ell} \in \mathcal{R}[\mathcal{H}_{\ell}, \mathcal{M}_{\ell}, \Phi_{\ell}]\}_{\ell=1}^{L}$, let $\Pi[P^{L}, \nu]$ be the ν -mixture of distributions $P_{1}, ..., P_{L}$ defined as the distribution of random vector $\omega \sim \Omega$ generated as follows: we draw at random, from probability distribution ν on $\{1, ..., L\}$, index ℓ , and then draw ω at random from the distribution P_{ℓ} . Finally, let \mathcal{P} be the set of all probability distributions on Ω which can be obtained as $\Pi[P^{L}, \nu]$ from the outlined tuples P^{L} and vectors ν running through the probabilistic simplex $\Delta_{L} = \{\mu \in \mathbf{R}^{L} : \nu \geq 0, \sum_{\ell} \nu_{\ell} = 1\}$.

Let us set

$$\begin{aligned}
\mathcal{H} &= \bigcap_{\ell=1}^{L} \mathcal{H}_{\ell}, \\
\Psi_{\ell}(h) &= \max_{\mu_{\ell} \in \mathcal{M}_{\ell}} \Phi_{\ell}(h; \mu_{\ell}) : \mathcal{H}_{\ell} \to \mathbf{R}, \\
\Phi(h; \nu) &= \ln\left(\sum_{\ell=1}^{L} \nu_{\ell} \exp\{\Psi_{\ell}(h)\}\right) : \mathcal{H} \times \mathbf{\Delta}_{L} \to \mathbf{R}.
\end{aligned}$$
(2.135)

Then $\mathcal{H}, \mathbf{\Delta}_L, \Phi$ clearly is a regular data (recall that all \mathcal{M}_ℓ are compact sets), and for every $\nu \in \mathbf{\Delta}_L$ and tuple P^L of the above type one has

$$P = \Pi[P^L, \nu] \Rightarrow \ln\left(\int_{\Omega} e^{h^T \omega} P(d\omega)\right) \le \Phi(h; \nu) \ \forall h \in \mathcal{H},$$
(2.136)

implying that $\mathcal{P} \subset \mathcal{S}[\mathcal{H}, \mathbf{\Delta}_L, \Phi]$, ν being a parameter of a distribution $P = \Pi[P^L, \nu] \in \mathcal{P}$.

Indeed,(2.136) is readily given by the fact that for $P = \Pi[P^L, \nu] \in \mathcal{P}$ and $h \in \mathcal{H}$ it holds

$$\ln\left(\mathbf{E}_{\omega\sim P}\left\{\mathbf{e}^{h^{T}\omega}\right\}\right) = \ln\left(\sum_{\ell=1}^{L}\nu_{\ell}\mathbf{E}_{\omega\sim P_{\ell}}\left\{\mathbf{e}^{h^{T}\omega}\right\}\right) \le \ln\left(\sum_{\ell=1}^{L}\nu_{\ell}\exp\{\Psi_{\ell}(h)\}\right) = \Phi(h;\nu),$$

with the concluding inequality given by $h \in \mathcal{H} \subset \mathcal{H}_{\ell}$ and $P_{\ell} \in \mathcal{R}[\mathcal{H}_{\ell}, \mathcal{M}_{\ell}, \Phi_{\ell}], 1 \leq \ell \leq L.$

We have build a simple family of distributions $S := S[\mathcal{H}, \Delta_L, \Phi]$ which contains all mixtures of distributions from given regular families $\mathcal{R}_{\ell} := \mathcal{R}[\mathcal{H}_{\ell}, \mathcal{M}_{\ell}, \Phi_{\ell}], 1 \leq \ell \leq L$, which makes S a simple outer approximation of mixtures of distributions from the simple families $S_{\ell} := S[\mathcal{H}_{\ell}, \mathcal{M}_{\ell}, \Phi_{\ell}], 1 \leq \ell \leq L$. In this latter capacity, Shas a drawback – the only parameter of the mixture $P = \Pi[P^L, \nu]$ of distributions $P_{\ell} \in S_{\ell}$ is ν , while the parameters of P_{ℓ} 's disappear. In some situations, this makes outer approximation S of \mathcal{P} too conservative. We are about to get rid, to come extent, of this drawback.

A modification. In the situation described in the beginning of 2.8.1.3.B, let a vector $\bar{\nu} \in \Delta_L$ be given, and let

$$\bar{\Phi}(h;\mu_1,...,\mu_L) = \sum_{\ell=1}^L \bar{\nu}_\ell \Phi_\ell(h;\mu_\ell) : \mathcal{H} \times (\mathcal{M}_1 \times ... \times \mathcal{M}_L) \to \mathbf{R}.$$
(2.137)

Let $d \times d$ matrix $Q \succeq 0$ satisfy

$$\left(\Phi_{\ell}(h;\mu_{\ell}) - \bar{\Phi}(h;\mu_{1},...,\mu_{L})\right)^{2} \leq h^{T}Qh \; \forall (h \in \mathcal{H}, \ell \leq L, \mu \in \mathcal{M}_{1} \times ... \times \mathcal{M}_{L}),$$
(2.138)

and let

136

$$\Phi(h;\mu_1,...,\mu_L) = \frac{3}{5}h^T Q h + \bar{\Phi}(h;\mu_1,...,\mu_L) : \mathcal{H} \times (\mathcal{M}_1 \times ... \times \mathcal{M}_L) \to \mathbf{R}.$$
(2.139)

 Φ clearly is convex-concave and continuous on its domain, whence $\mathcal{H} = \bigcap_{\ell} \mathcal{H}_{\ell}, \mathcal{M}_1 \times ... \times \mathcal{M}_L, \Phi$ is regular data.

Proposition 2.41. In the just defined situation, denoting by $\mathcal{P}_{\bar{\nu}}$ the family of all probability distributions $P = \Pi[P^L, \bar{\nu}]$, stemming from tuples

$$P^{L} = \{ P_{\ell} \in \mathcal{S}[\mathcal{H}_{\ell}, \mathcal{M}_{\ell}, \Phi_{\ell}] \}_{\ell=1}^{L}, \qquad (2.140)$$

one has

$$\mathcal{P}_{\bar{\nu}} \subset \mathcal{S}[\mathcal{H}, \mathcal{M}_1 \times ... \times \mathcal{M}_L, \Phi].$$
(2.141)

As a parameter of distribution $P = \Pi[P^L, \bar{\nu}] \in \mathcal{P}_{\bar{\nu}}$ with P^L as in (2.140), one can take $\mu^L = [\mu_1; ...; \mu_L]$.

Proof. It is easily seen that

$$\mathbf{e}^a \le a + \mathbf{e}^{\frac{3}{5}a^2}, \,\forall a.$$

As a result, when $a_{\ell}, \, \ell = 1, ..., L$, satisfy $\sum_{\ell} \bar{\nu}_{\ell} a_{\ell} = 0$, we have

$$\sum_{\ell} \bar{\nu}_{\ell} e^{a_{\ell}} \le \sum_{\ell} \bar{\nu}_{\ell} a_{\ell} + \sum_{\ell} \bar{\nu}_{\ell} e^{\frac{3}{5}a_{\ell}^2} \le e^{\frac{3}{5}\max_{\ell}a_{\ell}^2}.$$
 (2.142)

Now let P^L be as in (2.140), and let $h \in \mathcal{H} = \bigcap_L \mathcal{H}_\ell$. Setting $P = \prod [P^L, \bar{\nu}]$, we have

$$\ln\left(\int_{\Omega} e^{h^{T}\omega} P(d\omega)\right) = \ln\left(\sum_{\ell} \bar{\nu}_{\ell} \int_{\Omega} e^{h^{T}\omega} P_{\ell}(d\omega)\right) = \ln\left(\sum_{\ell} \bar{\nu}_{\ell} \exp\{\Phi_{\ell}(h,\mu_{\ell})\}\right)$$

$$= \bar{\Phi}(h;\mu_{1},...\mu_{L}) + \ln\left(\sum_{\ell} \bar{\nu}_{\ell} \exp\{\Phi_{\ell}(h,\mu_{\ell}) - \bar{\Phi}(h;\mu_{1},...\mu_{L})\}\right)$$

$$\leq a^{\bar{\Phi}}(h;\mu_{1},...\mu_{L}) + \frac{3}{5} \max_{\ell} [\Phi_{\ell}(h,\mu_{\ell}) - \bar{\Phi}(h;\mu_{1},...\mu_{L})]^{2} \leq b^{\bar{\Phi}}(h;\mu_{1},...,\mu_{L}),$$

where a is given by (2.142) as applied to $a_{\ell} = \Phi_{\ell}(h, \mu_{\ell}) - \overline{\Phi}(h; \mu_1, \dots, \mu_L)$, and b is due to (2.138), (2.139). The resulting inequality, which holds true for all $h \in \mathcal{H}$, is all we need.

2.8.1.3.C. I.I.D summation. Let $\Omega = \mathbf{R}^d$ be an observation space, $(\mathcal{H}, \mathcal{M}, \Phi)$ be regular data on this space, and let $\lambda = \{\lambda_\ell\}_{\ell=1}^K$ be a collection of reals. We can associate with the outlined entities a new data $(\mathcal{H}_\lambda, \mathcal{M}, \Phi_\lambda)$ on Ω by setting

$$\mathcal{H}_{\lambda} = \{h \in \Omega : \|\lambda\|_{\infty} h \in \mathcal{H}\}, \ \Phi_{\lambda}(h;\mu) = \sum_{\ell=1}^{L} \Phi(\lambda_{\ell}h;\mu) : \mathcal{H}_{\lambda} \times \mathcal{M} \to \mathbf{R}.$$

Now, given a probability distribution P on Ω , we can associate with it and with the above λ a new probability distribution P^{λ} on Ω as follows: P^{λ} is the distribution

of $\sum_{\ell} \lambda_{\ell} \omega_{\ell}$, where $\omega_1, \omega_2, ..., \omega_L$ are drawn, independently of each other, from *P*. An immediate observation is that the data $(\mathcal{H}_{\lambda}, \mathcal{M}, \Phi_{\lambda})$ is regular, and

• whenever a probability distribution P belongs to $\mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$, the distribution P^{λ} belongs to $\mathcal{S}[\mathcal{H}_{\lambda}, \mathcal{M}, \Phi_{\lambda}]$. In particular, when $\omega \sim P \in \mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$, then the distribution P^{L} of the sum of L independent copies of ω belongs to $\mathcal{S}[\mathcal{H}, \mathcal{M}, L\Phi]$.

2.8.1.3.D. Semi-direct summation. For $1 \le \ell \le L$, let regular data $\mathcal{H}_{\ell} \subset \Omega_{\ell} = \mathbf{R}^{d_{\ell}}$, \mathcal{M}_{ℓ} , Φ_{ℓ} be given. To avoid complications, we assume that for every ℓ ,

- $\mathcal{H}_{\ell} = \Omega_{\ell},$
- \mathcal{M}_{ℓ} is bounded.

Let also an $\epsilon > 0$ be given. We assume that ϵ is small, namely, $L\epsilon < 1$. Let us aggregate the given regular data into a new one by setting

$$\mathcal{H} = \Omega := \Omega_1 \times \ldots \times \Omega_L = \mathbf{R}^d, \ d = d_1 + \ldots + d_L, \ \mathcal{M} = \mathcal{M}_1 \times \ldots \times \mathcal{M}_L,$$

and let us define function $\Phi(h;\mu): \Omega^d \times \mathcal{M} \to \mathbf{R}$ as follows:

$$\Phi(h = [h^1; ...; h^L]; \mu = [\mu^1; ...; \mu^L]) = \inf_{\lambda \in \mathbf{\Delta}^\epsilon} \sum_{\ell=1}^d \lambda_\ell \Phi_\ell(h^\ell / \lambda_\ell; \mu^\ell),$$

$$\mathbf{\Delta}^\epsilon = \{\lambda \in \mathbf{R}^d : \lambda_\ell \ge \epsilon \,\forall \ell \& \, \sum_{\ell=1}^L \lambda_\ell = 1\}.$$
(2.143)

By evident reasons, the infimum in the description of Φ is achieved, and Φ is continuous. In addition, Φ is convex in $h \in \mathbf{R}^d$ and concave in $\mu \in \mathcal{M}$. Postponing for a moment verification, the consequences are that $\mathcal{H} = \Omega = \mathbf{R}^d$, \mathcal{M} and Φ form a regular data. We claim that

Whenever $\omega = [\omega^1; ...; \omega^L]$ is a random variable taking values in $\Omega = \mathbf{R}^{d_1} \times ... \times \mathbf{R}^{d_L}$, and the marginal distributions P_{ℓ} , $1 \leq \ell \leq L$, of ω belong to the families $S[\mathbf{R}^{d_{\ell}}, \mathcal{M}_{\ell}, \Phi_{\ell}]$ for all $1 \leq \ell \leq L$, the distribution P of ω belongs to $S[\mathbf{R}^d, \mathcal{M}, \Phi]$.

Indeed, since $P_{\ell} \in \mathcal{S}[\mathbf{R}^{d_{\ell}}, \mathcal{M}_{\ell}, \Phi_{\ell}]$, there exists $\hat{\mu}^{\ell} \in \mathcal{M}_{\ell}$ such that

$$\ln(\mathbf{E}_{\omega^{\ell} \sim P_{\ell}} \{ \exp\{g^{T} \omega^{\ell}\} \}) \leq \Phi_{\ell}(g; \widehat{\mu}^{\ell}) \; \forall g \in \mathbf{R}^{d_{\ell}}.$$

Let us set $\hat{\mu} = [\hat{\mu}^1; ...; \hat{\mu}^L]$, and let $h = [h^1; ...; h^L] \in \Omega$ be given. We can find $\lambda \in \Delta^{\epsilon}$ such that

$$\Phi(h;\widehat{\mu}) = \sum_{\ell=1}^{L} \lambda_{\ell} \Phi_{\ell}(h^{\ell}/\lambda_{\ell};\widehat{\mu}^{\ell}).$$

Applying Hölder inequality, we get

$$\mathbf{E}_{[\omega^1;\ldots;\omega^L]\sim P}\left\{\exp\{\sum_{\ell}[h^\ell]^T\omega^\ell\}\right\} \leq \prod_{\ell=1}^L \left(\mathbf{E}_{\omega^\ell\sim P_\ell}\left\{[h^\ell]^T\omega^\ell/\lambda_\ell\right\}\right)^{\lambda_\ell},$$

whence

$$\ln\left(\mathbf{E}_{[\omega^1;\ldots;\omega^L]\sim P}\left\{\exp\{\sum_{\ell}[h^{\ell}]^T\omega^{\ell}\}\right\}\right) \leq \sum_{\ell=1}^L \lambda_{\ell} \Phi_{\ell}(h^{\ell}/\lambda_{\ell};\widehat{\mu}^{\ell}) = \Phi(h;\widehat{\mu}).$$

We see that

$$\ln\left(\mathbf{E}_{[\omega^1;\ldots;\omega^L]\sim P}\left\{\exp\{\sum_{\ell}[h^\ell]^T\omega^\ell\}\right\}\right) \leq \Phi(h;\widehat{\mu}) \; \forall h \in \mathcal{H} = \mathbf{R}^d,$$

and thus $P \in \mathcal{S}[\mathbf{R}^d, \mathcal{M}, \Phi]$, as claimed.

It remains to verify that the function Φ defined by (2.143) indeed is convex in $h \in \mathbf{R}^d$ and concave in $\mu \in \mathcal{M}$. Concavity in μ is evident. Further, functions $\lambda_{\ell} \Phi_{\ell}(h^{\ell}/\lambda_{\ell};\mu)$ (as perspective transformation of convex functions $\Phi_{\ell}(\cdot;\mu)$) are jointly convex in λ and h^{ℓ} , and so is $\Psi(\lambda,h;\mu) = \sum_{\ell=1}^{L} \lambda_{\ell} \Phi_{\ell}(h^{\ell}/\lambda_{\ell},\mu)$. Thus $\Phi(\cdot;\mu)$, obtained by partial minimization of Ψ in λ , indeed is convex.

2.8.1.3.E. Affine image. Let \mathcal{H} , \mathcal{M} , Φ be regular data, Ω be the embedding space of \mathcal{H} , and $x \mapsto Ax + a$ be an affine mapping from Ω to $\overline{\Omega} = \mathbf{R}^{\overline{d}}$, and let us set

$$\bar{\mathcal{H}} = \{\bar{h} \in \mathbf{R}^d : A^T \bar{h} \in \mathcal{H}\}, \ \bar{\mathcal{M}} = \mathcal{M}, \ \bar{\Phi}(\bar{h};\mu) = \Phi(A^T \bar{h};\mu) + a^T \bar{h} : \ \bar{\mathcal{H}} \times \bar{\mathcal{M}} \to \mathbf{R}.$$

Note that $\overline{\mathcal{H}}$, $\overline{\mathcal{M}}$ and $\overline{\Phi}$ are regular data. It is immediately seen that

Whenever the probability distribution of a random variable ω belongs to $\mathcal{R}[\mathcal{H}, \mathcal{M}, \Phi]$ (or belongs to $\mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$), the distribution $\bar{P}[P]$ of the random variable $\bar{\omega} = A\omega + a$ belongs to $\mathcal{R}[\bar{\mathcal{H}}, \bar{\mathcal{M}}, \bar{\Phi}]$ (respectively, belongs to $\mathcal{S}[\bar{\mathcal{H}}, \bar{\mathcal{M}}, \bar{\Phi}]$).

2.8.1.3.F. Incorporating support information. Consider the situation as follows. We are given regular data $\mathcal{H} \subset \Omega = \mathbf{R}^d, \mathcal{M}, \Phi$ and are interested in a family \mathcal{P} of distributions known to belong to $\mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$. In addition, we know that all distributions P from \mathcal{P} are supported on a given closed convex set $X \subset \mathbf{R}^d$. How could we incorporate this domain information to pass from the family $\mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$ containing \mathcal{P} to a smaller family of the same type still containing \mathcal{P} ? We are about to give an answer in the simplest case of $\mathcal{H} = \Omega$. Specifically, denoting by $\phi_X(\cdot)$ the support function of X and selecting somehow a closed convex set $G \subset \mathbf{R}^d$ containing the origin, let us set

$$\widehat{\Phi}(h;\mu) = \inf_{g \in G} \left[\Phi^+(h,g;\mu) := \Phi(h-g;\mu) + \phi_X(g) \right],$$

where $\Phi(h; \mu) : \mathbf{R}^d \times \mathcal{M} \to \mathbf{R}$ is the continuous convex-concave function participating in the original regular data. Assuming that $\widehat{\Phi}$ is real-valued and continuous on the domain $\mathbf{R}^d \times \mathcal{M}$ (which definitely is the case when G is a compact set such that ϕ_X is finite and continuous on G), note that $\widehat{\Phi}$ is convex-concave on this domain, so that $\mathbf{R}^d, \mathcal{M}, \widehat{\Phi}$ is a regular data. We claim that

The family $\mathcal{S}[\mathbf{R}^d, \mathcal{M}, \widehat{\Phi}]$ contains \mathcal{P} , provided the family $\mathcal{S}[\mathbf{R}^d, \mathcal{M}, \Phi]$ does so, and the first of these two families is smaller than the second one.

Verification of the claim is immediate. Let $P \in \mathcal{P}$, so that for properly selected

 $\mu = \mu_P \in \mathcal{M}$ and for all $e \in \mathbf{R}^d$ it holds

$$\ln\left(\int_{\mathbf{R}^d} \exp\{e^T \omega\} P(d\omega)\right) \le \Phi(e;\mu_P).$$

Besides this, for every $g \in G$ we have $\phi_X(g) - g^T \omega \ge 0$ on the support of P, whence for every $h \in \mathbf{R}^d$ one has

$$\ln\left(\int_{\mathbf{R}^d} \exp\{h^T \omega\} P(d\omega)\right) \le \ln\left(\int_{\mathbf{R}^d} \exp\{h^T \omega + \phi_X(g) - g^T \omega\} P(d\omega)\right)$$

$$\le \phi_X(g) + \Phi(h - g; \mu_P).$$

Since the resulting inequality holds true for all $g \in G$, we get

$$\ln\left(\int_{\mathbf{R}^d} \exp\{h^T \omega\} P(d\omega)\right) \le \widehat{\Phi}(h; \mu_P) \; \forall h \in \mathbf{R}^d,$$

implying that $P \in \mathcal{S}[\mathbf{R}^d, \mathcal{M}, \widehat{\Phi}]$; since $P \in \mathcal{P}$ is arbitrary, the first part of the claim is justified. The inclusion $\mathcal{S}[\mathbf{R}^d, \mathcal{M}, \widehat{\Phi}] \subset \mathcal{S}[\mathbf{R}^d, \mathcal{M}, \Phi]$ is readily given by the inequality $\widehat{\Phi} \leq \Phi$, and the latter is due to $\widehat{\Phi}(h, \mu) \leq \Phi(h - 0, \mu) + \phi_X(0)$.

Illustration: distributions with bounded support revisited. In Section 2.8.1.2, given a convex compact set $X \subset \mathbf{R}^d$ with support function ϕ_X , we checked that the data $\mathcal{H} = \mathbf{R}^d$, $\mathcal{M} = X$, $\Phi(h;\mu) = h^T \mu + \frac{1}{8} [\phi_X(h) + \phi_X(-h)]^2$ are regular and the family $\mathcal{S}[\mathbf{R}^d, \mathcal{M}, \Phi]$ contains the family $\mathcal{P}[X]$ of all probability distributions supported on X. Moreover, for every $\mu \in \mathcal{M} = X$, the family $\mathcal{S}[\mathbf{R}^d, \{\mu\}, \Phi|_{\mathbf{R}^d \times \{\mu\}}]$ contains all supported on X distributions with the expectations $e[P] = \mu$. Note that $\Phi(h; e[P])$ describes well the behaviour of the logarithm $F_P(h) = \ln \left(\int_{\mathbf{R}^d} e^{h^T \omega} P(d\omega) \right)$ of the moment-generating function of $P \in \mathcal{P}[X]$ when h is small (indeed, $F_P(h) = h^T e[P] + O(||h||^2)$ as $h \to 0$), and by far overestimates $F_P(h)$ when h is large. Utilizing the above construction, we replace Φ with the real-valued, convex-concave and continuous on $\mathbf{R}^d \times \mathcal{M} = \mathbf{R}^d \times X$ (see Exercise 2.75) function

$$\widehat{\Phi}(h;\mu) = \inf_{g} \left[\widehat{\Psi}(h,g;\mu) := (h-g)^{T} \mu + \frac{1}{8} [\phi_{X}(h-g) + \phi_{X}(-h+g)]^{2} + \phi_{X}(g) \right] \\
\leq \Phi(h;\mu).$$
(2.144)

It is easy to see that $\widehat{\Phi}(\cdot; \cdot)$ still ensures the inclusion $P \in \mathcal{S}[\mathbf{R}^d, \{e[P]\}, \widehat{\Phi}|_{\mathbf{R}^d \times \{e[P]\}}]$ for every distribution $P \in \mathcal{P}[X]$ and "reproduces $F_P(h)$ reasonably well" for both small and large h. Indeed, since $F_P(h) \leq \widehat{\Phi}(h; e[P]) \leq \Phi(h; e[P])$, for small h $\widehat{\Phi}(h; e[P])$ reproduces $F_P(h)$ even better than $\Phi(h; e[P])$, and we clearly have

$$\widehat{\Phi}(h;\mu) \le \left[(h-h)^T \mu + \frac{1}{8} [\phi_X(h-h) + \phi_X(-h+h)]^2 + \phi_X(h) \right] = \phi_X(h) \ \forall \mu,$$

and $\phi_X(h)$ is a correct description of $F_P(h)$ for large h.

140

LECTURE 2

2.8.2 Main result

2.8.2.1 Situation & Construction

Assume we are given two collections of regular data with common $\Omega = \mathbf{R}^d$ and common \mathcal{H} , specifically, the collections $(\mathcal{H}, \mathcal{M}_{\chi}, \Phi_{\chi}), \chi = 1, 2$. We start with constructing a specific detector for the associated families of regular probability distributions

$$\mathcal{P}_{\chi} = \mathcal{R}[\mathcal{H}, \mathcal{M}_{\chi}, \Phi_{\chi}], \ \chi = 1, 2.$$

When building the detector, we impose on the regular data in question the following

Assumption I: The regular data $(\mathcal{H}, \mathcal{M}_{\chi}, \Phi_{\chi}), \chi = 1, 2$, are such that the convex-concave function

$$\Psi(h;\mu_1,\mu_2) = \frac{1}{2} \left[\Phi_1(-h;\mu_1) + \Phi_2(h;\mu_2) \right] : \mathcal{H} \times (\mathcal{M}_1 \times \mathcal{M}_2) \to \mathbf{R} \quad (2.145)$$

has a saddle point (min in $h \in \mathcal{H}$, max in $(\mu_1, \mu_2) \in \mathcal{M}_1 \times \mathcal{M}_2$).

A simple sufficient condition for existence of a saddle point of (2.145) is

Condition A: The sets \mathcal{M}_1 and \mathcal{M}_2 are compact, and the function

$$\overline{\Phi}(h) = \max_{\mu_1 \in \mathcal{M}_1, \mu_2 \in \mathcal{M}_2} \Phi(h; \mu_1, \mu_2)$$

is coercive on \mathcal{H} , meaning that $\overline{\Phi}(h_i) \to \infty$ along every sequence $h_i \in \mathcal{H}$ with $||h_i||_2 \to \infty$ as $i \to \infty$.

Indeed, under Condition A by Sion-Kakutani Theorem (Theorem 2.24) it holds

$$\operatorname{SadVal}[\Phi] := \inf_{h \in \mathcal{H}} \underbrace{\max_{\mu_1 \in M_1, \mu_2 \in \mathcal{M}_2} \Phi(h; \mu_1, \mu_2)}_{\overline{\Phi}(h)} = \sup_{\mu_1 \in M_1, \mu_2 \in \mathcal{M}_2} \underbrace{\inf_{h \in \mathcal{H}} \Phi(h; \mu_1, \mu_2)}_{\underline{\Phi}(\mu_1, \mu_2)},$$

so that the optimization problems

$$(P): \quad \operatorname{Opt}(P) = \min_{h \in \mathcal{H}} \Phi(h)$$

$$(D): \quad \operatorname{Opt}(D) = \max_{\mu_1 \in \mathcal{M}_1, \mu_2 \in \mathcal{M}_2} \underline{\Phi}(\mu_1, \mu_2)$$

have equal optimal values. Under Condition A, problem (P) clearly is a problem of minimizing a continuous coercive function over a closed set and as such is solvable; thus, Opt(P) = Opt(D) is a real. Problem (D) clearly is the problem of maximizing over a compact set an upper semi-continuous (since Φ is continuous) function taking real values and, perhaps, value $-\infty$, and not identically equal to $-\infty$ (since Opt(D) is a real), and thus (D) is solvable. Thus, (P) and (D) are solvable with common optimal values, and therefore Φ has a saddle point.

2.8.2.2 Main Result

An immediate (and crucial!) observation is as follows:

Proposition 2.42. In the situation of Section 2.8.2.1, let $h \in \mathcal{H}$ be such that the quantities

$$\Psi_1(h) = \sup_{\mu_1 \in \mathcal{M}_1} \Phi_1(-h;\mu_1), \ \Psi_2(h) = \sup_{\mu_2 \in \mathcal{M}_2} \Phi_2(h;\mu_2)$$

are finite. Consider the affine detector

$$\phi_h(\omega) = h^T \omega + \underbrace{\frac{1}{2} [\Psi_1(h) - \Psi_2(h)]}_{\varkappa}.$$

Then

$$\operatorname{Risk}[\phi_h | \mathcal{R}[\mathcal{H}, \mathcal{M}_1, \Phi_1], \mathcal{R}[\mathcal{H}, \mathcal{M}_2, \Phi_2]] \le \exp\{\frac{1}{2}[\Psi_1(h) + \Psi_2(h)]\}.$$
(2.146)

Proof. Let *h* satisfy the premise of Proposition. For every $\mu_1 \in \mathcal{M}_1$, we have $\Phi_1(-h; \mu_1) \leq \Psi_1(h)$, and for every $P \in \mathcal{R}[\mathcal{H}, \mathcal{M}_1, \Phi_1]$, we have

$$\int_{\Omega} \exp\{-h^T \omega\} P(d\omega) \le \exp\{\Phi_1(-h;\mu_1)\}$$

for properly selected $\mu_1 \in \mathcal{M}_1$. Thus,

$$\int_{\Omega} \exp\{-h^T \omega\} P(d\omega) \le \exp\{\Psi_1(h)\} \ \forall P \in \mathcal{R}[\mathcal{H}, \mathcal{M}_1, \Phi_1],$$

whence also

 $\int_{\Omega} \exp\{-h^{T}\omega - \varkappa\} P(d\omega) \leq \exp\{\Psi_{1}(h) - \varkappa\} = \exp\{\frac{1}{2}[\Psi_{1}(h) + \Psi_{2}(h)]\} \ \forall P \in \mathcal{R}[\mathcal{H}, \mathcal{M}_{1}, \Phi_{1}].$ Similarly, for every $\mu_{2} \in \mathcal{M}_{2}$, we have $\Phi_{2}(h; \mu_{2}) \leq \Psi_{2}(h)$, and for every $P \in \mathcal{R}[\mathcal{H}, \mathcal{M}_{2}, \Phi_{2}]$, we have

$$\int_{\Omega} \exp\{h^T \omega\} P(d\omega) \le \exp\{\Phi_2(h;\mu_2)\}$$

for properly selected $\mu_2 \in \mathcal{M}_2$. Thus,

$$\int_{\Omega} \exp\{h^T \omega\} P(d\omega) \le \exp\{\Psi_2(h)\} \ \forall P \in \mathcal{R}[\mathcal{H}, \mathcal{M}_2, \Phi_2],$$

whence also

$$\int_{\Omega} \exp\{h^{T}\omega + \varkappa\}P(d\omega) \le \exp\{\Psi_{2}(h) + \varkappa\}\exp\{\frac{1}{2}[\Psi_{1}(h) + \Psi_{2}(h)]\} \ \forall P \in \mathcal{R}[\mathcal{H}, \mathcal{M}_{2}, \Phi_{2}] \ \Box$$

An immediate corollary is as follows:

Proposition 2.43. In the situation of Section 2.8.2.1 and under Assumption I, let us associate with a saddle point $(h_*; \mu_1^*, \mu_2^*)$ of the convex-concave function (2.145) the following entities:

• the risk

$$\epsilon_{\star} := \exp\{\Psi(h_{\star}; \mu_1^{\star}, \mu_2^{\star})\}; \qquad (2.147)$$

this quantity is uniquely defined by the saddle point value of Ψ and thus is independent of how we select a saddle point;

• the detector $\phi_*(\omega)$ – the affine function of $\omega \in \mathbf{R}^d$ given by

$$\phi_*(\omega) = h_*^T \omega + a, \ a = \frac{1}{2} \left[\Phi_1(-h_*;\mu_1^*) - \Phi_2(h_*;\mu_2^*) \right].$$
(2.148)

Then

$$\operatorname{Risk}[\phi_* | \mathcal{R}[\mathcal{H}, \mathcal{M}_1, \Phi_1], \mathcal{R}[\mathcal{H}, \mathcal{M}_2, \Phi_2]] \le \epsilon_{\star}.$$
(2.149)

Consequences. Assume we are given L collections $(\mathcal{H}, \mathcal{M}_{\ell}, \Phi_{\ell})$ of regular data on a common observation space $\Omega = \mathbf{R}^d$ and with common \mathcal{H} , and let

$$\mathcal{P}_{\ell} = \mathcal{R}[\mathcal{H}, \mathcal{M}_{\ell}, \Phi_{\ell}]$$

be the corresponding families of regular distributions. Assume also that for every pair (ℓ, ℓ') , $1 \leq \ell < \ell' \leq L$, the pair of regular data $(\mathcal{H}, \mathcal{M}_{\ell}, \Phi_{\ell})$, $(\mathcal{H}, \mathcal{M}_{\ell'}, \Phi_{\ell'})$ satisfies Assumption I, so that the convex-concave functions

$$\Psi_{\ell\ell'}(h;\mu_{\ell},\mu_{\ell'}) = \frac{1}{2} \left[\Phi_{\ell}(-h;\mu_{\ell}) + \Phi_{\ell'}(h;\mu_{\ell'}) \right] : \mathcal{H} \times (\mathcal{M}_{\ell} \times \mathcal{M}_{\ell'}) \to \mathbf{R}$$

$$[1 < \ell < \ell' < L]$$

have saddle points $(h_{\ell\ell'}^*; (\mu_{\ell}^*, \mu_{\ell'}^*))$ (min in $h \in \mathcal{H}$, max in $(\mu_{\ell}, \mu_{\ell'}) \in \mathcal{M}_{\ell} \times \mathcal{M}_{\ell'}$). These saddle points give rise to affine detectors

$$\phi_{\ell\ell'}(\omega) = [h_{\ell\ell'}^*]^T \omega + \frac{1}{2} \left[\Phi_\ell(-h_{\ell\ell'}^*;\mu_\ell^*) - \Phi_{\ell'}(h_*;\mu_{\ell'}^*) \right] \qquad [1 \le \ell < \ell' \le L]$$

and the quantities

$$\epsilon_{\ell\ell'} = \exp\left\{\frac{1}{2} \left[\Phi_{\ell}(-h_{\ell\ell'}^*; \mu_{\ell}^*) + \Phi_{\ell'}(h_*; \mu_{\ell'}^*)\right]\right\}; \qquad [1 \le \ell < \ell' \le L]$$

by Proposition 2.43, $\epsilon_{\ell\ell'}$ are upper bounds on the risks, taken w.r.t. $\mathcal{P}_{\ell}, \mathcal{P}_{\ell'}$, of the detectors $\phi_{\ell\ell'}$:

$$\int_{\Omega} e^{-\phi_{\ell\ell'}(\omega)} P(d\omega) \le \epsilon_{\ell\ell'} \ \forall P \in \mathcal{P}_{\ell} \ \& \ \int_{\Omega} e^{\phi_{\ell\ell'}(\omega)} P(d\omega) \le \epsilon_{\ell\ell'} \ \forall P \in \mathcal{P}_{\ell'}.$$

 $[1 \leq \ell < \ell' \leq L].$ Setting $\phi_{\ell\ell'}(\cdot) = -\phi_{\ell'\ell}(\cdot)$ and $\epsilon_{\ell\ell'} = \epsilon_{\ell'\ell}$ when $L \geq \ell > \ell' \geq 1$ and $\phi_{\ell\ell}(\cdot) \equiv 0$, $\epsilon_{\ell\ell} = 1, 1 \leq \ell \leq L$, we get a system of detectors and risks satisfying (2.98) and, consequently, can use these "building blocks" in the developed so far machinery for pairwise- and multiple hypothesis testing from single and repeated observations (stationary, semi-stationary, and quasi-stationary).

Numerical example. To get some impression of how Proposition 2.43 extends the grasp of our computation-friendly test design machinery. consider a toy problem as follows:

We are given observation

$$\omega = Ax + \sigma A \operatorname{Diag} \left\{ \sqrt{x_1}, \dots, \sqrt{x_n} \right\} \xi, \qquad (2.150)$$

where

- unknown signal x is known to belong to a given convex compact subset M of the *interior* of Rⁿ₊;
- A is a given $n \times n$ matrix of rank $n, \sigma > 0$ is a given noise intensity, and $\xi \sim \mathcal{N}(0, I_n)$.

Our goal is to decide via K-repeated version of observations (2.150) on the pair of hypotheses $x \in X_{\chi}$, $\chi = 1, 2$, where X_1 , X_2 are given nonempty convex compact subsets of M.

Note that an essential novelty, as compared to the standard Gaussian o.s., is that now we deal with zero mean Gaussian noise with covariance matrix

$$\Theta(x) = \sigma^2 A \operatorname{Diag}\{x\} A^T$$

depending on the true signal – the larger signal, the larger noise.

We can easily process the situation in question via the machinery developed in this Section. Specifically, let us set

$$\mathcal{H}_{\chi} = \mathbf{R}^{n}, \ \mathcal{M}_{\chi} = \{(x, \operatorname{Diag}\{x\}) : x \in X_{\chi}\} \subset \mathbf{R}^{n} \times \mathbf{S}^{n}_{+}, \\ \Phi_{\chi}(h; x, \Xi) = h^{T}x + \frac{\sigma^{2}}{2}h^{T}[A\Xi A^{T}]h : \mathcal{M}_{\chi} \to \mathbf{R}$$
 $[\chi = 1, 2]$

It is immediately seen that for $\chi = 1, 2, \mathcal{H}, \mathcal{M}_{\chi}, \Phi_{\chi}$ is regular data, and that the distribution P of observation (2.150) stemming from a signal $x \in X_{\chi}$ belongs to $\mathcal{S}[\mathcal{H}, \mathcal{M}_{\chi}, \Phi_{\chi}]$, so that we can use Proposition 2.43 to build an affine detector for the families $\mathcal{P}_{\chi}, \chi = 1, 2$, of distributions of observations (2.150) stemming from signals $x \in X_{\chi}$. The corresponding recipe boils down to the necessity to find a saddle point $(h_*; x_*, y_*)$ of the simple convex-concave function

$$\Psi(h; x, y) = \frac{1}{2} \left[h^T(y - x) + \frac{\sigma^2}{2} h^T A \text{Diag}\{x + y\} A^T h \right]$$
(2.151)

(min in $h \in \mathbf{R}^n$, max in $(x, y) \in X_1 \times X_2$); such a point clearly exists and is easily found, and gives rise to affine detector

$$\phi_*(\omega) = h_*^T \omega + \underbrace{\frac{\sigma^2}{4} h_*^T A \text{Diag}\{x_* - y_*\} A^T h_* - \frac{1}{2} h_*^T [x_* + y_*]}_{a}$$

such that

$$\operatorname{Risk}[\phi_*|\mathcal{P}_1, \mathcal{P}_2] \le \exp\left\{\frac{1}{2}\left[h_*^T[y_* - x_*] + \frac{\sigma^2}{2}h_*^TA\operatorname{Diag}\{x_* + y_*\}A^Th_*\right]\right\}. \quad (2.152)$$

Note that we could also process the situation when defining the regular data as $\mathcal{H}, \mathcal{M}_{\chi}^+ = X_{\chi}, \Phi_{\chi}^+, \chi = 1, 2$, where

$$\Phi_{\chi}^{+}(h;x) = h^{T}x + \frac{\sigma^{2}\theta}{2}h^{T}AA^{T}h \qquad \qquad [\theta = \max_{x \in X_{1} \cup X_{2}} \|x\|_{\infty}]$$

which, basically, means passing from our actual observations (2.150) to the "more noisy" observations given by the Gaussian o.s.

$$\omega = Ax + \eta, \ \eta \sim \mathcal{N}(0, \sigma^2 \theta A A^T).$$
(2.153)

It is easily seen that the risk $\operatorname{Risk}[\phi_{\#}|\mathcal{P}_1, \mathcal{P}_2]$ of the optimal, for this Gaussian o.s., detector $\phi_{\#}$, can be upper-bounded by the known to us risk $\operatorname{Risk}[\phi_{\#}|\mathcal{P}_1^+, \mathcal{P}_2^+]$, where \mathcal{P}_{χ}^+ is the family of distributions of observations (2.153) induced by signals $x \in X_{\chi}$. Note that were we staying within the realm of detector-based tests in simple o.s.'s, $\operatorname{Risk}[\phi_{\#}|\mathcal{P}_1^+, \mathcal{P}_2^+]$ would be seemingly the best risk bound available for us. The goal of the small numerical experiment we are about to report was to understand how our new risk bound (2.152) compares to the "old" bound $\operatorname{Risk}[\phi_{\#}|\mathcal{P}_1^+, \mathcal{P}_2^+]$. We used

$$n = 16, X_1 = \left\{ x \in \mathbf{R}^{16} : \begin{array}{l} 0.001 \le x_1 \le \delta \\ 0.001 \le x_i \le 1, \ 2 \le i \le 16 \end{array} \right\}, \\ X_2 = \left\{ x \in \mathbf{R}^{16} : \begin{array}{l} 2\delta \le x_1 \le 1 \\ 0.001 \le x_i \le 1, \ 2 \le i \le 16 \end{array} \right\}$$

and $\sigma = 0.1$. The "separation parameter" δ was set to 0.1. Finally, 16×16 matrix A was generated to have condition number 100 (singular values $0.01^{(i-1)/15}$, $1 \leq i \leq 16$) with randomly oriented systems of left- and right singular vectors. With this setup, a typical numerical result is as follows:

- the right hand side in (2.152) is 0.4346, implying that with detector ϕ_* , 6-repeated observation is sufficient to decide on our two hypotheses with risk ≤ 0.01 ;
- the quantity $\operatorname{Risk}[\phi_{\#}|\mathcal{P}_1^+, \mathcal{P}_2^+]$ is 0.8825, meaning that with detector $\phi_{\#}$, we need at least 37-repeated observation to guarantee risk ≤ 0.01 .

When the separation parameter δ participating in the descriptions of X_1 , X_2 was reduced to 0.01, the risks in question grew to 0.9201 and 0.9988, respectively (56repeated observation to decide on the hypotheses with risk 0.01 when ϕ_* is used vs. 3685-repeated observation needed when $\phi_{\#}$ is used). The bottom line is that our new developments could improve quite significantly the performance of our inferences.

2.8.2.3 Illustration: sub-Gaussian and Gaussian cases

For $\chi = 1, 2$, let U_{χ} be nonempty closed convex set in \mathbf{R}^d , and \mathcal{V}_{χ} be a compact convex subset of the interior of the positive semidefinite cone \mathbf{S}^d_+ . We assume that U_1 is compact. Setting

$$\mathcal{H}_{\chi} = \Omega = \mathbf{R}^{d}, \ \mathcal{M}_{\chi} = U_{\chi} \times \mathcal{V}_{\chi}, \Phi_{\chi}(h;\theta,\Theta) = \theta^{T}h + \frac{1}{2}h^{T}\Theta h : \mathcal{H}_{\chi} \times \mathcal{M}_{\chi} \to \mathbf{R}, \ \chi = 1, 2,$$
(2.154)

we get two collections $(\mathcal{H}, \mathcal{M}_{\chi}, \Phi_{\chi}), \chi = 1, 2$, of regular data. As we know from Section 2.8.1.2, for $\chi = 1, 2$, the families of distributions $\mathcal{S}[\mathbf{R}^d, \mathcal{M}_{\chi}, \Phi_{\chi}]$ contain the families $\mathcal{S}G[U_{\chi}, \mathcal{V}_{\chi}]$ of sub-Gaussian distributions on \mathbf{R}^d with sub-Gaussianity parameters $(\theta, \Theta) \in U_{\chi} \times \mathcal{V}_{\chi}$ (see (2.133)), as well as families $\mathcal{G}[U_{\chi}, \mathcal{V}_{\chi}]$ of Gaussian distributions on \mathbf{R}^d with parameters (θ, Θ) (expectation and covariance matrix) running through $U_{\chi} \times \mathcal{V}_{\chi}$. Besides this, the pair of regular data in question clearly satisfies Condition A. Consequently, the test \mathcal{T}_*^K given by the above construction as applied to the collections of regular data (2.154) is well defined and allows to decide on hypotheses

$$H_{\chi}: P \in \mathcal{R}[\mathbf{R}^d, U_{\chi}, \mathcal{V}_{\chi}], \ \chi = 1, 2,$$

on the distribution P underlying K-repeated observation ω^K . The same test can be also used to decide on stricter hypotheses H_{χ}^G , $\chi = 1, 2$, stating that the observations $\omega_1, ..., \omega_K$ are i.i.d. and drawn from a Gaussian distribution P belonging to $\mathcal{G}[U_{\chi}, \mathcal{V}_{\chi}]$. Our goal now is to process in detail the situation in question and to refine our conclusions on the risk of the test \mathcal{T}_*^1 when the *Gaussian* hypotheses H_{χ}^G are considered and the situation is symmetric, that is, when $\mathcal{V}_1 = \mathcal{V}_2$.

Observe, first, that the convex-concave function Ψ from (2.145) in the situation under consideration becomes

$$\Psi(h;\theta_1,\Theta_1,\theta_2,\Theta_2) = \frac{1}{2}h^T[\theta_2 - \theta_1] + \frac{1}{4}h^T\Theta_1h + \frac{1}{4}h^T\Theta_2h.$$
(2.155)

We are interested in solutions to the saddle point problem

$$\min_{h \in \mathbf{R}^d} \max_{\substack{\theta_1 \in U_1, \theta_2 \in U_2\\\Theta_1 \in \mathcal{V}_1, \Theta_2 \in \mathcal{V}_2}} \Psi(h; \theta_1, \Theta_1, \theta_2, \Theta_2)$$
(2.156)

associated with the function (2.155). From the structure of Ψ and compactness of U_1 , \mathcal{V}_1 , \mathcal{V}_2 , combined with the fact that \mathcal{V}_{χ} , $\chi = 1, 2$, are comprised of positive definite matrices, it immediately follows that saddle points do exist, and a saddle point $(h_*; \theta_1^*, \Theta_1^*, \Theta_2^*, \Theta_2^*)$ satisfies the relations

(a)
$$h_* = [\Theta_1^* + \Theta_2^*]^{-1} [\theta_1^* - \theta_2^*],$$

 $\begin{array}{ll} (b) & h_*^T(\theta_1 - \theta_1^*) \ge 0 \ \forall \theta_1 \in U_1, \ h_*^T(\theta_2^* - \theta_2) \ge 0 \ \forall \theta_2 \in U_2, \\ (c) & h_*^T\Theta_1h_* \le h_*^T\Theta_1^*h_* \ \forall \Theta_1 \in \mathcal{V}_1, \ h_*^T\Theta_2^*h_* \le h_*\Theta_2^*h_* \ \forall \Theta_2 \in \mathcal{V}_2. \end{array}$ (2.157)

From (2.157.*a*) it immediately follows that the affine detector $\phi_*(\cdot)$ and risk ϵ_* , as given by (2.147) and (2.148), are

$$\begin{aligned}
\phi_*(\omega) &= h_*^T [\omega - w_*] + \frac{1}{2} h_*^T [\Theta_1^* - \Theta_2^*] h_*, \ w_* &= \frac{1}{2} [\theta_1^* + \theta_2^*]; \\
\epsilon_* &= \exp\{-\frac{1}{4} [\theta_1^* - \theta_2^*]^T [\Theta_1^* + \Theta_2^*]^{-1} [\theta_1^* - \theta_2^*]\} \\
&= \exp\{-\frac{1}{4} h_*^T [\Theta_1^* + \Theta_2^*] h_*\}.
\end{aligned}$$
(2.158)

Note that in the symmetric case (where $\mathcal{V}_1 = \mathcal{V}_2$), there always exists a saddle point of Ψ with $\Theta_1^* = \Theta_2^{*31}$, and the test \mathcal{T}_*^1 associated with such saddle point is quite transparent: it is the maximum likelihood test for two Gaussian distributions, $\mathcal{N}(\theta_1^*, \Theta_*), \mathcal{N}(\theta_2^*, \Theta_*)$, where Θ_* is the common value of Θ_1^* and Θ_2^* , and the bound ϵ_* on the risk of the test is nothing but the Hellinger affinity of these two Gaussian distributions, or, equivalently,

$$\epsilon_{\star} = \exp\left\{-\frac{1}{8}[\theta_1^* - \theta_2^*]^T \Theta_{\star}^{-1}[\theta_1^* - \theta_2^*]\right\}.$$
(2.159)

We arrive at the following result:

Proposition 2.44. In the symmetric sub-Gaussian case (i.e., in the case of (2.154) with $\mathcal{V}_1 = \mathcal{V}_2$), saddle point problem (2.155), (2.156) admits a saddle point of the

³¹Indeed, from (2.155) it follows that when $\mathcal{V}_1 = \mathcal{V}_2$, the function $\Psi(h;\theta_1,\Theta_1,\theta_2,\Theta_2)$ is symmetric w.r.t. Θ_1,Θ_2 , implying similar symmetry of the function $\underline{\Psi}(\theta_1,\Theta_1,\theta_2,\Theta_2) = \min_{h\in\mathcal{H}}\Psi(h;\theta_1,\Theta_1,\theta_2,\Theta_2)$. Since $\underline{\Psi}$ is concave, the set M of its maximizers over $\mathcal{M}_1 \times \mathcal{M}_2$ (which, as we know, is nonempty) is symmetric w.r.t. the swap of Θ_1 and Θ_2 and is convex, implying that if $(\theta_1,\Theta_1,\theta_2,\Theta_2) \in M$, then $(\theta_1,\frac{1}{2}[\Theta_1+\Theta_2],\theta_2,\frac{1}{2}[\Theta_1+\Theta_2]) \in M$ as well, and the latter point is the desired component of saddle point of Ψ with $\Theta_1 = \Theta_2$.

form $(h_*; \theta_1^*, \Theta_*, \theta_2^*, \Theta_*)$, and the associated affine detector and its risk are given by

$$\phi_*(\omega) = h_*^T[\omega - w_*], \ w_* = \frac{1}{2}[\theta_1^* + \theta_2^*];
\epsilon_* = \exp\{-\frac{1}{8}[\theta_1^* - \theta_2^*]^T\Theta_*^{-1}[\theta_1^* - \theta_2^*]\}.$$
(2.160)

As a result, when deciding, via ω^K , on "sub-Gaussian hypotheses" H_{χ} , $\chi = 1, 2$, the risk of the test \mathcal{T}^K_* associated with $\phi^{(K)}_*(\omega^K) := \sum_{t=1}^K \phi_*(\omega_t)$ is at most ϵ^K_* .

In the symmetric single-observation Gaussian case, that is, when $\mathcal{V}_1 = \mathcal{V}_2$ and we apply the test $\mathcal{T}_* = \mathcal{T}^1_*$ to observation $\omega \equiv \omega_1$ in order to decide on the hypotheses H^G_{χ} , $\chi = 1, 2$, the above risk bound can be improved:

Proposition 2.45. Consider symmetric case $\mathcal{V}_1 = \mathcal{V}_2 = \mathcal{V}$, let $(h_*; \theta_1^*; \Theta_1^*, \theta_2^*, \Theta_2^*)$ be "symmetric" – with $\Theta_1^* = \Theta_2^* = \Theta_*$ – saddle point of function Ψ given by (2.155), and let ϕ_* be the affine detector given by (2.157) and (2.158):

$$\phi_*(\omega) = h_*^T[\omega - w_*], \ h_* = \frac{1}{2}\Theta_*^{-1}[\theta_1^* - \theta_2^*], \ w_* = \frac{1}{2}[\theta_1^* + \theta_2^*].$$

Let also

$$\delta = \sqrt{h_*^T \Theta_* h_*} = \frac{1}{2} \sqrt{[\theta_1^* - \theta_2^*]^T \Theta_*^{-1} [\theta_1^* - \theta_2^*]}, \qquad (2.161)$$

so that

$$\delta^2 = h_*^T [\theta_1^* - w_*] = h_*^T [w_* - \theta_2^*] \quad and \ \epsilon_\star = \exp\{-\frac{1}{2}\delta^2\}.$$
(2.162)

Let, further, $\alpha \leq \delta^2$, $\beta \leq \delta^2$. Then

- (a) $\forall (\theta \in U_1, \Theta \in \mathcal{V}) : \operatorname{Prob}_{\omega \sim \mathcal{N}(\theta, \Theta)} \{ \phi_*(\omega) \le \alpha \} \le \operatorname{Erf}(\delta \alpha/\delta)$ (2.163)
- (b) $\forall (\theta \in U_2, \Theta \in \mathcal{V}) : \operatorname{Prob}_{\omega \sim \mathcal{N}(\theta, \Theta)} \{ \phi_*(\omega) \ge -\beta \} \le \operatorname{Erf}(\delta \beta/\delta),$ (2.163)

where

$$\operatorname{Erf}(s) = \frac{1}{\sqrt{2\pi}} \int_{s}^{\infty} \exp\{-r^{2}/2\} dr$$

is the normal error function. In particular, when deciding, via a single observation ω , on Gaussian hypotheses H_{χ}^G , $\chi = 1, 2$, with H_{χ}^G stating that $\omega \sim \mathcal{N}(\theta, \Theta)$ with $(\theta, \Theta) \in U_{\chi} \times \mathcal{V}$, the risk of the test \mathcal{T}_*^1 associated with ϕ_* is at most $\mathrm{Erf}(\delta)$.

Proof. Let us prove (a) (the proof of (b) is completely similar). For $\theta \in U_1, \Theta \in \mathcal{V}$ we have

$$\begin{aligned} \operatorname{Prob}_{\omega \sim \mathcal{N}(\theta,\Theta)} \{ \phi_*(\omega) \leq \alpha \} &= \operatorname{Prob}_{\omega \sim \mathcal{N}(\theta,\Theta)} \{ h_*^T[\omega - w_*] \leq \alpha \} \\ &= \operatorname{Prob}_{\xi \sim \mathcal{N}(0,I)} \{ h_*^T[\theta + \Theta^{1/2}\xi - w_*] \leq \alpha \} \\ &= \operatorname{Prob}_{\xi \sim \mathcal{N}(0,I)} \{ [\Theta^{1/2}h_*]^T \xi \leq \alpha - \underbrace{h_*^T[\theta - w_*]}_{\text{by } (2.157.b), (2.162)}^2 \} \\ &\leq \operatorname{Prob}_{\xi \sim \mathcal{N}(0,I)} \{ [\Theta^{1/2}h_*]^T \xi \leq \alpha - \delta^2 \} \\ &= \operatorname{Erf}([\delta^2 - \alpha] / \|\Theta^{1/2}h_*\|_2) \\ &\leq \operatorname{Erf}([\delta^2 - \alpha] / \|\Theta^{1/2}h_*\|_2) \text{ [since } \delta^2 - \alpha \geq 0 \text{ and } h_*^T \Theta h_* \leq h_*^T \Theta_* h_* \text{ by } (2.157.c) \\ &= \operatorname{Erf}([\delta^2 - \alpha] / \delta). \end{aligned}$$

The "in particular" part of Proposition is readily given by (2.163) as applied with

Note that the progress, as compared to our results on the minimum risk detectors for convex hypotheses in Gaussian o.s. is that we do *not* assume anymore that the covariance matrix is once for ever fixed. Now both the mean *and* the covariance matrix of Gaussian random variable we are observing are not known in advance, the mean is allowed to run through a closed convex set (depending on the hypothesis), the covariance is allowed to run, independently of the mean, through a given convex compact subset of the interior of the positive definite cone, and this subset should be common for both hypotheses we are deciding upon.

2.9 BEYOND THE SCOPE OF AFFINE DETECTORS: LIFTING OBSERVATIONS

2.9.1 Motivation

The detectors considered so far in this Section were affine functions of observations. Note, however, that what is an observation, it to some extent depends on us. To give an instructive example, consider the Gaussian observation

$$\zeta = A[u;1] + \xi \in \mathbf{R}^n,$$

where u is unknown signal known to belong to a given set $U \subset \mathbf{R}^n, u \mapsto A[u; 1]$ is a given affine mapping from \mathbf{R}^n into the observation space \mathbf{R}^d , and ξ is zero mean Gaussian observation noise with covariance matrix Θ known to belong to a given convex compact subset \mathcal{V} of the interior of the positive semidefinite cone \mathbf{S}_{\perp}^{d} . Treating observation "as is", affine in observation detector is affine in $[u;\xi]$. On the other hand, we can treat as our observation the image of the actual observation ζ under a whatever deterministic mapping, e.g., the "quadratic lift" $\zeta \mapsto (\zeta, \zeta \zeta^T)$. A detector affine in the new observation is quadratic in u and ξ – we get access to a wider set of detectors as compared to those affine in ζ ! At the first glance, applying our "affine detectors" machinery to appropriate "nonlinear lifts" of actual observations we can handle quite complicated detectors, e.g., polynomial, of arbitrary degree, in ζ . The bottleneck here stems from the fact that in general it is really difficult to "cover" the distribution of "nonlinearly lifted" actual observation ζ (even as simple as the above Gaussian observation) by an explicitly defined family of regular distributions; and such a "covering" is what we need in order to apply to the lifted observation our affine detector machinery. It turns out, however, that in some important cases the desired covering is achievable. We are about to demonstrate that this favorable situation indeed takes place when speaking about the quadratic lifting $\zeta \mapsto (\zeta, \zeta \zeta^T)$ of (sub)Gaussian observation ζ , and the resulting quadratic detectors allow to handle some important inference problems which are far beyond the grasp of "genuinely affine" detectors.

2.9.2 Quadratic lifting: Gaussian case

Given positive integer d, we define \mathcal{E}^d as the linear space $\mathbf{R}^d \times \mathbf{S}^d$ equipped with the inner product

$$\langle (z,S), (z',S') \rangle = s^T z' + \frac{1}{2} \operatorname{Tr}(SS').$$

Note that the quadratic lifting $z \mapsto (z, zz^T)$ maps the space \mathbb{R}^d into \mathcal{E}^d . In the sequel, an instrumental role is played by the following result.

Proposition 2.46.

(i) Assume we are given

- a nonempty and bounded subset U of \mathbf{R}^n ,
- a convex compact set \mathcal{V} contained in the interior of the cone \mathbf{S}^d_+ of positive semidefinite $d \times d$ matrices
- $a d \times (n+1)$ matrix A.

These data specify the family $\mathcal{G}_A[U, \mathcal{V}]$ of distributions of quadratic lifts $(\zeta, \zeta\zeta^T)$ of Gaussian random vectors $\zeta \sim \mathcal{N}(A[u; 1], \Theta)$ stemming from $u \in U$ and $\Theta \in \mathcal{V}$. Let us select somehow

- 1. $\gamma \in (0,1)$,
- 2. convex compact subset \mathcal{Z} of the set $\mathcal{Z}^n = \{Z \in \mathbf{S}^{n+1} : Z \succeq 0, Z_{n+1,n+1} = 1\}$ such that

$$Z(u) := [u; 1][u; 1]^T \in \mathcal{Z} \ \forall u \in U,$$
(2.164)

3. positive definite $d \times d$ matrix $\Theta_* \in \mathbf{S}^d_+$ and $\delta \in [0,2]$ such that

$$\Theta_* \succeq \Theta \ \forall \Theta \in \mathcal{V} \ \& \ \|\Theta^{1/2} \Theta_*^{-1/2} - I_d\| \le \delta \ \forall \Theta \in \mathcal{V}, \tag{2.165}$$

where $\|\cdot\|$ is the spectral norm,³²

and set

$$\mathcal{H} = \mathcal{H}^{\gamma} := \{ (h, H) \in \mathbf{R}^d \times \mathbf{S}^d : -\gamma \Theta_*^{-1} \preceq H \preceq \gamma \Theta_*^{-1} \},$$
(2.166)

$$\Phi_{A,\mathcal{Z}}(h,H;\Theta) = -\frac{1}{2}\ln \operatorname{Det}(I - \Theta_{*}^{1/2}H\Theta_{*}^{1/2}) + \frac{1}{2}\operatorname{Tr}([\Theta - \Theta_{*}]H) \\
+ \frac{\delta(2+\delta)}{2(1-||\Theta_{*}^{1/2}H\Theta_{*}^{1/2}||)} ||\Theta_{*}^{1/2}H\Theta_{*}^{1/2}||_{F}^{2} \\
+ \frac{1}{2}\phi_{\mathcal{Z}}\left(B^{T}\left[\left[\frac{H}{h^{T}}\right]^{h}\right] + [H,h]^{T}[\Theta_{*}^{-1} - H]^{-1}[H,h]\right]B\right): \\
\mathcal{H} \times \mathcal{V} \to \mathbf{R},$$
(2.167)

where B is given by

$$B = \begin{bmatrix} A\\ [0,...,0,1] \end{bmatrix}, \qquad (2.168)$$

the function

$$\phi_{\mathcal{Z}}(Y) := \max_{Z \in \mathcal{Z}} \operatorname{Tr}(ZY) \tag{2.169}$$

is the support function of \mathcal{Z} , and $\|\cdot\|_F$ is the Frobenius norm.

Function $\Phi_{A,\mathcal{Z}}$ is continuous on its domain, convex in $(h, H) \in \mathcal{H}$ and concave in $\Theta \in \mathcal{V}$, so that $(\mathcal{H}, \mathcal{V}, \Phi_{A,\mathcal{Z}})$ is a regular data. Besides this,

(#) Whenever $u \in \mathbf{R}^n$ is such that $[u; 1][u; 1]^T \in \mathcal{Z}$ and $\Theta \in \mathcal{V}$, the Gaussian

³²It is easily seen that with $\delta = 2$, the second relation in (2.165) is satisfied for all Θ such that $0 \leq \Theta \leq \Theta_*$, so that the restriction $\delta \leq 2$ is w.l.o.g.

random vector $\zeta \sim \mathcal{N}(A[u; 1], \Theta)$ satisfies the relation

$$\forall (h,H) \in \mathcal{H}: \ln\left(\mathbf{E}_{\zeta \sim \mathcal{N}(A[u;1],\Theta)}\left\{e^{\frac{1}{2}\zeta^{T}H\zeta+h^{T}\zeta}\right\}\right) \leq \Phi_{A,\mathcal{Z}}(h,H;\Theta).$$
(2.170)

which combines with (2.164) to imply that

$$\mathcal{G}_A[U,\mathcal{V}] \subset \mathcal{S}[\mathcal{H},\mathcal{V},\Phi_{A,\mathcal{Z}}]. \tag{2.171}$$

In addition, $\Phi_{A,\mathcal{Z}}$ is coercive in (h, H): $\Phi_{A,\mathcal{Z}}(h_i, H_i; \Theta) \to +\infty$ as $i \to \infty$ whenever $\Theta \in \mathcal{V}$, $(h_i, H_i) \in \mathcal{H}$ and $||(h_i, H_i)|| \to \infty$, $i \to \infty$.

(ii) Let two collections of entities from (i), $(\mathcal{V}_{\chi}, \Theta_*^{(\chi)}, \delta_{\chi}, \gamma_{\chi}, A_{\chi}, \mathcal{Z}_{\chi}), \chi = 1, 2,$ with common d be given, giving rise to the sets \mathcal{H}_{χ} , matrices B_{χ} , and functions $\Phi_{A_{\chi}, \mathcal{Z}_{\chi}}(h, H; \Theta), \chi = 1, 2$. These collections specify the families of normal distributions

$$\mathcal{G}_{\chi} = \{ \mathcal{N}(v, \Theta) : \Theta \in \mathcal{V}_{\chi} \& \exists u \in U : v = A_{\chi}[u; 1] \}, \ \chi = 1, 2.$$

Consider the convex-concave saddle point problem

$$SV = \min_{(h,H)\in\mathcal{H}_{1}\cap\mathcal{H}_{2}} \max_{\Theta_{1}\in\mathcal{V}_{1},\Theta_{2}\in\mathcal{V}_{2}} \underbrace{\frac{1}{2} \left[\Phi_{A_{1},\mathcal{Z}_{1}}(-h,-H;\Theta_{1}) + \Phi_{A_{2},\mathcal{Z}_{2}}(h,H;\Theta_{2}) \right]}_{\Phi(h,H;\Theta_{1},\Theta_{2})}.$$
(2.172)

A saddle point $(H_*, h_*; \Theta_1^*, \Theta_2^*)$ in this problem does exist, and the induced quadratic detector

$$\phi_*(\omega) = \frac{1}{2}\omega^T H_*\omega + h_*^T \omega + \underbrace{\frac{1}{2} \left[\Phi_{A_1, \mathcal{Z}_1}(-h_*, -H_*; \Theta_1^*) - \Phi_{A_2, \mathcal{Z}_2}(h_*, H_*; \Theta_2^*) \right]}_{a},$$
(2.173)

when applied to the families of Gaussian distributions \mathcal{G}_{χ} , $\chi = 1, 2$, has the risk

$$\operatorname{Risk}[\phi_*|\mathcal{G}_1, \mathcal{G}_2] \le \epsilon_\star := e^{\mathcal{S}V}$$

that is,

$$\begin{array}{ll} (a) & \int_{\mathbf{R}^d} e^{-\phi_*(\omega)} P(d\omega) \leq \epsilon_\star & \forall P \in \mathcal{G}_1, \\ (b) & \int_{\mathbf{R}^d} e^{\phi_*(\omega)} P(d\omega) \leq \epsilon_\star & \forall P \in \mathcal{G}_2. \end{array}$$

$$(2.174)$$

For proof, see Section 2.11.4.

Remark 2.47. Note that the computational effort to solve (2.172) reduces dramatically in the "*easy case*" of the situation described in item (ii) of Proposition 2.46, specifically, in the case where

- the observations are *direct*, meaning that $A_{\chi}[u; 1] \equiv u, u \in \mathbf{R}^d, \chi = 1, 2;$
- the sets \mathcal{V}_{χ} are comprised of positive definite *diagonal* matrices, and matrices $\Theta_*^{(\chi)}$ are diagonal as well, $\chi = 1, 2$;
- the sets \mathcal{Z}_{χ} , $\chi = 1, 2$, are convex compact sets of the form

$$\mathcal{Z}_{\chi} = \{ Z \in \mathbf{S}_{+}^{d+1} : Z \succeq 0, \operatorname{Tr}(ZQ_{j}^{\chi}) \le q_{j}^{\chi}, 1 \le j \le J_{\chi} \}$$

with diagonal matrices Q_j^{χ} , ³³ and these sets intersect the interior of the positive semidefinite cone \mathbf{S}_{+}^{d+1} .

In this case, the convex-concave saddle point problem (2.172) admits a saddle point $(h_*, H_*; \Theta_1^*, \Theta_2^*)$ where $h_* = 0$ and H_* is diagonal.

Justifying the remark. In the easy case, we have $B_{\chi} = I_{d+1}$ and therefore

$$M_{\chi}(h,H) := B_{\chi}^{T} \left[\left[\frac{H}{h^{T}} \right]^{-1} + [H,h]^{T} \left[[\Theta_{*}^{(\chi)}]^{-1} - H \right]^{-1} [H,h] \right] B_{\chi} \\ = \left[\frac{H + H \left[[\Theta_{*}^{(\chi)}]^{-1} - H \right]^{-1} H}{h^{T} + h^{T} \left[[\Theta_{*}^{(\chi)}]^{-1} - H \right]^{-1} H} \right]^{-1} H \left[h + H [[\Theta_{*}^{(\chi)}]^{-1} - H]^{-1} h \right]$$

and

$$\begin{split} \phi_{\mathcal{Z}_{\chi}}(Z) &= \max_{W} \left\{ \operatorname{Tr}(ZW) : W \succeq 0, \, \operatorname{Tr}(WQ_{j}^{\chi}) \leq q_{j}^{\chi}, \, 1 \leq j \leq J_{\chi} \right\} \\ &= \min_{\lambda} \left\{ \sum_{j} q_{j}^{\chi} \lambda_{j} : \, \lambda \geq 0, \, Z \preceq \sum_{j} \lambda_{j} Q_{j}^{\chi} \right\}, \end{split}$$

where the last equality is due to semidefinite duality³⁴. From the second representation of $\phi_{\mathcal{Z}_{\chi}}(\cdot)$ and the fact that all Q_j^{χ} are diagonal it follows that $\phi_{\mathcal{Z}_{\chi}}(M_{\chi}(0,H)) \leq \phi_{\mathcal{Z}_{\chi}}(M_{\chi}(h,H))$ (indeed, with diagonal Q_j^{χ} , if λ is feasible for the minimization problem participating in the representation when $Z = M_{\chi}(h,H)$, it clearly remains feasible when Z is replaced with $M_{\chi}(0,H)$). This, in turn, combines straightforwardly with (2.167) to imply that when replacing h_* with 0 in a saddle point $(h_*, H_*; \Theta_1^*, \Theta_2^*)$ of (2.172), we end up with another saddle point of (2.172). In other words, when solving (2.172), we can from the very beginning set h to 0, thus converting (2.172) into the convex-concave saddle point problem

$$SV = \min_{H:(0,H)\in\mathcal{H}_1\cap\mathcal{H}_2\,\Theta_1\in\mathcal{V}_1,\Theta_2\in\mathcal{V}_2} \Phi(0,H;\Theta_1,\Theta_2).$$
(2.175)

Taking into account the fact that we are in the case where all matrices from the sets \mathcal{V}_{χ} , same as the matrices $\Theta_*^{(\chi)}$ and all the matrices Q_j^{χ} , $\chi = 1, 2$, are diagonal, it is immediate to verify that if E is a $d \times d$ diagonal matrix with diagonal entries ± 1 , then $\Phi(0, H; \Theta_1, \Theta_2) = \Phi(0, EHE; \Theta_1, \Theta_2)$. Due to convexity-concavity of Φ this implies that (2.175) admits a saddle point $(0, H_*; \Theta_1^*, \Theta_2^*)$ with H_* invariant w.r.t. transformations $H_* \mapsto EH_*E$ with the above E, that is, with diagonal H_* , as claimed.

2.9.3 Quadratic lifting – does it help?

Assume that for $\chi = 1, 2$, we are given

- affine mappings $u \mapsto \mathcal{A}_{\chi}(u) = A_{\chi}[u;1] : \mathbf{R}^{n_{\chi}} \to \mathbf{R}^{d}$,
- nonempty convex compact sets $U_{\chi} \subset \mathbf{R}^{n_{\chi}}$,
- nonempty convex compact sets $\mathcal{V}_{\chi} \subset \operatorname{int} \mathbf{S}^d_+$.

 $^{^{33}\}text{In}$ terms of the sets $U_{\chi},$ this assumption means that the latter sets are given by linear inequalities on the squares of entries in u,

 $^{^{34}}$ see Section 4.1 (or [116, Section 7.1] for more details).

These data define families \mathcal{G}_{χ} of Gaussian distributions on \mathbf{R}^d : \mathcal{G}_{χ} is comprised of all distributions $\mathcal{N}(\mathcal{A}_{\chi}(u), \Theta)$ with $u \in U_{\chi}$ and $\Theta \in \mathcal{V}_{\chi}$. The data define also families \mathcal{SG}_{χ} of sub-Gaussian distributions on \mathbf{R}^d : \mathcal{SG}_{χ} is comprised of all sub-Gaussian distributions with parameters $(\mathcal{A}_{\chi}(u), \Theta)$ with $(u, \Theta) \in U_{\chi} \times \mathcal{V}_{\chi}$.

Assume we observe random variable $\zeta \in \mathbf{R}^d$ drawn from a distribution P known to belong to $\mathcal{G}_1 \cup \mathcal{G}_2$, and our goal is to decide from stationary K-repeated version of our observation on the pair of hypotheses $H_{\chi} : P \in \mathcal{G}_{\chi}, \chi = 1, 2$; we refer to this situation as to *Gaussian case*. We could also speak about *sub-Gaussian case*, where the hypotheses we would decide upon state that $P \in \mathcal{SG}_{\chi}$. In retrospect, all we are about to establish for the Gaussian case can be word by word repeated for the sub-Gaussian one, so that from now on, we assume we are in Gaussian case.

At present, we have developed two approaches to building detector-based tests for H_1, H_2 :

A. Utilizing the affine in ζ detector ϕ_{aff} given by solution to the saddle point problem (see (2.155), (2.156) and set $\theta_{\chi} = \mathcal{A}_{\chi}(u_{\chi})$ with u_{χ} running through U_{χ})

$$\operatorname{SadVal}_{\operatorname{aff}} = \min_{h \in \mathbf{R}^{d}} \max_{\substack{u_{1} \in U_{1}, u_{2} \in U_{2} \\ \Theta_{1} \in \mathcal{V}_{1}, \Theta_{2} \in \mathcal{V}_{2}}} \frac{1}{2} \left[h^{T} [\mathcal{A}_{2}(u_{2}) - \mathcal{A}_{1}(u_{1})] + \frac{1}{2} h^{T} [\Theta_{1} + \Theta_{2}] h \right];$$
(2.176)

this detector satisfies risk bound

$$\operatorname{Risk}[\phi_{\operatorname{aff}}|\mathcal{G}_1, \mathcal{G}_2] \le \exp\{\operatorname{SadVal}_{\operatorname{aff}}\}.$$
(2.177)

Q. Utilizing the quadratic in ζ detector ϕ_{lift} given by Proposition 2.46.ii, with the risk bound

$$\operatorname{Risk}[\phi_{\operatorname{lift}}|\mathcal{G}_1, \mathcal{G}_2] \le \exp\{\operatorname{SadVal}_{\operatorname{lift}}\},\tag{2.178}$$

with $\text{SadVal}_{\text{lift}}$ given by (2.172).

A natural question is, which one of these options results in a better risk bound. Note that we cannot just say "clearly, the second option is better, since there are more quadratic detectors than affine ones" – the difficulty is that the key, in the context of Proposition 2.46, relation (2.170) is inequality rather than equality³⁵. We are about to show that under reasonable assumptions, the second option indeed is better:

Proposition 2.48. In the situation in question, assume that the sets \mathcal{V}_{χ} , $\chi = 1, 2$, contain the \succeq -largest elements, and that these elements are taken as the matrices $\Theta_*^{(\chi)}$ participating in Proposition 2.46.ii. Let, further, the convex compact sets \mathcal{Z}_{χ} participating in Proposition 2.46.ii satisfy

$$\mathcal{Z}_{\chi} \subset \bar{\mathcal{Z}}_{\chi} := \{ Z = \begin{bmatrix} W & u \\ \hline u^T & 1 \end{bmatrix} \succeq 0, u \in U_{\chi} \}$$
(2.179)

(this assumption does not restrict generality, since \bar{Z}_{χ} is, along with U_{χ} , a closed

 $^{^{35}}$ One cannot make (2.170) an equality by redefining the right hand side function – it will lose the required in our context convexity-concavity properties.

ρ	σ_1	σ_2	$\begin{array}{c} \text{unrestricted} \\ H \text{ and } h \end{array}$	H = 0	h = 0
0.5	2	2	0.31	0.31	1.00
0.5	1	4	0.24	0.39	0.62
0.01	1	4	0.41	1.00	0.41

Table 2.2: Risk of quadratic detector $\phi(\zeta) = h^T \zeta + \frac{1}{2} \zeta^T H \zeta + \varkappa$

convex set which clearly contains all matrices $[u; 1][u; 1]^T$ with $u \in U_{\chi}$). Then

$$\operatorname{SadVal}_{lift} \leq \operatorname{SadVal}_{aff},$$
 (2.180)

that is, option Q is at least as efficient as option A.

Proof. Let $A_{\chi} = [\bar{A}_{\chi}, a_{\chi}]$. Looking at (2.155), where one should substitute $\theta_{\chi} = \mathcal{A}_{\chi}(u_{\chi})$ with u_{χ} running through U_{χ}) and taking into account that $\Theta_{\chi} \leq \Theta_{*}^{(\chi)} \in \mathcal{V}_{\chi}$ when $\Theta_{\chi} \in \mathcal{V}_{\chi}$, we conclude that

$$\operatorname{SadVal}_{\operatorname{aff}} = \min_{h} \max_{u_1 \in U_1, u_2 \in U_2} \frac{1}{2} \left[h^T [\bar{A}_2 u_2 - \bar{A}_1 u_1 + a_2 - a_1] + \frac{1}{2} h^T \left[\Theta_*^{(1)} + \Theta_*^{(2)} \right] h \right].$$
(2.181)

At the same time, we have by Proposition 2.46.ii:

where the concluding equality is due to (2.181).

Numerical illustration. To get an impression of the performance of quadratic detectors as compared to affine ones under the premise of Proposition 2.48, we present here the results of experiment where $U_1 = U_1^{\rho} = \{u \in \mathbf{R}^{12} : u_i \ge \rho, 1 \le i \le 12\}$, $U_2 = U_2^{\rho} = -U_1^{\rho}$, $A_1 = A_2 \in \mathbf{R}^{8 \times 13}$, and $\mathcal{V}_{\chi} = \{\Theta_*^{(\chi)} = \sigma_{\chi}^2 I_8\}$ are singletons. The risks of affine, quadratic and "purely quadratic" (with h set to 0) detectors on the associated families $\mathcal{G}_1, \mathcal{G}_2$ are given in Table 2.2.

We see that

• when deciding on families of Gaussian distributions with common covariance matrix and expectations varying in associated with the families convex sets, passing from affine detectors described by Proposition 2.44 to quadratic detectors

does not affect the risk (first row in the table). This should be expected: we are in the scope of Gaussian o.s., where minimum risk affine detectors are optimal among all possible detectors.

- when deciding on families of Gaussian distributions in the case where distributions from different families can have close expectations (third row in the table), affine detectors are useless, while the quadratic ones are not, provided that $\Theta_*^{(1)}$ differs from $\Theta_*^{(2)}$. This is how it should be – we are in the case where the first moments of the distribution of observation bear no definitive information on the family this distribution belongs to, which makes affine detectors useless. In contrast, quadratic detectors are able to utilize information (valuable when $\Theta_*^{(1)} \neq \Theta_*^{(2)}$) "stored" in the second moments of the observation.
- "in general" (second row in the table), both affine and purely quadratic components in a quadratic detector are useful; suppressing one of them can increase significantly the attainable risk.

2.9.4 Quadratic lifting: sub-Gaussian case

Sub-Gaussian version of Proposition 2.46 is as follows:

Proposition 2.49.

(i) Assume we are given

- a nonempty and bounded subset U of \mathbf{R}^n ,
- a convex compact set \mathcal{V} contained in the interior of the cone \mathbf{S}^d_+ of positive semidefinite $d \times d$ matrices
- $a d \times (n+1)$ matrix A.

These data specify the family $SG_A[U, \mathcal{V}]$ of distributions of quadratic lifts $(\zeta, \zeta\zeta^T)$ of sub-Gaussian random vectors ζ with sub-Gaussianity parameters $A[u; 1], \Theta$ stemming from $u \in U$ and $\Theta \in \mathcal{V}$.

 $Let \ us \ select \ somehow$

- 1. reals γ, γ^+ such that $0 < \gamma < \gamma^+ < 1$,
- 2. convex compact subset Z of the set $Z^n = \{Z \in \mathbf{S}^{n+1} : Z \succeq 0, Z_{n+1,n+1} = 1\}$ such that relation (2.164) takes place,
- 3. positive definite $d \times d$ matrix $\Theta_* \in \mathbf{S}^d_+$ and $\delta \in [0,2]$ such that (2.165) takes place.

These data specify the closed convex sets

$$\mathcal{H} = \mathcal{H}^{\gamma} := \{ (h, H) \in \mathbf{R}^{d} \times \mathbf{S}^{d} : -\gamma \Theta_{*}^{-1} \preceq H \preceq \gamma \Theta_{*}^{-1} \},$$

$$\widehat{\mathcal{H}} = \widehat{\mathcal{H}}^{\gamma, \gamma^{+}} = \left\{ (h, H, G) \in \mathbf{R}^{d} \times \mathbf{S}^{d} \times \mathbf{S}^{d} : \left\{ \begin{array}{c} -\gamma \Theta_{*}^{-1} \preceq H \preceq \gamma \Theta_{*}^{-1} \\ 0 \preceq G \preceq \gamma^{+} \Theta_{*}^{-1}, H \preceq G \end{array} \right\}$$

$$(2.182)$$

and the functions

$$\begin{split} \Psi_{A,\mathcal{Z}}(h,H,G) &= -\frac{1}{2}\ln\operatorname{Det}(I-\Theta_{*}^{1/2}G\Theta_{*}^{1/2}) \\ &+ \frac{1}{2}\phi_{\mathcal{Z}}\left(B^{T}\left[\left[\frac{H}{h^{T}}\right]^{h}\right] + [H,h]^{T}[\Theta_{*}^{-1}-G]^{-1}[H,h]\right]B\right): \\ \mathcal{H} \times \mathcal{Z} \to \mathbf{R}, \\ \Psi_{A,\mathcal{Z}}^{\delta}(h,H,G;\Theta) &= -\frac{1}{2}\ln\operatorname{Det}(I-\Theta_{*}^{1/2}G\Theta_{*}^{1/2}) \\ &+ \frac{1}{2}\operatorname{Tr}([\Theta-\Theta_{*}]G) + \frac{\delta(2+\delta)}{2(1-||\Theta_{*}^{1/2}G\Theta_{*}^{1/2}||)}||\Theta_{*}^{1/2}G\Theta_{*}^{1/2}||_{F}^{2} \\ &+ \frac{1}{2}\phi_{\mathcal{Z}}\left(B^{T}\left[\left[\frac{H}{h^{T}}\right]^{h}\right] + [H,h]^{T}[\Theta_{*}^{-1}-G]^{-1}[H,h]\right]B\right): \\ &\hat{\mathcal{H}} \times \{0 \leq \Theta \leq \Theta_{*}\} \to \mathbf{R}, \\ \Phi_{A,\mathcal{Z}}(h,H) &= \min_{G}\left\{\Psi_{A,\mathcal{Z}}(h,H,G):(h,H,G)\in\hat{\mathcal{H}}\right\}: \mathcal{H} \times \{0 \leq \Theta \leq \Theta_{*}\} \to \mathbf{R}, \\ \Phi_{A,\mathcal{Z}}^{\delta}(h,H;\Theta) &= \min_{G}\left\{\Psi_{A,\mathcal{Z}}^{\delta}(h,H,G;\Theta):(h,H,G)\in\hat{\mathcal{H}}\right\}: \mathcal{H} \times \{0 \leq \Theta \leq \Theta_{*}\} \to \mathbf{R}, \\ \end{split}$$

where B is given by (2.168) and $\phi_{\mathcal{Z}}(\cdot)$ is the support function of \mathcal{Z} given by (2.169).

Function $\Phi_{A,\mathcal{Z}}(h,H)$ is convex and continuous on its domain, while function $\Phi_{A,\mathcal{Z}}^{\delta}(h,H;\Theta)$ is continuous on its domain, convex in $(h,H) \in \mathcal{H}$ and concave in $\Theta \in \{0 \leq \Theta \leq \Theta_*\}$. Besides this,

(##) Whenever $u \in \mathbf{R}^n$ is such that $[u; 1][u; 1]^T \in \mathbb{Z}$ and $\Theta \in \mathcal{V}$, the sub-Gaussian, with parameters $(A[u; 1], \Theta)$, random vector ζ satisfies the relation

$$\begin{aligned} \forall (h,H) \in \mathcal{H}: \\ (a) \quad \ln\left(\mathbf{E}_{\zeta}\left\{e^{\frac{1}{2}\zeta^{T}H\zeta+h^{T}\zeta}\right\}\right) \leq \Phi_{A,\mathcal{Z}}(h,H), \\ (b) \quad \ln\left(\mathbf{E}_{\zeta}\left\{e^{\frac{1}{2}\zeta^{T}H\zeta+h^{T}\zeta}\right\}\right) \leq \Phi_{A,\mathcal{Z}}^{\delta}(h,H;\Theta). \end{aligned}$$
(2.184)

which combines with (2.164) to imply that

$$\mathcal{S}G_A[U,\mathcal{V}] \subset \mathcal{S}[\mathcal{H},\mathcal{V},\Phi_{A,\mathcal{S}}].$$
 (2.185)

In addition, $\Phi_{A,\mathcal{Z}}$ and $\Phi_{A,\mathcal{Z}}^{\delta}$ are coercive in (h, H): $\Phi_{A,\mathcal{Z}}(h_i, H_i) \to +\infty$ and $\Phi_{A,\mathcal{Z}}^{\delta}(h_i, H_i; \Theta) \to +\infty$ as $i \to \infty$ whenever $\Theta \in \mathcal{V}$, $(h_i, H_i) \in \mathcal{H}$ and $||(h_i, H_i)|| \to \infty$, $i \to \infty$.

(ii) Let two collections of data from (i): $(\mathcal{V}_{\chi}, \Theta_*^{(\chi)}, \delta_{\chi}, \gamma_{\chi}, \gamma_{\chi}^+, A_{\chi}, \mathcal{Z}_{\chi}), \chi = 1, 2$, with common d be given, giving rise to the sets \mathcal{H}_{χ} , matrices B_{χ} , and functions $\Phi_{A_{\chi}, \mathcal{Z}_{\chi}}(h, H), \Phi_{A_{\chi}, \mathcal{Z}_{\chi}}(h, H; \Theta), \chi = 1, 2$. These collections specify the families $SG_{\chi} = SG_{A_{\chi}}[U_{\chi}, \mathcal{V}_{\chi}]$ of sub-Gaussian distributions.

Consider the convex-concave saddle point problem

$$SV = \min_{(h,H)\in\mathcal{H}_{1}\cap\mathcal{H}_{2}} \max_{\Theta_{1}\in\mathcal{V}_{1},\Theta_{2}\in\mathcal{V}_{2}} \underbrace{\frac{1}{2} \left[\Phi_{A_{1},\mathcal{Z}_{1}}^{\delta_{1}}(-h,-H;\Theta_{1}) + \Phi_{A_{2},\mathcal{Z}_{2}}^{\delta_{2}}(h,H;\Theta_{2}) \right]}_{\Phi^{\delta_{1},\delta_{2}}(h,H;\Theta_{1},\Theta_{2})}.$$
(2.186)

A saddle point $(H_*, h_*; \Theta_1^*, \Theta_2^*)$ in this problem does exist, and the induced quadratic detector

$$\phi_*(\omega) = \frac{1}{2}\omega^T H_*\omega + h_*^T \omega + \underbrace{\frac{1}{2} \left[\Phi_{A_1, \mathcal{Z}_1}^{\delta_1}(-h_*, -H_*; \Theta_1^*) - \Phi_{A_2, \mathcal{Z}_2}^{\delta_2}(h_*, H_*; \Theta_2^*) \right]}_{a},$$
(2.187)

when applied to the families of sub-Gaussian distributions SG_{χ} , $\chi = 1, 2$, has the risk

$$\operatorname{Risk}[\phi_*|\mathcal{S}G_1, \mathcal{S}G_2] \le \epsilon_\star := \mathrm{e}^{\mathcal{S}V}.$$

As a result,

(a)
$$\int_{\mathbf{R}^{d}} e^{-\phi_{*}(\omega)} P(d\omega) \leq \epsilon_{\star} \quad \forall P \in SG_{1},$$

(b)
$$\int_{\mathbf{R}^{d}} e^{\phi_{*}(\omega)} P(d\omega) \leq \epsilon_{\star} \quad \forall P \in SG_{2}.$$
 (2.188)

Similarly, the convex minimization problem

$$Opt = \min_{(h,H)\in\mathcal{H}_1\cap\mathcal{H}_2} \underbrace{\frac{1}{2} \left[\Phi_{A_1,\mathcal{Z}_1}(-h,-H) + \Phi_{A_2,\mathcal{Z}_2}(h,H) \right]}_{\Phi(h,H)}.$$
 (2.189)

is solvable, and the induced by its optimal solution (h_*, H_*) quadratic detector

$$\phi_*(\omega) = \frac{1}{2}\omega^T H_*\omega + h_*^T \omega + \underbrace{\frac{1}{2} \left[\Phi_{A_1, \mathcal{Z}_1}(-h_*, -H_*) - \Phi_{A_2, \mathcal{Z}_2}(h_*, H_*)\right]}_{a}, \quad (2.190)$$

when applied to the families of sub-Gaussian distributions SG_{χ} , $\chi = 1, 2$, has the risk

$$\operatorname{Risk}[\phi_*|\mathcal{S}G_1, \mathcal{S}G_2] \le \epsilon_\star := e^{Opt}$$

so that for just defined ϕ_* and ϵ_* relation (2.189) takes place.

For proof, see Section 2.11.5.

Remark 2.50. Proposition 2.49 offers two options for building quadratic detectors for the families SG_1 , SG_2 , those based on saddle point of (2.186) and on optimal solution to (2.189). Inspecting the proof, the number of options can be increased to 4: we can replace any one of the functions $\Phi_{A_{\chi},Z_{\chi}}^{\delta_{\chi}}$, $\chi = 1, 2$ (or both these functions simultaneously) with $\Phi_{A_{\chi},Z_{\chi}}$. The second of the original two options is exactly what we get when replacing both $\Phi_{A_{\chi},Z_{\chi}}^{\delta_{\chi}}$, $\chi = 1, 2$, with $\Phi_{A_{\chi},Z_{\chi}}$. It is easily seen that depending on the data, every one of these 4 options can be the best – result in the smallest risk bound. Thus, it makes sense to keep all these options in mind and to use the one which, under the circumstances, results in the best risk bound. Note that the risk bounds are efficiently computable, so that identifying the best option is easy.

2.9.5 Recovering quadratic form of discrete distribution

Lemma 2.51. Let ζ be zero mean random variable taking values in the n-dimensional ℓ_1 -ball of radius 2 centered at the origin, K be a positive integer, and let $\eta_K = \frac{1}{K} \sum_{k=1}^{K} \zeta_k$, where $\zeta_1, ..., \zeta_K$ are independent copies of ζ . Then

$$0 < \gamma < K/4 \Rightarrow \mathbf{E}\left\{\exp\{\frac{\gamma\eta_K^T\eta_K}{2}\}\right\} \le \frac{n}{\sqrt{1 - 4\gamma/K}}.$$
(2.191)

In particular,

$$\mathbf{E}\left\{\exp\{\frac{K\eta_K^T\eta_K}{9}\}\right\} \le 3n.$$
(2.192)

Proof. By Proposition 2.40 we have

$$\mathbf{E}\left\{\exp\{h^{T}\eta_{K}\}\right\} = \left[\mathbf{E}\left\{\exp\{h^{T}\zeta/K\}\right\}\right]^{K} \le \left[\exp\{\|h/K\|_{\infty}^{2}/2\}\right]^{K} = \exp\{\frac{2\|h\|_{\infty}^{2}}{K}\}.$$

Now let $z \sim \mathcal{N}(0, I_n)$ be independent of η_K . We have

$$\begin{aligned} \mathbf{E}_{\eta_{K}}\left\{\exp\{\gamma\eta_{K}^{T}\eta_{K}/2\}\right\} &= \mathbf{E}_{\eta_{K}}\left\{\mathbf{E}_{z}\left\{\exp\{\sqrt{\gamma}\eta_{K}^{T}z\}\right\}\right\} = \mathbf{E}_{z}\left\{\mathbf{E}_{\eta_{K}}\left\{\exp\{\sqrt{\gamma}\eta_{K}^{T}z\}\right\}\right\} \\ &\leq \mathbf{E}_{z}\left\{\exp\{2\gamma\|z\|_{\infty}^{2}/K)\}\right\} \leq \mathbf{E}_{z}\left\{\sum_{i=1}^{n}\exp\{(4\gamma/K)z_{i}^{2}/2\}\right\} \\ &= \frac{n}{\sqrt{2\pi}}\int\exp\{-\frac{[1-(4\gamma/K)]s^{2}}{2}\}ds = \frac{n}{\sqrt{1-4\gamma/K}} \end{aligned}$$

Corollary 2.52. In the notation and under the premise of Lemma 2.51, for every $p \in [1,2]$ and $\epsilon \in (0,1)$ one has

$$\operatorname{Prob}\left\{\|\eta_K\|_p > \frac{3n^{\frac{2-p}{2p}}\sqrt{\ln(3n/\epsilon)}}{\sqrt{K}}\right\} < \epsilon$$

Proof. Let $L = \frac{9}{K} \ln(3n/\epsilon)$, and $\Xi = \{\eta_K : \eta_K^T \eta_K \leq L\}$, so that $\operatorname{Prob}\{\eta_k \notin \Xi\} < \epsilon$ by (2.192). When $\eta_K \in \Xi$, we have

$$\|\eta_K\|_p \le n^{\frac{1}{p} - \frac{1}{2}} \|\eta_K\|_2 \le \frac{3n^{\frac{d-p}{2p}} \sqrt{\ln(3n/\epsilon)}}{\sqrt{K}}.$$

Application: recovering quadratic form of discrete distribution. Consider the situation as follows: we are given an i.i.d. sample

$$\omega^K = (\omega_1, ..., \omega_K)$$

with $\omega_i \sim Au$, where $u \in \Delta_n$ is an unknown probabilistic vector, and A is $m \times n$ stochastic matrix, so that ω_k takes value e_i $(e_1, ..., e_m$ are basic orths in \mathbf{R}^m) with probability $[Au]_i$. Our goal is to recover

$$F(u) = u^T Q u,$$

where $Q \in \mathbf{S}^n$ is given, and to this end we want to build a "presumably good" estimate of the form

$$\widehat{g}_{H}(\omega^{K}) = \left[\frac{1}{K}\sum_{k}\omega_{k}\right]^{T}H\left[\frac{1}{K}\sum_{k}\omega_{k}\right].$$

Let us set

$$p = Au, \zeta_k = \omega_k - p, \eta_K = \frac{1}{K} \sum_k \zeta_k, \xi_K = \frac{1}{K} \sum_k \omega_k = \eta_K + p.$$

Note that $\|\zeta_k\|_1 \leq 2$. We have

$$\begin{aligned} \widehat{g}_{H}(\omega^{K}) - F(u) &= \xi_{K}^{T} H \xi_{K} - u^{T} Q u = \xi_{K}^{T} H \eta_{K} + \xi_{K}^{T} H p \\ &= \xi_{K}^{T} H \eta_{K} + \eta_{K}^{T} H p + [p^{T} H p - u^{T} Q u] \\ \Rightarrow |\widehat{g}_{H}(\omega^{K} - u^{T} Q u)| \leq |p^{T} H p - u^{T} Q u| + 2 ||H||_{\infty} ||\eta_{K}||_{1}, \end{aligned}$$

where $||H||_{\infty} = \max_{i,j} |H_{ij}|$. Invoking Corollary 2.52, it follows that for all $\epsilon \in (0, 1)$ and all $u \in \Delta_n$ it holds

$$\operatorname{Prob}\left\{ \left| \widehat{g}_{H}(\omega^{K}) - F(u) \right| \geq \mathcal{R}[H] := \frac{6\sqrt{m\ln(3m/\epsilon)}}{\sqrt{K}} \|H\|_{\infty} + \|A^{T}HA - Q\|_{\infty} \right\} \leq \epsilon.$$

and we can build a "presumably good" estimate by minimizing $\mathcal{R}[H]$ over H.

2.9.6 Generic application: quadratically constrained hypotheses

Propositions 2.46, 2.49 operate with Gaussian/sub-Gaussian observations ζ with matrix parameters Θ running through convex compact subsets \mathcal{V} of int \mathbf{S}_{+}^{d} , and means of the form A[u; 1], with "signals" u running through given sets $U \subset \mathbf{R}^{n}$. The constructions, however, involved additional entities – convex compact sets $\mathcal{Z} \subset \mathcal{Z}^{n} := \{Z \in \mathbf{S}_{+}^{n+1} : Z_{n+1,n+1} = 1\}$ containing quadratic lifts $[u; 1][u; 1]^{T}$ of all signals $u \in U$; other things being equal, the smaller is \mathcal{Z} , the smaller is the associated function $\Phi_{A,\mathcal{Z}}$ (or $\Phi_{A,\mathcal{Z}}^{\delta}$), and consequently, the smaller are the (upper bounds on the) risks of quadratic in ζ detectors we end up with. In order to apply these propositions, we should understand how to build the required sets \mathcal{Z} in an "economical" way. There exists a relatively simple case when it is easy to get reasonable candidates to the role of \mathcal{Z} – the case of quadratically constrained signal set U:

$$U = \{ u \in \mathbf{R}^n : f_k(u) := u^T Q_k u + 2q_k^T u \le b_k, \ 1 \le k \le K \}.$$
 (2.193)

Indeed, the constraints $f_k(u) \leq b_k$ are just linear constraints on the quadratic lifting $[u;1][u;1]^T$ of u:

$$u^T Q_k u + 2q_k^T u \le b_k \Leftrightarrow \operatorname{Tr}(F_k[u;1][u;1]^T) \le b_k, \ F_k = \left[\frac{Q_k | q_k}{q_k^T}\right] \in \mathbf{S}^{n+1}.$$

Consequently, in the case of (2.193), the simplest candidate on the role of \mathcal{Z} is the set

$$\mathcal{Z} = \{ Z \in \mathbf{S}^n : Z \succeq 0, Z_{n+1,n+1} = 1, \operatorname{Tr}(F_k Z) \le b_k, \ 1 \le k \le K \}.$$
(2.194)

This set clearly is closed and convex (the latter – even when U itself is not convex), and indeed contains the quadratic lifts $[u; 1][u; 1]^T$ of all points $u \in U$. We need also the compactness of \mathcal{Z} ; the latter definitely takes place when the quadratic constraints describing U contain constraint of the form $u^T u \leq R^2$, which, in turn, can be ensured, basically "for free," when U is bounded. It should be stressed that the "ideal" choice of \mathcal{Z} would be the convex hull $\mathcal{Z}[U]$ of all rank 1 matrices $[u; 1][u; 1]^T$ with $u \in U$ – this definitely is the smallest convex set which contains the quadratic lifts of all points from U; moreover, $\mathcal{Z}[U]$ is closed and bounded, provided U is so. The difficulty is that $\mathcal{Z}[U]$ can be computationally intractable (and thus

useless in our context) already for pretty simple sets U of the form (2.193). The set (2.194) is a simple outer approximation of $\mathcal{Z}[U]$, and this approximation can be very loose; for example, when $U = \{u : -1 \leq u_k \leq 1, 1 \leq k \leq d\}$ is just the unit box in \mathbb{R}^d , the set (2.194) is

$$\{Z \in \mathbf{S}^{n+1} : Z \succeq 0, Z_{n+1,n+1} = 1, |Z_{k,n+1}| \le 1, 1 \le k \le n\};$$

this set even is not bounded, while $\mathcal{Z}[U]$ clearly is bounded. There is, essentially, just one generic case when the set (2.194) is *exactly equal* to $\mathcal{Z}[U]$ – the case where

$$U = \{u : u^T Q u \le c\}, Q \succ 0$$

is an ellipsoid centered at the origin; the fact that in this case the set given by (2.194) is *exactly* $\mathcal{Z}[U]$ is a consequence of what is called *S*-Lemma.

The fact that, in general, the set \mathcal{Z} could be a very loose outer approximation of $\mathcal{Z}[U]$ does not mean that we cannot improve this construction. As an instructive example, let $U = \{u \in \mathbf{R}^n : ||u||_{\infty} \leq 1\}$. We get a much better that above approximation of $\mathcal{Z}[U]$ when applying (2.194) to equivalent description of the box by quadratic constraints:

$$U := \{ u \in \mathbf{R}^n : \|u\|_{\infty} \le 1 \} = \{ u \in \mathbf{R}^n : u_k^2 \le 1, 1 \le k \le n \}.$$

Applying recipe (2.194) to the second description of U, we arrive at a significantly less conservative outer approximation of $\mathcal{Z}[U]$, specifically,

$$\mathcal{Z} = \{ Z \in \mathbf{S}^{n+1} : Z \succeq 0, Z_{n+1,n+1} = 1, Z_{kk} \le 1, 1 \le k \le n \}.$$

Not only the resulting set \mathcal{Z} is bounded; we can get a reasonable "upper bound" on the discrepancy between \mathcal{Z} and $\mathcal{Z}[U]$. Namely, denoting by Z^o the matrix obtained from a symmetric $n \times n$ matrix Z by zeroing out the South-Eastern entry (the one in the cell (n + 1, n + 1)) and keeping the remaining entries intact, we have

$$\mathcal{Z}^{o}[U] := \{Z^{o} : Z \in \mathcal{Z}[U]\} \subset \mathcal{Z}^{o} := \{Z^{o} : Z \in \mathcal{Z}\} \subset O(1)\ln(n+1)\mathcal{Z}^{o}.$$

This is a particular case of a general result (going back to [119]; we shall get this result as a byproduct of our forthcoming considerations, specifically, Proposition 4.6) as follows:

Let U be a bounded set given by a system of convex quadratic constraints without linear terms:

$$U = \{ u \in \mathbf{R}^n : u^T Q_k u \le c_k, \ 1 \le k \le K \}, \ Q_k \succeq 0, \ 1 \le k \le K,$$

and let \mathcal{Z} be the associated set (2.194):

$$\mathcal{Z} = \{ Z \in \mathbf{S}^{n+1} : Z \succeq 0, Z_{n+1,n+1} = 1, \operatorname{Tr}(Z\operatorname{Diag}\{Q_k, 1\}) \le c_k, 1 \le k \le K \}$$

Then

$$\mathcal{Z}^{o}[U] := \{ Z^{o} : Z \in \mathcal{Z}[U] \} \subset \mathcal{Z}^{o} := \{ Z^{o} : Z \in \mathcal{Z} \} \subset 4\ln(5(K+1))\mathcal{Z}^{o}[U].$$

Note that when K = 1 (i.e., U is an ellipsoid centered at the origin), the factor $4\ln(5(K+1))$, as it was already mentioned, can be replaced by 1.

Finally, we remark that, same as in the case of a box, we can try to reduce the conservatism of outer approximation (2.194) of $\mathcal{Z}[U]$ by passing from the initial description of U to an equivalent one. The standard recipe here is to replace linear constraints in the description of U by their quadratic consequences; for example, we can augment a pair of linear constraints $q_i^T u \leq c_i$, $q_j^T u \leq c_j$, assuming there is such a pair, with the quadratic constraint $(c_i - q_i^T u)(c_j - q_j^T u) \geq 0$. While this constraint is redundant, as far as the description of U itself is concerned, adding this constraint reduces, and sometimes significantly, the set given by (2.194). Informally speaking, transition from (2.193) to (2.194) is by itself "too stupid" to utilize the fact (known to every kid) that the product of two nonnegative quantities is nonnegative; when augmenting linear constraints in the description of U by their pairwise products, we somehow compensate for this stupidity. Unfortunately, "computationally tractable" assistance of this type perhaps allows to reduce the conservatism of (2.194), but usually does not allow to eliminate it completely: a grave "fact of life" is that even in the case of unit box U, the set $\mathcal{Z}[U]$ is computationally intractable. Scientifically speaking: maximizing quadratic forms over the unit box U provably is an NP-hard problem; were we able to get a computationally tractable description of $\mathcal{Z}[U]$, we would be able to solve this NP-hard problem efficiently, implying that P=NP. While we do not know for sure that the latter is not the case, "the informal odds" are strongly against this possibility.

The bottom line is that while the approach we are discussing in *some* situations could result in quite conservative tests, "some" is by far not the same as "always;" on the positive side, this approach allows to process some important problems. We are about to present a simple and instructive illustration.

2.9.6.1 Simple change detection

On Figure 2.5, you see a sample of frames from a "movie" where noisy picture of a gentleman gradually transforms into noisy picture of a lady; several initial frames differ just by realizations of noise, and starting with some instant, the "signal" (the deterministic component of the image) starts to drift from the gentleman towards the lady. What, in your opinion, is the change point – the first time instant where the signal component of the image differs from the signal component of the initial image?

A simple model of the situation is as follows: we observe, one by one, vectors (in fact, 2D arrays, but we can "vectorize" them)

$$\omega_t = x_t + \xi_t, \ t = 1, 2, \dots, K,\tag{2.195}$$

where x_t are deterministic components of the observations and ξ_t are random noises. It may happen that for some $\tau \in \{2, 3, ..., K\}$, the vectors x_t are independent of t when $t < \tau$, and x_{τ} differs from $x_{\tau-1}$ (" τ is a change point"); if it is the case, τ is uniquely defined by $x^K = x_1, ..., x_K$. An alternative is that x_t is independent of $t, 1 \leq t \leq K$ ("no change"). The goal is to decide, based on observation $\omega^K = (\omega_1, ..., \omega_K)$, whether there was a change point, and if yes, then, perhaps, to localize it.

The model we have just described is the simplest case of "change detection,"



Figure 2.5: Frames from a "movie"

where, given noisy observations on some time horizon, one is interested to detect a "change" in some time series underlying the observations. In our simple model, this time series is comprised by deterministic components x_t of observations, and "change at time τ " is understood in the most straightforward way - as the fact that x_{τ} differs from equal to each other preceding x_t 's. In more complicated situations, our observations are obtained from the underlying time series $\{x_t\}$ by a non-anticipative transformation, like

$$\omega_t = \sum_{s=1}^t A_{ts} x_s + \xi_t, \ t = 1, ..., K,$$

and we still want to detect the change, if any, in the time series $\{x_t\}$. As an instructive example, consider observations, taken along equidistant time grid, of the positions of an aircraft which "normally" flies with constant velocity, but at some time instant can start to maneuver. In this situation, the underlying time series is comprised of the velocities of the aircraft at consecutive time instants, observations are obtained from this time series by integration, and to detect a maneuver means to detect that on the observation horizon, there was a change in the series of velocities.

Change detection is the subject of huge literature dealing with a wide range of models differing from each other in

- whether we deal with direct observations of the time series of interest, as in (2.195), or with indirect ones (in the latter case, there is a wide spectrum of options related to how the observations depend on the underlying time series),
- what are assumptions on noise,
- what happens with x_t 's after the change takes place do they jump from their common value prior to time τ to a new common value starting with this time, or start to depend on time (and if yes, then how), etc., etc.

A significant role in change detection is played by hypothesis testing; as far as affine/quadratic-detector-based techniques developed in this Section are concerned, their applications in the context of change detection are discussed in [71]. In what follows, we focus on the simplest of these applications.

Situation and goal. We consider the situation as follows:

- 1. Our observations are given by (2.195) with independent across t = 1, ..., K noises $\xi_t \sim \mathcal{N}(0, \sigma^2 I_d)$. We do not known σ a priori, what we know is that σ is independent of t and belongs to a given segment $[\underline{\sigma}, \overline{\sigma}]$, with $0 < \underline{\sigma} \leq \overline{\sigma}$;
- 2. Observations (2.195) arrive one by one, so that at time $t, 2 \le t \le K$ we have at our disposal observation $\omega^t = \omega_1, ..., \omega_t$. Our goal is to build a system of inferences $\mathcal{T}_t, 2 \le t \le K$, such that \mathcal{T}_t as applied to ω^t either infers that there was a change at time t or earlier, in which case we terminate, or infers that so far there was no change, in which case we either proceed to time t+1 (if t < K), or terminate (if t = K) with "no change" conclusion.

We are given $\epsilon \in (0, 1)$ and want from our collection of inferences to make the probability of *false alarm* (i.e., terminating somewhere on time horizon 2, 3, ..., K with "there was a change" conclusion in the situation when there was no change: $x_1 = \ldots = x_K$) at most ϵ . Under this restriction, we want to make as small as possible the probability of a *miss* (of not detecting the change at all in the

situation where there was a change).

The "small probability of a miss" desire should be clarified. When the noise is nontrivial, we have no chances to detect very small changes *and* respect the bound on the probability of false alarm. A realistic goal is to make as small as possible the probability of missing a *not too small* change, which can be formalized as follows. Given $\rho > 0$, and tolerances $\epsilon, \epsilon \in (0, 1)$, let us look for a system of inferences $\{\mathcal{T}_t : 2 \leq t \leq K\}$ such that

- the probability of false alarm is at most ϵ , and
- the probability of " ρ -miss" the probability to detect no change when there was a change of energy $\geq \rho^2$ (i.e., when there was a change point τ , and, moreover, at this point it holds $||x_{\tau} x_1||_2^2 \geq \rho^2$) is at most ε .

What we are interested in, is to achieve the just formulated goal with as small ρ as possible.

Construction. Let us select a large "safety parameter" R, like $R = 10^8$ or even $R = 10^{80}$, so that we can assume that for all time series we are interested in it holds $||x_t - x_\tau||_2^2 \leq R^{2-36}$. Let us associate with $\rho > 0$ "signal hypotheses" H_t^{ρ} , t = 2, 3, ..., K, on the distribution of observation ω^K given by (2.195), with H_t^{ρ} stating that in the time series $\{x_t\}_{t=1}^K$ underlying observation ω^K there is a change, of energy at least ρ^2 , at time t:

$$x_1 = x_2 = \dots = x_{t-1} \& ||x_t - x_{t-1}||_2^2 = ||x_t - x_1||_2^2 \ge \rho^2$$

(and on the top of it, $||x_t - x_\tau||_2^2 \leq R^2$ for all t, τ). Let us augment these hypotheses by the null hypothesis H_0 stating that there is no change at all – the observation ω^K stems from a stationary time series $x_1 = x_2 = \ldots = x_K$. We are about to use our machinery of detector-based tests in order to build a system of tests deciding, with partial risks ϵ, ε , on the null hypothesis vs. the "signal alternative" $\bigcup_t H_t^{\rho}$ for as small ρ as possible.

The implementation is as follows. Given $\rho > 0$ such that $\rho^2 < R^2$, consider two hypotheses, G_1 and G_2^{ρ} , on the distribution of observation

$$\zeta = x + \xi \in \mathbf{R}^d. \tag{2.196}$$

Both hypotheses state that $\xi \sim \mathcal{N}(0, \sigma^2 I_d)$ with unknown σ known to belong to a given segment $\Delta := [\sqrt{2}\sigma, \sqrt{2}\overline{\sigma}]$. In addition, G_1 states that x = 0, and G_2^{ρ} - that $\rho^2 \leq ||x||_2^2 \leq R^2$. We can use the result of Proposition 2.46.ii to build a quadratic in ζ detector for the families of distributions \mathcal{P}_1 , \mathcal{P}_2^{ρ} obeying the hypotheses G_1, G_2^{ρ} , respectively. To this end it suffices to apply Proposition to the collections of data

$$\mathcal{V}_{\chi} = \{\sigma^2 I_d : \sigma \in \Delta\}, \Theta_*^{(\chi)} = 2\overline{\sigma}^2 I_d, \delta_{\chi} = 1 - \underline{\sigma}/\overline{\sigma}, \gamma_{\chi} = 0.999, A_{\chi} = I_d, \mathcal{Z}_{\chi}, [\chi = 1, 2]$$

 $^{^{36}}R$ is needed by the only reason – to make the domains we are working with bounded, thus allowing to apply the theory we have developed so far. The actual value of R does *not* enter our constructions and conclusions.

where

$$\begin{aligned} \mathcal{Z}_1 &= \{ [0; ...; 0; 1] [0; ...; 0; 1]^T \} \subset \mathbf{S}_+^{d+1}, \\ \mathcal{Z}_2 &= \mathcal{Z}_2^{\rho} = \{ Z \in \mathbf{S}_+^{d+1} : Z_{d+1, d+1} = 1, 1 + R^2 \ge \operatorname{Tr}(Z) \ge 1 + \rho^2 \}. \end{aligned}$$

The (upper bound on the) risk of the quadratic in ζ detector yielded by a saddle point of function (2.172), as given by Proposition 2.46.ii, is immediate: by the same argument as used when justifying Remark 2.47, in the situation in question one can look for saddle point with h = 0, $H = \eta I_d$, and identifying the required η reduces to solving univariate convex problem

$$\begin{aligned} \operatorname{Opt}(\rho) &= \min_{\eta} \frac{1}{2} \left\{ -\frac{d}{2} \ln(1 - \widehat{\sigma}^4 \eta^2) - \frac{d}{2} \widehat{\sigma}^2 (1 - \underline{\sigma}^2 / \overline{\sigma}^2) \eta + \frac{d\delta(2 + \delta) \widehat{\sigma}^4 \eta^2}{1 + \widehat{\sigma}^2 \eta} \right. \\ &+ \frac{\rho^2 \eta}{2(1 - \widehat{\sigma}^2 \eta)} : -\gamma \leq \widehat{\sigma}^2 \eta \leq 0 \right\} \\ &\left[\widehat{\sigma} &= \sqrt{2} \overline{\sigma}, \ \delta &= 1 - \underline{\sigma} / \overline{\sigma} \right] \end{aligned}$$

which can be done in no time by Bisection. The resulting detector and the upper bound on its risk are given by optimal solution $\eta(\rho)$ to the latter problem according to

$$\begin{split} \phi_{\rho}^{*}(\zeta) &= \frac{1}{2}\eta(\rho)\zeta^{T}\zeta \\ &+ \underbrace{\frac{d}{4} \left[\ln \left(\frac{1 - \widehat{\sigma}^{2}\eta(\rho)}{1 + \widehat{\sigma}^{2}\eta(\rho)} \right) - \widehat{\sigma}^{2}(1 - \underline{\sigma}^{2}/\overline{\sigma}^{2})\eta(\rho) - \frac{\rho^{2}\eta(\rho)}{1 - \widehat{\sigma}^{2}\eta(\rho)} \right]}_{a(\rho)}, \\ \mathcal{P}_{1}, \mathcal{P}_{2}] &\leq \operatorname{Risk}(\rho) := e^{\operatorname{Opt}(\rho)}. \end{split}$$

 $\operatorname{Risk}[\phi_{\rho}^*|\mathcal{P}_1, \mathcal{P}_2] \leq \operatorname{Risk}(\rho) := e^{\operatorname{Opt}(\rho)}$

(2.197) Observe that R does not appear in (2.197) at all. Now, it is immediately seen that $Opt(\rho) \to 0$ as $\rho \to +0$ and $Opt(\rho) \to -\infty$ as $\rho \to +\infty$, implying that given $\kappa \in (0, 1)$, we can easily find by bisection $\rho = \rho(\kappa)$ such that $Risk(\rho) = \kappa$; in what follows, we assume w.l.o.g. that $R > \rho(\kappa)$ for the value of κ we end with, see below. Next, let us pass from the detector $\phi^*_{\rho(\kappa)}(\cdot)$ to its shift

$$\phi^{*,\kappa}(\zeta) = \phi^{*}_{\rho(\kappa)}(\zeta) + \ln(\varepsilon/\kappa),$$

so that for the simple test \mathcal{T}^{κ} which, given observation ζ , accepts G_1 and rejects $G_2^{\rho(\kappa)}$ whenever $\phi^{*,\kappa}(\zeta) \geq 0$, and accepts $G_2^{\rho(\kappa)}$ and rejects G_1 otherwise, it holds

$$\operatorname{Risk}_{1}(\mathcal{T}^{\kappa}|G_{1}, G_{2}^{\rho(\kappa)}) \leq \frac{\kappa^{2}}{\varepsilon}, \operatorname{Risk}_{2}(\mathcal{T}^{\kappa}|G_{1}, G_{2}^{\rho(\kappa)}) \leq \varepsilon,$$
(2.198)

see Proposition 2.16 and (2.52).

We are nearly done. Given $\kappa \in (0, 1)$, consider the system of tests \mathcal{T}_t^{κ} , t = 2, 3, ..., K, as follows. At time $t \in \{2, 3, ..., K\}$, given observations $\omega_1, ..., \omega_t$ stemming from (2.195), let us form the vector

$$\zeta_t = \omega_t - \omega_1$$

and compute the quantity $\phi^{*,\kappa}(\zeta_t)$. If this quantity is negative, we claim that the change has already taken place and terminate, otherwise we claim that so far, there was no change, and proceed to time t + 1 (if t < K) or terminate (if t = K).

The risk analysis for the resulting system of inferences is immediate. Observe that

- (!) For every t = 2, 3, ..., K:
- A. if there is no change on time horizon $1, ..., t: x_1 = x_2 = ... = x_t$, then the probability for \mathcal{T}_t^{κ} to conclude that there was a change is at most κ^2/ε ;
- B. if, on the other hand, $||x_t x_1||_2^2 \ge \rho^2$, then the probability for \mathcal{T}_t^{κ} to conclude that so far there was no change is at most ε .

Indeed, we clearly have

$$\zeta_t = [x_t - x_1] + \xi^t,$$

where $\xi^t = \xi_t - \xi_1 \sim \mathcal{N}(0, \sigma^2 I_d)$ with $\sigma \in [\sqrt{2\sigma}, \sqrt{2\sigma}]$. Our actions at time t are nothing but application of the test \mathcal{T}^{κ} to the observation ζ_t . In the case of A the distribution of this observation obeys the hypothesis G_1 , and the probability for \mathcal{T}_t^{κ} to claim that there was a change is at most κ^2/ε by the first inequality in (2.198). In the case of B, the distribution of ζ_t obeys the hypothesis $G_2^{\rho(\kappa)}$, and thus the probability for \mathcal{T}_t^{κ} to claim that there was no change on time horizon 1, ..., t is $\leq \varepsilon$ by the second inequality in (2.198).

In view of (!), the probability of false alarm for the system of inferences $\{\mathcal{T}_t^\kappa\}_{t=2}^K$ is at most $(K-1)\kappa^2/\varepsilon$, and specifying κ as

$$\kappa = \sqrt{\epsilon \varepsilon / (K - 1)},$$

we make this probability $\leq \epsilon$. The resulting procedure, by the same (!), detects a change at time $t \in \{2, 3, ..., K\}$ with probability at least $1 - \varepsilon$, provided that the energy of this change is at least ρ_*^2 , with

$$\rho_* = \rho\left(\sqrt{\epsilon\varepsilon/(K-1)}\right),\tag{2.199}$$

In fact we can say a bit more:

Proposition 2.53. Let the deterministic sequence $x_1, ..., x_K$ underlying observations (2.195) be such that for some t it holds $||x_t - x_1||_2^2 \ge \rho_*^2$, with ρ_* given by (2.199). Then the probability for the system of inferences we have built to detect a change at time t or earlier is at least $1 - \varepsilon$.

Indeed, under the premise of Proposition, the probability for \mathcal{T}_t^{κ} to claim that a change already took place is at least $1 - \varepsilon$, and this probability can be only smaller than the probability to detect change on time horizon 2, 3, ..., t.

How it works. As applied to the "movie" story we started with, the outlined procedure works as follows. The images in question are of the size 256×256 , so that we are in the case of $d = 256^2 = 65536$. The images are represented by 2D arrays in gray scale, that is, as 256×256 matrices with entries in the range [0, 255]. In the experiment to be reported (same as in the movie) we assumed the maximal noise intensity $\overline{\sigma}$ to be 10, and used $\underline{\sigma} = \overline{\sigma}/\sqrt{2}$. The reliability tolerances ϵ , ε were set to 0.01, and K was set to 9, resulting in

$$\rho_*^2 = 7.38 \cdot 10^6,$$
which corresponds to the per pixel energy $\rho_*^2/65536 = 112.68$ – just by 12% above the allowed expected per pixel energy of noise (the latter is $\overline{\sigma}^2 = 100$). The resulting detector is

$$\phi_*(\zeta) = -2.7138 \frac{\zeta^T \zeta}{10^5} + 366.9548;$$

in other words, test \mathcal{T}_t^{κ} claims that the change took place when the average, over pixels, per pixel energy in the difference $\omega_t - \omega_1$ is at least 206.33, which is pretty close to the expected per pixel energy (200.0) in the noise $\xi_t - \xi_1$ affecting the difference $\omega_t - \omega_1$.

Finally, this is how the just described system of inferences worked in simulations. The underlying sequence of images was obtained from the "basic sequence"

$$\bar{x}_t = G + 0.0357(t-1)(L-G), t = 1, 2, ...^{37}$$
 (2.200)

where G is the image of the gentlemen and L is the image of the lady (up to noise, these are the first and the last frames on Figure 2.5). To get the observations in a particular simulation, we augmented this sequence from the left by a random number of images G, took the first 9 images in the resulting sequence, and added to them independent across the images observation noises drawn at random from $\mathcal{N}(0, 100I_{65536})$. Augmentation was carried out in such a way that with probability 1/2, there was no change on the time horizon 1, 2, ..., 9, and with probability 0.5there was a change at time instant τ chosen at random according to uniform distribution on $\{2, 3, ..., 9\}$. In 3000 simulations of this type, not a *single* false alarm was observed, while the empirical probability of a miss was 0.0553. It should be added that the actual energy of a change, if any, that is, $0.0357^2 ||L - G||_F^2$, was "just" $3.37 \cdot 10^5$, that is, it was by factor ≈ 21 less than the energy of change ρ_*^2 which our inferences are bound to detect with probability at least 0.99. And in the series of 3000 experiments we have reported, there was no "no detection" simulations where $\max_{t \leq K} \|x_t - x_1\|_2^2$ was above ρ_*^2 (that is, no simulations where Proposition 2.53) ensures detectability with probability at least 0.99, and in fact the change is not detected). Thus, all misses came from simulations which are *not* covered by our risk guarantees³⁸. Moreover, the change at time t, if detected, never was detected with a delay more than 1.

Finally, in the particular "movie" we started with, the change takes place at time t = 3, and the system of inferences we have just developed discovered the change at time 4. How this compares to the time at which you managed to detect the change?

"Numerical near-optimality." Beyond the realm of simple o.s.'s we have no theoretical guarantees of near-optimality for the inferences we are developing; this does not mean, however, that we cannot quantify conservatism of our techniques

³⁸A reader can be surprised – how happens that with actual energy of change 20 times less than the "theoretical threshold" ρ_*^2 , in our experiments, the empirical probability of a miss was as low as 5%, instead of being 50% or 100%. A plausible explanation is as follows: our performance guarantees are given by worst-case oriented theoretical analysis , and in random simulations we usually do not generate the "worst case" situations. For example, with model (2.200), the change, when happens, is of energy 20 times below the threshold; however, 3 time units after the change, the quantity $||x_t - x_1||_2^2$ becomes 16 times larger the energy of change, so that by Proposition 2.53, already worst-case analysis shows that there are good chances to detect the change when it happens "deeply inside" the observation horizon.

numerically. To give an example, let us forget, for the sake of simplicity, about change detection *per se* and focus on the auxiliary problem we have introduced above, the one on deciding upon hypotheses G_1 and G_2^{ρ} via observation (2.196), and let our goal be to decide on these two hypotheses from a single observation with risk $\leq \epsilon$, for a given $\epsilon \in (0, 1)$. Whether this is possible or not, it depends on ρ ; let us denote by ρ^+ the smallest ρ for which we can meet the risk specification with our detector-based approach (ρ^+ is nothing but what was above called $\rho(\epsilon)$), and by ρ – the smallest ρ for which "in the nature" there exists a simple test deciding on \overline{G}_1 vs. G_2^{ρ} with risk $\leq \epsilon$. We can look at the ratio ρ^+/ρ as at the "index of conservatism" of our approach. Now, ρ^+ is given by an efficient computation; what about ρ ? Well, there is a simple way to get a *lower bound* on ρ , namely, as follows. Observe that if the two composite hypotheses G_1, G_2^{ρ} can be decided upon with risk $\leq \epsilon$, the same holds true for two simple hypotheses stating that the distribution of observation (2.196) is P_1 , respectively, P_2 , where P_1, P_2 correspond to the cases when

- (P₁): ζ is drawn from $\mathcal{N}(0, 2\overline{\sigma}^2 I_d)$
- (P_2) : ζ is obtained by adding $\mathcal{N}(0, 2\underline{\sigma}^2 I_d)$ -noise to a random, independent of the noise, signal u uniformly distributed on the sphere $\{\|u\|_2 = \rho\}$.

Indeed, P_1 obeys hypothesis G_1 , and P_2 is a mixture of distributions obeying G_2^{ρ} ; as a result, a simple test \mathcal{T} deciding $(1 - \epsilon)$ -reliably on G_1 vs. G_2^{ρ} would induce a test deciding equally reliably on P_1 vs. P_2 , specifically, the test which, given observation ζ , accepts P_1 if \mathcal{T} on the same observation accepts G_1 , and accepts P_2 otherwise.

Now, we can use two-point lower bound (Proposition 2.2) to lower-bound the risk of deciding on P_1 vs. P_2 ; since both distributions are spherically symmetric, computing this bound reduces to computing similar bound for the univariate distributions of $\zeta^T \zeta$ induced by P_1 and P_2 , and these univariate distributions are easy to compute. The resulting lower risk bound depends on ρ , and we can find the smallest ρ for which the bound is ≥ 0.01 , and use this ρ in the role of ρ ; the associated indexes of conservatism can be only larger than the true ones. Let us look what are these indexes for the data used in our change detection experiment, that is, $\epsilon = 0.01$, $d = 256^2 = 65536$, $\overline{\sigma} = 10$, $\underline{\sigma} = \overline{\sigma}/\sqrt{2}$. Computation shows that in this case we have

$$\rho^+ = 2702.4, \ \rho^+/\rho \le 1.04$$

– nearly no conservatism at all! When eliminating the uncertainty in the intensity of noise by increasing $\underline{\sigma}$ from $\overline{\sigma}/\sqrt{2}$ to $\overline{\sigma}$, we get

$$\rho^+ = 668.46, \ \rho^+/\rho \le 1.15$$

- still not that much of conservatism!

2.10 EXERCISES FOR LECTURE 2

[†] marks more difficult exercises.

2.10.1 Two-point lower risk bound

Exercise 2.54. Let p and q be two distinct from each other probability distributions on d-element observation space $\Omega = \{1, ..., d\}$, and consider two simple hypotheses on the distribution of observation $\omega \in \Omega$, $H_1 : \omega \sim p$, and $H_2 : \omega \sim q$.

- 1. Is it true that there always exists a simple deterministic test deciding on H_1, H_2 with risk < 1/2?
- 2. Is it true that there always exists a simple randomized test deciding on H_1, H_2 with risk < 1/2?
- 3. Is it true that when quasi-stationary K-repeated observations are allowed, one can decide on H_1 , H_2 with a whatever small risk, provided K is large enough?

2.10.2 Around Euclidean Separation

Exercise 2.55. Justify the "immediate observation" in Section 2.2.2.3.B.

Exercise 2.56. 1. Prove Proposition 2.11.

<u>Hint:</u> You can find useful the following simple observation (prove it, provided you indeed use it):

Let $f(\omega)$, $g(\omega)$ be probability densities taken w.r.t. a reference measure P on an observation space Ω , and let $\epsilon \in (0, 1/2]$ be such that

$$\int_{\Omega} \min[f(\omega), g(\omega)] P(d\omega) \le 2\epsilon.$$

Then

$$2\bar{\epsilon} := \int_{\Omega} \sqrt{f(\omega)g(\omega)} P(d\omega) \le 2\sqrt{\epsilon(1-\epsilon)}.$$

2. Justify Illustration in Section 2.2.3.2.C.

2.10.3 Hypothesis testing via ℓ_1 -separation

Let d be a positive integer, and the observation space Ω be the finite set $\{1, ..., d\}$ equipped with the counting reference measure³⁹. Probability distributions on Ω can be identified with points p of d-dimensional probabilistic simplex

$$\mathbf{\Delta}_d = \{ p \in \mathbf{R}^d : p \ge 0, \sum_i p_i = 1 \};$$

i-th entry p_i in $p \in \Delta_d$ is the probability for the distributed according to p random variable to take value $i \in \{1, ..., d\}$. With this interpretation, p is the probability density taken w.r.t. the counting measure on Ω .

Assume B and W are two nonintersecting nonempty closed convex subsets of Δ_d ; we interpret B and W as black and white probability distributions on Ω , and our goal is to find optimal, in terms of its total risk, test deciding on the hypotheses

$$H_1: p \in B, H_2: p \in W$$

³⁹Counting measure is the measure on a discrete (finite or countable) set Ω which assigns every point of Ω with mass 1, so that the measure of a subset of Ω is the cardinality of the subset when it is finite and is $+\infty$ otherwise.

via a single observation $\omega \sim p$.

Warning: Everywhere in this Section, "test" means "simple test."

Exercise 2.57. Consider the convex optimization problem

Opt =
$$\min_{p \in B, q \in W} \left[f(p,q) := \sum_{i=1}^{d} |p_i - q_i| \right]$$
 (2.201)

and let (p^*, q^*) be an optimal solution to this problem (it clearly exists).

1. Extract from optimality conditions that there exist reals $\rho_i \in [-1, 1], 1 \leq i \leq n$, such that

$$\rho_i = \begin{cases}
1, & p_i^* > q_i^* \\
-1, & p_i^* < q_i^*
\end{cases}$$
(2.202)

and

$$\rho^{T}(p-p^{*}) \ge 0 \,\forall p \in B \& \rho^{T}(q-q^{*}) \le 0 \,\forall q \in W.$$
(2.203)

2. Extract from the previous item that the test \mathcal{T} which, given an observation $\omega \in \{1, ..., d\}$, accepts H_1 with probability $\pi_{\omega} = (1 + \rho_{\omega})/2$ and accepts H_2 with complementary probability, has total risk equal to

$$\sum_{\omega \in \Omega} \min[p_{\omega}^*, q_{\omega}^*] \tag{2.204}$$

and thus is minimax optimal in terms of the total risk.

Comments. Exercise 2.57 describes an efficiently computable and *optimal in* terms of worst-case total risk simple test deciding on a pair of "convex" composite hypotheses on the distribution of a discrete random variable. While it seems an attractive result, we believe by itself this result is useless, since usually in the testing problem in question a single observation by far is not enough for a reasonable inference; such an inference requires observing several independent realizations $\omega_1, ..., \omega_K$ of the random variable in question. And construction presented in Exercise 2.57 says nothing on how to adjust the test to the case of repeated observation. Of course, when $\omega^K = (\omega_1, ..., \omega_K)$ is K-element i.i.d. sample drawn from a probability distribution p on $\Omega = \{1, ..., d\}, \omega^K$ can be thought of as a single observation of discrete random variable taking value in the set $\Omega^K = \underbrace{\Omega \times ... \times \Omega}_{K}$, the prob-

ability distribution p^K of ω^K being readily given by p; so why not to apply the construction from Exercise 2.57 to ω^K in the role of ω ? On a close inspection, this idea fails. One serious reason for this failure is that the cardinality of Ω^K (which, among other factors, is responsible for the computational complexity of building the test in Exercise 2.57) blows up exponentially as K grows. Another, even more serious, complication is that p^K depends on p nonlinearly, so that the family of distributions p^K of ω^K induced by a convex family of distributions p of ω , convexity meaning that p's in question fill a *convex* subset of the probabilistic simplex, is not convex; and convexity of the sets B, W in the context of Exercise 2.57 is crucial. Thus, passing from single realization of discrete random variable to the sample of K > 1 independent realizations of the variable results in severe structural and quantitative complications "killing," at least at the first glance, the

approach undertaken in Exercise 2.57.

In spite of the above pessimistic conclusions, the single-observation test from Exercise 2.57 admits a meaningful multi-observation modification, which is the subject of our next Exercise.

Exercise 2.58. There is a straightforward way to use the optimal, in terms of its total risk, single-observation test built in Exercise 2.57 in the "multi-observation" environment. Specifically, following the notation from Exercise 2.57, let $\rho \in \mathbf{R}^d, p^*, q^*$ be the entities built in this Exercise, so that $p^* \in B, q^* \in W$, all entries in ρ belong to [-1, 1], and

$$\{ \rho^T p \ge \alpha := \rho^T p^* \ \forall p \in B \} \& \{ \rho^T q \le \beta := \rho^T q^* \ \forall q \in W \} \\ \& \alpha - \beta = \rho^T [p^* - q^*] = \| p^* - q^* \|_1.$$

Given an i.i.d. sample $\omega^K = (\omega_1, ..., \omega_K)$ with $\omega_t \sim p$, where $p \in B \cup W$, we could try to decide on the hypotheses $H_1 : p \in B$, $H_2 : p \in W$ as follows. Let us set $\zeta_t = \rho_{\omega_t}$. For large K the observable, given ω^K , quantity $\zeta^K := \frac{1}{K} \sum_{t=1}^K \zeta_t$, by the Law of Large Numbers, will be with overwhelming probability close to $\mathbf{E}_{\omega \sim p} \{\rho_\omega\} = \rho^T p$, and the latter quantity is $\geq \alpha$ when $p \in B$ and is $\leq \beta < \alpha$ when $p \in W$. Consequently, selecting a "comparison level" $\ell \in (\beta, \alpha)$, we can decide on the hypotheses $p \in B$ vs. $p \in W$ by computing ζ^K , comparing the result with ℓ , and accepting the hypothesis $p \in B$ when $\zeta^K \geq \ell$, otherwise accepting the alternative $p \in W$. The goal of this Exercise is to quantify the above qualitative considerations. To this end let us fix $\ell \in (\beta, \alpha)$ and K and ask ourselves the following questions:

A. For $p \in B$, how to upper-bound the probability $\operatorname{Prob}_{p_K} \{\zeta^K \leq \ell\}$? B. For $p \in W$, how to upper-bound the probability $\operatorname{Prob}_{p_K} \{\zeta^K \geq \ell\}$?

Here p_K is the probability distribution of the i.i.d. sample $\omega^K = (\omega_1, ..., \omega_K)$ with $\omega_t \sim p$.

The simplest way to answer these questions is to use Bernstein's bounding scheme. Specifically, to answer question A, let us select $\gamma \ge 0$ and observe that for every probability distribution p on $\{1, 2, ..., d\}$ it holds

$$\underbrace{\operatorname{Prob}_{p_{K}}\left\{\zeta^{K} \leq \ell\right\}}_{\pi_{K,-}[p]} \exp\{-\gamma\ell\} \leq \mathbf{E}_{p_{K}}\left\{\exp\{-\gamma\zeta^{K}\}\right\} = \left[\sum_{i=1}^{d} p_{i} \exp\{-\frac{1}{K}\gamma\rho_{i}\}\right]^{K},$$

whence

$$\ln(\pi_{K,-}[p]) \le K \ln\left(\sum_{i=1}^{d} p_i \exp\{-\frac{1}{K}\gamma \rho_i\}\right) + \gamma \ell$$

implying, via substitution $\gamma = \mu K$, that

$$\forall \mu \ge 0 : \ln(\pi_{K,-}[p]) \le K\psi_{-}(\mu,p), \ \psi_{-}(\mu,p) = \ln\left(\sum_{i=1}^{d} p_{i} \exp\{-\mu\rho_{i}\}\right) + \mu\ell.$$

Similarly, setting $\pi_{K,+}[p] = \operatorname{Prob}_{p_K} \left\{ \zeta^K \ge \ell \right\}$, we get

$$\forall \nu \ge 0 : \ln(\pi_{K,+}[p]) \le K\psi_+(\nu,p), \ \psi_+(\nu,p) = \ln\left(\sum_{i=1}^d p_i \exp\{\nu\rho_i\}\right) - \nu\ell.$$

Now goes the exercise:

1. Extract from the above observations that

$$\operatorname{Risk}(\mathcal{T}^{K,\ell}|H_1,H_2) \le \exp\{K\varkappa\}, \ \varkappa = \max\left[\max_{p\in B} \inf_{\mu\ge 0}\psi_-(\mu,p), \max_{q\in W} \inf_{\nu\ge 0}\psi_+(\nu,q)\right],$$

where $\mathcal{T}^{K,\ell}$ is the K-observation test which accepts the hypothesis $H_1: p \in B$ when $\zeta^K \geq \ell$ and accepts the hypothesis $H_2: p \in W$ otherwise.

2. Verify that $\psi_{-}(\mu, p)$ is convex in μ and concave in p, and similarly for $\psi_{+}(\nu, q)$, so that

$$\max_{p \in B} \inf_{\mu \ge 0} \psi_{-}(\mu, p) = \inf_{\mu \ge 0} \max_{p \in B} \psi_{-}(\mu, p), \ \max_{q \in W} \inf_{\nu \ge 0} \psi_{+}(\nu, q) = \inf_{\nu \ge 0} \max_{q \in W} \psi_{+}(\nu, q)$$

Thus, computing \varkappa reduces to minimizing on the nonnegative ray the convex functions $\phi_{-}(\mu) = \max_{p \in B} \psi_{+}(\mu, p)$ and $\phi_{+}(\nu) = \max_{q \in W} \psi_{+}(\nu, q)$.

3. Prove that when $\ell = \frac{1}{2}[\alpha + \beta]$, one has

$$\varkappa \le -\frac{1}{12}\Delta^2, \ \Delta = \alpha - \beta = \|p^* - q^*\|_1.$$
(2.205)

Note that the above test and the quantity \varkappa responsible for the upper bound on its risk depend, as on a parameter, on the "acceptance level" $\ell \in (\beta, \alpha)$. The simplest way to select a reasonable value of ℓ is to minimize \varkappa over an equidistant grid $\Gamma \subset (\beta, \alpha)$, of small cardinality, of values of ℓ .

Now let us consider an alternative way to pass from a "good" single-observation test to its multi-observation version. Our "building block" now is the minimum risk randomized single-observation test⁴⁰, and its multi-observation modification is just the majority version of this building block. Our first observation is that building the minimum risk single-observation test reduces to solving a *convex* optimization problem:

Exercise 2.59. Let, as above, B and W be nonempty nonintersecting closed convex subsets of probabilistic simplex Δ_d . Demonstrate that the problem of finding the best, in terms of its risk, randomized single-observation test deciding on $H_1 : p \in B$ vs. $H_2 : p \in W$ via observation $\omega \sim p$ reduces to solving a convex optimization problem. Write down this problem as an explicit LO program when B and W are polyhedral sets given by polyhedral representations:

$$B = \{p : \exists u : P_B p + Q_B u \le a_B\},\$$

$$W = \{p : \exists u : P_W p + Q_W u \le a_W\},\$$

We see that the "ideal building block" – the minimum-risk single-observation test – can be built efficiently. What is at this point unclear, is whether this block is of any use for majority modifications, that is, whether it is true that the risk of this test is < 1/2 – this is what we need to get the possibility for low-risk testing from repeated observations via majority version of the minimum-risk single-observation test.

 $^{^{40}}$ this test can differ from the one built in Exercise 2.57 – the latter test is optimal in terms of the sum, rather than the maximum, of its partial risks.

Exercise 2.60. Extract from Exercise 2.57 that in the situation of this Section, denoting by Δ the optimal value in the optimization problem (2.201), one has

- 1. The risk of any single-observation test, deterministic or randomized alike, is $\geq \frac{1}{2} \frac{\Delta}{4}$
- 2. There exists a single-observation randomized test with risk $\leq \frac{1}{2} \frac{\Delta}{8}$, and thus the risk of the minimum risk single-observation test given by Exercise 2.59 does not exceed $\frac{1}{2} \frac{\Delta}{8} < 1/2$ as well.

Pay attention to the fact that $\Delta > 0$ (since, by assumption, B and W do not intersect).

The bottom line is that in the situation of this Section, given a target value ϵ of risk and assuming stationary repeated observations are allowed, we have (at least) three options to meet the risk specifications:

- 1. To start with the optimal, in terms of its total risk, single-observation detector as explained in Exercise 2.57, and the to pass to its multi-observation version built in Exercise 2.58;
- 2. To use the majority version of the minimum-risk randomized single-observation test built in Exercise 2.59;
- 3. To use the test based on the minimum risk detector for B, W, as explained in the main body of Lecture 2.

In all cases, the number K of observations should be specified as "presumably the smallest" K ensuring that the risk of the resulting multi-observation test is at most a given target ϵ ; this K can be easily specified by utilizing the results on the risk of a detector-based test in a Discrete o.s. from the main body of Lecture 2 along with risk-related results of Exercises 2.58, 2.59.

Exercise 2.61. Run numerical experimentation to get an idea whether one of the three options above always dominates other options (that is, requires smaller sample of observations to ensure the same risk).

Now let us focus on theoretical comparison of the detector-based test and the majority version of the minimum-risk single-observation test (options 1 and 2 above) in the general situation described in the beginning of Section 2.10.3. Given $\epsilon \in (0, 1)$, the corresponding sample sizes, K_d and K_m , are completely specified each by its own "measure of closeness" between B and W. Specifically,

• For K_d , the closeness measure is

$$\rho_d(B,W) = 1 - \max_{p \in B, q \in W} \sum_{\omega} \sqrt{p_{\omega} q_{\omega}}; \qquad (2.206)$$

 $1 - \rho_d(B, W)$ is the minimal risk of a detector for B, W, and for $\rho_d(B, W)$ and ϵ small, we have $K_d \approx \ln(1/\epsilon)/\rho_d(B, W)$ (why?).

• Given ϵ , K_m is fully specified by the minimal risk ρ of simple randomized singleobservation test \mathcal{T} deciding on the associated with B, W hypotheses; by Exercise 2.60, we have $\rho = \frac{1}{2} - \delta$, where δ is within absolute constant factor of the optimal value $\Delta = \min_{p \in B, q \in W} ||p - q||_1$ of (2.201). The risk bound for the Kobservation majority version of \mathcal{T} is the probability to get at least K/2 heads in K independent tosses of coin with probability to get head in a single toss

equal to $\rho = 1/2 - \delta$. When ρ is not close to 0 and ϵ is small, the $(1 - \epsilon)$ quantile of the number of heads in our K coin tosses is $K\rho + O(1)\sqrt{K \ln(1/\epsilon)} = K/2 - \delta K + O(1)\sqrt{K \ln(1/\epsilon)}$ (why?). K_m is the smallest K for which this quantile is $\langle K/2 \rangle$, so that K_m is of order of $\ln(1/\epsilon)/\delta^2$, or, which is the same, of order of $\ln(1/\epsilon)/\Delta^2$. We see that the "responsible for K_m " closeness between B and W is

$$\rho_m(B, W) = \Delta^2 = \left[\min_{p \in B, q \in W} \|p - q\|_1\right]^2, \qquad (2.207)$$

and K_m is of order of $\ln(1/\epsilon)/\rho_m(B, W)$.

The goal of the next exercise is to compare ρ_b and ρ_m .

Exercise 2.62. Prove that in the situation of this Section one has

$$\frac{1}{8}\rho_m(B,W) \le \rho_d(B,W) \le \frac{1}{2}\sqrt{\rho_m(B,W)}.$$
(2.208)

Relation (2.208) suggests that while K_d never is "much larger" than K_m (this we know in advance: in repeated version of Discrete o.s., properly built detectorbased test provably is nearly optimal), but K_m could be much larger than K_d . This indeed is the case:

Exercise 2.63. Given $\delta \in (0, 1/2)$, let $B = \{[\delta; 0; 1 - \delta]\}$ and $W = \{[0; \delta; 1 - \delta]\}$. Verify that in this case, the numbers of observations K_d and K_m resulting in a given risk $\epsilon \ll 1$ of multi-observation tests, as functions of δ are proportional to $1/\delta$ and $1/\delta^2$, respectively. Compare the numbers when $\epsilon = 0.01$ and $\delta \in \{0.01; 0.05; 0.1\}$.

2.10.4 Miscellaneous exercises

Exercise 2.64. Prove that the conclusion in Proposition 2.20 remains true when the test \mathcal{T} in the premise of Proposition is randomized.

Exercise 2.65. Let $p_1(\omega), p_2(\omega)$ be two positive probability densities, taken w.r.t. a reference measure Π , on an observation space Ω , and let $\mathcal{P}_{\chi} = \{p_{\chi}\}, \chi = 1, 2$. Find the optimal, in terms of its risk, balanced detector for $\mathcal{P}_{\chi}, \chi = 1, 2$.

Exercise 2.66. Recall that the exponential, with parameter $\mu > 0$, distribution on $\Omega = \mathbf{R}_+$ is the distribution with the density $p_{\mu}(\omega) = \mu e^{-\mu\omega}$, $\omega \ge 0$. Given positive reals $\alpha < \beta$, consider two families of exponential distributions, $\mathcal{P}_1 = \{p_{\mu} : 0 < \mu \le \alpha\}$, and $\mathcal{P}_2 = \{p_{\mu} : \mu \ge \beta\}$. Build the optimal, in terms of its risk, balanced detector for $\mathcal{P}_1, \mathcal{P}_2$. What happens with the risk of the detector you have built when the families $\mathcal{P}_{\chi}, \chi = 1, 2$, are replaced with their convex hulls?

Exercise 2.67. [Follow-up to Exercise 2.66] Assume that the "lifetime" ζ of a lightbulb is a realization of random variable with exponential distribution (i.e., the density $p_{\mu}(\zeta) = \mu e^{-\mu\zeta}$, $\zeta \geq 0$; in particular, the expected lifetime of a lightbulb in this model is $1/\mu$)⁴¹. Given a lot of lightbulbs, you should decide whether they were

⁴¹In Reliability, probability distribution of the lifetime ζ of an organism or a technical device is characterized by the failure rate $\lambda(t) = \lim_{dt\to+0} \frac{\operatorname{Prob}\{t \leq \zeta \leq t+dt\}}{dt \cdot \operatorname{Prob}\{\zeta \geq t\}}$ (so that for small dt, $\lambda(t)dt$ is the conditional probability to "die" in the time interval [t, t+dt] provided the lifetime is at least

produced under normal conditions (resulting in $\mu \leq \alpha = 1$) or under abnormal ones (resulting in $\mu \geq \beta = 1.5$). To this end, you can select at random K lightbulbs and test them. How many lightbulbs should you test in order to make a 0.99-reliable conclusion? Answer this question in the situations when the observation ω in a test is

- 1. the lifetime of a lightbulb (i.e., $\omega \sim p_{\mu}(\cdot)$)
- 2. the minimum $\omega = \min[\zeta, \delta]$ of the lifetime $\zeta \sim p_{\mu}(\cdot)$ of a lightbulb and the allowed duration $\delta > 0$ of your test (i.e., if the lightbulb you are testing does not "die" on time horizon δ , you terminate the test)
- 3. $\omega = \chi_{\zeta < \delta}$, that is, $\omega = 1$ when $\zeta < \delta$, and $\omega = 0$ otherwise; here, as above, $\zeta \sim p_{\mu}(\cdot)$ is the random lifetime of a lightbulb, and $\delta > 0$ is the allowed test duration (i.e., you observe whether or not a lightbulb "dies" on time horizon δ , but do not register the lifetime when it is $< \delta$).

Consider the values 0.25, 0.5, 1, 2, 4 of δ .

Exercise 2.68. [Follow-up to Exercise 2.67] In the situation of Exercise 2.67, build a sequential test for deciding on Null hypothesis "the lifetime of a lightbulb from a given lot is $\zeta \sim p_{\mu}(\cdot)$ with $\mu \leq 1$ " (recall that $p_{\mu}(z)$ is the exponential density $\mu e^{-\mu z}$ on the ray $\{z \geq 0\}$) vs. the alternative "the lifetime is $\zeta \sim p_{\mu}(\cdot)$ with $\mu > 1$." In this test, you can select a number K of lightbulbs from the lot, switch them on at time 0 and record the actual lifetimes of the lightbulbs you are testing. As a result at the end of (any) observation interval $\Delta = [0, \delta]$, you observe K independent realizations of r.v. $\min[\zeta, \delta]$, where $\zeta \sim p_{\mu}(\cdot)$ with some unknown μ . In your sequential test, you are welcome to make conclusions at the endpoints $\delta_1 < \delta_2 < ... < \delta_S$ of several observation intervals.

Note: We deliberately skip details of problem's setting; how you decide on these missing details, is part of your solution to Exercise.

Exercise 2.69. In Section 2.6, we considered a model of elections where every member of population was supposed to cast a vote. Enrich the model by incorporating the option for a voter not to participate in the elections at all. Implement Sequential test for the resulting model and run simulations.

Exercise 2.70. Work out the following extension of the DOP problem. You are given two finite sets, $\Omega_1 = \{1, ..., I\}$ and $\Omega_2 = \{1, ..., M\}$, along with L nonempty closed convex subsets Y_{ℓ} of the set

$$\boldsymbol{\Delta}_{IM} = \{ [y_{im} > 0]_{i,m} : \sum_{i=1}^{I} \sum_{m=1}^{M} y_{im} = 1 \}$$

of all non-vanishing probability distributions on $\Omega = \Omega_1 \times \Omega_2 = \{(i, m) : 1 \leq i \leq I, 1 \leq m \leq M\}$. The sets Y_{ℓ} are such that all distributions from Y_{ℓ} have a common

t). The exponential distribution corresponds to the case of failure rate independent of t; usually, this indeed is nearly so except for "very small" and "very large" values of t.

marginal distribution $\theta^{\ell} > 0$ of *i*:

$$\sum_{m=1}^{M} y_{im} = \theta_i^{\ell}, \, 1 \leq i \leq I, \, \forall y \in Y_{\ell}, \, 1 \leq \ell \leq L.$$

Your observations $\omega_1, \omega_2, \ldots$ are sampled, independently of each other, from a distribution partly selected "by nature," and partly – by you. Specifically, the nature selects $\ell \leq L$ and a distribution $y \in Y_{\ell}$, and you select a positive *I*-dimensional probabilistic vector q from a given convex compact subset Q of the positive part of *I*-dimensional probabilistic simplex. Let $y_{|i}$ be the conditional, i being given, distribution of $m \in \Omega_2$ induced by y, so that $y_{|i}$ is the *M*-dimensional probabilistic vector with entries

$$[y_{|i}]_m = \frac{y_{im}}{\sum_{\mu \le M} y_{i\mu}} = \frac{y_{im}}{\theta_i^\ell}$$

In order to generate $\omega_t = (i_t, m_t) \in \Omega$, you draw i_t at random from the distribution q, and then the nature draws m_t at random from the distribution $y_{|i_t}$.

Given closeness relation C, your goal is to decide, up to closeness C, on the hypotheses H_1, \ldots, H_L , with H_ℓ stating that the distribution y selected by the nature belongs to Y_ℓ . Given "observation budget" (a number K of observations ω_k you can use), you want to find a probabilistic vector q which results in the test with as small C-risk as possible. Pose this Measurement Design problem as an efficiently solvable convex optimization problem.

Exercise 2.71. [probabilities of deviations from the mean]

The goal of what follows is to present the most straightforward application of simple families of distributions – bounds on probabilities of deviations of random vectors from their means.

Let $\mathcal{H} \subset \Omega = \mathbf{R}^d$, \mathcal{M}, Φ be a regular data such that $0 \in \operatorname{int} \mathcal{H}$, \mathcal{M} is compact, $\Phi(0; \mu) = 0 \forall \mu \in \mathcal{M}$, and $\Phi(h; \mu)$ is differentiable at h = 0 for every $\mu \in \mathcal{M}$. Let, further, $\overline{P} \in \mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$ and let $\overline{\mu} \in \mathcal{M}$ be a parameter of \overline{P} . Prove that

1. \bar{P} possesses expectation $e[\bar{P}]$, and

$$e[\bar{P}] = \nabla_h \Phi(0; \bar{\mu})$$

2. For every linear form $e^T \omega$ on Ω it holds

$$\pi := \bar{P}\{\omega : e^T(\omega - e[\bar{P}]) \ge 1\} \le \exp\left\{\sup_{\mu \in \mathcal{M}} \inf_{t \ge 0: t \in \mathcal{H}} \left[\Phi(te;\bar{\mu}) - te^T \nabla_h \Phi(0;\bar{\mu}) - t\right]\right\}.$$
(2.209)

Exercise 2.72. [testing convex hypotheses on mixtures] Consider the situation as follows. For given positive integers K, L and for $\chi = 1, 2$, given are

- nonempty convex compact signal sets $U_{\chi} \subset \mathbf{R}^{n_{\chi}}$
- regular data $\mathcal{H}_{k\ell}^{\chi} \subset \mathbf{R}^{d_k}, \mathcal{M}_{k\ell}^{\chi}, \Phi_{k\ell}^{\chi}$, and affine mappings

$$u_{\chi} \mapsto A_{k\ell}^{\chi}[u_{\chi}; 1] : \mathbf{R}^{n_{\chi}} \to \mathbf{R}^{d_{k}}$$

such that

$$u_{\chi} \in U_{\chi} \Rightarrow A_{k\ell}^{\chi}[u_{\chi}; 1] \in \mathcal{M}_{k\ell}^{\chi}$$

 $\label{eq:klassical} \begin{array}{l} 1 \leq k \leq K, \ 1 \leq \ell \leq L, \\ \bullet \ \mbox{probability vectors} \ \mu^k = [\mu_1^k;...;\mu_L^k], \ 1 \leq k \leq K. \end{array}$

We can associate with the outlined data families of probability distributions \mathcal{P}_{χ} on the observation space $\Omega = \mathbf{R}^{d_1} \times ... \times \mathbf{R}^{d_K}$ as follows. For $\chi = 1, 2, \mathcal{P}_{\chi}$ is comprised of all probability distributions P of random vectors $\omega^K = [\omega_1; ...; \omega_K] \in \Omega$ generated as follows:

We select

- a signal $u \in U_{\chi}$,
- a collection of probability distributions $P_{k\ell} \in \mathcal{S}[\mathcal{H}_{k\ell}^{\chi}, \mathcal{M}_{k\ell}^{\chi}, \Phi_{k\ell}^{\chi}], 1 \leq k \leq K, 1 \leq \ell \leq L$, in such a way that $A_{k\ell}^{\chi}[u_{\chi}; 1]$ is a parameter of $P_{k\ell}$:

$$\forall h \in \mathcal{H}_{k\ell}^{\chi} : \ln\left(\mathbf{E}_{\omega_k \sim P_{k\ell}} e^{h^T \omega_k}\right) \leq \Phi_{k\ell}^{\chi}(h_k; A_{k\ell}^{\chi}[u_{\chi}; 1]);$$

• we generate the components ω_k , k = 1, ..., K, independently across k, from μ^k mixture $\Pi[\{P_{k\ell}\}_{\ell=1}^L, \mu]$ of distributions $P_{k\ell}, \ell = 1, ..., L$, that is, draw at random, from distribution μ^k on $\{1, ..., L\}$, index ℓ , and then draw ω_k from the distribution $P_{k\ell}$.

Prove that setting

$$\begin{aligned} \mathcal{H}_{\chi} &= \{h = [h_1; ...; h_K] \in \mathbf{R}^{d = d_1 + ... + d_K} : h_k \in \bigcap_{\ell=1}^L \mathcal{H}_{k\ell}^{\chi}, 1 \le k \le K\}, \\ \mathcal{M}_{\chi} &= \{0\} \subset \mathbf{R}, \\ \Phi_{\chi}(h; \mu) &= \sum_{k=1}^K \ln\left(\sum_{\ell=1}^L \mu_{\ell}^k \max_{u_{\chi} \in U_{\chi}} \Phi_{k\ell}^{\chi}(h_k; A_{k\ell}^{\chi}[u_{\chi}; 1])\right) : \mathcal{H}_{\chi} \times \mathcal{M}_{\chi} \to \mathbf{R}, \end{aligned}$$

we get regular data such that

$$\mathcal{P}_{\chi} \subset \mathcal{S}[\mathcal{H}_{\chi}, \mathcal{M}_{\chi}, \Phi_{\chi}].$$

Explain how to use this observation to compute via Convex Programming affine detector, along with its risk, for the families of distributions $\mathcal{P}_1, \mathcal{P}_2$.

Exercise 2.73. [mixture of sub-Gaussian distributions] Let P_{ℓ} be sub-Gaussian distributions on \mathbb{R}^d with sub-Gaussianity parameters θ_{ℓ} , Θ , $1 \leq \ell \leq L$, with a common Θ -parameter, and let $\nu = [\nu_1; ...; \nu_L]$ be a probabilistic vector. Consider the ν -mixture $P = \prod[P^L, \nu]$ of distributions P_{ℓ} , so that $\omega \sim P$ is generated as follows: we draw at random from distribution ν index ℓ and then draw ω at random from distribution P_{ℓ} . Prove that P is sub-Gaussian with sub-Gaussianity parameters $\bar{\theta} = \sum_{\ell} \nu_{\ell} \theta_{\ell}$ and $\bar{\Theta}$, with (any) $\bar{\Theta}$ chosen to satisfy

$$\bar{\Theta} \succeq \Theta + \frac{6}{5} [\theta_{\ell} - \bar{\theta}] [\theta_{\ell} - \bar{\theta}]^T \,\forall \ell,$$

in particular, according to any one of the following rules:

 $\begin{array}{ll} 1. \ \bar{\Theta} = \Theta + \left(\frac{6}{5} \max_{\ell} \| \theta_{\ell} - \bar{\theta} \|_2^2 \right) I_d, \\ 2. \ \bar{\Theta} = \Theta + \frac{6}{5} \sum_{\ell} (\theta_{\ell} - \bar{\theta}) (\theta_{\ell} - \bar{\theta})^T, \\ 3. \ \bar{\Theta} = \Theta + \frac{8}{5} \sum_{\ell} \theta_{\ell} \theta_{\ell} \theta_{\ell}^T, \text{ provided that } \nu_1 = \ldots = \nu_L = 1/L, \end{array}$

Exercise 2.74. The goal of this Exercise is to give a simple sufficient condition for quadratic lift "to work" in the Gaussian case. Specifically, let \mathcal{A}_{χ} , U_{χ} , \mathcal{V}_{χ} , \mathcal{G}_{χ} , $\chi = 1, 2$, be as in Section 2.9.3, with the only difference that now we do *not* assume the compact sets \mathcal{U}_{χ} to be convex, and let \mathcal{Z}_{χ} be convex compact subsets of the sets $\mathcal{Z}^{n_{\chi}}$, see (2.164), such that

$$[u_{\chi};1][u_{\chi};1]^T \in \mathcal{Z}_{\chi} \ \forall u_{\chi} \in U_{\chi}, \ \chi = 1,2$$

Augmenting the above data with $\Theta_{\chi}^{(*)}$, δ_{χ} such that $\mathcal{V} = \mathcal{V}_{\chi}$, $\Theta_* = \Theta_*^{(\chi)}$, $\delta = \delta_{\chi}$ satisfy (2.165), $\chi = 1, 2$, and invoking Proposition 2.46.ii, we get at our disposal a quadratic detector ϕ_{lift} such that

$$\operatorname{Risk}[\phi_{\operatorname{lift}}|\mathcal{G}_1, \mathcal{G}_2] \le \exp\{-\operatorname{SadVal_{\operatorname{lift}}}\},\$$

with SadVal_{lift} given by (2.172). A natural question is, when SadVal_{lift} is negative, meaning that our quadratic detector indeed "is working" – its risk is < 1, implying that when repeated observations are allowed, tests based upon this detector allow to decide on the hypotheses $H_{\chi}: P \in \mathcal{G}_{\chi}, \chi = 1, 2$, on the distribution of observation $\zeta \sim P$ with a whatever small desired risk $\epsilon \in (0, 1)$. With our computationoriented ideology, this is not too important question, since we can answer it via efficient computation. This being said, there is no harm in a "theoretical" answer which could provide us with an additional insight. The goal of the Exercise is to justify a simple result on the subject. Here is the Exercise:

In the situation in question, assume that $\mathcal{V}_1 = \mathcal{V}_2 = \{\Theta_*\}$, which allows to set $\Theta_*^{(\chi)} = \Theta_*$, $\delta_{\chi} = 0$, $\chi = 1, 2$. Prove that in the case in question a necessary and sufficient condition for SadVal_{lift} to be negative is that the convex compact sets

$$\mathcal{U}_{\chi} = \{ B_{\chi} Z B_{\chi}^T : Z \in \mathcal{Z}_{\chi} \} \subset \mathbf{S}_{+}^{d+1}, \, \chi = 1, 2$$

do not intersect with each other.

Exercise 2.75. Prove if X is a nonempty convex compact set in \mathbf{R}^d , then the function $\widehat{\Phi}(h;\mu)$ given by (2.144) is real-valued and continuous on $\mathbf{R}^d \times X$ and is convex in h and concave in μ .

2.11 PROOFS

2.11.1 Proof of Claim in Remark 2.10

What we should prove is that is $p = [p_1; ...; p_K] \in B = [0, 1]^K$, then the probability $P_M(p)$ of the event

The total number of heads in K independent coin tosses, with probability p_k to get head in k-th toss, is at least M

is a nondecreasing function of p: if $p' \leq p''$, $p', p'' \in B$, then $P_M(p') \leq P_M(p'')$. To see it, let us associate with $p \in B$ a subset of B, specifically, $B_p = \{x \in B : 0 \leq x_k \leq p_k, 1 \leq k \leq K\}$, and a function $\chi_p(x) : B \to \{0,1\}$ which is equal to 0 at every point $x \in B$ where the number of entries x_k satisfying $x_k \leq p_k$ is less than

177

M, and is equal to 1 otherwise. It is immediately seen that

$$P_M(p) \equiv \int_B \chi_p(x) dx \tag{2.210}$$

(since with respect to the uniform distribution on B, the events $E_k = \{x \in B : x_k \leq p_k\}$ are independent across k and have probabilities p_k , and the right hand side in (2.210) is exactly the probability, taken w.r.t. the uniform distribution on B, of the event "at least M of the events E_1, \ldots, E_K take place"). But the right hand side in (2.33) clearly is nondecreasing in $p \in B$, since χ_p , by construction, is the characteristic function of the set

$$B[p] = \{x : \text{ at least } M \text{ of the entries } x_k \text{ in } x \text{ satisfy } x_k \leq p_k\},\$$

and these sets clearly grow when p is entrywise increased.

2.11.2 Proof of Proposition 2.8 in the case of quasi-stationary K-repeated observations

2.11.2.A Situation and goal. We are in the case **QS**, see Section 2.2.3.1, of the situation described in the beginning of Section 2.2.3; it suffices to verify that if \mathcal{H}_{ℓ} , $\ell \in \{1, 2\}$, is true, then the probability for $\mathcal{T}_{K}^{\text{maj}}$ to reject \mathcal{H}_{ℓ} is at most the quantity ϵ_{K} defined in (2.25). Let us verify this statement in the case of $\ell = 1$; the reasoning for $\ell = 2$ "mirrors" the one to follow.

It is clear that our situation and goal can be formulated as follows:

- "In the nature" there exists a random sequence $\zeta^K = (\zeta_1, ..., \zeta_K)$ of driving factors and collection of deterministic functions $\theta_k(\zeta^k = (\zeta_1, ..., \zeta_k))^{42}$ such that our k-th observation is $\omega_k = \theta_k(\zeta^k)$. Besides this, the conditional, ζ^{k-1} given, distribution $P_{\omega_k|\zeta^{k-1}}$ of ω_k always belongs to the family \mathcal{P}_1 comprised of distributions of random vectors of the form $x + \xi$, where deterministic x belongs to X_1 and the distribution of ξ belongs to \mathcal{P}_{γ}^d .
- There exist deterministic functions $\chi_k : \Omega \to \{0, 1\}$ and integer $M, 1 \leq M \leq K$, such that the test $\mathcal{T}_K^{\text{maj}}$, as applied to observation $\omega^K = (\omega_1, ..., \omega_K)$, rejects \mathcal{H}_1 if and only if the number of ones among the quantities $\chi_k(\omega_k), 1 \leq k \leq K$, is at least M.

In the situation of Proposition 2.8, M = |K/2| and $\chi_k(\cdot)$ are in fact independent of k: $\chi_k(\omega) = 1$ if and only if $\phi(\omega) \leq 0^{43}$.

• What we know is that the conditional, ζ^{k-1} being given, probability of the event $\chi_k(\omega_k = \theta_k(\zeta^k)) = 1$ is at most ϵ_* :

 $P_{\omega_k|\zeta^{k-1}}\{\omega_k:\chi_k(\omega_k)=1\}\leq \epsilon_\star\,\forall\zeta^{k-1}.$

⁴²as always, given a K-element sequence, say, $\zeta_1, ..., \zeta_K$, we write $\zeta^t, t \leq K$, as a shorthand for the fragment $\zeta_1, ..., \zeta_t$ of this sequence.

⁴³in fact, we need to write $\phi(\omega) < 0$ instead of $\phi(\omega) \leq 0$; we replace the strict inequality with its nonstrict version in order to make our reasoning applicable to the case of $\ell = 2$, where nonstrict inequalities do arise. Clearly, replacing in the definition of χ_k strict inequality with the nonstrict one, we only increase the "rejection domain" of \mathcal{H}_1 , so that upper bound on the probability of this domain we are about to get automatically is valid for the true rejection domain.

Indeed, $P_{\omega_k|\zeta^{k-1}} \in \mathcal{P}_1$, that is $P_{\omega_k|\zeta^{k-1}} \in \mathcal{P}_1$. As a result,

$$\begin{aligned} P_{\omega_k|\zeta^{k-1}}\{\omega_k:\phi_k(\omega_k)=1\} &= P_{\omega_k|\zeta^{k-1}}\{\omega_k:\phi(\omega_k)\leq 0\} \\ &= P_{\omega_k|\zeta^{k-1}}\{\omega_k:\phi(\omega_k)<0\}\leq \epsilon_\star, \end{aligned}$$

where the second equality is due to the fact that $\phi(\omega)$ is a nonconstant affine function and $P_{\omega_k|\zeta^{k-1}}$, along with all distributions from \mathcal{P}_1 , has density, and the inequality is given by the origin of ϵ_{\star} which upper-bounds the risk of the single-observation test underlying $\mathcal{T}_K^{\text{maj}}$.

What we want to prove is that under the circumstances we have just summarized, we have

$$P_{\omega^{K}}\{\omega^{K} = (\omega_{1}, ..., \omega_{K}) : \operatorname{Card}\{k \leq K : \chi_{k}(\omega_{k}) = 1\} \geq M\}$$

$$\leq \epsilon_{M} = \sum_{M \leq k \leq K} {K \choose k} \epsilon_{\star}^{k} (1 - \epsilon_{\star})^{K-k}, \qquad (2.211)$$

where $P_{\omega^{K}}$ is the distribution of $\omega^{K} = \{\omega_{k} = \theta_{k}(\zeta^{k-1})\}_{k=1}^{K}$ induced by the distribution of hidden factors. There is nothing to prove when $\epsilon_{\star} = 1$, since in this case $\epsilon_{M} = 1$. Thus, we assume from now on that $\epsilon_{\star} < 1$.

2.11.2.B Achieving the goal, step 1. Our reasoning, inspired by the one we used to justify Remark 2.10, is as follows. Consider a sequence of random variables η_k , $1 \leq k \leq K$, uniformly distributed on [0, 1] and independent of each other and of ζ^K , and consider new driving factors $\lambda_k = [\zeta_k; \eta_k]$ and new observations

$$\mu_k = [\omega_k = \theta_k(\zeta^k); \eta_k] = \Theta_k(\lambda^k = (\lambda_1, ..., \lambda_k))$$
(2.212)

driven by these new driving factors, and let

$$\psi_k(\mu_k = [\omega_k; \eta_k]) = \chi_k(\omega_k).$$

It is immediately seen that

• $\mu_k = [\omega_k = \theta_k(\zeta^k); \eta_k]$ is a deterministic function, $\Theta_k(\lambda^k)$, of λ^k , and the conditional, $\lambda^{k-1} = [\zeta^{k-1}; \eta^{k-1}]$ given, distribution $P_{\mu_k|\lambda^{k-1}}$ of μ_k is the product distribution $P_{\omega_k|\zeta^{k-1}} \times U$ on $\Omega \times [0, 1]$, where U is the uniform distribution on [0, 1]. In particular,

$$\pi_k(\lambda^{k-1}) := P_{\mu_k|\lambda^{k-1}}\{\mu_k = [\omega_k; \eta_k] : \chi_k(\omega_k) = 1\} = P_{\omega_k|\zeta^{k-1}}\{\omega_k : \chi_k(\omega_k) = 1\}$$

$$\leq \epsilon_\star.$$
(2.213)

• We have

$$P_{\lambda^{K}}\{\lambda^{K}: \operatorname{Card}\{k \leq K: \psi_{k}(\mu_{k} = \Theta_{k}(\lambda^{k})) = 1\} \geq M\}$$

= $P_{\omega^{K}}\{\omega^{K} = (\omega_{1}, ..., \omega_{K}): \operatorname{Card}\{k \leq K: \chi_{k}(\omega_{k}) = 1\} \geq M\},$
(2.214)

where $P_{\omega^{\kappa}}$ is as in (2.211), and $\Theta_k(\cdot)$ is defined in (2.212).

Now let us define $\psi_k^+(\lambda^k)$ as follows:

- when $\psi_k(\Theta_k(\lambda^k)) = 1$, or, which is the same, $\chi_k(\omega_k = \theta_k(\zeta^k)) = 1$, we set $\psi_k^+(\lambda^k) = 1$ as well;
- when $\psi_k(\Theta_k(\lambda^k)) = 0$, or, which is the same, $\chi_k(\omega_k = \theta_k(\zeta^k)) = 0$, we set

 $\psi_k^+(\lambda^k) = 1$ whenever

$$\eta_k \le \gamma_k(\lambda^{k-1}) := \frac{\epsilon_\star - \pi_k(\lambda^{k-1})}{1 - \pi_k(\lambda^{k-1})}$$

and $\psi_k^+(\lambda^k) = 0$ otherwise.

Let us make the following immediate observations:

(A) Whenever λ^k is such that $\psi_k(\mu_k = \Theta_k(\lambda^k)) = 1$, we have also $\psi_k^+(\lambda^k) = 1$; (B) The conditional, $\lambda^{k-1} = [\zeta^{k-1}; \eta^{k-1}]$ being fixed, probability of the event

$$\psi_k^+(\lambda^k) = 1$$

is exactly ϵ_{\star} .

Indeed, let $P_{\lambda_k|\lambda^{k-1}}$ be the conditional, λ^{k-1} being fixed, distribution of λ_k . Let us fix λ^{k-1} . The event $E = \{\lambda_k : \psi_k^+(\lambda^k) = 1\}$, by construction, is the union of two nonoverlapping events:

$$E_1 = \{\lambda_k = [\zeta_k; \eta_k] : \chi_k(\theta_k(\zeta^k)) = 1\}; \\ E_2 = \{\lambda_k = [\zeta_k; \eta_k] : \chi_k(\theta_k(\zeta^k)) = 0, \eta_k \le \gamma_k(\lambda^k)\}$$

Taking into account that the conditional, λ^{k-1} fixed, distribution of $\mu_k = [\omega_k =$ $\theta_k(\zeta^k); \eta_k$ is the product distribution $P_{\omega_k|\zeta^{k-1}} \times U$, we conclude in view of (2.213) that

$$P_{\lambda_{k}|\lambda^{k-1}}\{E_{1}\} = P_{\omega_{k}|\zeta^{k-1}}\{\omega_{k}:\chi_{k}(\omega_{k})=1\} = \pi_{k}(\lambda^{k-1}), P_{\lambda_{k}|\lambda^{k-1}}\{E_{2}\} = P_{\omega_{k}|\zeta^{k-1}}\{\omega_{k}:\chi_{k}(\omega_{k})=0\}U\{\eta \leq \gamma_{k}(\lambda^{k-1})\} = (1 - \pi_{k}(\lambda^{k-1}))\gamma_{k}(\lambda^{k-1}),$$

which combines with the definition of $\gamma_k(\cdot)$ to imply (B).

2.11.2.C Achieving the goal, step 2. By (A) combined with (2.214) we have

$$P_{\omega^{K}}\{\omega^{K} : \operatorname{Card}\{k \leq K : \chi_{k}(\omega_{k}) = 1\} \geq M\}$$

= $P_{\lambda^{K}}\{\lambda^{K} : \operatorname{Card}\{k \leq K : \psi_{k}(\mu_{k} = \Theta_{k}(\lambda^{k})) = 1\} \geq M\}$
 $\leq P_{\lambda^{K}}\{\lambda^{K} : \operatorname{Card}\{k \leq K : \psi_{k}^{+}(\lambda^{k}) = 1\} \geq M\},$

and all we need to verify is that the first quantity in this chain is upper-bounded by the quantity ϵ_M given by (2.211). Invoking the chain and (B), it is enough to justify the following claim:

(!) Let $\lambda^K = (\lambda_1, ..., \lambda_K)$ be a random sequence with probability distribution P, let $\psi_k(\lambda^k)$ take values 0 and 1 only, and let for every $k \leq K$ the conditional, λ^{k-1} being fixed, probability for $\psi_k^+(\lambda^k)$ to take value 1 is, for all λ^{k-1} , equal to ϵ_{\star} . Then the P-probability of the event

$$\{\lambda^K : \operatorname{Card}\{k \le K : \psi_k^+(\lambda_k) = 1\} \ge M\}$$

is exactly equal to ϵ_M given by (2.211).

This is immediate. For integers k, m, $1 \le k \le K$, $m \ge 0$, Let $\chi_m^k(\lambda^k)$ be the

characteristic function of the event

$$\{\lambda^k : \operatorname{Card}\{t \le k : \psi_t^+(\lambda^t) = 1\} = m\},\$$

and let

$$\pi_m^k = P\{\lambda^K : \chi_m^k(\lambda^k) = 1\}.$$

We have the following evident recurrence:

$$\chi_m^k(\lambda^k) = \chi_m^{k-1}(\lambda^{k-1})(1 - \psi_k^+(\lambda^k)) + \chi_{m-1}^{k-1}(\lambda^{k-1})\psi_k^+(\lambda^k), \ k = 1, 2, \dots$$

augmented by the "boundary conditions" $\chi_m^0 = 0, m > 0, \chi_0^0 = 1, \chi_{-1}^{k-1} = 0$ for all $k \ge 1$. Taking expectation w.r.t. P and utilizing the fact that conditional, λ^{k-1} being given, expectation of $\psi_k^+(\lambda^k)$ is, identically in λ^{k-1} , equal to ϵ_{\star} , we get

$$\pi_m^k = \pi_m^{k-1}(1-\epsilon_\star) + \pi_{m-1}^{k-1}\epsilon_\star, \ k = 1, \dots, K, \ \pi_m^0 = \left\{ \begin{array}{cc} 1, & m = 0 \\ 0, & m > 0 \end{array}, \ \pi_{-1}^{k-1} = 0, \ k = 1, 2, \dots \right. \right.$$

whence

$$\pi_m^k = \begin{cases} \binom{k}{m} \epsilon_\star^m (1 - \epsilon_\star)^{k-m}, & m \le k \\ 0, & m > k \end{cases}$$

and therefore

$$P\{\lambda^{K}: \operatorname{Card}\{k \le K: \psi_{k}^{+}(\lambda^{k}) = 1\} \ge M\} = \sum_{M \le k \le K} \pi_{k}^{K} = \epsilon_{M},$$

as required.

2.11.3 Proof of Proposition 2.40

All we need is to verify (2.134) and to check that the right hand side function in this relation is convex. The latter is evident, since $\phi_X(h) + \phi_X(-h) \ge 2\phi_X(0) = 0$ and $\phi_X(h) + \phi_X(-h)$ is convex. To verify (2.134), let us fix $P \in \mathcal{P}[X]$ and $h \in \mathbf{R}^d$ and set

$$\nu = h^T e[P],$$

so that ν is the expectation of $h^T \omega$ with $\omega \sim P$. Note that $-\phi_X(-h) \leq \nu \leq \phi_X(h)$, so that (2.134) definitely holds true when $\phi_X(h) + \phi_X(-h) = 0$. Now let

$$\eta := \frac{1}{2} \left[\phi_X(h) + \phi_X(-h) \right] > 0,$$

and let

$$a = \frac{1}{2} \left[\phi_X(h) - \phi_X(-h) \right], \ \beta = (\nu - a)/\eta.$$

Denoting by P_h the distribution of $h^T \omega$ induced by the distribution P of ω and noting that this distribution is supported on $[-\phi_X(-h), \phi_X(h)] = [a - \eta, a + \eta]$ and has expectation ν , we get

$$\beta \in [-1,1]$$

and

$$\gamma := \int \exp\{h^T \omega\} P(d\omega) = \int_{a-\eta}^{a+\eta} [e^s - \lambda(s-\nu)] P_h(ds)$$

180

for all $\lambda \in \mathbf{R}$. Hence,

$$\begin{aligned} \ln(\gamma) &\leq \inf_{\lambda} \ln\left(\max_{a-\eta \leq s \leq a+\eta} [e^s - \lambda(s-\nu)]\right) \\ &= a + \inf_{\rho} \ln\left(\max_{-\eta \leq t \leq \eta} [e^t - \rho(t-[\nu-a])]\right) \\ &= a + \inf_{\rho} \ln\left(\max_{-\eta \leq t \leq \eta} [e^t - \rho(t-\eta\beta)]\right) \leq a + \ln\left(\max_{-\eta \leq t \leq \eta} [e^t - \bar{\rho}(t-\eta\beta)\right) \end{aligned}$$

with $\bar{\rho} = (2\eta)^{-1}(e^{\eta} - e^{-\eta})$. The function $g(t) = e^t - \bar{\rho}(t - \eta\beta)$ is convex on $[-\eta, \eta]$, and

$$g(-\eta) = g(\eta) = \cosh(\eta) + \beta \sinh(\eta),$$

which combines with the above computation to yield the relation

$$\ln(\gamma) \le a + \ln(\cosh(\eta) + \beta \sinh(\eta)), \qquad (2.215)$$

and all we need to verify is that

$$\forall (\eta > 0, \beta \in [-1, 1]) : \ \beta \eta + \frac{1}{2} \eta^2 - \ln(\cosh(\eta) + \beta \sinh(\eta)) \ge 0.$$
 (2.216)

Indeed, if (2.216) holds true (2.215) implies that

$$\ln(\gamma) \le a + \beta \eta + \frac{1}{2}\eta^2 = \nu + \frac{1}{2}\eta^2$$

which, recalling what γ , ν and η are, is exactly what we want to prove.

Verification of (2.216) is as follows. The left hand side in (2.216) is convex in β for $\beta > -\frac{\cosh(\eta)}{\sinh(\eta)}$ containing, due to $\eta > 0$, the range of β in (2.216). Furthermore, the minimum of the left hand side of (2.216) over $\beta > -\coth(\eta)$ is attained when $\beta = \frac{\sinh(\eta) - \eta \cosh(\eta)}{\eta \sinh(\eta)}$ and is equal to

$$r(\eta) = \frac{1}{2}\eta^2 + 1 - \eta \coth(\eta) - \ln(\sinh(\eta)/\eta).$$

All we need to prove is that the latter quantity is nonnegative whenever $\eta > 0$. We have

$$r'(\eta) = \eta - \coth(\eta) - \eta(1 - \coth^2(\eta)) - \coth(\eta) + \eta^{-1} = (\eta \coth(\eta) - 1)^2 \eta^{-1} \ge 0,$$

and since $r(+0) = 0$, we get $r(\eta) \ge 0$ when $\eta > 0$.

2.11.4 Proof of Proposition 2.46

2.11.4.A Proof of Proposition 2.46.i

1⁰. Let $b = [0; ...; 0; 1] \in \mathbf{R}^{n+1}$, so that $B = \begin{bmatrix} A \\ b^T \end{bmatrix}$, and let $\mathcal{A}(u) = A[u; 1]$. For any $u \in \mathbf{R}^n$, $h \in \mathbf{R}^d$, $\Theta \in \mathbf{S}^d_+$ and $H \in \mathbf{S}^d$ such that $-I \prec \Theta^{1/2} H \Theta^{1/2} \prec I$ we have

due to

$$h^{T}\mathcal{A}(u) = [u;1]^{T}bh^{T}A[u;1] = [u;1]^{T}A^{T}hb^{T}[u;1], \ H\mathcal{A}(u) + h = [H,h]B[u;1].$$

Observe that when $(h, H) \in \mathcal{H}^{\gamma}$, we have

$$\Theta^{1/2}[I - \Theta^{1/2}H\Theta^{1/2}]^{-1}\Theta^{1/2} = [\Theta^{-1} - H]^{-1} \preceq [\Theta_*^{-1} - H]^{-1},$$

so that (2.217) implies that for all $u \in \mathbf{R}^n$, $\Theta \in \mathcal{V}$, and $(h, H) \in \mathcal{H}^{\gamma}$,

$$\begin{split} \Psi(h,H;u,\Theta) &\leq -\frac{1}{2} \ln \operatorname{Det}(I-\Theta^{1/2}H\Theta^{1/2}) \\ +\frac{1}{2}[u;1]^T \underbrace{\left[bh^TA + A^Thb^T + A^THA + B^T[H,h]^T[\Theta_*^{-1} - H]^{-1}[H,h]B\right]}_{Q[H,h]}[u;1] \\ &= -\frac{1}{2} \ln \operatorname{Det}(I-\Theta^{1/2}H\Theta^{1/2}) + \frac{1}{2}\operatorname{Tr}(Q[H,h]Z(u)) \\ &\leq -\frac{1}{2} \ln \operatorname{Det}(I-\Theta^{1/2}H\Theta^{1/2}) + \Gamma_{\mathcal{Z}}(h,H), \\ \Gamma_{\mathcal{Z}}(h,H) &= \frac{1}{2}\phi_{\mathcal{Z}}(Q[H,h]) \end{split}$$
(2.218)

(we have taken into account that $Z(u) \in \mathcal{Z}$ when $u \in U$ (premise of the proposition) and therefore $\operatorname{Tr}(Q[H,h]Z(u)) \leq \phi_{\mathcal{Z}}(Q[H,h])$).

2^0 . We need the following

Lemma 2.76. Let Θ_* be a $d \times d$ symmetric positive definite matrix, let $\delta \in [0, 2]$, and let \mathcal{V} be a closed convex subset of \mathbf{S}^d_+ such that

$$\Theta \in \mathcal{V} \Rightarrow \{\Theta \preceq \Theta_*\} \& \{ \|\Theta^{1/2}\Theta_*^{-1/2} - I\| \le \delta \}$$

$$(2.219)$$

(cf. (2.165)). Let also $\mathcal{H}^o := \{ H \in \mathbf{S}^d : -\Theta_*^{-1} \prec H \prec \Theta_*^{-1} \}$. Then

$$\begin{aligned} \forall (H, \Theta) \in \mathcal{H}^{o} \times \mathcal{V} : \\ G(H; \Theta) &:= -\frac{1}{2} \ln \operatorname{Det}(I - \Theta^{1/2} H \Theta^{1/2}) \\ &\leq G^{+}(H; \Theta) := -\frac{1}{2} \ln \operatorname{Det}(I - \Theta^{1/2}_{*} H \Theta^{1/2}_{*}) + \frac{1}{2} \operatorname{Tr}([\Theta - \Theta_{*}] H) \\ &\quad + \frac{\delta(2 + \delta)}{2(1 - ||\Theta^{1/2}_{*} H \Theta^{1/2}_{*}||)} ||\Theta^{1/2}_{*} H \Theta^{1/2}_{*}||_{F}^{2}, \end{aligned}$$

$$(2.220)$$

where $\|\cdot\|$ is the spectral, and $\|\cdot\|_F$ - the Frobenius norm of a matrix. In addition, $G^+(H,\Theta)$ is continuous function on $\mathcal{H}^o \times \mathcal{V}$ which is convex in $H \in H^o$ and concave (in fact, affine) in $\Theta \in \mathcal{V}$

Proof. Let us set

$$d(H) = \|\Theta_*^{1/2} H \Theta_*^{1/2} \|,$$

so that d(H) < 1 for $H \in \mathcal{H}^o$. For $H \in \mathcal{H}^o$ and $\Theta \in \mathcal{V}$ fixed we have

$$\begin{aligned} \|\Theta^{1/2}H\Theta^{1/2}\| &= \|[\Theta^{1/2}\Theta_{*}^{-1/2}][\Theta_{*}^{1/2}H\Theta_{*}^{1/2}][\Theta^{1/2}\Theta_{*}^{-1/2}]^{T}\| \\ &\leq \|\Theta^{1/2}\Theta_{*}^{-1/2}\|^{2}\|\Theta_{*}^{1/2}H\Theta_{*}^{1/2}\| \leq \|\Theta_{*}^{1/2}H\Theta_{*}^{1/2}\| = d(H) \end{aligned}$$

$$(2.221)$$

(we have used the fact that $0 \leq \Theta \leq \Theta_*$ implies $\|\Theta^{1/2}\Theta_*^{-1/2}\| \leq 1$). Noting that $\|AB\|_F \leq \|A\| \|B\|_F$, computation completely similar to the one in (2.221) yields

$$\|\Theta^{1/2}H\Theta^{1/2}\|_F \le \|\Theta_*^{1/2}H\Theta_*^{1/2}\|_F =: D(H)$$
(2.222)

Besides this, setting $F(X) = -\ln \text{Det}(X) : \text{int } \mathbf{S}^d_+ \to \mathbf{R}$ and equipping \mathbf{S}^d with the Frobenius inner product, we have $\nabla F(X) = -X^{-1}$, so that with $R_0 = \Theta^{1/2}_* H \Theta^{1/2}_*$, $R_1 = \Theta^{1/2} H \Theta^{1/2}$, and $\Delta = R_1 - R_0$, we have for properly selected $\lambda \in (0, 1)$ and $R_\lambda = \lambda R_0 + (1 - \lambda) R_1$:

$$F(I-R_1) = F(I-R_0-\Delta) = F(I-R_0) + \langle \nabla F(I-R_\lambda), -\Delta \rangle$$

= $F(I-R_0) + \langle (I-R_\lambda)^{-1}, \Delta \rangle$
= $F(I-R_0) + \langle I, \Delta \rangle + \langle (I-R_\lambda)^{-1} - I, \Delta \rangle.$

We conclude that

$$F(I - R_1) \le F(I - R_0) + \text{Tr}(\Delta) + \|I - (I - R_\lambda)^{-1}\|_F \|\Delta\|_F.$$
 (2.223)

Denoting by μ_i the eigenvalues of R_{λ} and noting that $||R_{\lambda}|| \leq \max[||R_0||, ||R_1||] = d(H)$ (see (2.221)), we have $|\mu_i| \leq d(H)$, and therefore eigenvalues $\nu_i = 1 - \frac{1}{1-\mu_i} = -\frac{\mu_i}{1-\mu_i}$ of $I - (I - R_{\lambda})^{-1}$ satisfy $|\nu_i| \leq |\mu_i|/(1-\mu_i) \leq |\mu_i|/(1-d(H))$, whence

$$||I - (I - R_{\lambda})^{-1}||_F \le ||R_{\lambda}||_F / (1 - d(H)).$$

Noting that $||R_{\lambda}||_F \leq \max[||R_0||_F, ||R_1||_F] \leq D(H)$, see (2.222), we conclude that $||I - (I - R_{\lambda})^{-1}||_F \leq D(H)/(1 - d(H))$, so that (2.223) yields

$$F(I - R_1) \le F(I - R_0) + \operatorname{Tr}(\Delta) + D(H) \|\Delta\|_F / (1 - d(H)).$$
(2.224)

Further, by (2.165) the matrix $D = \Theta^{1/2} \Theta_*^{-1/2} - I$ satisfies $\|D\| \le \delta$, whence

$$\Delta = \underbrace{\Theta^{1/2} H \Theta^{1/2}}_{R_1} - \underbrace{\Theta^{1/2}_* H \Theta^{1/2}_*}_{R_0} = (I+D) R_0 (I+D^T) - R_0 = DR_0 + R_0 D^T + DR_0 D^T.$$

Consequently,

$$\begin{aligned} \|\Delta\|_F &\leq \|DR_0\|_F + \|R_0D^T\|_F + \|DR_0D^T\|_F \leq [2\|D\| + \|D\|^2] \|R_0\|_F \\ &\leq \delta(2+\delta) \|R_0\|_F = \delta(2+\delta)D(H). \end{aligned}$$

This combines with (2.224) and the relation

$$\operatorname{Tr}(\Delta) = \operatorname{Tr}(\Theta^{1/2}H\Theta^{1/2} - \Theta_*^{1/2}H\Theta_*^{1/2}) = \operatorname{Tr}([\Theta - \Theta_*]H)$$

to yield

$$F(I - R_1) \le F(I - R_0) + \operatorname{Tr}([\Theta - \Theta_*]H) + \frac{\delta(2 + \delta)}{1 - d(H)} \|\Theta_*^{1/2} H \Theta_*^{1/2} \|_F^2,$$

and we arrive at (2.220). It remains to prove that $G^+(H;\Theta)$ is convex-concave and continuous on $\mathcal{H}^o \times \mathcal{V}$. The only component of this claim which is not completely evident is convexity of the function in $H \in \mathcal{H}^o$; to see that it is the case, note that $\ln \operatorname{Det}(S)$ is concave on the interior of the semidefinite cone, the function $f(u,v) = \frac{u^2}{1-v}$ is convex and nondecreasing in u, v in the convex domain $\Pi =$ $\{(u,v) : u \geq 0, v < 1\}$, and the function $\frac{\|\Theta_*^{1/2}H\Theta_*^{1/2}\|_F^2}{1-\|\Theta_*^{1/2}H\Theta_*^{1/2}\|}$ is obtained from f by convex substitution of variables $H \mapsto (\|\Theta_*^{1/2}H\Theta_*^{1/2}\|_F, \|\Theta_*^{1/2}H\Theta_*^{1/2}\|)$ mapping \mathcal{H}^o into Π . \Box

3⁰. Combining (2.220), (2.218), (2.167) and the origin of Ψ , see (2.217), we arrive at

$$\forall ((u, \Theta) \in U \times \mathcal{V}, (h, H) \in \mathcal{H}^{\gamma} = \mathcal{H}) : \\ \ln \left(\mathbf{E}_{\zeta \sim \mathcal{N}(A[u; 1], \Theta)} \left\{ \exp\{h^{T} \zeta + \frac{1}{2} \zeta^{T} H \zeta\} \right\} \right) \leq \Phi_{A, \mathcal{Z}}(h, H; \Theta),$$

as claimed in (2.170).

4⁰. Now let us check that $\Phi_{A,\mathcal{Z}}(h, H; \Theta) : \mathcal{H} \times \mathcal{V} \to \mathbf{R}$ is continuous and convexconcave. Recalling that the function $G^+(H; \Theta)$ from (2.220) is convex-concave and continuous on $\mathcal{H}^o \times \mathcal{V}$, all we need to verify is that $\Gamma_{\mathcal{Z}}(h, H)$ is convex and continuous on \mathcal{H} . Recalling that \mathcal{Z} is nonempty compact set, the function $\phi_{\mathcal{Z}}(\cdot) : \mathbf{S}^{d+1} \to \mathbf{R}$ is continuous, implying the continuity of $\Gamma_{\mathcal{Z}}(h, H) = \frac{1}{2}\phi_{\mathcal{Z}}(Q[H, h])$ on $\mathcal{H} = \mathcal{H}^{\gamma}$ (Q[H, h] is defined in (2.218)). To prove convexity of $\Gamma_{\mathcal{Z}}$, note that \mathcal{Z} is contained in \mathbf{S}^{n+1}_+ , implying that $\phi_{\mathcal{Z}}(\cdot)$ is convex and \succeq -monotone. On the other hand, by Schur Complement Lemma, we have

$$\begin{split} S &:= & \left\{ (h,H,G) : G \succeq Q[H,h], (h,H) \in \mathcal{H}^{\gamma} \right\} \\ &= & \left\{ (h,H,G) : \left[\begin{array}{c|c} G - [bh^T A + A^T h b^T + A^T H A] & B^T [H,h]^T \\ \hline & [H,h] B & \Theta_*^{-1} - H \end{array} \right] \succeq 0, \\ & (h,H) \in \mathcal{H}^{\gamma} \right\}, \end{split}$$

implying that S is convex. Since $\phi_{\mathcal{Z}}(\cdot)$ is \succeq -monotone, we have

$$\begin{aligned} &\{(h,H,\tau):(h,H)\in\mathcal{H}^{\gamma},\ \tau\geq\Gamma_{\mathcal{Z}}(h,H)\}\\ &=\{(h,H,\tau):\ \exists G:G\succeq Q[H,h],\ 2\tau\geq\phi_{\mathcal{Z}}(G),\ (h,H)\in\mathcal{H}^{\gamma}\},\end{aligned}$$

and we see that the epigraph of $\Gamma_{\mathcal{Z}}$ is convex (since the set S and the epigraph of $\phi_{\mathcal{Z}}$ are so), as claimed.

5⁰. It remains to prove that $\Phi_{\mathcal{A},\mathcal{Z}}$ is coercive in H, h. Let $\Theta \in \mathcal{V}$ and $(h_i, H_i) \in \mathcal{H}^{\gamma}$ with $\|(h_i, H_i)\| \to \infty$ as $i \to \infty$, and let us prove that $\Phi_{\mathcal{A},\mathcal{Z}}(h_i, H_i; \Theta) \to \infty$. Looking at the expression for $\Phi_{\mathcal{A},\mathcal{Z}}(h_i, H_i; \Theta)$, it is immediately seen that all terms in this expression, except for the terms coming from $\phi_{\mathcal{Z}}(\cdot)$, remain bounded as i grows, so that all we need to verify is that the $\phi_{\mathcal{Z}}(\cdot)$ -term goes to ∞ as $i \to \infty$. Observe that H_i are uniformly bounded due to $(h_i, H_i) \in \mathcal{H}^{\gamma}$, implying that

 $\|h_i\|_2 \to \infty$ as $i \to \infty$. Denoting $e = [0; ...; 0; 1] \in \mathbf{R}^{d+1}$ and, as before, $b = [0; ...; 0; 1] \in \mathbf{R}^{n+1}$, note that, by construction, $B^T e = b$. Now let $W \in \mathcal{Z}$, so that $W_{n+1,n+1} = 1$. Taking into account that the matrices $[\Theta_*^{-1} - H_i]^{-1}$ satisfy $\alpha I_d \preceq [\Theta_*^{-1} - H_i]^{-1} \preceq \beta I_d$ for some positive α, β due to $H_i \in \mathcal{H}^{\gamma}$, observe that

$$\underbrace{\left[\left[\begin{array}{c|c} H_i & h_i \\ \hline h_i^T & \end{array}\right] + [H_i, h_i]^T \left[\Theta_*^{-1} - H_i\right]^{-1} [H_i, h_i]\right]}_{Q_i} = \underbrace{\left[\begin{array}{c} h_i^T [\Theta_*^{-1} - H_i]^{-1} h_i \right]}_{\alpha_i \|h_i\|_2^2} ee^T + R_i,$$

where $\alpha_i \geq \alpha > 0$ and $||R_i||_F \leq C(1 + ||h_i||_2)$. As a result,

$$\phi_{\mathcal{Z}}(B^{T}Q_{i}B) \geq \operatorname{Tr}(WB^{T}Q_{i}B) = \operatorname{Tr}(WB^{T}[\alpha_{i}\|h_{i}\|_{2}^{2}ee^{T} + R_{i}]B)$$

$$= \alpha_{i}\|h_{i}\|_{2}^{2} \underbrace{\operatorname{Tr}(Wbb^{T})}_{=W_{n+1,n+1}=1} - \|BWB^{T}\|_{F}\|R_{i}\|_{F}$$

$$\geq \alpha\|h_{i}\|_{2}^{2} - C(1 + \|h_{i}\|_{2})\|BWB^{T}\|_{F},$$

and the concluding quantity tends to ∞ as $i \to \infty$ due to $||h_i||_2 \to \infty$, $i \to \infty$. Part (i) is proved.

2.11.4.B Proof of Proposition 2.46.ii

By (i) the function $\Phi(h, H; \Theta_1, \Theta_2)$ is continuous and convex-concave on the domain $(\mathcal{H}_1 \cap \mathcal{H}_2) \times (\mathcal{V}_1 \times \mathcal{V}_2)$ and are coercive in (h, H), while \mathcal{H} and \mathcal{V} are closed and \mathcal{V}

convex, and \mathcal{V} in addition is compact, saddle point problem (2.172) is solvable (Sion-Kakutani Theorem, a.k.a. Theorem 2.24). Now let $(h_*, H_*; \Theta_1^*, \Theta_2^*)$ be a saddle point. To prove (2.174), let $P \in \mathcal{G}_1$, that is, $P = \mathcal{N}(A_1[u; 1], \Theta_1)$ for some $\Theta_1 \in \mathcal{V}_1$ and some u with $[u; 1][u; 1]^T \in \mathcal{Z}_1$. Applying (2.170) to the first collection of data, with a given by (2.173), we get the first \leq in the following chain:

$$\ln\left(\int e^{-\frac{1}{2}\omega^{T}H_{*}\omega-\omega^{T}h_{*}-a}P(d\omega)\right) \leq \Phi_{A_{1},\mathcal{Z}_{1}}(-h_{*},-H_{*};\Theta_{1})-a \underbrace{\leq}_{(a)} \Phi_{A_{1},\mathcal{Z}_{1}}(-h_{*},-H_{*};\Theta_{1}^{*})-a \underbrace{=}_{(b)} \mathcal{S}V,$$

where (a) is due to the fact that $\Phi_{A_1,\mathcal{Z}_1}(-h_*,-H_*;\Theta_1)+\Phi_{A_2,\mathcal{Z}_2}(h_*,H_*;\Theta_2)$ attains its maximum over $(\Theta_1,\Theta_2) \in \mathcal{V}_1 \times \mathcal{V}_2$ at the point (Θ_1^*,Θ_2^*) , and (b) is due to the origin of a and the relation $\mathcal{S}V = \frac{1}{2}[\Phi_{A_1,\mathcal{Z}_1}(-h_*,-H_*;\Theta_1^*)+\Phi_{A_2,\mathcal{Z}_2}(h_*,H_*;\Theta_2^*)]$. The bound in (2.174.a) is proved. Similarly, let $P \in \mathcal{G}_2$, that is, $P = \mathcal{N}(A_2[u;1],\Theta_2)$ for some $\Theta_2 \in \mathcal{V}_2$ and some u with $[u;1][u;1]^T \in \mathcal{Z}_2$. Applying (2.170) to the second collection of data, with the same a as above, we get the first \leq in the following chain:

$$\ln\left(\int e^{\frac{1}{2}\omega^{T}H_{*}\omega+\omega^{T}h_{*}+a}P(d\omega)\right) \leq \Phi_{A_{2},\mathcal{Z}_{2}}(h_{*},H_{*};\Theta_{2})+a$$

$$\leq \Phi_{A_{2},\mathcal{Z}_{2}}(h_{*},H_{*};\Theta_{2}^{*})+a = \mathcal{S}V,$$
(b)

with exactly the same as above justification of (a) and (b). The bound in (2.174.b) is proved.

2.11.5 Proof of Proposition 2.49

2.11.5.A Preliminaries

We start with the following result:

Lemma 2.77. Let Θ be a positive definite $d \times d$ matrix, and let

 $u \mapsto \mathcal{C}(u) = A[u; 1]$

be an affine mapping from \mathbf{R}^n into \mathbf{R}^d . Finally, let $h \in \mathbf{R}^d$, $H \in \mathbf{S}^d$ and $P \in \mathbf{S}^d$ satisfy the relations

$$0 \leq P \prec I_d \& P \succeq \overline{\Theta}^{1/2} H \overline{\Theta}^{1/2}.$$
 (2.225)

Then, setting $B = \begin{bmatrix} A \\ 0, ..., 0, 1 \end{bmatrix}$, for every $u \in \mathbf{R}^n$ it holds

$$\zeta \sim \mathcal{S}G(\mathcal{C}(u), \bar{\Theta}) \Rightarrow \ln\left(\mathbf{E}_{\zeta}\left\{e^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta}\right\}\right) \leq -\frac{1}{2}\ln\operatorname{Det}(I-P) + \frac{1}{2}[u; 1]^{T}B^{T}\left[\left[\frac{H}{h^{T}}\right] + [H, h]^{T}\bar{\Theta}^{1/2}[I-P]^{-1}\bar{\Theta}^{1/2}[H, h]\right]B[u; 1]$$

$$(2.226)$$

Equivalently (set $G = \overline{\Theta}^{-1/2} P \overline{\Theta}^{-1/2}$): Whenever $h \in \mathbf{R}^d$, $H \in \mathbf{S}^d$ and $G \in \mathbf{S}^d$ satisfy the relations

$$0 \preceq G \prec \bar{\Theta}^{-1} \& G \succeq H, \tag{2.227}$$

one has for every for every $u \in \mathbf{R}^n$:

$$\begin{aligned} \zeta \sim \mathcal{S}G(\mathcal{C}(u), \bar{\Theta}) \Rightarrow \ln\left(\mathbf{E}_{\zeta}\left\{\mathrm{e}^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta}\right\}\right) &\leq -\frac{1}{2}\ln\operatorname{Det}(I - \bar{\Theta}^{1/2}G\bar{\Theta}^{1/2}) \\ &+ \frac{1}{2}[u; 1]^{T}B^{T}\left[\left[\frac{H}{h^{T}}\right]^{h}\right] + [H, h]^{T}\left[\bar{\Theta}^{-1} - G\right]^{-1}[H, h]\right]B[u; 1] \end{aligned}$$

$$(2.228)$$

Proof. 1^0 . Let us start with the following observation:

Lemma 2.78. Let $\Theta \in \mathbf{S}^d_+$ and $S \in \mathbf{R}^{d \times d}$ be such that $S \Theta S^T \prec I_d$. Then for every $\nu \in \mathbf{R}^d$ one has

$$\ln\left(\mathbf{E}_{\xi\sim\mathcal{S}G(0,\Theta)}\left\{\mathrm{e}^{\nu^{T}S\xi+\frac{1}{2}\xi^{T}S^{T}S\xi}\right\}\right) \leq \ln\left(\mathbf{E}_{x\sim\mathcal{N}(\nu,I_{d})}\left\{\mathrm{e}^{\frac{1}{2}x^{T}S\Theta S^{T}x}\right\}\right)$$

= $-\frac{1}{2}\ln\operatorname{Det}(I_{d}-S\Theta S^{T})+\frac{1}{2}\nu^{T}\left[S\Theta S^{T}(I_{d}-S\Theta S^{T})^{-1}\right]\nu.$ (2.229)

Indeed, let $\xi \sim SG(0, \Theta)$ and $x \sim \mathcal{N}(\nu, I_d)$ be independent. We have

$$\mathbf{E}_{\xi} \left\{ e^{\nu^{T} S\xi + \frac{1}{2}\xi^{T} S^{T} S\xi} \right\} \underbrace{=}_{a} \mathbf{E}_{\xi} \left\{ \mathbf{E}_{x} \left\{ e^{[S\xi]^{T} x} \right\} \right\} = \mathbf{E}_{x} \left\{ \mathbf{E}_{\xi} \left\{ e^{[S^{T} x]^{T} \xi} \right\} \right\}$$
$$\underbrace{\leq}_{b} \mathbf{E}_{x} \left\{ e^{\frac{1}{2}x^{T} S \Theta S^{T} x} \right\},$$

where a is due to $x \sim \mathcal{N}(\nu, I_d)$ and b is due to $\xi \sim \mathcal{S}G(0, \Theta)$. We have verified the inequality in (2.229); the equality in (2.229) is given by direct computation.

2⁰. Now, in the situation described in Lemma 2.77, by continuity it suffices to prove (2.226) in the case when $P \succeq 0$ in (2.225) is replaced with $P \succ 0$. Under the premise of Lemma, given $u \in \mathbf{R}^n$ and assuming $P \succ 0$, let us set $\mu = \mathcal{C}(u) = A[u; 1]$, $\nu = P^{-1/2}\overline{\Theta}^{1/2}[H\mu + h], S = P^{1/2}\overline{\Theta}^{-1/2}$, so that $S\overline{\Theta}S^T = P \prec I_d$, and let

 $G = \overline{\Theta}^{-1/2} P \overline{\Theta}^{-1/2}$, so that $G \succeq H$. Let $\zeta \sim SG(\mu, \overline{\Theta})$. Representing ζ as $\zeta = \mu + \xi$ with $\xi \sim SG(0, \overline{\Theta})$, we have

It is immediately seen that the concluding quantity in this chain is nothing but the right hand side quantity in (2.226).

2.11.5.B Completing proof of Proposition 2.49.

1⁰. Let us prove (2.184.*a*). By Lemma 2.77 (see (2.228)) applied with $\overline{\Theta} = \Theta_*$, setting $\mathcal{C}(u) = A[u; 1]$, we have

$$\forall \left((h, H) \in \mathcal{H}, G : 0 \leq G \leq \gamma^{+} \Theta_{*}^{-1}, G \succeq H, u \in \mathbf{R}^{n} : [u; 1][u; 1]^{T} \in \mathcal{Z} \right) : \\ \ln \left(\mathbf{E}_{\zeta \sim \mathcal{S}G(\mathcal{C}(u),\Theta_{*})} \left\{ e^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta} \right\} \right) \leq -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2}G\Theta_{*}^{1/2}) \\ + \frac{1}{2}[u; 1]^{T}B^{T} \left[\left[\frac{H}{h^{T}} \right]^{h} + [H, h]^{T} [\Theta_{*}^{-1} - G]^{-1} [H, h] \right] B[u; 1] \\ \leq -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2}G\Theta_{*}^{1/2}) \\ + \frac{1}{2}\phi_{\mathcal{Z}} \left(B^{T} \left[\left[\frac{H}{h^{T}} \right]^{h} + [H, h]^{T} [\Theta_{*}^{-1} - G]^{-1} [H, h] \right] B \right) = \Psi_{A, \mathcal{Z}}(h, H, G),$$

$$(2.230)$$

implying, due to the origin of $\Phi_{A,\mathcal{Z}}$, that under the premise of (2.230) we have

$$\ln\left(\mathbf{E}_{\zeta\sim\mathcal{S}G(\mathcal{C}(u),\Theta_*)}\left\{\mathrm{e}^{h^T\zeta+\frac{1}{2}\zeta^TH\zeta}\right\}\right) \leq \Phi_{A,\mathcal{Z}}(h,H), \,\forall (h,H)\in\mathcal{H}.$$

Taking into account that when $\zeta \sim SG(\mathcal{C}(u), \Theta)$ with $\Theta \in \mathcal{V}$, we have also $\zeta \sim SG(\mathcal{C}(u), \Theta_*)$, (2.184.*a*) follows.

 2^{0} . Now let us prove (2.184.b). All we need is to verify the relation

$$\forall \left((h, H) \in \mathcal{H}, G : 0 \leq G \leq \gamma^+ \Theta_*^{-1}, G \geq H, u \in \mathbf{R}^n : [u; 1][u; 1]^T \in \mathcal{Z}, \Theta \in \mathcal{V} \right) : \\ \ln \left(\mathbf{E}_{\zeta \sim \mathcal{S}G(\mathcal{C}(u), \Theta)} \left\{ e^{h^T \zeta + \frac{1}{2} \zeta^T H \zeta} \right\} \right) \leq \Psi_{A, \mathcal{Z}}^{\delta}(h, H, G; \Theta);$$

$$(2.231)$$

with this relation at our disposal (2.184.b) can be obtained by the same argument as the one we used in item 1⁰ to derive (2.184.a).

To establish (2.231), let us fix h, H, G, u, Θ satisfying the premise of (2.231); note that under the premise of Proposition 2.49.i, we have $0 \leq \Theta \leq \Theta_*$. Now let $\lambda \in (0, 1)$, and let $\Theta_{\lambda} = \Theta + \lambda(\Theta_* - \Theta)$, so that $0 \prec \Theta_{\lambda} \leq \Theta_*$, and let $\delta_{\lambda} = \|\Theta_{\lambda}^{1/2} \Theta_*^{-1/2} - I\|$, so that $\delta_{\lambda} \in [0, 2]$. We have $0 \leq G \leq \gamma^+ \Theta_*^{-1} \leq \gamma^+ \Theta_{\lambda}^{-1}$ that is, H, G satisfy (2.227) w.r.t. $\overline{\Theta} = \Theta_{\lambda}$. As a result, for our h, G, H, u and the just

defined $\overline{\Theta}$ relation (2.228) holds true:

$$\begin{split} &\zeta \sim \mathcal{S}G(\mathcal{C}(u), \Theta_{\lambda}) \Rightarrow \\ &\ln\left(\mathbf{E}_{\zeta}\left\{e^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta}\right\}\right) \leq -\frac{1}{2}\ln\operatorname{Det}(I - \Theta_{\lambda}^{1/2}G\Theta_{\lambda}^{1/2}) \\ &+ \frac{1}{2}[u;1]^{T}B^{T}\left[\left[\frac{H}{h^{T}}\right] + [H,h]^{T}\left[\Theta_{\lambda}^{-1} - G\right]^{-1}[H,h]\right]B[u;1] \\ &\leq -\frac{1}{2}\ln\operatorname{Det}(I - \Theta_{\lambda}^{1/2}G\Theta_{\lambda}^{1/2}) \\ &+ \frac{1}{2}\phi_{\mathcal{Z}}\left(B^{T}\left[\left[\frac{H}{h^{T}}\right] + [H,h]^{T}\left[\Theta_{\lambda}^{-1} - G\right]^{-1}[H,h]\right]B\right) \end{split}$$
(2.232)

(recall that $[u; 1][u; 1]^T \in \mathcal{Z}$). As a result,

$$\begin{aligned} \zeta \sim \mathcal{S}G(\mathcal{C}(u),\Theta) \Rightarrow \ln\left(\mathbf{E}_{\zeta}\left\{e^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta}\right\}\right) &\leq -\frac{1}{2}\ln\operatorname{Det}(I - \Theta_{\lambda}^{1/2}G\Theta_{\lambda}^{1/2}) \\ &+ \frac{1}{2}\phi_{\mathcal{Z}}\left(B^{T}\left[\left[\frac{H}{h^{T}}\right] + \left[H,h\right]^{T}\left[\Theta_{*}^{-1} - G\right]^{-1}\left[H,h\right]\right]B\right) \end{aligned}$$

$$(2.233)$$

When deriving (2.233) from (2.232), we have used that

 $\begin{array}{l} -\Theta \preceq \Theta_{\lambda} \text{, so that when } \zeta \sim \mathcal{S}G(\mathcal{C}(u), \Theta) \text{, we have also } \zeta \sim \mathcal{S}G(\mathcal{C}(u), \Theta_{\lambda}), \\ - 0 \preceq \Theta_{\lambda} \preceq \Theta_{*} \text{ and } G \prec \Theta_{*}^{-1} \text{, whence } [\Theta_{\lambda}^{-1} - G]^{-1} \preceq [\Theta_{*}^{-1} - G]^{-1}, \\ - \mathcal{Z} \subset \mathbf{S}_{+}^{n+1} \text{, whence } \phi_{\mathcal{Z}} \text{ is } \succeq \text{-monotone: } \phi_{\mathcal{Z}}(M) \leq \phi_{\mathcal{Z}}(N) \text{ whenever } M \preceq N. \end{array}$

By Lemma 2.76 applied with Θ_{λ} in the role of Θ and δ_{λ} in the role of δ , we have

$$\begin{aligned} &-\frac{1}{2}\ln \operatorname{Det}(I - \Theta_{\lambda}^{1/2}G\Theta_{\lambda}^{1/2}) \\ &\leq -\frac{1}{2}\ln \operatorname{Det}(I - \Theta_{*}^{1/2}G\Theta_{*}^{1/2}) + \frac{1}{2}\operatorname{Tr}([\Theta_{\lambda} - \Theta_{*}]G) + \frac{\delta_{\lambda}(2 + \delta_{\lambda})}{2(1 - \|\Theta_{*}^{1/2}G\Theta_{*}^{1/2}\|)} \|\Theta_{*}^{1/2}G\Theta_{*}^{1/2}\|_{F}^{2} \end{aligned}$$

Consequently, (2.233) implies that

$$\begin{split} \zeta &\sim \mathcal{S}G(\mathcal{C}(u), \Theta) \Rightarrow \\ \ln\left(\mathbf{E}_{\zeta}\left\{\mathbf{e}^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta}\right\}\right) &\leq -\frac{1}{2}\ln\operatorname{Det}(I - \Theta_{*}^{1/2}G\Theta_{*}^{1/2}) + \frac{1}{2}\operatorname{Tr}([\Theta_{\lambda} - \Theta_{*}]G) \\ &+ \frac{\delta_{\lambda}(2 + \delta_{\lambda})}{2(1 - ||\Theta_{*}^{1/2}G\Theta_{*}^{1/2}||)} ||\Theta_{*}^{1/2}G\Theta_{*}^{1/2}||_{F}^{2} \\ &+ \frac{1}{2}\phi_{\mathcal{Z}}\left(B^{T}\left[\left[\frac{H}{h^{T}}\right] + H, h\right]^{T}[\Theta_{*}^{-1} - G]^{-1}[H, h]\right]B\right). \end{split}$$

The resulting inequality holds true for all small positive λ ; taking limit of the right hand side as $\lambda \to +0$, and recalling that $\Theta_0 = \Theta$, we get

$$\begin{split} \zeta &\sim \mathcal{S}G(\mathcal{C}(u), \Theta) \Rightarrow \\ \ln\left(\mathbf{E}_{\zeta} \left\{ e^{h^{T}\zeta + \frac{1}{2}\zeta^{T}H\zeta} \right\} \right) &\leq -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2}G\Theta_{*}^{1/2}) + \frac{1}{2}\operatorname{Tr}([\Theta - \Theta_{*}]G) \\ &+ \frac{\delta(2+\delta)}{2(1 - ||\Theta_{*}^{1/2}G\Theta_{*}^{1/2}||)} ||\Theta_{*}^{1/2}G\Theta_{*}^{1/2}||_{F}^{2} \\ &+ \frac{1}{2}\phi_{\mathcal{Z}} \left(B^{T} \left[\left[\frac{H}{h^{T}} \right] + [H, h]^{T} \left[\Theta_{*}^{-1} - G \right]^{-1} [H, h] \right] B \right) \end{split}$$

(note that under the premise of Proposition 2.49.i we clearly have $\liminf_{A\to+0} \delta_A \leq \delta$). The right hand side of the resulting inequality is nothing but $\Psi_{A,\mathcal{Z}}^{\delta}(h, H, G; \Theta)$, see (2.183), and we arrive at the inequality required in the conclusion of (2.231).

3⁰. To complete the proof of Proposition 2.49.i, it remains to prove that the functions $\Phi_{A,\mathcal{Z}}$, $\Phi_{A,\mathcal{Z}}^{\delta}$ possess the announced in Proposition continuity, convexity-concavity, and coerciveness properties. Let us verify that this indeed is so for $\Phi_{A,\mathcal{Z}}^{\delta}$;

reasoning to follow, with evident simplifications, is applicable to $\Phi_{A,\mathcal{Z}}$ as well.

Observe, first, that by exactly the same reasons as in item 4^0 of the proof of Proposition 2.46, the function $\Psi^{\delta}_{A,\mathcal{Z}}(h, H, G; \Theta)$ is real valued, continuous and convex-concave on the domain

$$\widehat{\mathcal{H}} \times \mathcal{V} = \{(h, H, G) : -\gamma^+ \Theta_*^{-1} \preceq H \preceq \gamma^+ \Theta_*^{-1}, 0 \preceq G \preceq \gamma^+ \Theta_*^{-1}, H \preceq G\} \times \mathcal{V}.$$

The function $\Phi_{A,\mathbb{Z}}^{\delta}(h, H; \Theta) : \mathcal{H} \times \mathcal{V} \to \mathbf{R}$ is obtained from $\Psi^{\delta}(h, H, G; \Theta)$ by the following two operations: we first minimize $\Psi_{A,\mathbb{Z}}^{\delta}(h, H, G; \Theta)$ over G linked to (h, H) by the convex constraints $0 \leq G \leq \gamma^+ \Theta_*^{-1}$ and $G \succeq H$, thus obtaining a function

$$\bar{\Phi}(h,H;\Theta):\underbrace{\{(h,H):-\gamma^+\Theta_*^{-1} \leq H \leq \gamma^+\Theta_*^{-1}\}}_{\bar{\mathcal{H}}} \times \mathcal{V} \to \mathbf{R} \cup \{+\infty\} \cup \{-\infty\}.$$

Second, we restrict the function $\overline{\Phi}(h, H; \Theta)$ from $\overline{\mathcal{H}} \times \mathcal{V}$ onto $\mathcal{H} \times \mathcal{V}$. For $(h, H) \in \overline{\mathcal{H}}$, the set of G's linked to (h, H) by the above convex constraints clearly is a nonempty compact set; as a result, $\overline{\Phi}$ is real-valued convex-concave function on $\overline{\mathcal{H}} \times \mathcal{V}$. From continuity of $\Psi_{A,Z}^{\delta}$ on its domain it immediately follows that $\Psi_{A,Z}^{\delta}$ is bounded and uniformly continuous on every bounded subset of this domain, implying by evident reasons that $\overline{\Phi}(h, H; \Theta)$ is bounded in every domain of the form $\overline{B} \times \mathcal{V}$, where \overline{B} is a bounded subset of $\overline{\mathcal{H}}$, and is continuous on $\overline{B} \times \mathcal{V}$ in $\Theta \in \mathcal{V}$ with properly selected modulus of continuity independent of $(h, H) \in \overline{B}$. Besides this, by construction, $\mathcal{H} \subset \operatorname{int} \overline{\mathcal{H}}$, implying that if B is a convex compact subset of \mathcal{H} , it belongs to the interior of a properly selected convex compact subset \overline{B} of $\overline{\mathcal{H}}$. Since $\overline{\Phi}$ is bounded on $\overline{B} \times \mathcal{V}$ and is convex in (h, H), the function $\overline{\Phi}$ is Lipschitz continuous in $(h, H) \in B$ with Lipschitz constant which can be selected to be independent of $\Theta \in \mathcal{V}$. Taking into account that \mathcal{H} is convex and closed, the bottom line is that $\Phi_{A,Z}^{\delta}$ is not just real-valued convex-concave function on the domain $\mathcal{H} \times \mathcal{V}$, it is also continuous on this domain.

Coerciveness of $\Phi_{A,\mathbb{Z}}^{\delta}(h, H; \Theta)$ in (h, H) is proved in exactly the same fashion as the similar property of function (2.167), see item 5⁰ in the proof of Proposition 2.46. The proof of item (i) of Proposition 2.49 is complete.

 4^{0} . Item (ii) of Proposition 2.49 can be derived from item (i) of Proposition in exactly the same fashion as when proving Proposition 2.46.

Lecture Three

Estimating Functions via Hypothesis Testing

In this Lecture we apply the hypothesis testing techniques developed in Lecture 2 to estimating properly structured scalar functionals in simple o.s.'s (Section 3.2) and beyond (Section 3.4).

3.1 ESTIMATING LINEAR FORMS ON UNIONS OF CONVEX SETS

3.1.1 The problem

Let $\mathcal{O} = ((\Omega, \Pi), \{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}, \mathcal{F})$ be a simple observation scheme. The problem we are interested in this section is as follows:

We are given a positive integer K and I nonempty convex compact sets $X_j \subset \mathbf{R}^n$, along with affine mappings $A_j(\cdot) : \mathbf{R}^n \to \mathbf{R}^M$ such that $A_j(x) \in \mathcal{M}$ whenever $x \in X_j$, $1 \leq j \leq I$. In addition, we are given a linear function $g^T x$ on \mathbf{R}^n .

Given random observation

$$\omega^K = (\omega_1, ..., \omega_K)$$

with ω_k drawn, independently across k, from $p_{A_j(x)}$ with $j \leq I$ and $x \in X_j$, we want to recover $g^T x$. It should be stressed that we do *not* know neither j nor x underlying our observation.

Given reliability tolerance $\epsilon \in (0, 1)$, we quantify the performance of a candidate estimate – a Borel function $\hat{g}(\cdot) : \Omega \to \mathbf{R}$ – by the worst case, over j and x, width of $(1 - \epsilon)$ -confidence interval, Specifically, we say that $\hat{g}(\cdot)$ is (ρ, ϵ) -reliable, if

$$\forall (j \leq I, x \in X_j) : \operatorname{Prob}_{\omega \sim p_{A,(x)}} \{ |\widehat{g}(\omega) - g^T x| > \rho \} \leq \epsilon.$$

We define ϵ -risk of the estimate as

$$\operatorname{Risk}_{\epsilon}[\widehat{g}] = \inf \{ \rho : \widehat{g} \text{ is } (\rho, \epsilon) \text{-reliable} \};$$

note that \hat{g} is the smallest ρ such that \hat{g} is (ρ, ϵ) -reliable.

We remark that the technique we are about to use originates from [86] where recovery, in a simple o.s., of a linear form on a convex compact set (i.e., the case I = 1 of the estimation problem at hand) was considered; it was proved that in this situation the estimate

$$\widehat{g}(\omega^K) = \sum_k \phi(\omega_k) + \varkappa$$

with properly selected $\phi \in \mathcal{F}$ and $\kappa \in \mathbf{R}$ is near-optimal; for Gaussian o.s. similar

ESTIMATING FUNCTIONS VIA HYPOTHESIS TESTING

191

fact was discovered, by different technique, by D. Donoho [44] as early as in 1994.

3.1.2 The estimate

In the sequel, we associate with the simple o.s. $\mathcal{O} = ((\Omega, \Pi), \{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}, \mathcal{F})$ in question the function

$$\Phi_{\mathcal{O}}(\phi;\mu) = \ln\left(\int e^{\phi(\omega)} p_{\mu}(\omega) \Pi(d\omega)\right) : \mathcal{F} \times \mathcal{M} \to \mathbf{R}.$$

Recall that by definition of a simple o.s., this function is real-valued on its domain and is concave in $\mu \in \mathcal{M}$, convex in $\phi \in \mathcal{F}$, and continuous on $\mathcal{F} \times \mathcal{M}$ (the latter follows from convexity-concavity and relative openness of \mathcal{M} and \mathcal{F}).

Let us associate with a pair $(i, j), 1 \leq i, j \leq I$, the functions

$$\begin{split} \Phi_{ij}(\alpha,\phi;x,y) &= \frac{1}{2} \begin{bmatrix} K\alpha \Phi_{\mathcal{O}}(\phi/\alpha;A_{i}(x)) + K\alpha \Phi_{\mathcal{O}}(-\phi/\alpha;A_{j}(y)) \\ &+ g^{T}(y-x) + 2\alpha \ln(2I/\epsilon) \end{bmatrix} : \{\alpha > 0, \phi \in \mathcal{F}\} \times [X_{i} \times X_{j}] \to \mathbf{R}, \\ \Psi_{ij}(\alpha,\phi) &= \max_{x \in X_{i}, y \in X_{j}} \Phi_{ij}(\alpha,\phi;x,y) \\ &= \frac{1}{2} \left[\Psi_{i,+}(\alpha,\phi) + \Psi_{j,-}(\alpha,\phi) \right] : \{\alpha > 0\} \times \mathcal{F} \to \mathbf{R}, \\ \Psi_{\ell,+}(\beta,\psi) &= \max_{x \in X_{\ell}} \begin{bmatrix} K\beta \Phi_{\mathcal{O}}(\psi/\beta;A_{\ell}(x)) - g^{T}x + \beta \ln(2I/\epsilon) \end{bmatrix} : \\ &\{\beta > 0, \psi \in \mathcal{F}\} \to \mathbf{R}, \\ \Psi_{\ell,-}(\beta,\psi) &= \max_{x \in X_{\ell}} \begin{bmatrix} K\beta \Phi_{\mathcal{O}}(-\psi/\beta;A_{\ell}(x)) + g^{T}x + \beta \ln(2I/\epsilon) \end{bmatrix} : \\ &\{\beta > 0, \psi \in \mathcal{F}\} \to \mathbf{R}. \\ &\{\beta > 0, \psi \in \mathcal{F}\} \to \mathbf{R}. \end{split}$$

$$\end{split}$$

Note that the function $\alpha \Phi_{\mathcal{O}}(\phi/\alpha; A_i(x))$ is obtained from continuous convex-concave function $\Phi_{\mathcal{O}}(\cdot, \cdot)$ by projective transformation in the convex argument, and affine substitution in the concave argument, so that the former function is convex-concave and continuous on the domain $\{\alpha > 0, \phi \in \mathcal{X}\} \times X_i$. By similar argument, the function $\alpha \Phi_{\mathcal{O}}(-\phi/\alpha; A_j(y))$ is convex-concave and continuous on the domain $\{\alpha > 0, \phi \in \mathcal{F}\} \times X_j$. These observations combine with compactness of X_i, X_j to imply that $\Psi_{ij}(\alpha, \phi)$ is real-valued continuous convex function on the domain

$$\mathcal{F}^+ = \{\alpha > 0\} \times \mathcal{F}.$$

Observe that functions $\Psi_{ii}(\alpha, \phi)$ are nonnegative on \mathcal{F}^+ . Indeed, selecting somehow $\bar{x} \in X_i$, and setting $\mu = A_i(\bar{x})$, we have

$$\begin{split} \Psi_{ii}(\alpha,\phi) &\geq \Phi_{ii}(\alpha,\phi;\bar{x},\bar{x}) = \frac{\alpha}{2} \left[K[\Phi_{\mathcal{O}}(\phi/\alpha;\mu) + \Phi_{\mathcal{O}}(-\phi/\alpha;\mu)] + 2\ln(2I/\epsilon) \right] \\ &= \frac{\alpha}{2} \left[K \ln\left(\underbrace{\left[\int \exp\{\phi(\omega)/\alpha\} p_{\mu}(\omega) \Pi(d\omega)\right] \left[\int \exp\{-\phi(\omega)/\alpha\} p_{\mu}(\omega) \Pi(d\omega)\right]}_{\geq \left[\int \exp\{\frac{1}{2}\phi(\omega)/\alpha\} \exp\{-\frac{1}{2}\phi(\omega)/\alpha\} p_{\mu}(\omega) \Pi(d\omega)\right]^2 = 1} \right] \\ &+ 2\ln(2I/\epsilon) \end{split}$$

 $\geq \alpha \ln(2I/\epsilon) > 0$

(we have used Cauchy inequality).

Functions Ψ_{ij} give rise to convex and feasible optimization problems

$$Opt_{ij} = Opt_{ij}(K) = \min_{(\alpha,\phi)\in\mathcal{F}^+} \Psi_{ij}(\alpha,\phi).$$
(3.2)

By its origin, Opt_{ij} is either a real, or $-\infty$; by the observation above, Opt_{ii} are nonnegative. Our estimate is as follows.

1. For $1 \leq i, j \leq I$, we select somehow feasible solutions α_{ij}, ϕ_{ij} to problems (3.2) (the less the values of the corresponding objectives, the better) and set

$$\begin{aligned}
\rho_{ij} &= \Psi_{ij}(\alpha_{ij}, \phi_{ij}) = \frac{1}{2} \left[\Psi_{i,+}(\alpha_{ij}, \phi_{ij}) + \Psi_{j,-}(\alpha_{ij}, \phi_{ij}) \right] \\
\varkappa_{ij} &= \frac{1}{2} \left[\Psi_{j,-}(\alpha_{ij}, \phi_{ij}) - \Psi_{i,+}(\alpha_{ij}, \phi_{ij}) \right] \\
g_{ij}(\omega^K) &= \sum_{k=1}^{K} \phi_{ij}(\omega_k) + \varkappa_{ij} \\
\rho &= \max_{1 \le i, j \le I} \rho_{ij}
\end{aligned} (3.3)$$

2. Given observation ω^{K} , we specify the estimate $\hat{g}(\omega^{K})$ as follows:

$$r_{i} = \max_{j \leq I} g_{ij}(\omega^{K})$$

$$c_{j} = \min_{i \leq I} g_{ij}(\omega^{K})$$

$$\widehat{g}(\omega^{K}) = \frac{1}{2} [\min_{i \leq I} r_{i} + \max_{j \leq I} c_{j}].$$
(3.4)

3.1.3 Main result

Proposition 3.1. The ϵ -risk of the estimate we have built can be upper-bounded as follows:

$$\operatorname{Risk}_{\epsilon}[\widehat{g}] \le \rho. \tag{3.5}$$

Proof. Let the common distribution p of independent across k components ω_k in observation ω^K be $p_{A_\ell(u)}$ for some $\ell \leq I$ and $u \in X_\ell$. Let us fix these ℓ and u, let $\mu = A_\ell(u)$, and let p^K stand for the distribution of ω^K .

$\mathbf{1}^0$. We have

and we arrive at

$$\operatorname{Prob}_{\omega^{K} \sim p^{K}} \left\{ g_{\ell j}(\omega^{K}) > g^{T} u = \rho_{\ell j} \right\} \leq \frac{\epsilon}{2I}.$$
(3.6)

ESTIMATING FUNCTIONS VIA HYPOTHESIS TESTING

Similarly,

and we arrive at

$$\operatorname{Prob}_{\omega^{K} \sim p^{K}} \left\{ g_{i\ell}(\omega^{K}) < g^{T}u - \rho_{i\ell} \right\} \leq \frac{\epsilon}{2I}.$$
(3.7)

 2^{0} . Let

$$\mathcal{E} = \{ \omega^K : g_{\ell j}(\omega^K) \le g^T u + \rho_{\ell j}, g_{i\ell}(\omega^K) \ge g^T u - \rho_{i\ell}, 1 \le i, j \le I \}.$$

From (3.6), (3.7) and the union bound it follows that p^{K} -probability of the event \mathcal{E} is $\geq 1 - \epsilon$. As a result, all we need to complete the proof of Proposition is to verify that

$$\omega^K \in \mathcal{E} \Rightarrow |\widehat{g}(\omega^K) - g^T u| \le \rho.$$
(3.8)

Indeed, let us fix $\omega^K \in \mathcal{E}$, and let E be the $I \times I$ matrix with entries $E_{ij} = g_{ij}(\omega^K)$, $1 \leq i, j \leq I$. The quantity r_i , see (3.4), is the maximum of entries in *i*-th row of E, and the quantity c_j is the minimum of entries in *j*-th column of E; in particular, $r_i \geq E_{ij} \geq c_j$ for all i, j, implying that $r_i \geq c_j$ for all i, j. Now, since $\omega^K \in \mathcal{E}$, we have $E_{\ell\ell} = g_{\ell\ell}(\omega^K) \geq g^T u - \rho_{\ell\ell} \geq g^T u - \rho$ and $E_{\ell j} = g_{\ell j}(\omega^K) \leq g^T u + \rho_{\ell j} \leq g^T u + \rho$ for all j, implying that $r_{\ell} = \max_j E_{\ell j} \in \Delta = [g^T u - \rho, g^T u + \rho]$. Similarly, $\omega \in \mathcal{E}$ implies that $E_{\ell\ell} = g_{\ell\ell}(\omega^K) \leq g^T u + \rho$ and $E_{i\ell} = g_{i\ell}(\omega^K) \geq g^T u - \rho_{i\ell} \geq g^T u - \rho$ for all i, implying that $c_{\ell} = \min_i E_{i\ell} \in \Delta$. We see that both r_{ℓ} and c_{ℓ} belong to Δ ; since $r_* := \min_i r_i \leq r_{\ell}$ and, as have already seen, $r_i \geq c_{\ell}$ for all i, we conclude that $r_* \in \Delta$. By similar argument, $c_* := \max_j c_j \in \Delta$ as well. By construction, $\widehat{g}(\omega^K) = \frac{1}{2}[r_* + c_*]$, that is, $\widehat{g}(\omega^K) \in \Delta$, and the conclusion in (3.8) indeed takes place. \Box

3.1.4 Near-optimality

Observe that properly selecting ϕ_{ij} and α_{ij} we can make, in a computationally efficient manner, the upper bound ρ on the ϵ -risk of the above estimate arbitrarily close to

$$\operatorname{Opt}(K) = \max_{1 \le i, j \le I} \operatorname{Opt}_{ij}(K).$$

We are about to demonstrate that the quantity Opt(K) "nearly lower-bounds" the minimax optimal ϵ -risk

$$\operatorname{Risk}_{\epsilon}^{*}(K) = \inf_{\widehat{g}(\cdot)} \operatorname{Risk}_{\epsilon}[\widehat{g}],$$

the infimum being taken over all K-observation Borel estimates. The precise statement is as follows:

Proposition 3.2. In the situation of this Section, let $\epsilon \in (0, 1/2)$ and \overline{K} be a positive integer. Then for every integer K satisfying

$$K/\bar{K} > \frac{2\ln(2I/\epsilon)}{\ln(\frac{1}{4\epsilon(1-\epsilon)})}$$

one has

$$\operatorname{Opt}(K) \le \operatorname{Risk}^*_{\epsilon}(K).$$
 (3.9)

In addition, in the special case where for every i, j there exists $x_{ij} \in X_i \cap X_j$ such that $A_i(x_{ij}) = A_j(x_{ij})$ one has

$$K \ge \bar{K} \Rightarrow \operatorname{Opt}(K) \le \frac{2\ln(2I/\epsilon)}{\ln(\frac{1}{4\epsilon(1-\epsilon)})} \operatorname{Risk}_{\epsilon}^{*}(\bar{K}).$$
(3.10)

Proof. 1⁰. Observe that $Opt_{ij}(K)$ is the saddle point value in the convexconcave saddle point problem:

$$\begin{aligned}
\operatorname{Opt}_{ij}(K) &= \inf_{\alpha > 0, \phi \in \mathcal{F}} \max_{x \in X_i, y \in X_j} \left[\frac{1}{2} K \alpha \left\{ \Phi_{\mathcal{O}}(\phi/\alpha; A_i(x)) + \Phi_{\mathcal{O}}(-\phi/\alpha; A_j(y)) \right\} \\
&+ \frac{1}{2} g^T [y - x] + \alpha \ln(2I/\epsilon) \right].
\end{aligned}$$

The domain of the maximization variable is compact and the cost function is continuous on its domain, whence, by Sion-Kakutani Theorem, we have also

$$\begin{aligned}
\operatorname{Opt}_{ij}(K) &= \max_{x \in X_i, y \in X_j} \Theta_{ij}(x, y), \\
\Theta_{ij}(x, y) &= \inf_{\alpha > 0, \phi \in \mathcal{F}} \left[\frac{1}{2} K \alpha \left\{ \Phi_{\mathcal{O}}(\phi/\alpha; A_i(x)) + \Phi_{\mathcal{O}}(-\phi/\alpha; A_j(y)) \right\} + \alpha \ln(2I/\epsilon) \right] + \frac{1}{2} g^T [y - x].
\end{aligned}$$
(3.11)

We have

$$\Theta_{ij}(x,y) = \inf_{\alpha>0,\psi\in\mathcal{F}} \left[\frac{1}{2} K \alpha \left\{ \Phi_{\mathcal{O}}(\psi;A_i(x)) + \Phi_{\mathcal{O}}(-\psi;A_j(y)) \right\} + \frac{1}{2} g^T[y-x] + \alpha \ln(2I/\epsilon) \right]$$
$$= \inf_{\substack{\alpha>0\\ +\frac{1}{2} g^T[y-x]}} \left[\frac{1}{2} \alpha K \inf_{\substack{\psi\in\mathcal{F}\\ \psi\in\mathcal{F}}} \left\{ \Phi_{\mathcal{O}}(\psi;A_i(x)) + \Phi_{\mathcal{O}}(-\psi;A_j(y)) \right\} + \alpha \ln(2I/\epsilon) \right]$$

Given $x \in X_i$, $y \in X_j$ and setting $\mu = A_i(x)$, $\nu = A_j(y)$, we obtain

$$\inf_{\psi \in \mathcal{F}} [\Phi_{\mathcal{O}}(\psi; A_i(x)) + \Phi_{\mathcal{O}}(-\psi; A_j(y))] = \inf_{\psi \in \mathcal{F}} \left[\ln \left(\int \exp\{\psi(\omega)\} p_\mu(\omega) P(d\omega) \right) + \ln \left(\int \exp\{-\psi(\omega)\} p_\nu(\omega) P(d\omega) \right) \right]$$

ESTIMATING FUNCTIONS VIA HYPOTHESIS TESTING

Since
$$\mathcal{O}$$
 is a good o.s., the function $\psi(\omega) = \frac{1}{2} \ln(p_{\nu}(\omega)/p_{\mu}(\omega))$ belongs to \mathcal{F} , and

$$\inf_{\psi \in \mathcal{F}} \left[\ln\left(\int \exp\{\psi(\omega)\}p_{\mu}(\omega)P(d\omega)\right) + \ln\left(\int \exp\{-\psi(\omega)\}p_{\nu}(\omega)P(d\omega)\right) \right]$$

$$= \inf_{\delta \in \mathcal{F}} \left[\ln\left(\int \exp\{\bar{\psi}(\omega) + \delta(\omega)\}p_{\mu}(\omega)P(d\omega)\right) + \ln\left(\int \exp\{-\bar{\psi}(\omega) - \delta(\omega)\}p_{\nu}(\omega)P(d\omega)\right) \right]$$

$$= \inf_{\delta \in \mathcal{F}} \underbrace{\left[\ln\left(\int \exp\{\delta(\omega)\}\sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right) + \ln\left(\int \exp\{-\delta(\omega)\}\sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right) \right]}_{f(\delta)}.$$

Observe that $f(\delta)$ clearly is a convex and even function of $\delta \in \mathcal{F}$; as such, it attains its minimum over $\delta \in \mathcal{F}$ when $\delta = 0$. The bottom line is that

$$\inf_{\psi \in \mathcal{F}} [\Phi_{\mathcal{O}}(\psi; A_i(x)) + \Phi_{\mathcal{O}}(-\psi; A_j(y))] = 2 \ln \left(\int \sqrt{p_{A_i(x)}(\omega) p_{A_j(y)}(\omega)} P(d\omega) \right),$$
(3.12)

and

$$\begin{aligned} \Theta_{ij}(x,y) &= \inf_{\alpha>0} \alpha \left[K \ln\left(\int \sqrt{p_{A_i(x)}(\omega)p_{A_j(y)}(\omega)}P(d\omega) \right) + \ln(2I/\epsilon) \right] + \frac{1}{2}g^T[y-x] \\ &= \begin{cases} \frac{1}{2}g^T[y-x] &, K \ln\left(\int \sqrt{p_{A_i(x)}(\omega)p_{A_j(y)}(\omega)}P(d\omega) \right) + \ln(2I/\epsilon) \ge 0, \\ -\infty, & \text{otherwise.} \end{cases} \end{aligned}$$

This combines with (3.11) to imply that

$$\operatorname{Opt}_{ij}(K) = \max_{x,y} \left\{ \frac{1}{2} g^T[y-x] : x \in X_i, y \in X_j, \left[\int \sqrt{p_{A_i(x)}(\omega)p_{A_j(y)}(\omega)} P(d\omega) \right]^K \ge \frac{\epsilon}{2I} \right\}.$$
(3.13)

 $\mathbf{2^{0}}.$ We claim that under the premise of Proposition, for all $i,j,\,1\leq i,j\leq I,$ one has

$$\operatorname{Opt}_{ij}(K) \le \operatorname{Risk}^*_{\epsilon}(K)$$

implying the validity of (3.9). Indeed, assume that for some pair i, j the opposite inequality holds true:

$$\operatorname{Opt}_{ij}(K) > \operatorname{Risk}^*_{\epsilon}(\bar{K}),$$

and let us lead this assumption to a contradiction. Under our assumption optimization problem in (3.13) has a feasible solution (\bar{x}, \bar{y}) such that

$$r := \frac{1}{2}g^T[\bar{y} - \bar{x}] > \operatorname{Risk}^*_{\epsilon}(\bar{K}), \qquad (3.14)$$

implying, due to the origin of $\operatorname{Risk}^*_{\epsilon}(\bar{K})$, that there exists an estimate $\widehat{g}(\omega^{\bar{K}})$ such that for $\mu = A_i(\bar{x}), \ \nu = A_j(\bar{y})$ it holds

$$\begin{aligned} &\operatorname{Prob}_{\omega^{\bar{K}} \sim p_{\nu}^{\bar{K}}} \left\{ \widehat{g}(\omega^{\bar{K}}) \leq \frac{1}{2} g^{T}[\bar{x} + \bar{y}] \right\} &\leq \operatorname{Prob}_{\omega^{\bar{K}} \sim p_{\nu}^{\bar{K}}} \left\{ |\widehat{g}(\omega^{\bar{K}}) - g^{T}\bar{y}| \geq r \right\} \leq \epsilon \\ &\operatorname{Prob}_{\omega^{\bar{K}} \sim p_{\mu}^{\bar{K}}} \left\{ \widehat{g}(\omega^{\bar{K}}) \geq \frac{1}{2} g^{T}[\bar{x} + \bar{y}] \right\} &\leq \operatorname{Prob}_{\omega^{\bar{K}} \sim p_{\mu}^{\bar{K}}} \left\{ |\widehat{g}(\omega^{\bar{K}}) - g^{T}\bar{x}| \geq r \right\} \leq \epsilon, \end{aligned}$$

so that we can decide on two simple hypotheses stating that observation $\omega^{\vec{K}}$ obeys distribution $p_{\mu}^{\vec{K}}$, resp., $p_{\nu}^{\vec{K}}$, with risk $\leq \epsilon$. Therefore,

$$\int \min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right] P^{\bar{K}}(d\omega^{\bar{K}}) \le 2\epsilon.$$

Hence, when setting $p_{\theta}^{\bar{K}}(\omega^{\bar{K}}) = \prod_{k} p_{\theta}(\omega_{k})$ and $P^{\bar{K}} = \underbrace{P \times \dots \times P}_{\bar{K}}$, we have

$$\begin{split} &\left[\int \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right]^{\bar{K}} = \int \sqrt{p_{\mu}^{\bar{K}}(\omega^{\bar{K}})p_{\nu}^{\bar{K}}(\omega^{\bar{K}})}P^{\bar{K}}(d\omega^{\bar{K}}) \\ &= \int \sqrt{\min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]}\sqrt{\max\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]}P^{\bar{K}}(d\omega^{\bar{K}}) \\ &\leq \left[\int \min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]P^{\bar{K}}(d\omega^{\bar{K}})\right]^{1/2} \left[\int \max\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]P^{\bar{K}}(d\omega^{\bar{K}})\right]^{1/2} \\ &= \left[\int \min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]P^{\bar{K}}(d\omega^{\bar{K}})\right]^{1/2} \\ &\times \left[\int \left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}}) - \min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]\right]P^{\bar{K}}(d\omega^{\bar{K}})\right]^{1/2} \\ &= \left[\int \min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]P^{\bar{K}}(d\omega^{\bar{K}})\right]^{1/2} \left[2 - \int \min\left[p_{\mu}^{\bar{K}}(\omega^{\bar{K}}), p_{\nu}^{\bar{K}}(\omega^{\bar{K}})\right]P^{\bar{K}}(d\omega^{\bar{K}})\right]^{1/2} \\ &\leq 2\sqrt{\epsilon(1-\epsilon)}. \end{split}$$

Consequently,

$$\left[\int \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right]^{K} \leq \left[2\sqrt{\epsilon(1-\epsilon)}\right]^{K/\bar{K}} < \frac{\epsilon}{2I},$$

which is the desired contradiction (recall that $\mu = A_i(\bar{x}), \nu = A_j(\bar{y})$ and (\bar{x}, \bar{y}) is feasible for (3.13)).

 3^0 . Now let us prove that under the premise of Proposition, (3.10) takes place. To this end let us set

$$w_{ij}(s) = \max_{x \in X_j, y \in X_j} \left\{ \frac{1}{2} g^T[y - x] : \bar{K} \underbrace{\ln\left(\int \sqrt{p_{A_i(x)}(\omega)p_{A_j(y)}(\omega)}P(d\omega)\right)}_{H(x,y)} + s \ge 0 \right\}.$$
(3.15)

As we have seen in item 1^0 , see (3.12), one has

$$H(x,y) = \inf_{\psi \in \mathcal{F}} \frac{1}{2} \left[\Phi_{\mathcal{O}}(\psi; A_i(x)) + \Phi_{\mathcal{O}}(-\psi, A_j(y)) \right],$$

that is, H(x, y) is the infimum of a parametric family of concave functions of $(x, y) \in X_i \times X_j$ and as such is concave. Besides this, the optimization problem in (3.15) is feasible whenever $s \ge 0$, a feasible solution being $y = x = x_{ij}$. At this feasible solution we have $g^T[y - x] = 0$, implying that $w_{ij}(s) \ge 0$ for $s \ge 0$. Observe also that from concavity of H(x, y) it follows that $w_{ij}(s)$ is concave on the ray $\{s \ge 0\}$. Finally, we claim that

$$w_{ij}(\bar{s}) \le \operatorname{Risk}^*_{\epsilon}(\bar{K}), \ \bar{s} = -\ln(2\sqrt{\epsilon(1-\epsilon)}).$$
 (3.16)

Indeed, $w_{ij}(s)$ is nonnegative, concave and bounded (since X_i, X_j are compact) on \mathbf{R}_+ , implying that $w_{ij}(s)$ is continuous on $\{s > 0\}$. Assuming, on the contrary to our claim, that $w_{ij}(\bar{s}) > \operatorname{Risk}^*_{\epsilon}(\bar{K})$, there exists $s' \in (0, \bar{s})$ such that $w_{ij}(s') > \operatorname{Risk}^*_{\epsilon}(\bar{K})$ and thus there exist $\bar{x} \in X_i, \bar{y} \in X_j$ such that (\bar{x}, \bar{y}) is feasible for the optimization problem specifying $w_{ij}(s')$ and (3.14) takes place. We have seen in item 2⁰ that the latter relation implies that for $\mu = A_i(\bar{x}), \nu = A_j(\bar{y})$ it holds

$$\left[\int \sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right]^{\bar{K}} \leq 2\sqrt{\epsilon(1-\epsilon)},$$

ESTIMATING FUNCTIONS VIA HYPOTHESIS TESTING

that is,

$$\bar{K}\ln\left(\int\sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right) + \bar{s} \le 0,$$

whence

$$\bar{K}\ln\left(\int\sqrt{p_{\mu}(\omega)p_{\nu}(\omega)}P(d\omega)\right)+s'<0,$$

contradicting the fact that (\bar{x}, \bar{y}) is feasible for the optimization problem specifying $w_{ij}(s')$.

It remains to note that (3.16) combines with concavity of $w_{ij}(\cdot)$ and the relation $w_{ij}(0) \ge 0$ to imply that

$$w_{ij}(\ln(2I/\epsilon)) \le \vartheta w_{ij}(\bar{s}) \le \vartheta \operatorname{Risk}^*_{\epsilon}(\bar{K}), \ \ \vartheta = \ln(2I/\epsilon)/\bar{s} = \frac{2\ln(2I/\epsilon)}{\ln([4\epsilon(1-\epsilon)]^{-1})}.$$

Invoking (3.13), we conclude that

$$\operatorname{Opt}_{ij}(\bar{K}) = w_{ij}(\ln(2I/\epsilon)) \le \vartheta \operatorname{Risk}^*_{\epsilon}(\bar{K}) \,\forall i, j.$$

Finally, from (3.13) it immediately follows that $\operatorname{Opt}_{ij}(K)$ is nonincreasing in K (since as K grows, the feasible set of the right hand side optimization problem in (3.13) shrinks), that is,

$$K \ge \bar{K} \Rightarrow \operatorname{Opt}(K) \le \operatorname{Opt}(\bar{K}) = \max_{i,j} \operatorname{Opt}_{ij}(\bar{K}) \le \vartheta \operatorname{Risk}^*_{\epsilon}(\bar{K}),$$

and (3.10) follows.

3.1.5 Illustration

We illustrate our construction on the simplest possible example – one where $X_i = \{x_i\}$ are singletons in \mathbb{R}^n , i = 1, ..., I, the observation scheme is Gaussian. Thus, setting $y_i = A_i(x_i) \in \mathbb{R}^m$, the observation's components ω_k , $1 \le k \le K$, stemming from signal x_i , are drawn, independently of each other, from the normal distribution $\mathcal{N}(y_i, I_m)$. The family \mathcal{F} of functions ϕ associated with Gaussian o.s. is the family of all affine functions $\phi(\omega) = \phi_0 + \varphi^T \omega$ on the observation space (which at present is \mathbb{R}^m); we identify $\phi \in \mathcal{F}$ with the pair (ϕ_0, φ) . The function $\Psi_{\mathcal{O}}$ associated with the Gaussian observation scheme with *m*-dimensional observations is

$$\Phi_{\mathcal{O}}(\phi;\mu) = \phi_0 + \varphi^T \mu + \frac{1}{2} \varphi^T \varphi : (\mathbf{R} \times \mathbf{R}^m) \times \mathbf{R}^m \to \mathbf{R},$$

a straightforward computation shows that in the case in question, setting

$$\theta = \ln(2I/\epsilon),$$

(3.17)

we have

$$\begin{split} \Psi_{i,+}(\alpha,\phi) &= K\alpha \left[\phi_0 + \varphi^T y_i/\alpha + \frac{1}{2}\varphi^T \varphi/\alpha^2\right] + \alpha\theta - g^T x_i \\ &= K\phi_0 + K\varphi^T y_i - g^T x_i + \frac{K}{2\alpha}\varphi^T \varphi + \alpha\theta \\ \Psi_{j,-}(\alpha,\phi) &= -K\phi_0 - K\varphi^T y_j + g^T x_j + \frac{K}{2\alpha}\varphi^T \varphi + \alpha\theta \\ \operatorname{Opt}_{ij} &= \inf_{\alpha>0,\phi} \frac{1}{2} \left[\Psi_{i,+}(\alpha,\phi) + \Psi_{j,-}(\alpha,\phi)\right] \\ &= \frac{1}{2}g^T [x_j - x_i] + \inf_{\varphi} \left[\frac{K}{2}\varphi^T [y_i - y_j] + \inf_{\alpha>0} \left[\frac{K}{2\alpha}\varphi^T \varphi + \alpha\theta\right]\right] \\ &= \frac{1}{2}g^T [x_j - x_i] + \inf_{\varphi} \left[\frac{K}{2}\varphi^T [y_i - y_j] + \sqrt{2K\theta} \|\varphi\|_2\right] \\ &= \begin{cases} \frac{1}{2}g^T [x_j - x_i], & \|y_i - y_j\|_2 \le 2\sqrt{2\theta/K} \\ -\infty, & \|y_i - y_j\|_2 > 2\sqrt{2\theta/K}. \end{cases} \end{split}$$

We see that we can safely set $\phi_0 = 0$ and that setting

$$\mathcal{I} = \{(i, j) : \|y_i - y_j\|_2 \le 2\sqrt{2\theta/K}\},\$$

 $\operatorname{Opt}_{ij}(K)$ is finite iff $(i, j) \in \mathcal{I}$ and is $-\infty$ otherwise; in both cases, the optimization problem specifying Opt_{ij} has no optimal solution. Indeed, this clearly is the case when $(i, j) \notin \mathcal{I}$; when $(i, j) \in \mathcal{I}$, a minimizing sequence is, e.g., $\phi_0 \equiv 0, \varphi \equiv 0, \alpha_i \to 0$, but its limits is not in the minimization domain (on this domain, α should be positive). Coping with this case was exactly the reason why in our construction we required from ϕ_{ij}, α_{ij} to be feasible, and not necessary optimal, solutions to the optimization problems in question). In the illustration under consideration, the simplest way to overcome the difficulty is to restrict the optimization domain \mathcal{F}^+ in (3.2) with its compact subset $\{\alpha \geq 1/R, \phi_0 = 0, \|\varphi\|_2 \leq R\}$ with large R, like $R = 10^{10}$ or 10^{20} . With this approach, we specify the entities participating in (3.3) as

$$\phi_{ij}(\omega) = \varphi_{ij}^T \omega, \ \varphi_{ij} = \begin{cases} 0, & (i,j) \in \mathcal{I} \\ -R[y_i - y_j]/\|y_i - y_j\|_2, & (i,j) \notin \mathcal{I} \end{cases}$$

$$\alpha_{ij} = \begin{cases} 1/R, & (i,j) \in \mathcal{I} \\ \sqrt{\frac{K}{2\theta}}R, & (i,j) \notin \mathcal{I} \end{cases}$$
(3.18)

resulting in

$$\begin{aligned} \varkappa_{ij} &= \frac{1}{2} \left[\Psi_{j,-}(\alpha_{ij},\phi_{ij}) - \Psi_{i,+}(\alpha_{ij},\phi_{ij}) \right] \\ &= \frac{1}{2} \left[-K\varphi_{ij}^{T}y_{j} + g^{T}x_{j} + \frac{K}{2\alpha_{ij}}\varphi_{ij}^{T}\varphi_{ij} + \alpha_{ij}\theta - K\varphi_{ij}^{T}y_{i} + g^{T}x_{i} - \frac{K}{2\alpha_{ij}}\varphi_{ij}^{T}\varphi_{ij} - \alpha_{ij}\theta \right] \\ &= \frac{1}{2} g^{T}[x_{i} + x_{j}] - \frac{K}{2}\varphi_{ij}^{T}[y_{i} + y_{j}] \\ \rho_{ij} &= \frac{1}{2} \left[\Psi_{i,+}(\alpha_{ij},\phi_{ij}) + \Psi_{j,-}(\alpha_{ij},\phi_{ij}) \right] \\ &= \frac{1}{2} \left[K\varphi_{ij}^{T}y_{i} - g^{T}x_{i} + \frac{K}{2\alpha_{ij}}\varphi_{ij}^{T}\varphi_{ij} + \alpha_{ij}\theta - K\varphi_{ij}^{T}y_{j} + g^{T}x_{j} + \frac{K}{2\alpha_{ij}}\varphi_{ij}^{T}\varphi_{ij} + \alpha_{ij}\theta \right] \\ &= \frac{K}{2\alpha_{ij}}\varphi_{ij}^{T}\phi_{ij} + \alpha_{ij}\theta + \frac{1}{2}g^{T}[x_{j} - x_{i}] + \frac{K}{2}\varphi_{ij}^{T}[y_{i} - y_{j}] \\ &= \begin{cases} \frac{1}{2}g^{T}[x_{j} - x_{i}] + R^{-1}\theta, & (i,j) \in \mathcal{I} \\ \frac{1}{2}g^{T}[x_{j} - x_{i}] + [\sqrt{2K\theta} - \frac{K}{2}||y_{i} - y_{j}||_{2}]R, & (i,j) \notin \mathcal{I} \end{cases} \end{aligned}$$

$$(3.19)$$

In the numerical experiments we are about to report we used n = 20, m = 10, and I = 100, with x_i , $i \leq I$, drawn independently of each other from $\mathcal{N}(0, I_n)$, and $y_i = Ax_i$ with randomly generated matrix A (specifically, matrix with independent $\mathcal{N}(0, 1)$ entries normalized to have unit spectral norm), and used $R = 10^{20}$; the linear form to be recovered was just the first coordinate of x. The results of typical

ESTIMATING FUNCTIONS VIA HYPOTHESIS TESTING

experiment are as follows:

K	$\max_{i,j} \rho_{ij}$	Empirical recovery error [mean/median/max]
2	2.541	0.9243/0.8292/2.541
4	2541	0.9859/0.9066/2.541
8	2.541	0.8057/0.7316/2.541
16	2.541	0.6807/0.6567/2.115
32	1.758	0.3630/0.2845/1.758
64	0.954	0.0860/0.0000/0.954
128	0.000	0.0000/0.0000/0.000
256	0.000	0.0000/0.0000/0.000

For every K, the empirical recovery errors shown in the table stem from 20 experiments, with the signal underlying an experiment selected at random among $x_1, ..., x_{100}$.

3.2 ESTIMATING *N*-CONVEX FUNCTIONS ON UNIONS OF CONVEX SETS

In this Section, we apply our testing machinery to the estimation problem as follows.

Given are:

- a simple o.s. $\mathcal{O} = (\Omega, \Pi; \{p_{\mu} : \mu \in \mathcal{M}\}; \mathcal{F}),$
- a signal space $X \subset \mathbf{R}^n$ along with affine encoding $x \mapsto A(x) : X \to \mathcal{M}$,
- a real-valued function f on X.

Given observation $\omega \sim p_{A(x_*)}$ stemming from unknown signal x_* known to belong to X, we want to recover $f(x_*)$.

Our approach imposes severe restrictions on f (satisfied, e.g., when f is linear, or linear-fractional, or is the maximum of several linear functions); as a compensation, we allow for rather "complex" X – finite unions of convex sets.

3.2.1 Outline

The approach we intend to develop is, in nutshell, extremely simple; its formal description, however, turns to be lengthy and obscures, to some extent, the simple ideas underlying the construction. By this reason, it makes sense to start with informal outline of the strategy underlying the forthcoming developments. Consider the situation where the signal space X is the 2D rectangle depicted on the top of Figure 3.1.(a), and let the function to be recovered be $f(u) = u_1$. Thus, "the nature" has somehow selected x in the rectangle, and we observe, say, Gaussian random variable with the mean A(x) and known covariance matrix, where $A(\cdot)$ is a given affine mapping. Note that hypotheses $f(x) \ge b$ and $f(x) \le a$ translate into convex hypotheses on the expectation of the observed Gaussian r.v., so that we can use out hypothesis testing machinery to decide on hypotheses of this type and to localize f(x) in a (hopefully, small) segment by a Bisection-type process.



Figure 3.1: Bisection via Hypothesis Testing

Before describing the process, let us make a terminological agreement. In the sequel we shall use pairwise hypothesis testing in the situation where it may happen the *neither one* of the hypotheses H_1 , H_2 we are deciding upon is true. In this case, we will say that the outcome of a test is correct, if the rejected hypothesis indeed is wrong (the accepted hypothesis can be wrong as well, but the latter can happen only in the case when both our hypotheses are wrong).

This is how our Bisection could look like.

1. Were we able to decide reliably on the Blue and Red hypotheses on Figure 3.1.(a), that is, to understand via observations whether x belongs to the left or to the right half of the original rectangle, our course of actions would be clear: depending on this decision, we would replace our original rectangle with a smaller rectangle localizing x, as shown on Figure 3.1.(a), and then iterate this process. The difficulty, of course, is that our Red and Blue hypotheses intersect, so that is impossible to decide on them reliably.

2. In order to make Red and Blue hypotheses distinguishable from each other, we could act as shown on Figure 3.1.(b), by shrinking a little bit the blue and the red rectangles and inserting between the resulting rectangles the green "no-man land." Assuming that the width of the green rectangle allows to decide reliably on our new Blue and Red hypotheses and utilizing available observation, we can localize x either in the blue, or in the red rectangles as shown on Figure 3.1.(b). Specifically, assume that our "Red vs. Blue" test rejected correctly the red hypothesis. Then x can be located either in blue, or in green rectangles shown on the top of the figure, and thus x is in the new blue localizer which is the union of the blue and the green original rectangles. Similarly, if our test rejects correctly the blue hypothesis, then we can take, as the new localizer of x, the union of the original red and green rectangles, as shown on Figure 3.1.(b). Note that our localization is as reliable as our test is, and that it reduces the width of localizer by factor close to 2, provided the width of the green rectangle is small as compared to the width of the original "tricolor" localizer of x. We can iterate this process, with the new – smaller –
localizer in the role of the old till arriving at a localizer so narrow that "no-man land" part of it (this part cannot be too narrow, since it should allow for reliable decision on the current blue and red hypotheses) becomes too large to allow for significant progress in localizer's width.

The bottleneck of this approach is where to take observations to be used in our subsequent tests. In principle, we could use in all of them the initial observation; the difficulty with this approach is, that the hypotheses we need to decide upon depend on the observations (e.g., when x belongs to the green part of the "tricolor" rectangle on Figure 3.1, deciding on Blue vs. Red can, depending on observation, lead to accepting either red or blue hypothesis, thus leading to different updated localizers), and we arrive at the situation when we should decide on random hypotheses via observation statistically depending on these hypotheses – a mess we have no idea how to analyze. To circumvent this difficulty, we could use in every one of the tests its own observation drawn, independently of the previous observations, from the distribution $p_{A(x)}$. However, to do this, we need repeated observations to be allowed, and the number of observations we will use will be proportional to the number of tests we intend to run.

3. Finally, there is a theoretically sound way to implement Bisection based on a *single* observation, and this is what we intend to do. The policy we use now is as follows: given current localizer for x (at the first step - our initial rectangle), we consider two "tricolor" partitions of it depicted at the top of Figure 3.1.(c). In the first partition, the blue rectangle is the left half of the original rectangle, in the second the red rectangle is the right half of the original rectangle. We then run *two* Blue vs. Red tests, the first on the pair of Blue and Red hypotheses stemming from the first partition, and the second on the pair of Blue and Red hypotheses stemming from the second partition. Assuming that in both tests the rejected hypotheses indeed were wrong, the results of these tests allow us to make conclusions as follows:

- when both tests reject red hypotheses from the corresponding pairs, x is located in the left half of the initial rectangle (since otherwise in the second test the rejected hypothesis were in fact true, contradicting to the assumption that both tests make no wrong rejections);
- when both tests reject blue hypotheses from the corresponding pairs, x is located in the right half of the original rectangle (by the same reasons as in the previous case);
- when the tests "disagree," rejecting hypotheses of different colors, x is located in the union of the two green rectangles we deal with. Indeed, otherwise x should be either in the blue rectangles of both our "tricolors," or in the red rectangles of both of them. Since we have assumed that in both tests no wrong rejections took place, in the first case both tests must reject red hypotheses, and in the second both should reject blue ones, while in fact neither one of these two options took place.

Now, in the first two cases we can safely say to which one of "halves" – left or right – of the initial rectangle x belongs, and take this half as our new localizer. In the third case, we take as a new localizer for x the green rectangle shown on the bottom of Figure 3.1 and terminate our estimation process – the new localizer already is narrow! Now, in the proposed algorithm, unless we terminate at the very first step, we carry out the second step exactly in the same fashion as the first one,

with the localizer of x yielded by the first step in the role of the initial localizer, then carry out, in the same fashion, the third step, etc., until termination either due to running into a disagreement, or due to reaching a prescribed number of steps. Upon termination, we return the last localizer for x which we have built, and claim that $f(x) = x_1$ belongs to the projection of this localizer onto the x_1 -axis. In all tests from the above process, we use the same observation. Note that in our current situation, in contrast to the one we have discussed earlier, re-utilizing a single observation creates no difficulties, since with no wrong rejections in the pairwise tests we use, the pairs of hypotheses participating in the tests are not random at all – they are uniquely defined by $f(x) = x_1!$ Indeed, with no wrong rejections, prior to termination everything is as if we were running perfect Bisection, that is, were updating subsequent rectangles Δ_t containing x according to the rules

- Δ_1 is a given in advance rectangle containing x,
- Δ_{t+1} is either the left, or the right half of Δ_t , depending on which one of these two halves contains x.

Thus, given x and with no wrong rejections, the situation is as if a single observation were used in a number L of tests "in parallel" rather than sequentially, and the only elaboration caused by the sequential nature of our process is in "risk accumulation" – we want the probability of error *in one or more of our* L *tests* to be less than the desired risk ϵ of wrong "bracketing" of f(x), implying, for absence of something better, that the risks of the individual tests should be at most ϵ/L . These risks, in turn, define the allowed width of "no man land" zones, and thus – the accuracy to which f(x) can be estimated. It should be noted that the number L of steps of Bisection always is a moderate integer (since otherwise the width of "no-man land" zone, which at the concluding Bisection steps is of order of 2^{-L} , will be by far too small to allow for deciding on the concluding pairs of our hypotheses with risk ϵ/L , at least when our observations possess non-negligible volatility). As a result, "the price" of Bisection turns out to be low as compared to the case where every test uses its own observation.

We have outlined the strategy we are about to implement. From the outline it is clear that all what matters is our ability to decide on the pairs of hypotheses $\{x \in X : f(x) \leq a\}$ and $\{x \in X : f(X) \geq b\}$, with a and b given, via observation drawn from $p_{A(x)}$. In our outline, these were convex hypotheses in Gaussian o.s., and in this case we can use detector-based pairwise tests yielded by Theorem 2.25. Applying the machinery developed in Section 2.5.1, we could also handle the case when the sets $\{x \in X : f(x) \leq a\}$ and $\{x \in X : f(X) \geq b\}$ are unions of a moderate number of convex sets (e.g., f is affine, and X is the union of a number of convex sets), the o.s. in question still being simple, and this is the situation we intend to consider.

3.2.2 Estimating N-convex functions: problem's setting

In the rest of this Section, we consider the situation as follows. Given are:

- 1. simple o.s. $\mathcal{O} = ((\Omega, P), \{p_{\mu}(\cdot) : \mu \in \mathcal{M}\}, \mathcal{F}),$
- 2. convex compact set $\mathcal{X} \subset \mathbf{R}^n$ along with a collection of I convex compact sets $X_i \subset \mathcal{X}$,
- 3. affine "encoding" $x \mapsto A(x) : \mathcal{X} \to \mathcal{M},$
- 4. a continuous function $f(x): \mathcal{X} \to \mathbf{R}$ which is *N*-convex, meaning that for every

 $a \in \mathbf{R}$ the sets $\mathcal{X}^{a,\geq} = \{x \in \mathcal{X} : f(x) \geq a\}$ and $\mathcal{X}^{a,\leq} = \{x \in \mathcal{X} : f(x) \leq a\}$ can be represented as the unions of at most N closed convex sets $\mathcal{X}^{a,\geq}_{\nu}, \mathcal{X}^{a,\leq}_{\nu}$:

$$\mathcal{X}^{a,\geq} = \bigcup_{\nu=1}^{N} \mathcal{X}^{a,\geq}_{\nu}, \ \mathcal{X}^{a,\leq} = \bigcup_{\nu=1}^{N} \mathcal{X}^{a,\leq}_{\nu}.$$
 (3.20)

For some unknown x known to belong to $X = \bigcup_{i=1}^{I} X_i$, we have at our disposal observation $\omega^K = (\omega_1, ..., \omega_K)$ with i.i.d. $\omega_t \sim p_{A(x)}(\cdot)$, and our goal is to estimate from this observation the quantity f(x).

Given tolerances $\rho > 0$, $\epsilon \in (0,1)$, let us call a candidate estimate $\widehat{f}(\omega^K)$ (ρ, ϵ) -reliable, if for every $x \in X$, with the $p_{A(x)}$ -probability at least $1 - \epsilon$ it holds $|\widehat{f}(\omega^K) - f(x)| \leq \rho$ or, which is the same, if

$$\forall (x \in X) : \operatorname{Prob}_{\omega^{K} \sim p_{A(x)} \times \ldots \times p_{A(x)}} \left\{ |\widehat{f}(\omega^{K}) - f(x)| > \rho \right\} \le \epsilon.$$
(3.21)

3.2.2.1 Examples of N-convex functions

Example 3.3. [Minima and Maxima of linear-fractional functions] Every function which can be obtained from linear-fractional functions $\frac{g_{\nu}(x)}{h_{\nu}(x)}$ (g_{ν} , h_{ν} are affine functions on \mathcal{X} and h_{ν} are positive on \mathcal{X}) by taking maxima and minima is *N*-convex for appropriately selected *N* due to the following immediate observations:

- linear-fractional function $\frac{g(x)}{h(x)}$ with positive on \mathcal{X} denominator is 1-convex on \mathcal{X} ;
- if f(x) is N-convex, so is -f(x);
- if $f_i(x)$ is N_i -convex, i = 1, 2, ..., I, then $f(x) = \max_i f_i(x)$ is N-convex with

$$N = \max[\prod_{i} N_i, \sum_{i} N_i],$$

due to

$$\{x \in \mathcal{X} : f(x) \le a\} = \bigcap_{\substack{i=1\\I}}^{I} \{x : f_i(x) \le a\},$$
$$\{x \in \mathcal{X} : f(x) \ge a\} = \bigcup_{\substack{i=1\\i=1}}^{I} \{x : f_i(x) \ge a\}.$$

The first right hand side set is the intersection of I unions of convex sets with N_i components in *i*-th union, and thus is the union of $\prod_i N_i$ convex sets; the second right hand side set is the union of I unions, N_i components in *i*-th of them, of convex sets, and thus is the union of $\sum_i N_i$ convex sets.

Example 3.4. [Conditional quantile]. Let $S = \{s_1 < s_2 < ... < s_N\} \subset \mathbf{R}$ and T be a finite set, and let \mathcal{X} be a convex compact set in the space of nonvanishing probability distributions on $S \times T$. Given $\tau \in T$, consider the conditional, by the condition $t = \tau$, distribution $q_{\tau}[p]$ of $s \in S$ induced by a distribution $p(\cdot, \cdot) \in \mathcal{X}$:

$$(q_{\tau}[p])_{\mu} = \frac{p(\mu, \tau)}{\sum_{\nu=1}^{N} p(\nu, \tau)}.$$

For a nonvanishing probability distribution q on S and $\alpha \in (0, 1)$, let $\chi_{\alpha}(q)$ be the regularized α -quantile of q defined as follows: we pass from q to the distribution on $[s_1, s_N]$ by spreading uniformly the mass q_{ν} , $1 \leq \nu < N$, over $[s_{\nu}, s_{\nu+1}]$, and assigning mass q_N to the point s_N ; $\chi_{\alpha}(q)$ is the usual α -quantile of the resulting distribution \bar{q} :

$$\chi_{\alpha}(q) = \min \{ s \in [s_1, s_N] : \bar{q}\{(s, s_N]\} \le \alpha \}$$

The function $\chi_{\alpha}(q_{\tau}[p]): \mathcal{X} \to \mathbf{R}$ turns out to be 2-convex, see Section 3.6.2.

3.2.3 Bisection Estimate: Construction

While the construction to be presented admits numerous refinements, we focus here on its simplest version as follows.

3.2.3.1 Preliminaries

Upper and lower feasibility/infeasibility, sets $Z_i^{a,\geq}$ and $Z_i^{a,\leq}$. Let a be a real. We associate with a the collection of upper a-sets defined as follows: we look at the sets $X_i \cap \mathcal{X}_{\nu}^{a,\geq}$, $1 \leq i \leq I$, $1 \leq \nu \leq N$, and arrange the nonempty sets from this family into a sequence $Z_i^{a,\geq}$, $1 \leq i \leq I_{a,\geq}$, where $I_{a,\geq} = 0$ if all sets in the family are empty; in the latter case, we call a upper-infeasible, otherwise upper-feasible. Similarly, we associate with a the collection of lower a-sets $Z_i^{a,\leq}$, $1 \leq i \leq I_{a,\leq}$ by arranging into a sequence all nonempty sets from the family $X_i \cap \mathcal{X}_{\nu}^{a,\leq}$, and call a lower-feasible or lower-infeasible depending on whether $I_{a,\leq}$ is positive or zero. Note that upper and lower a-sets are nonempty convex compact sets, and

$$\begin{array}{ll} X^{a,\geq} &:= & \{x \in X : f(x) \geq a\} = \bigcup_{1 \leq i \leq I_{a,\geq}} Z_i^{a,\geq}, \\ X^{a,\leq} &:= & \{x \in X : f(x) \leq a\} = \bigcup_{1 \leq i \leq I_{a,\leq}} Z_i^{a,\leq}. \end{array}$$
(3.22)

Right-side tests. Given a segment $\Delta = [a, b]$ of positive length with lower-feasible a, we associate with this segment *right-side test* – a function $\mathcal{T}_{\Delta,r}^{K}(\omega^{K})$ taking values red and blue, and risk $\sigma_{\Delta,r} \geq 0$ – as follows:

- 1. if b is upper-infeasible, $\mathcal{T}_{\Delta,\mathbf{r}}^{K}(\cdot) \equiv$ blue and $\sigma_{\Delta,\mathbf{r}} = 0$.
- 2. if b is upper-feasible, the collections $\{A(Z_i^{b,\geq})\}_{i\leq I_{b,\geq}}$ ("red sets"), $\{A(Z_j^{a,\leq})\}_{j\leq I_{a,\leq}}$ ("blue sets"), are nonempty, and the test is given by the construction from Section 2.5.1 as applied to these sets and the stationary K-repeated version of \mathcal{O} in the role of \mathcal{O} , specifically,
 - for $1 \leq i \leq I_{b,\geq}$, $1 \leq j \leq I_{a,\leq}$, we build the detectors

$$\phi_{ij\Delta}^{K}(\omega^{K}) = \sum_{t=1}^{K} \phi_{ij\Delta}(\omega_{t}),$$

with
$$\phi_{ij\Delta}(\omega)$$
 given by

 set

$$\epsilon_{ij\Delta} = \int_{\Omega} \sqrt{p_{A(r_{ij\Delta})}(\omega)p_{A(s_{ij\Delta})}(\omega)}\Pi(d\omega)$$
(3.24)

and build the $I_{b,\geq} \times I_{a,\leq}$ matrix $E_{\Delta,\mathbf{r}} = [\epsilon_{ij\Delta}^K]_{\substack{1 \leq i \leq I_{b,\geq} \\ 1 \leq j \leq I_{a,\leq}}}$;

• $\sigma_{\Delta,\mathbf{r}}$ is defined as the spectral norm of $E_{\Delta,\mathbf{r}}$. We compute the Perron-Frobenius eigenvector $[g^{\Delta,\mathbf{r}};h^{\Delta,\mathbf{r}}]$ of the matrix $\left[\frac{|E_{\Delta,\mathbf{r}}|}{|E_{\Delta,\mathbf{r}}|}\right]$, so that (see Section 2.5.1.2)

$$\begin{split} g^{\Delta,\mathbf{r}} > 0, \, h^{\Delta,\mathbf{r}} > 0, \sigma_{\Delta,\mathbf{r}} g^{\Delta,\mathbf{r}} = E_{\Delta,\mathbf{r}} h^{\Delta,\mathbf{r}}, \\ \sigma_{\Delta,\mathbf{r}} h^{\Delta,\mathbf{r}} = E_{\Delta,\mathbf{r}}^T g^{\Delta,\mathbf{r}}. \end{split}$$

Finally, we define the matrix-valued function

$$D_{\Delta,\mathbf{r}}(\omega^K) = [\phi_{ij\Delta}^K(\omega^K) + \ln(h_j^{\Delta,\mathbf{r}}) - \ln(g_i^{\Delta,\mathbf{r}})]_{\substack{1 \le i \le I_{b,\geq}\\1 \le j \le I_{a,\leq}}}$$

Test $\mathcal{T}_{\Delta,\mathbf{r}}^{K}(\omega^{K})$ takes value red iff the matrix $D_{\Delta,\mathbf{r}}(\omega^{K})$ has a nonnegative row, and takes value blue otherwise.

Given $\delta > 0$, $\varkappa > 0$, we call segment $\Delta = [a, b] \delta$ -good (right), if a is lower-feasible, b > a, and $\sigma_{\Delta, \mathbf{r}} \leq \delta$. We call a δ -good (right) segment $\Delta = [a, b] \varkappa$ -maximal, if the segment $[a, b - \varkappa]$ is not δ -good (right).

Left-side tests. The "mirror" version of the above is as follows. Given a segment $\Delta = [a, b]$ of positive length with upper-feasible b, we associate with this segment *left-side test* – a function $\mathcal{T}_{\Delta,l}^{K}(\omega^{K})$ taking values red and blue, and risk $\sigma_{\Delta,l} \geq 0$ – as follows:

1. if a is lower-infeasible, $\mathcal{T}_{\Delta,l}^{K}(\cdot) \equiv \text{red and } \sigma_{\Delta,l} = 0.$

2. if a is lower-feasible, we set $\mathcal{T}_{\Delta,l}^K \equiv \mathcal{T}_{\Delta,r}^K$, $\sigma_{\Delta,l} = \sigma_{\Delta,r}$.

Given $\delta > 0$, $\varkappa > 0$, we call segment $\Delta = [a, b] \delta$ -good (left), if b is upper-feasible, b > a, and $\sigma_{\Delta,1} \leq \delta$. We call a δ -good (left) segment $\Delta = [a, b] \varkappa$ -maximal, if the segment $[a + \varkappa, b]$ is not δ -good (left).

Explanation: When a < b and a is lower-feasible, b is upper-feasible, so that the sets

$$X^{a,\leq} = \{x \in X : f(x) \le a\}, X^{b,\geq} = \{x \in X : f(x) \ge b\}$$

are nonempty, the right-side and the left-side tests $\mathcal{T}_{\Delta,l}^{K}$, $\mathcal{T}_{\Delta,r}^{K}$ are identical to each other and coincide with the minimal risk test, built as explained in Section 2.5.1, deciding, via stationary K-repeated observations, on the "color" of the distribution $p_{A(x)}$ underlying observations – whether this color is blue ("blue" hypothesis stating that $x \in X$ and $f(x) \leq a$, whence $A(x) \in \bigcup_{1 \leq i \leq I_{a,\leq}} A(Z_i^{a,\leq}))$, or red ("red" hypothesis, stating that $x \in X$ and $f(x) \geq b$, whence $A(x) \in \bigcup_{1 \leq i \leq I_{b,>}} A(Z_i^{b,\geq}))$. When a is

lower-feasible and b is *not* upper-feasible, the red one of the above two hypotheses is empty, and the left-side test associated with [a, b], naturally, always accepts the blue hypothesis; similarly, when a is lower-infeasible and b is upper-feasible, the right-side test associated with [a, b] always accepts the red hypothesis.

A segment [a, b] with a < b is δ -good (left), if the corresponding to the segment "red" hypothesis is nonempty, and the left-hand side test $\mathcal{T}_{\Delta \ell}^{K}$ associated with [a, b]decides on the "red" and the "blue" hypotheses with risk $\leq \delta$, and similarly for δ -good (right) segment [a, b].

3.2.4 Building the Bisection estimate

3.2.4.1 Control parameters

The control parameters of our would-be Bisection estimate are

- 1. positive integer L the maximum allowed number of bisection steps,
- 2. tolerances $\delta \in (0,1)$ and $\varkappa > 0$.

3.2.4.2 Bisection estimate: construction

The estimate of f(x) (x is the signal underlying our observations: $\omega_t \sim p_{A(x)}$) is given by the following recurrence run on the observation $\bar{\omega}^K = (\bar{\omega}_1, ..., \bar{\omega}_K)$ which we have at our disposal:

1. Initialization. We find a valid upper bound b_0 on $\max_{u \in X} f(u)$ and valid lower bound a_0 on $\min_{u \in X} f(u)$ and set $\Delta_0 = [a_0, b_0]$. We assume w.l.o.g. that $a_0 < b_0$, otherwise the estimation is trivial. Note: $f(a) \in \Delta_0$.

2. Bisection Step ℓ , $1 \leq \ell \leq L$. Given *localizer* $\Delta_{\ell-1} = [a_{\ell-1}, b_{\ell-1}]$ with $a_{\ell-1} < b_{\ell-1}$, we act as follows:

a) We set $c_{\ell} = \frac{1}{2}[a_{\ell-1} + b_{\ell-1}]$. If c_{ℓ} is not upper-feasible, we set $\Delta_{\ell} = [a_{\ell-1}, c_{\ell}]$ and pass to 2e, and if c_{ℓ} is not lower-feasible, we set $\Delta_{\ell} = [c_{\ell}, b_{\ell-1}]$ and pass to 2e.

<u>Note</u>: In the latter two cases, $\Delta_{\ell} \setminus \Delta_{\ell-1}$ does not intersect with f(X); in particular, in these cases $f(x) \in \Delta_{\ell}$ provided that $f(x) \in \Delta_{\ell-1}$.

b) When c_{ℓ} is both upper- and lower-feasible, we check whether the segment $[c_{\ell}, b_{\ell-1}]$ is δ -good (right). If it is not the case, we terminate and claim that $f(x) \in \overline{\Delta} := \Delta_{\ell-1}$, otherwise find $v_{\ell}, c_{\ell} < v_{\ell} \leq b_{\ell-1}$, such that the segment $\Delta_{\ell}^{rg} = [c_{\ell}, v_{\ell}]$ is δ -good (right) \varkappa -maximal.

<u>Note</u>: In terms of the outline of our strategy presented in Section 3.2.1, termination when the segment $[c_{\ell}, b_{\ell-1}]$ is not δ -good (right) corresponds to the case when the current localizer is too small to allow for "no-man land" wide enough to ensure low-risk decision on the blue and the red hypotheses.

<u>Note</u>: To find v_{ℓ} , we look one by one at the candidates with $v_{\ell}^{k} = b_{\ell-1} - k\varkappa$, k = 0, 1, ... until arriving for the first time at segment $[c_{\ell}, v_{\ell}^{k}]$ which is not δ -good (right), and take, as v_{ℓ} , the quantity v^{k-1} (when v_{ℓ} indeed is sought, we clearly have $k \geq 1$, so that our recipe for building v_{ℓ} is well-defined and clearly meets the above requirements on v_{ℓ}).

c) Similarly, we check whether the segment $[a_{\ell-1}, c_{\ell}]$ is δ -good (left). If it is not the case, we terminate and claim that $f(x) \in \overline{\Delta} := \Delta_{\ell-1}$, otherwise find $u_{\ell}, a_{\ell-1} \leq u_{\ell} < c_{\ell}$, such that the segment $\Delta_{\ell, \text{lf}} = [u_{\ell}, c_{\ell}]$ is δ -good (left)

207

 \varkappa -maximal.

- <u>Note</u>: The rules for building u_{ℓ} are completely similar to those for v_{ℓ} .
- d) We compute $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}},\mathrm{r}}^{K}(\bar{\omega}^{K})$ and $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}},\mathrm{l}}^{K}(\bar{\omega}^{K})$. If $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}},\mathrm{r}}^{K}(\bar{\omega}^{K}) = \mathcal{T}_{\Delta_{\ell,\mathrm{lf}},\mathrm{l}}^{K}(\bar{\omega}^{K})$ ("consensus"), we set

$$\Delta_{\ell} = [a_{\ell}, b_{\ell}] = \begin{cases} [c_{\ell}, b_{\ell-1}], & \mathcal{T}_{\Delta_{\ell, \mathrm{rg}}, \mathrm{r}}^{K}(\bar{\omega}^{K}) = \mathrm{red}, \\ [a_{\ell-1}, c_{\ell}], & \mathcal{T}_{\Delta_{\ell, \mathrm{rg}}, \mathrm{r}}^{K}(\bar{\omega}^{K}) = \mathrm{blue} \end{cases}$$
(3.25)

and pass to 2e. Otherwise ("disagreement") we terminate and claim that $f(x) \in \overline{\Delta} = [u_{\ell}, v_{\ell}].$

- e) When arriving at this rule, Δ_{ℓ} is already built. When $\ell < L$, we pass to step $\ell + 1$, otherwise we terminate with the claim that $f(x) \in \overline{\Delta} := \Delta_L$.
- 3. Output of the estimation procedure is the segment $\overline{\Delta}$ built upon termination and claimed to contain f(x), see rules 2b–2e; the midpoint of this segment is the estimate of f(x) yielded by our procedure.

3.2.5 Bisection estimate: Main result

Our main result on Bisection is as follows:

Proposition 3.5. Consider the situation described in the beginning of Section 3.2.2, and let $\epsilon \in (0, 1/2)$ be given. Then

(i) [reliability of Bisection] For every positive integer L and every $\kappa > 0$, Bisection with control parameters L,

$$\delta = \frac{\epsilon}{2L},$$

 κ is $(1-\epsilon)$ -reliable: for every $x \in X$, the $p_{A(x)}$ -probability of the event

$$f(x) \in \bar{\Delta}$$

 $(\Delta \text{ is the output of Bisection as defined above}) \text{ is at least } 1 - \epsilon.$

(ii) [near-optimality] Let $\rho > 0$ and positive integer \bar{K} be such that in the nature there exists a (ρ, ϵ) -reliable estimate $\hat{f}(\cdot)$ of $f(x), x \in X := \bigcup_{i \leq I} X_i$, via stationary \bar{K} -repeated observation $\omega^{\bar{K}}$ with $\omega_k \sim p_{A(x)}, 1 \leq k \leq \bar{K}$. Given $\hat{\rho} > 2\rho$, the Bisection estimate utilizing stationary K-repeated observations, with

$$K = \rfloor \frac{2\ln(2LNI/\epsilon)}{\ln(1/\epsilon) - \ln(4(1-\epsilon))} \bar{K} \lfloor, \qquad (3.26)$$

the control parameters of the estimate being

$$L = \lfloor \log_2\left(\frac{b_0 - a_0}{2\widehat{\rho}}\right) \lfloor, \ \delta = \frac{\epsilon}{2L}, \ \varkappa = \widehat{\rho} - 2\rho, \tag{3.27}$$

is $(\widehat{\rho}, \epsilon)$ -reliable. Not that K is only "slightly larger" than \overline{K} .

For proof, see Section 3.6.1.

Note that the running time K of Bisection estimate as given by (3.26) is just by (at most) logarithmic in N, I, L, $1/\epsilon$ factor larger than \bar{K} ; note also that L is just logarithmic in $1/\rho$. Assume, e.g., that for some $\gamma > 0$ "in the nature" there exist

 $(\epsilon^{\gamma}, \epsilon)$ reliable estimates, parameterized by $\epsilon \in (0, 1/2)$, with $\bar{K} = \bar{K}(\epsilon)$. Bisection with the volume of observation and control parameters given by (3.26) (3.27), where $\bar{\rho} = 3\rho = 3\epsilon^{\gamma}$ and $\bar{K} = \bar{K}(\epsilon)$, is $(3\epsilon^{\gamma}, \epsilon)$ -reliable and requires $K = K(\epsilon)$ -repeated observations with $\overline{\lim_{\epsilon \to +0} K(\epsilon)}/\bar{K}(\epsilon) \leq 2$.

3.2.6 Illustration

To illustrate bisection-based estimation of N-convex functional, consider the situation as follows⁴⁵. There are M devices ("receivers") recording a signal u known to belong to a given convex compact and nonempty set $U \subset \mathbf{R}^n$; the output of *i*-th receiver is the vector

$$y_i = A_i u + \sigma \xi \in \mathbf{R}^m \qquad [\xi \sim \mathcal{N}(0, I_m)]$$

where A_i are given $m \times n$ matrices; you may think about M allowed positions of a single receiver, and on y_i – as on the output of receiver when the latter is in position i. Our observation ω is one of the vectors y_i , $1 \le i \le M$ with unknown to us index i ("we observe a noisy record of signal, but do not know the position in which this record was taken"). Given ω , we want to recover a given linear function $g(x) = e^T u$ of the signal.

The problem can be modeled as follows. Consider the sets

$$X_i = \{x = [x^1; \dots; x^M] \in \mathbf{R}^{Mn} = \underbrace{\mathbf{R}^n \times \dots \times \mathbf{R}^n}_M : x^j = 0, j \neq i; x^i \in U\}$$

along with the linear mapping

$$A[x^1;...;x^M] = \sum_{i=1}^M A_i x^i : \mathbf{R}^{Mn} \to \mathbf{R}^m$$

and linear function

$$f([x^1;...;x^M]) = e^T \sum_i x^i : \mathbf{R}^{Mn} \to \mathbf{R},$$

and let \mathcal{X} be a convex compact set in \mathbf{R}^{Mn} containing all the sets X_i , $1 \leq i \leq m$. Observe that the problem we are interested in is nothing but the problem of recovering f(x) via observation

$$\omega = Ax + \sigma\xi, \ \xi \sim \mathcal{N}(0, I_m), \tag{3.28}$$

where the unknown signal x is known to belong to the union $\bigcup_{i=1}^{M} X_i$ of known convex compact sets X_i . As a result, our problem can be solved via the machinery we have developed.

Numerical illustration. In the numerical results to be reported, we used n = 128, m = 64 and M = 2. The data was generated as follows:

• The set $U \subset \mathbf{R}^{128}$ of candidate signals was comprised by restrictions onto equidistant (n = 128)-point grid in [0, 1] of twice differentiable functions h(t) of continu-

⁴⁵Our local goal is to illustrate a mathematical construction rather than to work out a particular application; the reader is welcome to invent a plausible "covering story" for this construction.

Characteristic	min	median	mean	max
error bound	0.008	0.015	0.014	0.015
actual error	0.001	0.002	0.002	0.005
# of Bisection steps	5	7.00	6.60	8

Table 3.1: Experiments with Bisection, data over 10 experiments, $\sigma = 0.01$. In the table, "error bound" is half-length of final localizer, which is an 0.99-reliable upper bound on the estimation error, and "actual error" is the actual estimation error.

ous argument $t \in [0, 1]$ satisfying the relations $|h(0)| \le 1$, $|h'(0)| \le 1$, $|h''(t)| \le 1$, $0 \le t \le 1$, which for the discretized signal u = [h(0); h(1/n); h(2/n); ...; h(1 - 1/n)] translates to the system of convex constraints

$$|u_1| \le 1, n|u_2 - u_1| \le 1, n^2|u_{i+1} - 2u_i + u_{i-1}| \le 1, 2 \le i \le n - 1.$$

- We were interested to recover the discretized counterpart of the integral $\int_0^1 h(t)dt$, specifically, dealt with $e = \bar{e}$, $\bar{e}^T u = \alpha \sum_{i=1}^n u_i$. The normalizing constant α was selected to ensure $\max_{u \in U} \bar{e}^T u = 1$, $\min_{u \in U} \bar{e}^T u = -1$, allowing to run Bisection with $\Delta_0 = [-1; 1]$.
- We generated A_1 as $(m = 64) \times (n = 128)$ matrix with singular values $\sigma_i = \theta^{i-1}$, $1 \leq i \leq m$, with θ selected from the requirement $\sigma_m = 0.1$. The system of left singular vectors of A_1 was obtained from the system of basic orths in \mathbb{R}^n by random rotation.

Matrix A_2 was selected as $A_2 = A_1 S$, where S was "reflection w.r.t. the axis \bar{e} ", that is,

$$S\bar{e} = \bar{e} \& Sh = -h$$
 whenever h is orthogonal to \bar{e} . (3.29)

Signals u underlying the observations were selected in U at random.

• The reliability $1 - \epsilon$ of our estimate was set to 0.99, and the maximal allowed number L of Bisection steps was set to 8. We used single observation (3.28) (i.e., used K = 1 in our general scheme) with σ set to 0.01.

The results of our experiments are presented in Table 3.1. Note that in the problem we are considering, there exists an intrinsic obstacle for high accuracy estimation even in the case of noiseless observations and invertible matrices A_i , i = 1, 2 (recall that we are in the case of M = 2). Indeed, assume that there exist $u \in U$, $u' \in U$ such that $A_1u = A_2u'$ and $e^T u \neq e^T u'$. In this case, when the signal is u and the (noiseless) observation is A_1u , the true quantity to be estimated is $e^T u$, and when the signal is u' and the observation is A_2u' , the true quantity to be estimated is $e^T u' \neq e^T u$. Since we do not know which of the matrices, A_1 or A_2 , underlies the observation and $A_1u = A_2u'$, there is no way to distinguish between the two cases we have described, implying that the quantity

$$\rho = \max_{u,u' \in U} \left\{ \frac{1}{2} |e^T(u - u')| : A_1 u = A_2 u' \right\}$$
(3.30)

is a lower bound on the worst-case, over signals from U, error of a reliable recovery of $e^T u$, independently of how small is the noise. In the reported experiments, we used $A_2 = A_1 S$ with S linked to $e = \bar{e}$, see (3.29); with this selection of S, $e = \bar{e}$

			Characteristic			m	in 1	median mean		m	iax			
			error bound		0.0	057	0.457 0.		0.441	1.(000			
			actual error		0.0	001	0.297		0.350	1.000				
			# of Bisection steps			1	1.00		2.20		5			
				"Difficult	" sign	als,	data ov	ver 1	10 exper	iments				
	ρ	0.0223	0.0281	0.1542	0.17	01	0.2130	0	0.2482	0.250)3	0.4999	0.6046	0.9238
;	error bound	0.0569	0.0625	0.2188	0.23	93	0.4063	3	0.5078	0.51	56	0.6250	0.7734	1.0000
		E	rror boun	d vs. ρ , e	xperin	nent	s sorted	d ac	cording	to the	valı	ues of ρ		
		ĺ	Char	acteristic		m	in n	ned	ian n	nean	m	ax		
	error bound		П	0.0	16	0.274 0.348 1.000		000						
	actual error			0.0	05	0.0	66 0	.127	0.5	556				
	# of Bisection steps		1		2.0	00 2	2.80	1	7					
Random signals, data over 10 experiments														
	ρ	0.0100	0.0853	0.1768	0.24	31	0.2940	0	0.3336	0.336	35	0.5535	0.6300	0.7616
1	error bound	0.0156	0.1816	0.3762	0.43	75	0.6016	6	0.0293	0.03	13	0.6875	0.1250	1.0000
	Error bound vs. ρ , experiments sorted according to the values of ρ													

Table 3.2: Experiments with randomly selected linear form, $\sigma = 0.01$

and A_2 , were A_1 invertible, the lower bound ρ would be just trivial – zero. In fact, our A_1 was not invertible, resulting in a positive ρ ; computation shows, however, that with our data, this positive ρ is negligibly small (about 2.0e-5). When we destroy the link between e and S, the estimation problem can become intrinsically more difficult, and the performance of our estimation procedure can deteriorate. Let us look what happens when we keep A_1 and $A_2 = A_1 S$ exactly as they are, but replace the linear form $\bar{e}^T u$ to be estimated with $e^T u$, e being randomly selected e⁴⁶. The corresponding data are presented in Table 3.2. The data in the top part of Table relate to the case of "difficult" signals u – those participating in forming the lower bound (3.30) on the recovery error, while the data in the bottom part of Table relate to randomly selected signals ⁴⁷. We see that when recovering the value of a randomly selected linear form, the error bounds indeed deteriorate, as compared to those in Table 3.1. We see also that the resulting error bounds are in reasonably good agreement with the lower bound ρ , illustrating the basic property of nearly optimal estimates: the guaranteed performance of an estimate can be bad or good, but it always is nearly as good as is possible under the circumstances. As about actual estimation errors, they in some experiments were essentially less than the error bounds, especially when random signals were used. This phenomenon, of course, should not be overestimated; remember that even a broken clock twice a day shows the correct time.

⁴⁶in the experiments to be reported, e was selected as follows: we start with a random unit vector drawn from the uniform distribution on the unit sphere in \mathbf{R}^n and then normalize it to make $\max_{u \in U} e^T u - \min_{u \in U} e^T u = 2$.

⁴⁷specifically, to generate a signal u, we drew a point \bar{u} at random, from the uniform distribution on the sphere of radius $10\sqrt{n}$, and took as u the $\|\cdot\|_2$ -closest to \bar{u} point of U.

3.2.7 Estimating N-convex functions: an alternative

Observe that the problem of estimating an N-convex function on the union of convex sets posed in Section 3.2.2 can be processed not only by Bisection. An alternative is as follows. In the notation from Section 3.2.2, we start with computing the range Δ of function f on the set $X = \bigcup_{i \leq I} X_i$, that is, we compute the quantities

$$\underline{f} = \min_{x \in X} f(x), \ \overline{f} = \max_{x \in X} f(x)$$

and set $\Delta = [\underline{f}, \overline{f}]$. We assume that this segment is not a singleton, otherwise estimating f is trivial. Further, we split Δ in a number L of consecutive bins – segments Δ_{ℓ} of equal length $\delta_L = (\overline{f} - \underline{f})/L$. δ_L will be the accuracy of our estimate; given a desired accuracy, we can select L accordingly. We now consider the sets

$$X_{i\ell} = \{ x \in X_i : f(x) \in \Delta_\ell \}, \ 1 \le i \le I, 1 \le \ell \le L.$$

Since f is N-convex, every one of these sets is the union of $M_{i\ell} \leq N^2$ convex compact sets $X_{i\ell j}$, $1 \leq j \leq M_{i\ell}$. Thus, we get at our disposal a collection of at most ILN^2 convex compact sets; let us eliminate from this collection empty sets and arrange the nonempty ones into a sequence $Y_1, ..., Y_M, M \leq ILN^2$. Note that $\bigcup_{s \leq M} Y_s = X$, so that the goal posed in Section 3.2.2 can be reformulated as follows:

For some unknown x known to belong to $X = \bigcup_{s=1}^{M} Y_s$, we have at our disposal observation $\omega^K = (\omega_1, ..., \omega_K)$ with i.i.d. $\omega_t \sim p_{A(x)}(\cdot)$; our goal is to estimate from this observation the quantity f(x).

The sets Y_s give rise to M hypotheses $H_1, ..., H_M$ on the distribution of our observations $\omega_t, 1 \le t \le K$; according to $H_s, \omega_t \sim p_{A(x)}(\cdot)$ with some $x \in Y_s$.

Let us define a closeness \mathcal{C} on the set of our M hypotheses as follows. Given $s \leq M$, the set Y_s is some $X_{i(s)\ell(s)j(s)}$; we say that two hypotheses, H_s and $H_{s'}$, are \mathcal{C} -close, if the segments $\Delta_{\ell(s)}$ and $\Delta_{\ell(s')}$ intersect. Observe that when H_s and $H_{s'}$ are not \mathcal{C} -close, the convex compact sets Y_s , Y'_s do not intersect, since the values of f on Y_s belong to $\Delta_{\ell(s)}$, the values of f on $Y_{s'}$ belong to $\Delta_{\ell(s')}$, and the segments $\Delta_{\ell(s)}$ and $\Delta_{\ell(s')}$ do not intersect.

Now let us apply to the hypotheses $H_1, ..., H_M$ our machinery for testing up to closeness \mathcal{C} , see Section 2.5.2. Assuming that whenever H_s and $H_{s'}$ are not \mathcal{C} -close, the risks $\epsilon_{ss'}$ defined in Section 2.5.2.2 are $< 1^{48}$, we, given tolerance $\epsilon \in (0, 1)$, can find $K = K(\epsilon)$ such that stationary K-repeated observation ω^K allows to decide $(1-\epsilon)$ -reliably on $H_1, ..., H_M$ up to closeness \mathcal{C} . As applied to ω^K , the corresponding test \mathcal{T}^K will accept some (perhaps, none) of the hypotheses, let the indexes of the accepted hypotheses form set $S = S(\omega^K)$. We convert S into an estimate $\hat{f}(\omega^K)$ of $f(x), x \in X = \bigcup_{s \le M} Y_s$ being the signal underlying our observation, as follows:

• when S is empty, the estimate is, say $(\overline{f} + f)/2$;

⁴⁸In our standard simple o.s.'s, this is the case whenever for s, s' in question the images of Y_s and $Y_{s'}$ under the mapping $x \mapsto A(x)$ do not intersect; this definitely is the case when $A(\cdot)$ is an embedding, since for our s, s', Y_s and $Y_{s'}$ do not intersect.

• when S is nonempty, we take the union $\Delta(S)$ of the segments $\Delta_{\ell(s)}$, $s \in S$, and our estimate is the average of the largest and the smallest elements of $\Delta(S)$.

It is immediately seen (check it!) that if the signal x underlying our stationary K-repeated observation ω^K belongs to some Y_{s_*} , so that the hypothesis H_{s_*} is true, and the outcome S of \mathcal{T}^K contains s_* and is such that for all $s \in S$ H_s and H_{s_*} are C-close to each other, we have $|f(x) - \hat{f}(\omega^K)| \leq \delta_L$. Note that since C-risk of \mathcal{T}^K is $\leq \epsilon$, the $p_{A(x)}$ -probability to get $|f(x) - \hat{f}(\omega^K)| \leq \delta_L$ is at least $1 - \epsilon$.

3.2.7.1 Numerical illustration

Our illustration deals with the situation when I = 1, $X = X_1$ is a convex compact set, and f(x) is fractional-linear: $f(x) = a^T x/c^T x$ with positive on X denominator. Specifically, assume we are given noisy measurements of voltages V_i at *some* nodes *i* and currents I_{ij} in *some* arcs (i, j) of an electric circuit, and want to recover the resistance of a particular arc (\hat{i}, \hat{j}) :

$$r_{\hat{i}\hat{j}} = \frac{V_{\hat{j}} - V_{\hat{i}}}{I_{\hat{i}\hat{j}}}.$$

In our experiment, we work with the data as follows:



We are in the situation N = 1, I = 1, implying M = L. When using L = 8, the projections of the sets Y_s , $1 \le s \le L = 8$ onto the 2D plane of variables $(V_{\hat{i}} - V_{\hat{i}}, I_{\hat{i}\hat{i}})$

213

are the "stripes" shown below:



The range of the unknown resistance turns out to be $\Delta = [1, 10]$.

In our experiment we worked with $\epsilon = 0.01$. Instead of looking for K such that K-repeated observation allows to recover 0.99-reliably the resistance in the arc of interest within accuracy $|\Delta|/L$, we looked for the largest observation noise σ allowing to achieve the desired recovery with single observation. The results for L = 8, 16, 32 are as follows

L	8	16	32
δ_L	$9/8 \approx 1.13$	$9/16 \approx 0.56$	$9/32 \approx 0.28$
σ	0.024	0.010	0.005
$\sigma_{ m opt}/\sigma \leq 1$	1.31	1.31	1.33
σ	0.031	0.013	0.006
$\sigma_{ m opt}/\sigma \leq 1$	1.01	1.06	1.08

In the table:

- σ_{opt} is the largest σ for which "in the nature" there exists a test deciding on $H_1, ..., H_L$ with C-risk ≤ 0.01 ;
- Red data: Risks $\epsilon_{ss'}$ of pairwise tests are bounded via risks of optimal detectors, C-risk of T is bounded by

$$\left| \left| \left[\epsilon_{ss'} \chi_{ss'} \right]_{s,s'=1}^{L} \right| \right|_{2,2}, \chi_{ss'} = \begin{cases} 1, & (s,s') \notin \mathcal{C} \\ 0, & (s,s') \in \mathcal{C} \end{cases}$$

see Proposition 2.33;

• Brown data: Risks $\epsilon_{ss'}$ of pairwise tests are bounded via error function, C-risk of \mathcal{T} is bounded by

$$\max_{s} \sum_{s':(s,s') \notin \mathcal{C}} \epsilon_{ss'}$$

(check that in the case of Gaussian o.s., this indeed is a legitimate risk bound).



Figure 3.2: A circuit (9 nodes, 16 arcs). Red: arc of interest; Green: arcs with measured currents and nodes with measured voltages.

3.2.7.2 Estimating dissipated power

The alternative approach to estimating N-convex functions proposed in Section 3.2.7 can be combined with quadratic lifting from Section 2.9 to yield, under favorable circumstances, estimates of quadratic and quadratic fractional functions. We are about to consider an instructive example of this sort. Figure 3.2 represent a DC electrical circuit. We have access to repeated noisy measurements of currents in green arcs and voltages at green nodes, with the voltage of the ground node equal to 0. The arcs are somehow oriented; this orientation, however, is of no relevance in our context and therefore is not displayed. Our goal is to use these observations to estimate the power dissipated in a given "arc of interest." Our a priori information is as follows:

- the (unknown) resistances of arcs are known to belong to a given range [r, R], with 0 < r < R < ∞;
- the currents and the voltages are linked by Kirchhoff Laws:
 - at every node, the sum of currents in the outgoing arcs is equal to the sum of currents in the incoming arcs plus the external current at the node.
 In our circuit, there are just two external currents, one at the ground node and one at the input node marked by dashed line.
 - the voltages and the currents are linked by Ohm's Law: for every (inner) arc γ , we have

$$I_{\gamma}r_{\gamma} = V_{j(\gamma)} - V_{i(\gamma)}$$

where I_{γ} is the current in the arc, r_{γ} is the arc's resistance, V_s is the voltage at node s, and $i(\gamma)$, $j(\gamma)$ are the initial and the final nodes linked by arc γ ;

• magnitudes of all currents and voltages are bounded by 1.

We assume that the measurements of observable currents and voltages are affected by zero mean Gaussian noise with scalar covariance matrix $\theta^2 I$, with unknown θ from a given range $[\underline{\sigma}, \overline{\sigma}]$.

Processing the problem. We specify the "signal" underlying our observation as the collection u of the voltages at our 9 nodes and currents I_{γ} in our 16 (inner) arcs γ , augmented by external current I_o at the input node (so that $-I_o$ is the external current at the ground node), so that our single-time observation is

$$\zeta = Au + \theta\xi, \tag{3.31}$$

where A extracts from u four entries, $\xi \sim \mathcal{N}(0, I_4)$, and $\theta \in [\underline{\sigma}, \overline{\sigma}]$. Our a priori information on u states that u belongs to the compact set U given by the quadratic constraints, namely, as follows:

$$U = \begin{cases} I_{\gamma}^{2} \leq 1, V_{i}^{2} \leq 1 \,\forall \gamma, i; u^{T} J^{T} J u = 0 \\ [V_{j(\gamma)} - V_{i(\gamma)}]^{2} / R - I_{\gamma} [V_{j(\gamma)} - V_{i(\gamma)}] \leq 0 \\ I_{\gamma} [V_{j(\gamma)} - V_{i(\gamma)}] - [V_{j(\gamma)} - V_{i(\gamma)}]^{2} / r \leq 0 \\ I_{\gamma} [V_{j(\gamma)} - V_{i(\gamma)}] - [V_{j(\gamma)} - V_{i(\gamma)}] \leq 0 \\ I_{\gamma} [V_{j(\gamma)} - V_{i(\gamma)}] - RI_{\gamma}^{2} \leq 0 \end{cases} \forall \gamma \quad (b) \end{cases}$$

$$(3.32)$$

where Ju = 0 expresses the first Kirchhoff's Law, and quadratic constraints (a), (b)account for the Ohm's Law in the situation when we do not know the resistances, just the range [r, R] of them. Note that the groups (a), (b) of constraints in (3.32)are "logical consequences" of each other, and thus one of groups seems to be redundant. However, on a closest inspection, valid on U quadratic inequalities are indeed redundant in our context, that is, do not tighten the outer approximation \mathcal{Z} of $\mathcal{Z}[U]$, only when these inequalities can be obtained from the inequalities we do include into the description of \mathcal{Z} "in a linear fashion" – by taking weighted sum with nonnegative coefficients; this is *not* how (b) is obtained from (a). As a result, to get a smaller \mathcal{Z} , it makes sense to keep both (a) and (b).

The dissipated power we are interested to estimate is the quadratic function

$$f(u) = I_{\gamma_*}[V_{j_*} - V_{i_*}] = [u; 1]^T G[u; 1]$$

where $\gamma_* = (i_*, j_*)$ is the arc of interest, and $G \in \mathbf{S}^{n+1}$, $n = \dim u$, is a properly built matrix.

In order to build an estimate, we "lift quadratically" the observations:

$$\zeta \mapsto \omega = (\zeta, \zeta \zeta^T)$$

and pass from the domain U of actual signals to the outer approximation \mathcal{Z} of the quadratic lifting of U:

$$\begin{aligned} \mathcal{Z} &:= & \{ Z \in \mathbf{S}^{n+1} : Z \succeq 0, Z_{n+1,n+1} = 1, \operatorname{Tr}(Q_s Z) \le c_s, \, 1 \le s \le S \} \\ &\supset & \{ [u; 1] [u; 1]^T : u \in \mathcal{V} \} \,, \end{aligned}$$

where the matrix $Q_s \in \mathbf{S}^{n+1}$ represents the left hand side $F_s(u)$ of s-th quadratic constraint participating in the description (3.32) of U: $F_s(u) \equiv [u; 1]^T Q_s[u; 1]$, and c_s is the right hand side of s-th constraint.

We process the problem similarly to what was done in Section 3.2.7.1, where our goal was to estimate a fractional-linear function. Specifically,

1. We compute the range of f on U; the smallest value \underline{f} of f on U clearly is zero, and an upper bound on the maximum of f(u) over $u \in U$, is the optimal value

in the convex optimization problem

$$\overline{f} = \max_{Z \in \mathcal{Z}} \operatorname{Tr}(GZ)$$

2. Given a positive integer L, we split the range $[\underline{f}, \overline{f}]$ into L segments $\Delta_{\ell} = [a_{\ell-1}, a_{\ell}]$ of equal length $\delta_L = (\overline{f} - f)/L$ and define convex compact sets

$$\mathcal{Z}_{\ell} = \{ Z \in \mathcal{Z} : a_{\ell-1} \le \operatorname{Tr}(GZ) \le a_{\ell} \}, \ 1 \le \ell \le L,$$

so that

$$u \in U, f(u) \in \Delta_{\ell} \Rightarrow [u; 1][u; 1]^T \in \mathcal{Z}_{\ell}, 1 \le \ell \le L_{\ell}$$

3. We specify L quadratically constrained hypotheses $H_1, ..., H_L$ on the distribution of observation (3.31), with H_ℓ stating that $\zeta \sim \mathcal{N}(Au, \theta^2 I_4)$ with some $u \in U$ satisfying $f(u) \in \Delta_\ell$ (so that $[u; 1][u; 1]^T \in \mathcal{Z}_\ell$), and θ belongs to the above segment $[\underline{\sigma}, \overline{\sigma}]]$.

We equip our hypotheses with closeness relation C, specifically, say that H_{ℓ} , $H_{\ell'}$ are C-close if and only if the segments Δ_{ℓ} and $\Delta_{\ell'}$ intersect.

4. We use Proposition 2.46.ii to build quadratic in ζ detectors $\phi_{\ell\ell'}$ for the families of distributions obeying H_{ℓ} and $H_{\ell'}$, respectively, along with upper bounds $\epsilon_{\ell\ell'}$ on the risks of these detectors, and then use the machinery from Section 2.5.2 to find the smallest K and a test $\mathcal{T}_{\mathcal{C}}^{K}$, based on stationary K-repeated version of observation (3.31), capable to decide on H_1, \ldots, H_L with \mathcal{C} -risk $\leq \epsilon$, where $\epsilon \in (0, 1)$ is a given tolerance.

Finally, given stationary K-repeated observation (3.31), we apply to it test $\mathcal{T}_{\mathcal{C}}^{K}$, look at the hypotheses, if any, accepted by the test, and build the union Δ of the corresponding segments Δ_{ℓ} . If $\Delta = \emptyset$, we estimate f(u) by the midpoint of the range $[\underline{f}, \overline{f}]$ of power, otherwise the estimate is the mean of the largest and the smallest points in Δ . It is easily seen (check it!) that for this estimate, the probability for the estimation error to be $> \delta_{\ell}$ is $\leq \epsilon$.

Numerical results we are about to report deal with the circuit presented on Figure 3.2; we used $\overline{\sigma} = 0.01$, $\underline{\sigma} = \overline{\sigma}/\sqrt{2}$, [r, R] = [1, 2], $\epsilon = 0.01$, and L = 8. The numerical results are as follows. The range $[\underline{f}, \overline{f}]$ of the dissipated power turned out to be [0, 0.821], so that the estimate we have built with reliability 0.99 recovers the dissipated power within accuracy 0.103. The resulting value of K was K = 95.

In a series of 500 simulations, the actual recovery error *all the time* was less than the bound 0.103, and the average error was as small as 0.041.

3.3 ESTIMATING LINEAR FORMS

We are about to demonstrate that the techniques developed in Section 2.8 can be applied to building estimates of linear and quadratic forms of the parameters of observed distributions. As compared to the machinery of Section 3.2, our new approach has somehow restricted scope: we cannot estimate anymore general Nconvex functions and/or handle domains which are unions of convex sets; now we need the function to be linear (perhaps, after quadratic lifting of observations) and the domain to be convex. As a compensation, the new approach, when applicable,

seems to be cheaper computationally: the estimate is yielded by solving a single convex problem, while the techniques developed so far require solving several (perhaps even few tens) of problems of similar structure and complexity. In this Section, we focus on estimating linear forms; estimating quadratic forms will be our subject in Section 3.4.

3.3.1 Situation and goal

Consider the situation as follows: given are Euclidean spaces $\Omega = \mathcal{E}_H, \mathcal{E}_M, \mathcal{E}_X$ along with

- regular data $\mathcal{H} \subset \mathcal{E}_H, \mathcal{M} \subset \mathcal{E}_M, \Phi(\cdot; \cdot) : \mathcal{H} \times \mathcal{M} \to \mathbf{R}$, with $0 \in \operatorname{int} \mathcal{H}$,
- a nonempty convex compact set $\mathcal{X} \subset \mathcal{E}_X$,
- an affine mapping $x \mapsto \mathcal{A}(x) : \mathcal{E}_X \to \mathcal{E}_M$ such that $\mathcal{A}(\mathcal{X}) \subset \mathcal{M}$,
- a continuous convex calibrating function $v(x) : \mathcal{X} \to \mathbf{R}$
- a vector $g \in \mathcal{E}_X$ and a constant c specifying the linear form $G(x) = \langle g, x \rangle + c : \mathcal{E}_X \to \mathbf{R}^{49}$,
- a tolerance $\epsilon \in (0, 1)$.

These data specify, in particular, the family

$$\mathcal{P} = \mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$$

of probability distributions on $\Omega = \mathcal{E}_H$, see Section 2.8.1.1. Given random observation

$$\omega \sim P(\cdot) \tag{3.33}$$

where $P \in \mathcal{P}$ is such that

$$\forall h \in \mathcal{H} : \ln\left(\int_{\mathcal{E}_H} e^{\langle h, \omega \rangle} P(d\omega)\right) \le \Phi(h; \mathcal{A}(x))$$
(3.34)

for some $x \in \mathcal{X}$ (that is, $\mathcal{A}(x)$ is a parameter, as defined in Section 2.8.1.1, of distribution P), we want to recover the quantity G(x).

 ϵ -risk. Given $\rho > 0$, we call an estimate $\widehat{g}(\cdot) : \mathcal{E}_H \to \mathbf{R} \ (\rho, \epsilon, v(\cdot))$ -accurate, if for all pairs $x \in \mathcal{X}, P \in \mathcal{P}$ satisfying (3.34) it holds

$$\operatorname{Prob}_{\omega \sim P}\left\{ \left| \widehat{g}(\omega) - G(x) \right| > \rho + \upsilon(x) \right\} \le \epsilon.$$

$$(3.35)$$

If ρ_* is the infimum of those ρ for which estimate \hat{g} is $(\rho, \epsilon, v(\cdot))$ -accurate, then clearly \hat{g} is $(\rho_*, \epsilon, v(\cdot))$ -accurate; we shall call ρ_* the ϵ -risk of the estimate \hat{g} taken w.r.t. the data $G(\cdot)$, \mathcal{X} , $v(\cdot)$ and $(\mathcal{A}, \mathcal{H}, \mathcal{M}, \Phi)$:

$$\operatorname{Risk}_{\epsilon}(\widehat{g}(\cdot)|G,\mathcal{X},\upsilon,\mathcal{A},\mathcal{H},\mathcal{M},\Phi) = \min\left\{\rho:\operatorname{Prob}_{\omega\sim P}\{\omega:|\widehat{g}(\omega)-G(x)|>\rho+\upsilon(x)\} \le \epsilon \\ \forall (x,P): \left\{ \begin{array}{c} P\in\mathcal{P}, x\in X\\ \ln\left(\int e^{h^{T}\omega}P(d\omega)\right) \le \Phi(h;\mathcal{A}(x))\forall h\in\mathcal{H} \end{array} \right\}.$$

$$(3.36)$$

⁴⁹from now on, $\langle u, v \rangle$ denotes the inner product of vectors u, v belonging to a Euclidean space; what is this space, it always will be clear from the context.

218

LECTURE 3

When $G, X, v, \mathcal{A}, \mathcal{H}, \mathcal{M}, \Phi$ are clear from the context, we shorten

$$\operatorname{Risk}_{\epsilon}(\widehat{g}(\cdot)|G, X, \upsilon, \mathcal{A}, \mathcal{H}, \mathcal{M}, \Phi)$$

to $\operatorname{Risk}_{\epsilon}(\widehat{g}(\cdot))$.

Given the data listed in the beginning of this section, we are about to build, in a computationally efficient fashion, an affine estimate $\hat{g}(\omega) = \langle h_*, \omega \rangle + \varkappa$ along with ρ_* such that the estimate is $(\rho_*, \epsilon, \upsilon(\cdot))$ -accurate.

3.3.2 Construction & Main results

Let us set

$$\mathcal{H}^+ = \{(h, \alpha) : h \in \mathcal{E}_H, \alpha > 0, h/\alpha \in \mathcal{H}\}$$

so that \mathcal{H}^+ is a nonempty convex set in $\mathcal{E}_H \times \mathbf{R}_+$, and let

(a)
$$\Psi_{+}(h,\alpha) = \sup_{x \in \mathcal{X}} \left[\alpha \Phi(h/\alpha, \mathcal{A}(x)) - G(x) - \upsilon(x) \right] : \mathcal{H}^{+} \to \mathbf{R},$$

(b) $\Psi_{-}(h,\beta) = \sup_{x \in \mathcal{X}} \left[\beta \Phi(-h/\beta, \mathcal{A}(x)) + G(x) - \upsilon(x) \right] : \mathcal{H}^{+} \to \mathbf{R},$
(3.37)

so that Ψ_{\pm} are convex real-valued functions on \mathcal{H}^+ (recall that Φ is convex-concave and continuous on $\mathcal{H} \times \mathcal{M}$, while $\mathcal{A}(\mathcal{X})$ is a compact subset of \mathcal{M}).

Our starting point is pretty simple:

Proposition 3.6. Given $\epsilon \in (0,1)$, let \bar{h} , $\bar{\alpha}$, $\bar{\beta}$, $\bar{\varkappa}$, $\bar{\rho}$ be a feasible solution to the system of convex constraints

in variables $h, \alpha, \beta, \rho, \varkappa$. Setting

$$\widehat{g}(\omega) = \langle \overline{h}, \omega \rangle + \overline{\varkappa},$$

we get an estimate with ϵ -risk at most $\bar{\rho}$.

Proof. Let $\epsilon \in (0, 1)$, \bar{h} , $\bar{\alpha}$, $\bar{\beta}$, $\bar{\varkappa}$, $\bar{\rho}$ satisfy the premise of Proposition, and let $x \in X, P$ satisfy (3.34). We have

$$\begin{aligned} \operatorname{Prob}_{\omega \sim P} \{ \widehat{g}(\omega) > G(x) + \bar{\rho} + \upsilon(x) \} &= \operatorname{Prob}_{\omega \sim P} \left\{ \frac{\langle \bar{h}, \omega \rangle}{\bar{\alpha}} > \frac{G(x) + \bar{\rho} - \bar{\varkappa} + \upsilon(x)}{\bar{\alpha}} \right\} \\ \Rightarrow \quad \operatorname{Prob}_{\omega \sim P} \{ \widehat{g}(\omega) > G(x) + \bar{\rho} + \upsilon(x) \} \leq \left[\int e^{\langle \bar{h}, \omega \rangle / \bar{\alpha}} P(d\omega) \right] e^{-\frac{G(x) + \bar{\rho} - \bar{\varkappa} + \upsilon(x)}{\bar{\alpha}}} \\ &< e^{\Phi(\bar{h}/\bar{\alpha}, \mathcal{A}(x))} e^{-\frac{G(x) + \bar{\rho} - \bar{\varkappa} + \upsilon(x)}{\bar{\alpha}}} \end{aligned}$$

- $\Rightarrow \quad \bar{\alpha} \ln \left(\operatorname{Prob}_{\omega \sim P} \{ \widehat{g}(\omega) > G(x) + \bar{\rho} + \upsilon(x) \} \right) \\ \leq \bar{\alpha} \Phi(\bar{h}/\bar{\alpha}, \mathcal{A}(x)) G(x) \bar{\rho} \upsilon(x) + \bar{\varkappa} \\ \leq \Psi_{+}(\bar{h}, \bar{\alpha}) \bar{\rho} + \bar{\varkappa} \text{ [by definition of } \Psi_{+} \text{ and due to } x \in X] \\ \leq \bar{\alpha} \ln(\epsilon/2) \text{ [by (b_1)]}$
- $\Rightarrow \operatorname{Prob}_{\omega \sim P} \{ \widehat{g}(\omega) > G(x) + \overline{\rho} + \upsilon(x) \} \le \epsilon/2,$

and similarly

$$\begin{aligned} \operatorname{Prob}_{\omega\sim P}\{\widehat{g}(\omega) < G(x) - \bar{\rho} - \upsilon(x)\} &= \operatorname{Prob}_{\omega\sim P}\left\{\frac{-\langle \bar{h}, \omega \rangle}{\bar{\beta}} > \frac{-G(x) + \bar{\rho} + \bar{\varkappa} + \upsilon(x)}{\bar{\beta}}\right\} \\ \Rightarrow & \operatorname{Prob}_{\omega\sim P}\{\widehat{g}(\omega) < G(x) - \bar{\rho} - \upsilon(x)\} \leq \left[\int e^{-\langle \bar{h}, \omega \rangle / \bar{\beta}} P(d\omega)\right] e^{-\frac{-G(x) + \bar{\rho} + \bar{\varkappa} + \upsilon(x)}{\bar{\beta}}} \\ & \leq e^{\Phi(-\bar{h}/\bar{\beta}, \mathcal{A}(x))} e^{\frac{G(x) - \bar{\rho} - \bar{\varkappa} - \upsilon(x)}{\bar{\beta}}} \\ \Rightarrow & \bar{\beta} \ln\left(\operatorname{Prob}_{\omega\sim P}\{\widehat{g}(\omega) < G(x) - \bar{\rho} - \upsilon(x)\}\right) \\ &\leq \bar{\beta} \Phi(-\bar{h}/\bar{\beta}, \mathcal{A}(x)) + G(x) - \bar{\rho} - \bar{\varkappa} - \upsilon(x) \\ &\leq \Psi_{-}(\bar{h}, \bar{\beta}) - \bar{\rho} - \bar{\varkappa} \text{ [by definition of } \Psi_{-} \text{ and due to } x \in X] \\ &\leq \bar{\beta} \ln(\epsilon/2) \text{ [by } (b_2)] \\ \Rightarrow & \operatorname{Prob}_{\omega\sim P}\{\widehat{g}(\omega) < G(x) - \bar{\rho} - \upsilon(x)\} \leq \epsilon/2. \end{aligned}$$

Corollary 3.7. In the situation described in Section 3.3.1, let Φ satisfy the relation

$$\Phi(0;\mu) \ge 0 \ \forall \mu \in \mathcal{M}. \tag{3.39}$$

Then

$$\begin{split} \Psi_{+}(h) &:= \inf_{\alpha} \left\{ \Psi_{+}(h,\alpha) + \alpha \ln(2/\epsilon) : \alpha > 0, (h,\alpha) \in \mathcal{H}^{+} \right\} \\ &= \sup_{x \in \mathcal{X}} \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^{+}} \left[\alpha \Phi(h/\alpha, \mathcal{A}(x)) - G(x) - \upsilon(x) + \alpha \ln(2/\epsilon) \right], \quad (a) \\ \widehat{\Psi}_{-}(h) &:= \inf_{\alpha} \left\{ \Psi_{-}(h,\alpha) + \alpha \ln(2/\epsilon) : \alpha > 0, (h,\alpha) \in \mathcal{H}^{+} \right\} \end{split}$$

$$= \sup_{x \in \mathcal{X}} \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^+} \left[\alpha \Phi(-h/\alpha, \mathcal{A}(x)) + G(x) - \upsilon(x) + \alpha \ln(2/\epsilon) \right].$$
(b)
(3.40)

and functions $\widehat{\Psi}_{\pm} : \mathcal{E}_H \to \mathbf{R}$ are convex. Furthermore, let $\overline{h}, \overline{\varkappa}, \widetilde{\rho}$ be a feasible solution to the system of convex constraints

$$\widehat{\Psi}_{+}(h) \le \rho - \varkappa, \ \widehat{\Psi}_{-}(h) \le \rho + \varkappa$$
(3.41)

in variables h, ρ, \varkappa . Then, setting

$$\widehat{g}(\omega) = \langle \overline{h}, \omega \rangle + \overline{\varkappa}$$

we get an estimate of G(x), $x \in X$, with ϵ -risk at most $\widehat{\Psi}(\overline{h})$:

$$\operatorname{Risk}_{\epsilon}(\widehat{g}(\cdot)|G, X, v, \mathcal{A}, \mathcal{H}, \mathcal{M}, \Phi) \leq \widehat{\Psi}(\overline{h}).$$
(3.42)

Relation (3.41) (and thus – the risk bound (3.42)) clearly holds true when \bar{h} is a candidate solution to the convex optimization problem

$$Opt = \min_{h} \left\{ \widehat{\Psi}(h) := \frac{1}{2} \left[\widehat{\Psi}_{+}(h) + \widehat{\Psi}_{-}(h) \right] \right\},$$
(3.43)

 $\bar{\rho} = \widehat{\Psi}(\bar{h}), and$

$$\bar{\varkappa} = \frac{\widehat{\Psi}_{-}(\bar{h}) - \widehat{\Psi}_{+}(\bar{h})}{2}$$

As a result, properly selecting \overline{h} , we can make (an upper bound on) the ϵ -risk of estimate $\widehat{g}(\cdot)$ arbitrarily close to Opt, and equal to Opt when optimization problem (3.43) is solvable.

Proof. Let us first verify the equalities in (3.40). The function

$$\Theta_{+}(h,\alpha;x) = \alpha \Phi(h/\alpha,\mathcal{A}(x)) - G(x) - \upsilon(x) + \alpha \ln(2/\epsilon) : \mathcal{H}^{+} \times \mathcal{X} \to \mathbf{R}$$

is convex-concave and continuous, and \mathcal{X} is compact, whence by Sion-Kakutani Theorem

$$\begin{split} \widehat{\Psi}_{+}(h) &:= \inf_{\alpha} \left\{ \Psi_{+}(h,\alpha) + \alpha \ln(2/\epsilon) : \alpha > 0, (h,\alpha) \in \mathcal{H}^{+} \right\} \\ &= \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^{+}} \max_{x \in \mathcal{X}} \Theta_{+}(h,\alpha;x) \\ &= \sup_{x \in \mathcal{X}} \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^{+}} \Theta_{+}(h,\alpha;x) \\ &= \sup_{x \in \mathcal{X}} \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^{+}} \left[\alpha \Phi(h/\alpha, \mathcal{A}(x)) - G(x) - \upsilon(x) + \alpha \ln(2/\epsilon) \right], \end{split}$$

as required in (3.40.a). As we know, $\Psi_+(h,\alpha)$ is real-valued continuous function on \mathcal{H}^+ , so that Ψ_+ is convex on \mathcal{E}_H , provided that the function is real-valued. Now, let $\bar{x} \in \mathcal{X}$, and let e be a subgradient of $\phi(h) = \Phi(h; \mathcal{A}(x))$ taken at h = 0. For $h \in \mathcal{E}_H$ and all $\alpha > 0$ such that $(h, \alpha) \in \mathcal{H}^+$ we have

$$\begin{split} \Psi_{+}(h,\alpha) &\geq \alpha \Phi(h/\alpha;\mathcal{A}(\bar{x})) - G(\bar{x}) - \upsilon(\bar{x}) + \alpha \ln(2/\epsilon) \\ &\geq \alpha [\Phi(0;\mathcal{A}(\bar{x})) + \langle e, h/\alpha \rangle] - G(\bar{x}) - \upsilon(\bar{x}) + \alpha \ln(2/\epsilon) \\ &\geq \langle e, h \rangle - G(\bar{x}) - \upsilon(\bar{x}) \end{split}$$

(we have used (3.39)), and therefore $\Psi_{+}(h, \alpha)$ is bounded below on the set $\{\alpha > 0 : \{\alpha > 0\}$ $h/\alpha \in \mathcal{H}$; in addition, this set is nonempty, since \mathcal{H} contains a neighbourhood of the origin. Thus, $\widehat{\Psi}_+$ is real-valued and convex on \mathcal{E}_H . Verification of (3.40.b) and of the fact that $\widehat{\Psi}_{-}(h)$ is real-valued convex function on \mathcal{E}_{H} is completely similar.

Now, given a feasible solution $(\bar{h}, \bar{\varkappa}, \tilde{\rho})$ to (3.41), let us select somehow $\bar{\rho} > \tilde{\rho}$. Taking into account the definition of $\widehat{\Psi}_{\pm}$, we can find $\overline{\alpha}$ and $\overline{\beta}$ such that

$$(\bar{h},\bar{\alpha}) \in \mathcal{H}^+ \& \Psi_+(\bar{h},\bar{\alpha}) + \bar{\alpha}\ln(2/\epsilon) \le \bar{\rho} - \bar{\varkappa}, (\bar{h},\bar{\beta}) \in \mathcal{H}^+ \& \Psi_-(\bar{h},\bar{\beta}) + \bar{\beta}\ln(2/\epsilon) \le \bar{\rho} + \bar{\varkappa},$$

implying that the collection $(\bar{h}, \bar{\alpha}, \bar{\beta}, \bar{\varkappa}, \bar{\rho})$ is a feasible solution to (3.38). Invoking Proposition 3.6, we get

$$\operatorname{Prob}_{\omega \sim P} \left\{ \omega : |\widehat{g}(\omega) - G(x)| > \overline{\rho} + \upsilon(x) \right\} \le \epsilon$$

for all $(x \in X, P \in \mathcal{P})$ satisfying (3.34). Since $\bar{\rho}$ can be selected arbitrarily close to $\widetilde{\rho}, \, \widehat{q}(\cdot)$ indeed is a $(\widetilde{\rho}, \epsilon, v(\cdot))$ -accurate estimate.

3.3.3 Estimation from repeated observations

Assume that in the situation described in section 3.3.1 we have access to K observations $\omega_1, \ldots, \omega_K$ sampled, independently of each other, from a probability distribution P, and are allowed to build our estimate based on these K observations rather than on a single observation. We can immediately reduce this new situation to the previous one, just by redefining the data. Specifically, given initial data $\mathcal{H} \subset \mathcal{E}_H$, $\mathcal{M} \subset \mathcal{E}_M, \Phi(\cdot; \cdot) : \mathcal{H} \times \mathcal{M} \to \mathbf{R}, X \subset \mathcal{X} \subset \mathcal{E}_X, \mathcal{A}(\cdot), G(x) = g^T x + c$, see section 3.3.1 and a positive integer K, let us update part of the data, specifically, replace $\mathcal{H} \subset \mathcal{E}_{\mathcal{H}}$ with $\mathcal{H}^{K} := \underbrace{\mathcal{H} \times ... \times \mathcal{H}}_{K} \subset \mathcal{E}_{H}^{K} := \underbrace{\mathcal{E}_{H} \times ... \times \mathcal{E}_{H}}_{K}$ and replace $\Phi(\cdot, \cdot) : \mathcal{H} \times \mathcal{M} \to \mathbf{R}$ with $\Phi^{K}(h^{K} = (h_{1}, ..., h_{K}); \mu) = \sum_{i=1}^{K} \Phi(h_{i}; \mu) : \mathcal{H}^{K} \times \mathcal{M} \to \mathbf{R}$. It is immediately

seen that the updated data satisfy all requirements imposed on the data in section 3.3.1, and that whenever a Borel probability distribution \mathcal{P} on $\mathcal{E}_{\mathcal{H}}$ and $x \in X$ are linked by (3.34), the distribution P^K of K-element i.i.d. sample $\omega^K = (\omega_1, ..., \omega_K)$ drawn from P and x are linked by the relation

$$\forall h^{K} = (h_{1}, ..., h_{K}) \in \mathcal{H}^{K} : \ln \left(\int_{\mathcal{E}_{H}^{K}} e^{\langle h^{K}, \omega^{K} \rangle} P^{K}(d\omega^{K}) \right) = \sum_{i} \ln \left(\int_{\mathcal{E}_{H}} e^{\langle h_{i}, \omega_{i} \rangle} P(d\omega_{i}) \right)$$

$$\leq \Phi^{K}(h^{K}; \mathcal{A}(x)).$$

$$(3.44)$$

Applying to our new data the construction from section 3.3.2, we arrive at "repeated observations" versions of Proposition 3.6 and Corollary 3.7. Note that the resulting convex constraints/objectives are symmetric w.r.t. permutations functions of the components $h_1, ..., h_K$ of h^K , implying that we lose nothing when restricting ourselves with collections h^K with equal to each other components; it is convenient to denote the common value of these components h/K. With these observations, Proposition 3.6 and Corollary 3.7 become the statements as follows (we use the assumptions and the notation from the previous sections):

Proposition 3.8. Given $\epsilon \in (0,1)$ and positive integer K, let

(a)
$$\Psi_{+}(h,\alpha) = \sup_{x \in \mathcal{X}} \left[\alpha \Phi(h/\alpha, \mathcal{A}(x)) - G(x) - \upsilon(x) \right] : \mathcal{H}^{+} \to \mathbf{R},$$

(b) $\Psi_{-}(h,\beta) = \sup_{x \in \mathcal{X}} \left[\beta \Phi(-h/\beta, \mathcal{A}(x)) + G(x) - \upsilon(x) \right] : \mathcal{H}^{+} \to \mathbf{R},$

and let \bar{h} , $\bar{\alpha}$, $\bar{\beta}$, $\bar{\varkappa}$, $\bar{\rho}$ be a feasible solution to the system of convex constraints

in variables $h, \alpha, \beta, \rho, \varkappa$. Setting

$$\widehat{g}(\omega^{K}) = \langle \overline{h}, \frac{1}{K} \sum_{i=1}^{K} \omega_{i} \rangle + \overline{\varkappa}$$

we get an estimate of G(x) via independent K-repeated observations

$$\omega_i \sim P, i = 1, ..., K$$

with ϵ -risk on X not exceeding $\overline{\rho}$, meaning that whenever $x \in X$ and a Borel probability distribution P on $\mathcal{E}_{\mathcal{H}}$ are linked by (3.34), one has

$$\operatorname{Prob}_{\omega^{K} \sim P^{K}} \left\{ \omega^{K} : |\widehat{g}(\omega^{K}) - G(x)| > \overline{\rho} + \upsilon(x) \right\} \le \epsilon.$$

$$(3.46)$$

Corollary 3.9. In the situation described in the beginning of section 3.3.1, let Φ

satisfy the relation (3.39), and let a positive integer K be given. Then

$$\begin{split} \Psi_{+}(h) &:= \inf_{\alpha} \left\{ \Psi_{+}(h,\alpha) + K^{-1}\alpha\ln(2/\epsilon) : \alpha > 0, (h,\alpha) \in \mathcal{H}^{+} \right\} \\ &= \sup_{x \in \mathcal{X}} \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^{+}} \left[\alpha \Phi(h/\alpha, \mathcal{A}(x)) - G(x) - \upsilon(x) + K^{-1}\alpha\ln(2/\epsilon) \right], \quad (a) \\ \widehat{\Psi}_{-}(h) &:= \inf_{\alpha} \left\{ \Psi_{-}(h,\alpha) + K^{-1}\alpha\ln(2/\epsilon) : \alpha > 0, (h,\alpha) \in \mathcal{H}^{+} \right\} \\ &= \sup_{x \in \mathcal{X}} \inf_{\alpha > 0, (h,\alpha) \in \mathcal{H}^{+}} \left[\alpha \Phi(-h/\alpha, \mathcal{A}(x)) + G(x) - \upsilon(x) + K^{-1}\alpha\ln(2/\epsilon) \right]. \quad (b) \end{split}$$

$$(3.47)$$

and functions $\widehat{\Psi}_{\pm} : \mathcal{E}_H \to \mathbf{R}$ are convex. Furthermore, let $\overline{h}, \overline{\varkappa}, \widetilde{\rho}$ be a feasible solution to the system of convex constraints

$$\widehat{\Psi}_{+}(h) \le \rho - \varkappa, \ \widehat{\Psi}_{-}(h) \le \rho + \varkappa \tag{3.48}$$

in variables h, ρ, \varkappa . Then, setting

$$\widehat{g}(\omega^{K}) = \langle \overline{h}, \frac{1}{K} \sum_{i=1}^{K} \omega_{i} \rangle + \overline{\varkappa}$$

we get an estimate of G(x), $x \in X$, with ϵ -risk at most $\widehat{\Psi}(\overline{h})$, meaning that whenever $x \in X$ and a Borel probability distribution P on $\mathcal{E}_{\mathcal{H}}$ are linked by (3.34), relation (3.46) holds true.

Relation (3.48) clearly holds true when \bar{h} is a candidate solution to the convex optimization problem

$$Opt = \min_{h} \left\{ \widehat{\Psi}(h) := \frac{1}{2} \left[\widehat{\Psi}_{+}(h) + \widehat{\Psi}_{-}(h) \right] \right\}, \qquad (3.49)$$

 $\bar{\rho} = \widehat{\Psi}(\bar{h})$ and

$$\bar{\varkappa} = \frac{\widehat{\Psi}_{-}(\bar{h}) - \widehat{\Psi}_{+}(\bar{h})}{2}.$$

As a result, properly selecting \bar{h} , we can make (an upper bound on) the ϵ -risk of estimate $\hat{g}(\cdot)$ arbitrarily close to Opt, and equal to Opt when optimization problem (3.49) is solvable.

From now on, if otherwise is not explicitly stated, we deal with K-repeated observations; to get back to single-observation case, it suffices to set K = 1.

3.3.4 Application: Estimating linear form of sub-Gaussianity parameters

Consider the simplest case of the situation from sections 3.3.1, 3.3.3, where

- $\mathcal{H} = \mathcal{E}_H = \mathbf{R}^d$, $\mathcal{M} = \mathcal{E}_M = \mathbf{R}^d \times \mathbf{S}^d_+$, $\Phi(h; \mu, M) = h^T \mu + \frac{1}{2} h^T M h : \mathbf{R}^d \times (\mathbf{R}^d \times \mathbf{S}^d_+) \to \mathbf{R}$, so that $\mathcal{S}[\mathcal{H}, \mathcal{M}, \Phi]$ is the family of all sub-Gaussian distributions on \mathbf{R}^d ;
- $X = \mathcal{X} \subset \mathcal{E}_X = \mathbf{R}^{n_x}$ is a nonempty convex compact set, and
- $\mathcal{A}(x) = (Ax + a, M(x))$, where A is $d \times n_x$ matrix, and M(x) is affinely depending on x symmetric $d \times d$ matrix such that M(x) is $\succeq 0$ when $x \in X$,
- v(x) is a convex continuous function on \mathcal{X} ,
- G(x) is an affine function on \mathcal{E}_X .

In the case in question (3.39) clearly takes place, and the left hand sides in the constraints (3.48) are

$$\begin{split} \widehat{\Psi}_{+}(h) &= \sup_{x \in X} \inf_{\alpha > 0} \left\{ h^{T}[Ax + a] + \frac{1}{2\alpha} h^{T} M(x) h + K^{-1} \alpha \ln(2/\epsilon) - G(x) - \upsilon(x) \right\} \\ &= \max_{x \in X} \left\{ \sqrt{2K^{-1} \ln(2/\epsilon) [h^{T} M(x) h]} + h^{T}[Ax + a] - G(x) - \upsilon(x) \right\}, \\ \widehat{\Psi}_{-}(h) &= \sup_{x \in X} \inf_{\alpha > 0} \left\{ -h^{T}[Ax + a] + \frac{1}{2\alpha} h^{T} M(x) h + K^{-1} \alpha \ln(2/\epsilon) + G(x) - \upsilon(x) \right\} \\ &= \max_{x \in X} \left\{ \sqrt{2K^{-1} \ln(2/\epsilon) [h^{T} M(x) h]} - h^{T}[Ax + a] + G(x) - \upsilon(x) \right\}. \end{split}$$

Thus, system (3.48) reads

$$\begin{aligned} a^T h + \max_{x \in X} \left[\sqrt{2K^{-1} \ln(2/\epsilon) [h^T M(x)h]} + h^T A x - G(x) - \upsilon(x) \right] &\leq \rho - \varkappa, \\ -a^T h + \max_{x \in X} \left[\sqrt{2K^{-1} \ln(2/\epsilon) [h^T M(x)h]} - h^T A x + G(x) - \upsilon(x) \right] &\leq \rho + \varkappa. \end{aligned}$$

We arrive at the following version of Corollary 3.9:

Proposition 3.10. In the situation described in the beginning of section 3.3.4, given $\epsilon \in (0, 1)$, let \overline{h} be a feasible solution to the convex optimization problem

Then, setting

$$\bar{\varkappa} = \frac{1}{2} \left[\widehat{\Psi}_{-}(\bar{h}) - \widehat{\Psi}_{+}(\bar{h}) \right], \ \bar{\rho} = \widehat{\Psi}(\bar{h}), \tag{3.51}$$

the affine estimate

$$\widehat{g}(\omega^K) = \frac{1}{K} \sum_{i=1}^K \overline{h}^T \omega_i + \overline{\varkappa}$$

has ϵ -risk, taken w.r.t. the data listed in the beginning of this section, at most $\bar{\rho}$.

It is immediately seen that optimization problem (3.50) is solvable, provided that

$$\bigcap_{x \in X} \operatorname{Ker}(M(x)) = \{0\},\$$

and an optimal solution h_* to the problem, taken along with

$$\varkappa_{*} = \frac{1}{2} \left[\widehat{\Psi}_{-}(h_{*}) - \widehat{\Psi}_{+}(h_{*}) \right], \qquad (3.52)$$

yields the affine estimate

$$\widehat{g}_*(\omega) = \frac{1}{K} \sum_{i=1}^K h_*^T \omega_i + \varkappa_*$$

with $\epsilon\text{-risk},$ taken w.r.t. the data listed in the beginning of this section, at most Opt.

3.3.4.1 Consistency

Assuming $v(x) \equiv 0$, we can easily answer the natural question "when the proposed estimation scheme is consistent", meaning that for every $\epsilon \in (0,1)$, it allows to achieve arbitrarily small ϵ -risk, provided that K is large enough. Specifically, denoting by $g^T x$ the linear part of G(x): $G(x) = g^T x + c$, from Proposition 3.10 it is immediately seen that a sufficient condition for consistency is the existence of $\bar{h} \in \mathbf{R}^d$ such that $\bar{h}^T Ax = g^T x$ for all $x \in \mathcal{X} - \mathcal{X}$, or, equivalently, the condition that g is orthogonal to the intersection of the kernel of A with the linear span of $\mathcal{X} - \mathcal{X}$. Indeed, under this assumption, for every fixed $\epsilon \in (0, 1)$ we clearly have $\lim_{K\to\infty} \widehat{\Phi}(\bar{h}) = 0$, implying that $\lim_{K\to\infty} \operatorname{Opt} = 0$, with $\widehat{\Psi}$ and Opt given by (3.50). Still assuming $v(x) \equiv 0$, the condition in question is necessary for consistency as well, since when the condition is violated, we have Ax' = Ax'' for properly selected $x', x'' \in \mathcal{X}$ with $G(x') \neq G(x'')$, making low risk recovery of $G(x), x \in \mathcal{X}$, impossible already in the case of zero noisy component in observations⁵⁰.

3.3.4.2 Direct product case

Further simplifications are possible in the *direct product case*, where, in addition to what was assumed in the beginning of section 3.3.4,

- $\mathcal{E}_X = \mathcal{E}_U \times \mathcal{E}_V$ and $X = U \times V$, with convex compact sets $U \subset \mathcal{E}_U = \mathbf{R}^{n_u}$ and $V \subset E_V = \mathbf{R}^{n_v}$,
- $\mathcal{A}(x = (u, v)) = [Au + a, M(v)] : U \times V \to \mathbf{R}^d \times \mathbf{S}^d$, with $M(v) \succeq 0$ for $v \in V$,
- $G(x = (u, v)) = g^T u + c$ depends solely on u, and
- $v(x = (u, v)) = \varrho(u)$ depends solely on u.

It is immediately seen that in the direct product case problem (3.50) reads

$$Opt = \min_{h \in \mathbf{R}^d} \left\{ \frac{\phi_U(A^T h - g) + \phi_U(-A^T h + g)}{2} + \max_{v \in V} \sqrt{2K^{-1} \ln(2/\epsilon) h^T M(v) h} \right\},$$
(3.53)

where

$$\phi_U(f) = \max_{u \in U} \left[u^T f - \varrho(u) \right].$$
(3.54)

Assuming $\bigcap_{v \in V} \operatorname{Ker}(M(v)) = \{0\}$, the problem is solvable, and its optimal solution h_* produces affine estimate

$$\widehat{g}_{*}(\omega^{K}) = \frac{1}{K} \sum_{i} h_{*}^{T} \omega_{i} + \varkappa_{*}, \ \varkappa_{*} = \frac{1}{2} [\phi_{U}(-A^{T}h + g) - \phi_{U}(A^{T}h - g)] - a^{T}h_{*} + c$$

with ϵ -risk \leq Opt.

Near-optimality In addition to the assumption that we are in the direct product case, assume that $v(\cdot) \equiv 0$ and, for the sake of simplicity, that $M(v) \succ 0$ whenever

⁵⁰Note that in Gaussian case with M(x) depending on x the above condition is, in general, not necessary for consistency, since a nontrivial information on x (and thus on G(x)) can, in principle, be extracted from the covariance matrix M(x) which can be estimated from observations.

225

 $v \in V$. In this case (3.50) reads

$$Opt = \min_{h} \max_{v \in V} \left\{ \Theta(h, v) := \frac{1}{2} [\phi_U(A^T h - g) + \phi_U(-A^T h + g)] + \sqrt{2K^{-1} \ln(2/\epsilon) h^T M(v) h} \right\},$$

whence, taking into account that $\Theta(h, v)$ clearly is convex in h and concave in v, while V is a convex compact set, by Sion-Kakutani Theorem we get also

$$Opt = \max_{v \in V} \left[Opt(v) = \min_{h} \frac{1}{2} [\phi_U(A^T h - g) + \phi_U(-A^T h + g)] + \sqrt{2K^{-1} \ln(2/\epsilon) h^T M(v) h} \right]$$
(3.55)

Now consider the problem of recovering $g^T u$ from observation ω_i , $i \leq K$, independently of each other sampled from $\mathcal{N}(Au + a, M(v))$, where unknown u is known to belong to U and $v \in V$ is known. Let $\rho_{\epsilon}(v)$ be the minimax ϵ -risk of the recovery:

$$\rho_{\epsilon}(v) = \inf_{\widehat{g}(\cdot)} \left\{ \rho : \operatorname{Prob}_{\omega^{K} \sim [\mathcal{N}(Au+a, M(v))]^{K}} \{ \omega^{K} : |\widehat{g}(\omega^{K}) - g^{T}u| > \rho \} \le \epsilon \ \forall u \in U \right\},$$

where inf is taken over all Borel functions $\hat{g}(\cdot) : \mathbf{R}^{Kd} \to \mathbf{R}$. Invoking [86, Theorem 3.1], it is immediately seen that whenever $\epsilon < 1/4$, one has

$$\rho_{\epsilon}(v) \geq \left[\frac{2\ln(2/\epsilon)}{\ln\left(\frac{1}{4\epsilon}\right)}\right]^{-1} \operatorname{Opt}(v).$$

Since the family SG(U, V) of all sub-Gaussian, with parameters $(Au+a, M(v)), u \in U, v \in V$, distributions on \mathbb{R}^d contains all Gaussian distributions $\mathcal{N}(Au+a, M(v))$ induced by $(u, v) \in U \times V$, we arrive at the following conclusion:

Proposition 3.11. In the just described situation, the minimax optimal ϵ -risk

$$\operatorname{Risk}_{\epsilon}^{\operatorname{opt}}(K) = \inf_{\widehat{g}(\cdot)} \operatorname{Risk}_{\epsilon}(\widehat{g}(\cdot)),$$

of recovering $g^T u$ from K-repeated i.i.d. sub-Gaussian, with parameters $(Au + a, M(v)), (u, v) \in U \times V$, random observations is within a moderate factor of the upper bound Opt on the ϵ -risk, taken w.r.t. the same data, of the affine estimate $\hat{g}_*(\cdot)$ yielded by an optimal solution to (3.53), namely,

$$Opt \le \frac{2\ln(2/\epsilon)}{\ln\left(\frac{1}{4\epsilon}\right)} Risk_{\epsilon}^{opt}.$$

3.3.4.3 Numerical illustration

The numerical illustration we are about to discuss models the situation when we want to recover a linear form of a signal x known to belong to a given convex compact subset X via indirect observations Ax affected by sub-Gaussian "relative noise," meaning that the variance of observation is the larger the larger is the signal. Specifically, our observation is

$$\omega \sim \mathcal{S}G(Ax, M(x)),$$

where

$$x \in X = \left\{ x \in \mathbf{R}^n : 0 \le x_j \le j^{-\alpha}, 1 \le j \le n \right\}, \ M(x) = \sigma^2 \sum_{j=1}^n x_j \Theta_j \qquad (3.56)$$

where $A \in \mathbf{R}^{d \times n}$ and $\Theta_j \in \mathbf{S}_+^d$, j = 1, ..., n, are given matrices; the linear form to be recovered from observation ω is $G(x) = g^T x$. The entities $g, A, \{\Theta_j\}_{j=1}^n$ and reals $\alpha \ge 0$ ("degree of smoothness"), $\sigma > 0$ ("noise intensity") are parameters of the estimation problem we intend to process. The parameters g, A, Θ_j were generated as follows:

• $g \ge 0$ was selected at random and then normalized to have

$$\max_{x \in X} g^{T} x = \max_{x, y \in X} g^{T} [x - y] = 2;$$

- we dealt with n > d ("deficient observations"); the *d* nonzero singular values of *A* were set to $\theta^{-\frac{i-1}{d-1}}$, where "condition number" $\theta \ge 1$ is a parameter; the orthonormal systems *U* and *V* of the first *d* left, respectively, right singular vectors of *A* were drawn at random from rotationally invariant distributions;
- the positive semidefinite $d \times d$ matrices Θ_j were orthogonal projectors on randomly selected subspaces in \mathbf{R}^d of dimension |d/2|;
- in all our experiments, we dealt with single-observation case K = 1, and used $v(\cdot) \equiv 0$.

Note that X possesses \geq -largest point \bar{x} , whence $M(x) \leq M(\bar{x})$ whenever $x \in X$; as a result, sub-Gaussian distributions with matrix parameter M(x), $x \in X$, can be thought also to have matrix parameter $M(\bar{x})$. One of the goals of experiment to be reported was to understand how much would be lost were we replacing $M(\cdot)$ with $\widehat{M}(x) \equiv M(\bar{x})$, that is, were we ignoring the fact that small signals result in low-noise observations.

In the experiment to be reported, we use d = 32, m = 48, $\alpha = 2$, $\theta = 2$, and $\sigma = 0.01$. Utilizing these parameters, we generated at random, as described above, 10 collections $\{g, A, \Theta_j, j \leq d\}$, thus arriving at 10 estimation problems. For every one of these problems, we used the outlined machinery to build affine in ω estimate of $g^T x$ as yielded by optimal solution to (3.50), and computed upper bound Opt on $(\epsilon = 0.01)$ -risk of this estimate. In fact, for every one of the 10 generated estimation problems, we build two estimates and two risk bounds: the first – for the problem "as is," and the second – for the aforementioned "direct product envelope" of the problem, where the mapping $x \mapsto M(x)$ is replaced with $x \mapsto \widehat{M}(x) := M(\overline{x})$. The results are as follows:

\min	median	mean	max
0.138	0.190	0.212	0.299
0.150	0.210	0.227	0.320

0.01-Risk, data over 10 estimation problems $[d = 32, m = 48, \alpha = 2, \theta = 2, \sigma = 0.1]$ First row: $\omega \sim SG(Ax, M(x))$. Second row: $\omega \sim SG(Ax, M(\bar{x}))$

Pay attention to "amplification of noise" in the estimate (about 20 times the level σ of observation noise) and significant variability of risk across the experiments;

seemingly, both these phenomena stem from the fact that we have highly deficient observations (n/d = 1.5) combined with "random interplay" between the directions of coordinate axes in \mathbf{R}^m (along these directions, X becomes more and more thin) and the orientation of the 16-dimensional kernel of A.

3.4 ESTIMATING QUADRATIC FORMS VIA QUADRATIC LIFTING

In the situation of Section 3.33, passing from "original" observations (3.33) to their quadratic lifting: we can use the just developed machinery to estimate quadratic forms of the underlying parameters rather than linear ones. We are about to investigate the related possibilities in the cases of Gaussian and sub-Gaussian observations.

3.4.1 Estimating quadratic forms, Gaussian case

3.4.1.1 Preliminaries

Consider the situation where we are given

- a nonempty bounded set U in \mathbf{R}^m ;
- a nonempty convex compact subset \mathcal{V} of the positive semidefinite cone \mathbf{S}^{d}_{+} ;
- a matrix $\Theta_* \succ 0$ such that $\Theta_* \succeq \Theta$ for all $\Theta \in \mathcal{V}$;
- an affine mapping $u \mapsto A[u;1] : \mathbf{R}^m \to \Omega = \mathbf{R}^d$, where A is a given $d \times (m+1)$ matrix,
- a convex continuous function $\rho(\cdot)$ on \mathbf{S}^{m+1}_+ .

A pair $(u \in U, \Theta \in \mathcal{V})$ specifies Gaussian random vector $\zeta \sim \mathcal{N}(A[u; 1], \Theta)$ and thus specifies probability distribution $P[u, \Theta]$ of $(\zeta, \zeta\zeta^T)$. Let $\mathcal{Q}(U, \mathcal{V})$ be the family of probability distributions on $\Omega = \mathbf{R}^d \times \mathbf{S}^d$ stemming in this fashion from Gaussian distributions with parameters from $U \times \mathcal{V}$. Our goal is to cover the family $\mathcal{Q}(U, \mathcal{V})$ by a family of the type $\mathcal{S}[N, \mathcal{M}, \Phi]$.

It is convenient to represent a linear form on $\Omega = \mathbf{R}^d \times \mathbf{S}^d$ as

$$h^T z + \frac{1}{2} \operatorname{Tr}(HZ),$$

where $(h, H) \in \mathbf{R}^d \times \mathbf{S}^d$ is the "vector of coefficients" of the form, and $(z, Z) \in \mathbf{R}^d \times \mathbf{S}^d$ is the argument of the form.

We assume that for some $\delta \in [0, 2]$ it holds

$$\|\Theta^{1/2}\Theta_*^{-1/2} - I\| \le \delta \ \forall \Theta \in \mathcal{V},\tag{3.57}$$

where $\|\cdot\|$ is the spectral norm (cf. (2.165). Finally, we set $B = \begin{bmatrix} A \\ b^T \end{bmatrix}$ and

$$\mathcal{Z}^{+} = \{ W \in \mathbf{S}_{+}^{m+1} : W_{m+1,m+1} = 1 \}.$$
(3.58)

The statement below is a straightforward reformulation of Proposition 2.46.i:

Proposition 3.12. In the just described situation, let us select $\gamma \in (0,1)$ and set

$$\mathcal{H} = \mathcal{H}_{\gamma} := \{(h, H) \in \mathbf{R}^{d} \times \mathbf{S}^{d} : -\gamma \Theta_{*}^{-1} \leq H \leq \gamma \Theta_{*}^{-1} \},$$

$$\mathcal{M}^{+} = \mathcal{V} \times \mathcal{Z}^{+},$$

$$\Phi(h, H; \Theta, Z) = -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2} H \Theta_{*}^{1/2}) + \frac{1}{2} \operatorname{Tr}([\Theta - \Theta_{*}] H)$$

$$+ \frac{\delta(2 + \delta)}{2(1 - \|\Theta_{*}^{1/2} H \Theta_{*}^{1/2}\|)} \|\Theta_{*}^{1/2} H \Theta_{*}^{1/2}\|_{F}^{2}$$

$$+ \Gamma(h, H; Z) : \mathcal{H} \times \mathcal{M}^{+} \to \mathbf{R}$$

$$[\|\cdot\| \text{ is the spectral, and } \|\cdot\|_{F} \text{ is the Frobenius norm}],$$

$$\Gamma(h, H; Z) = \frac{1}{2} \operatorname{Tr} \left(Z \left(bh^{T} A + A^{T} hb^{T} + A^{T} HA + B^{T} [H, h]^{T} [\Theta_{*}^{-1} - H]^{-1} [H, h] B \right) \right).$$

$$= \frac{1}{2} \operatorname{Tr} \left(Z \left(B^{T} \left[\left[\frac{H + h}{h^{T}} \right] + [H, h]^{T} [\Theta_{*}^{-1} - H]^{-1} [H, h] \right] B \right) \right).$$

$$(3.59)$$

Then $\mathcal{H}, \mathcal{M}^+, \Phi$ form a regular data, and for every $(u, \Theta) \in \mathbf{R}^m \times \mathcal{V}$ it holds

$$\forall (h,H) \in \mathcal{H} : \ln\left(\mathbf{E}_{\zeta \sim \mathcal{N}(\mathcal{C}(u),\Theta)}\left\{e^{h^T\zeta + \frac{1}{2}\zeta^T H\zeta}\right\}\right) \leq \Phi(h,H;\Theta,[u;1][u;1]^T).$$
(3.60)

Besides this, function $\Phi(h, H; \Theta, Z)$ is coercive in the convex argument: whenever $(\Theta, Z) \in \mathcal{M}$ and $(h_i, H_i) \in \mathcal{H}$ and $||(h_i, H_i)|| \to \infty$ as $i \to \infty$, we have $\Phi(h_i, H_i; \Theta, Z) \to \infty$, $i \to \infty$.

3.4.1.2 Estimating quadratic form: Situation & goal

We are interested in the situation as follows: we are given a sample $\zeta^K = (\zeta_1, ..., \zeta_K)$ of independent across *i* and identically distributed random observations

$$\zeta_i \sim \mathcal{N}(A[u;1], M(v)), \ 1 \le i \le K, \tag{3.61}$$

where

• (u, v) is unknown "signal" known to belong to a given set $U \times V$, where

 $- U \subset \mathbf{R}^m$ is a compact set, and

- $V \subset \mathbf{R}^k$ is a compact convex set;
- A is a given $d \times (m+1)$ matrix, and $v \mapsto M(v) : \mathbf{R}^k \to \mathbf{S}^d$ is affine mapping such that $M(v) \succeq 0$ whenever $v \in V$.

We are also given a convex calibrating function $\varrho(Z):{\bf S}^{m+1}_+\to {\bf R}$ and "functional of interest"

$$F(u,v) = [u;1]^T Q[u;1] + q^T v, \qquad (3.62)$$

where Q and q are known $(m + 1) \times (m + 1)$ symmetric matrix and k-dimensional vector, respectively. Our goal is to recover F(u, v), for unknown (u, v) known to belong to $U \times V$, via observation (3.61). Given a tolerance $\epsilon \in (0, 1)$, we quantify the quality of a candidate estimate $\hat{g}(\zeta^K)$ of F(u, v) by the smallest ρ such that for all $(u, v) \in U \times V$ it holds

$$\operatorname{Prob}_{\zeta^{K} \sim \mathcal{N}(A[u;1],M(v)) \times \ldots \times \mathcal{N}(A[u;1],M(v))} \left\{ |\widehat{g}(\zeta^{K}) - F(u,v) \rho + \varrho([u;1][u;1]^{T}) \right\} \leq \epsilon.$$

$$(3.63)$$

3.4.1.3 Construction & Result

Let

$$\mathcal{V} = \{ M(v) : v \in V \},\$$

so that \mathcal{V} is a convex compact subset of the positive semidefinite cone \mathbf{S}^{d+1}_+ . Let us select somehow

- 1. a matrix $\Theta_* \succ 0$ such that $\Theta_* \succeq \Theta$, for all $\Theta \in \mathcal{V}$;
- 2. a convex compact subset \mathcal{Z} of the set $\mathcal{Z}^+ = \{Z \in \mathbf{S}^{m+1}_+ : Z_{m+1,m+1} = 1\}$ such that $[u; 1][u; 1]^T \in \mathcal{Z}$ for all $u \in U$;
- 3. a real $\gamma \in (0,1)$ and a nonnegative real δ such that (3.57) takes place.

We further set (cf. Proposition 3.12)

$$B = \begin{bmatrix} A \\ [0,...,0,1] \end{bmatrix} \in \mathbf{R}^{(d+1)\times(m+1)},$$

$$\mathcal{H} = \mathcal{H}_{\gamma} := \{(h,H) \in \mathbf{R}^{d} \times \mathbf{S}^{d} : -\gamma\Theta_{*}^{-1} \preceq H \preceq \gamma\Theta_{*}^{-1}\},$$

$$\mathcal{M} = \mathcal{V} \times \mathcal{Z},$$

$$\Phi(h,H;\Theta,Z) = -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2}H\Theta_{*}^{1/2}) + \frac{1}{2}\operatorname{Tr}([\Theta - \Theta_{*}]H) + \frac{\delta(2+\delta)}{2(1-\|\Theta_{*}^{1/2}H\Theta_{*}^{1/2}\|)} \|\Theta_{*}^{1/2}H\Theta_{*}^{1/2}\|_{F}^{2} + \Gamma(h,H;Z) : \mathcal{H} \times \mathcal{M} \to \mathbf{R}$$

$$[\|\cdot\| \text{ is the spectral, and } \|\cdot\|_{F} \text{ is the Frobenius norm}],$$

$$\Gamma(h,H;Z) = \frac{1}{2}\operatorname{Tr}\left(Z[bh^{T}A + A^{T}hb^{T} + A^{T}HA + B^{T}[H,h]^{T}[\Theta_{*}^{-1} - H]^{-1}[H,h]B]\right) = \frac{1}{2}\operatorname{Tr}\left(Z\left(B^{T}\left[\left[\frac{H}{h^{T}}\right] + [H,h]^{T}[\Theta_{*}^{-1} - H]^{-1}[H,h]\right]B\right)\right)\right)$$

$$(3.64)$$

and treat, as our observation, the quadratic lift of observation (3.61), that is, our observation is

$$\omega^{K} = \{\omega_{i} = (\zeta_{i}, \zeta_{i}\zeta_{i}^{T})\}_{i=1}^{K}, \zeta_{i} \sim \mathcal{N}(A[u; 1], M(v)) \text{ are independent across } i.$$
(3.65)

Note that by Proposition 3.12, function $\Phi(h, H; \Theta, Z) : \mathcal{H} \times \mathcal{M} \to \mathbf{R}$ is continuous convex-concave function which is coercive in convex argument and is such that

$$\forall (u \in U, v \in V, (h, H) \in \mathcal{H}) :$$

$$\ln \left(\mathbf{E}_{\zeta \sim \mathcal{N}(A[u;1], M(v))} \left\{ e^{\frac{1}{2} \zeta^T H \zeta + h^T \zeta} \right\} \right) \leq \Phi(h, H; M(v), [u;1][u;1]^T).$$
 (3.66)

We are about to demonstrate that as far as estimating the functional of interest (3.62) at a point $(u, v) \in U \times V$ via observation (3.65) is concerned, we are in the situation considered in Section 3.3 and can use the machinery developed there. Indeed, let us specify the data introduced in section 3.3.1 and participating in the constructions of section 3.3 as follows:

• $\mathcal{H} = \{f = (h, H) \in \mathcal{H}\} \subset \mathcal{E}_H = \mathbf{R}^d \times \mathbf{S}^d$, with \mathcal{H} defined in (3.64), and the inner product on \mathcal{E}_H defined as

$$\langle (h,H), (h',H') \rangle = h^T h' + \frac{1}{2} \operatorname{Tr}(HH'),$$

 $\mathcal{E}_M = \mathbf{S}^{d+1} \times \mathbf{S}^{m+1}, \text{ and } \mathcal{M}, \Phi \text{ are as defined in (3.64);}$ • $\mathcal{E}_X = \mathbf{R}^k \times \mathbf{S}^{d+1}, X := \{x = (v, Z) : v \in V, Z = [u; 1][u; 1]^T, u \in U\} \subset \mathcal{X} :=$

 $\{x = (v, Z) \in V \times \mathcal{Z}\};\$

• $\mathcal{A}(x) = \mathcal{A}(v, Z) = (M(v), Z)$; note that \mathcal{A} is affine mapping from \mathcal{E}_X into \mathcal{E}_M mapping \mathcal{X} into \mathcal{M} , as required in section 3.3. Observe that when $u \in U$ and $v \in V$, the distribution $P = P_{u,v}$ of observation ω defined by (3.65) satisfies the relation

$$\forall (f = (h, H) \in \mathcal{H}) : \ln \left(\mathbf{E}_{\omega \sim P} \left\{ e^{\langle f, \omega \rangle} \right\} \right) = \ln \left(\mathbf{E}_{\zeta \sim \mathcal{N}(A[u;1], M(v))} \left\{ e^{h^T \zeta + \frac{1}{2} \zeta^T H \zeta} \right\} \right)$$

$$\leq \Phi(h, H; M(v), [u;1][u;1]^T),$$

$$(3.67)$$

see (3.66);

- $v(x = (v, Z)) = \varrho(Z) : \mathcal{X} \to \mathbf{R},$
- we define affine functional G(x) on \mathcal{E}_X by the relation

$$\langle g, x := (v, Z) \rangle = q^T v + \operatorname{Tr}(QZ),$$

see (3.62). As a result, for $x \in X$, that is, for $x = (v, [u; 1][u; 1]^T)$ with $v \in V$ and $u \in U$ we have

$$F(u,v) = G(x).$$
 (3.68)

Applying to the just specified data Corollary 3.9 (which is legitimate, since our Φ clearly satisfies (3.39)), we arrive at the result as follows:

Proposition 3.13. In the just described situation, let us set

$$\begin{split} \widehat{\Psi}_{+}(h,H) &:= \inf_{\alpha} \left\{ \max_{(v,Z)\in V\times \mathcal{Z}} \left[\alpha \Phi(h/\alpha,H/\alpha;M(v),Z) - G(v,Z) - \varrho(Z) + K^{-1}\alpha \ln(2/\epsilon) \right] : \\ \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right\} \\ &= \max_{(v,Z)\in V\times \mathcal{Z}} \inf_{-\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1}} \left[\alpha \Phi(h/\alpha,H/\alpha;M(v),Z) - G(v,Z) - \varrho(Z) + K^{-1}\alpha \ln(2/\epsilon) \right], \\ \widehat{\Psi}_{-}(h,H) &:= \inf_{\alpha} \left\{ \max_{(v,Z)\in V\times \mathcal{Z}} \left[\alpha \Phi(-h/\alpha,-H/\alpha;M(v),Z) + G(v,Z) - \varrho(Z) + K^{-1}\alpha \ln(2/\epsilon) \right] : \\ \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right\} \\ &= \max_{(v,Z)\in V\times \mathcal{Z}} \inf_{-\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1}} \left[\alpha \Phi(-h/\alpha,-H/\alpha;M(v),Z) + G(v,Z) - \varrho(Z) + K^{-1}\alpha \ln(2/\epsilon) \right]. \end{split}$$

so that the functions $\Psi_{\pm}(h, H) : \mathbf{R}^d \times \mathbf{S}^d \to \mathbf{R}$ are convex. Furthermore, whenever $\bar{h}, \bar{H}, \bar{\rho}, \bar{\varkappa}$ form a feasible solution to the system of convex constraints

$$\widehat{\Psi}_{+}(h,H) \le \rho - \varkappa, \ \widehat{\Psi}_{-}(h,H) \le \rho + \varkappa$$
(3.70)

in variables $(h, H) \in \mathbf{R}^d \times \mathbf{S}^d$, $\rho \in \mathbf{R}$, $\varkappa \in \mathbf{R}$, setting

$$\widehat{g}(\zeta^{K} := (\zeta_{1}, ..., \zeta_{K})) = \frac{1}{K} \sum_{i=1}^{K} \left[h^{T} \zeta_{i} + \frac{1}{2} \zeta_{i}^{T} H \zeta_{i} \right] + \bar{\varkappa}, \qquad (3.71)$$

we get an estimate of the functional of interest $F(u, v) = [u; 1]^T Q[u; 1] + q^T v$ via K independent observations

$$\zeta_i \sim \mathcal{N}(A[u;1], M(v)), \ i = 1, \dots, K,$$

with the following property:

$$\forall (u,v) \in U \times V : Prob_{\zeta^{K} \sim [\mathcal{N}(A[u;1],M(v))]^{K}} \left\{ |F(u,v) - \widehat{g}(\zeta^{K})| > \bar{\rho} + \varrho([u;1][u;1]^{T}) \right\} \leq \epsilon.$$

$$(3.72)$$

Proof. Under the premise of Proposition, let us fix $u \in U$, $v \in V$, so that $x := (v, Z) := [u; 1][u; 1]^T \in X$. Denoting, as above, by $P = P_{u,v}$ the distribution of $\omega := (\zeta, \zeta\zeta^T)$ with $\zeta \sim \mathcal{N}(A[u; 1], M(v))$, and invoking (3.67), we see that for just defined (x, P), relation (3.34) takes place. Applying Corollary 3.9, we conclude that

$$\operatorname{Prob}_{\zeta^{K} \sim [\mathcal{N}(A[u;1],M(v))]^{K}} \left\{ |\widehat{g}(\zeta^{K}) - G(x)| > \bar{\rho} + \varrho([u;1][u;1]^{T}) \right\} \leq \epsilon.$$

It remains to note that by construction for x = (v, Z) in question it holds

$$G(x) = q^{T}v + \operatorname{Tr}(QZ) = q^{T}v + \operatorname{Tr}(Q[u;1][u;1]^{T}) = q^{T}v + [u;1]^{T}Q[u,1] = F(u,v).$$

An immediate consequence of Proposition 3.13 is as follows:

Corollary 3.14. Under the premise and in the notation of Proposition 3.13, let $(h, H) \in \mathbf{R}^d \times \mathbf{S}^d$. Setting

$$\rho = \frac{1}{2} \left| \widehat{\Psi}_{+}(h, H) + \widehat{\Psi}_{-}(h, H) \right|,
\varkappa = \frac{1}{2} \left| \widehat{\Psi}_{-}(h, H) - \widehat{\Psi}_{+}(h, H) \right|,$$
(3.73)

the ϵ -risk of estimate (3.71) does not exceed ρ .

Indeed, with ρ and \varkappa given by (3.73), h, H, ρ, \varkappa satisfy (3.70).

3.4.1.4 Consistency

We are about to present a simple sufficient condition for the estimator suggested by Proposition 3.13 to be consistent, in the sense of Section 3.3.4.1. Specifically, in the situation and with the notation from Sections 3.4.1.1, 3.4.1.3 assume that

A.1.
$$\varrho(\cdot) \equiv 0$$
,

A.2. $V = \{\bar{v}\}$ is a singleton, which allows to set $\Theta_* = M(\bar{v})$, to satisfy (3.57) with $\delta = 0$, and to assume w.l.o.g. that

$$F(u,v) = [u;1]^T Q[u;1], \ G(Z) = \operatorname{Tr}(QZ);$$

A.3. the first m columns of the $d \times (m+1)$ matrix A are linearly independent.

By A.3, the columns of $(d+1) \times (m+1)$ matrix B, see (3.64), are linearly independent, so that we can find $(m+1) \times (d+1)$ matrix C such that $CB = I_{m+1}$. Let us define $(\bar{h}, \bar{H}) \in \mathbf{R}^d \times \mathbf{S}^d$ from the relation

$$\begin{bmatrix} \bar{H} & \bar{h} \\ \bar{h}^T & \end{bmatrix} = 2(C^T Q C)^o, \qquad (3.74)$$

where for $(d + 1) \times (d + 1)$ matrix S, S^o is the matrix obtained from S by zeroing our the entry in the cell (d + 1, d + 1).

The consistency of our estimation machinery is given by the following simple statement:

Proposition 3.15. In the just described situation and under assumptions A.1-3, given $\epsilon \in (0, 1)$, consider the estimate

$$\widehat{g}_{K,\epsilon}(\zeta^K) = \frac{1}{K} \sum_{k=1}^K [\overline{h}^T \zeta_k + \frac{1}{2} \zeta^T \overline{H} \zeta_k] + \varkappa_{K,\epsilon},$$

where

$$\varkappa_{K,\epsilon} = \frac{1}{2} \left[\widehat{\Psi}_{-}(\bar{h},\bar{H}) - \widehat{\Psi}_{+}(\bar{h},\bar{H}) \right]$$

and $\widehat{\Psi}_{\pm} = \widehat{\Psi}_{\pm}^{K,\epsilon}$ are given by (3.69). Then the ϵ -risk of $\widehat{g}_{K,\epsilon}(\cdot)$ goes to 0 as $K \to \infty$.

For proof, see Section 3.6.3.

3.4.1.5 A modification

In the situation described in the beginning of this Section, let a set $W \subset U \times V$ be given, and assume we are interested in recovering functional of interest (3.62) at points $(u, v) \in W$ only. When reducing the "domain of interest" $U \times V$ to W, we hopefully can reduce the achievable ϵ -risk of recovery. To utilize for this purpose the machinery we have developed, assume that we can point our a convex compact set $W \subset V \times Z$ such that

$$(u, v) \in W \Rightarrow (v, [u; 1][u; 1]^T) \in \mathcal{W}$$

A straightforward inspection justifies the following

Remark 3.16. In the just described situation, the conclusion of Proposition 3.13 remains valid when the set $U \times V$ participating in (3.72) and in relations (3.69) is reduced to \mathcal{W} . This modification enlarges the feasible set of (3.70) and thus reduces the achievable values of risk bound $\bar{\rho}$.

3.4.2 Estimating quadratic form, sub-Gaussian case

3.4.2.1 Situation

In the rest of this Section we are interested in the situation is as follows: we are given i.i.d. random observations

$$\zeta_i \sim SG(A[u;1], M(v)), \ i = 1, ..., K,$$
(3.75)

where $\zeta \sim SG(\mu, \Theta)$ means that ζ is sub-Gaussian with parameters $\mu \in \mathbf{R}^d$, $\Theta \in S^d_+$, and

- (u, v) is unknown "signal" known to belong to a given set $U \times V$, where
 - $U \subset \mathbf{R}^m$ is a compact set, and
 - $V \subset \mathbf{R}^k$ is a compact convex set;
- A is a given $d \times (m+1)$ matrix, and $v \mapsto M(v) : \mathbf{R}^k \to \mathbf{S}^{d+1}$ is affine mapping such that $M(v) \succeq 0$ whenever $v \in V$.

We are also given a convex calibrating function $\varrho(Z): \mathbf{S}^{m+1}_+ \to \mathbf{R}$ and "functional of interest"

$$F(u,v) = [u;1]^T Q[u;1] + q^T v, \qquad (3.76)$$

where Q and q are known $(m + 1) \times (m + 1)$ symmetric matrix and k-dimensional vector, respectively. Our goal is to recover F(u, v), for unknown (u, v) known to belong to $U \times V$, via observation (3.75).

Note that the only difference of our present situation with the one considered in Section 3.4.1.1 is that now we allow for sub-Gaussian, and not necessary Gaussian, observations.

3.4.2.2 Construction & Result

Let

$$\mathcal{V} = \{ M(v) : v \in V \}.$$

so that \mathcal{V} is a convex compact subset of the positive semidefinite cone \mathbf{S}^d_+ . Let us select somehow

- 1. a matrix $\Theta_* \succ 0$ such that $\Theta_* \succeq \Theta$, for all $\Theta \in \mathcal{V}$;
- 2. a convex compact subset \mathcal{Z} of the set $\mathcal{Z}^+ = \{ Z \in \mathbf{S}^{m+1}_+ : Z_{m+1,m+1} = 1 \}$ such that $[u; 1][u; 1]^T \in \mathcal{Z}$ for all $u \in U$;
- 3. reals $\gamma, \gamma^+ \in (0, 1)$ with $\gamma < \gamma^+$ (say, $\gamma = 0.99, \gamma^+ = 0.999$).

Preliminaries Given the data of the above description and $\delta \in [0, 2]$, we set (cf.

Proposition 3.12)

$$\begin{aligned} \mathcal{H} &= \mathcal{H}_{\gamma} := \{(h, H) \in \mathbf{R}^{d} \times \mathbf{S}^{d} : -\gamma \Theta_{*}^{-1} \preceq H \preceq \gamma \Theta_{*}^{-1} \}, \\ B &= \begin{bmatrix} A \\ [0, ..., 0, 1] \end{bmatrix} \in \mathbf{R}^{(d+1) \times (m+1)}, \\ \mathcal{M} &= \mathcal{V} \times \mathcal{Z}, \\ \Psi(h, H, G; Z) &= -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2} G \Theta_{*}^{1/2}) \\ &+ \frac{1}{2} \operatorname{Tr} \left(Z B^{T} \left[\left[\frac{H}{h^{T}} \right]^{+} \right] + [H, h]^{T} [\Theta_{*}^{-1} - G]^{-1} [H, h] \right] B \right) : \\ \left(\mathcal{H} \times \{G : 0 \preceq G \preceq \gamma^{+} \Theta_{*}^{-1}\} \right) \times \mathcal{Z} \to \mathbf{R}, \\ \Psi_{\delta}(h, H, G; \Theta, Z) &= -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_{*}^{1/2} G \Theta_{*}^{1/2}) + \frac{1}{2} \operatorname{Tr}([\Theta - \Theta_{*}]G) \\ &+ \frac{\delta(2+\delta)}{2(1-||\Theta_{*}^{1/2} G \Theta_{*}^{1/2}||)} ||\Theta_{*}^{1/2} G \Theta_{*}^{1/2}||_{H}^{2} \\ &+ \frac{1}{2} \operatorname{Tr} \left(Z B^{T} \left[\left[\frac{H}{h^{T}} \right]^{+} \right] + [H, h]^{T} [\Theta_{*}^{-1} - G]^{-1} [H, h] \right] B \right) : \\ \left(\mathcal{H} \times \{G : 0 \preceq G \preceq \gamma^{+} \Theta_{*}^{-1}\} \right) \times (\{0 \preceq \Theta \preceq \Theta_{*}\} \times \mathcal{Z}) \to \mathbf{R}, \\ \Phi(h, H; Z) &= \min_{G} \left\{ \Psi(h, H, G; Z) : 0 \preceq G \preceq \gamma^{+} \Theta_{*}^{-1}, G \succeq H \right\} : \\ \mathcal{H} \times (\{0 \prec \Theta \prec \Theta_{*}\} \times \mathcal{Z}) \to \mathbf{R}. \end{aligned}$$

The following statement is straightforward reformulation of Proposition 2.49.i:

Proposition 3.17. In the situation described in Section 3.4.2.1, we have

(i) Φ is well-defined real-valued continuous function on the domain $\mathcal{H} \times \mathcal{Z}$; the function is convex in $(h, H) \in \mathcal{H}$, concave in $Z \in \mathcal{Z}$, and $\Phi(0; Z) \geq 0$. Furthermore, let $(h, H) \in \mathcal{H}$, $u \in U$, $v \in V$, and let $\zeta \sim \mathcal{S}G(A[u; 1], M(v))$. Then

$$\ln\left(\mathbf{E}_{\zeta}\left\{\exp\{h^{T}\zeta+\frac{1}{2}\zeta^{T}H\zeta\}\right\}\right) \leq \Phi(h,H;[u;1][u;1]^{T}).$$
(3.78)

(ii) Let \mathcal{V} be a convex compact subset of \mathbf{S}^d_+ such that $M(v) \in \mathcal{V}$ for all $v \in V$, and let $\delta \in [0,2]$ be such that

$$\Theta \in \mathcal{V} \Rightarrow \{\Theta \preceq \Theta_*\} \& \{ \|\Theta^{1/2}\Theta_*^{-1/2} - I\| \le \delta \}.$$
(3.79)

Then $\Phi_{\delta}(h, H; \Theta, Z)$ is well-defined real-valued continuous function on the domain $\mathcal{H} \times (\mathcal{V} \times Z)$; the function is convex in $(h, H) \in \mathcal{H}$, concave in $(\Theta, Z) \in \mathcal{V} \times Z$, and $\Phi_{\delta}(0; \Theta, Z) \geq 0$. Furthermore, let $(h, H) \in \mathcal{H}$, $u \in U$, $v \in V$, and let $\zeta \sim SG(A[u; 1], M(v))$. Then

$$\ln\left(\mathbf{E}_{\zeta}\left\{\exp\{h^{T}\zeta+\frac{1}{2}\zeta^{T}H\zeta\}\right\}\right) \leq \Phi_{\delta}(h,H;M(v),[u;1][u;1]^{T}).$$
(3.80)

The estimate. We estimate the functional of interest similarly to the case of Gaussian observations. Specifically, let us pass from observations (3.75) to their quadratic lifts, so that our observations become

$$\omega_i = (\zeta_i, \zeta_i \zeta_i^T), \ 1 \le i \le K, \zeta_i \sim \mathcal{S}G(A[u; 1], M(v)) \text{ are i.i.d.}$$
(3.81)

As in the Gaussian case, we find ourselves in the situation considered in section 3.3.3 and can use the machinery developed there. Indeed, let us specify the data introduced in section 3.3.1 and participating in the constructions of section 3.3 as

follows:

• $\mathcal{H} = \{f = (h, H) \in \mathcal{H}\} \subset \mathcal{E}_H = \mathbf{R}^d \times \mathbf{S}^d$, with \mathcal{H} defined in (3.77), and the inner product on \mathcal{E}_H defined as

$$\langle (h,H), (h',H') \rangle = h^T h' + \frac{1}{2} \operatorname{Tr}(HH'),$$

- $\mathcal{E}_M = \mathbf{S}^{m+1}, \text{ and } \mathcal{M}, \Phi \text{ are as defined in (3.77);}$ $\mathcal{E}_X = \mathbf{R}^k \times \mathbf{S}^{d+1}, X := \{(v, Z) : v \in V, Z = [u; 1][u; 1]^T, u \in U\} \subset \mathcal{X} := \{(v, Z) : v \in V, Z = [u; 1][u; 1]^T, u \in U\} \subset \mathcal{X} := \{(v, Z) : v \in V, Z = [u; 1][u; 1]^T, u \in U\} \subset \mathcal{X} := \{(v, Z) : v \in V, Z = [u; 1][u; 1]^T, u \in U\} \subset \mathcal{X} := \{(v, Z) : v \in V, Z = [u; 1][u; 1]^T, u \in U\}$ $v \in V, Z \in \mathcal{Z} \};$
- $\mathcal{A}(x) = \mathcal{A}(v, Z) = (M(v), Z)$; note that \mathcal{A} is affine mapping from \mathcal{E}_X into \mathcal{E}_M mapping \mathcal{X} into \mathcal{M} , as required in section 3.3. Observe that when $u \in U$ and $v \in V$, the distribution P = P of observation ω_i defined by (3.81) satisfies the relation

$$\forall (f = (h, H) \in \mathcal{H}) : \ln \left(\mathbf{E}_{\omega \sim P} \left\{ \mathbf{e}^{\langle f, \omega \rangle} \right\} \right) = \ln \left(\mathbf{E}_{\zeta \sim \mathcal{S}G(A[u;1],M(v))} \left\{ \mathbf{e}^{h^T \zeta + \frac{1}{2} \zeta^T H \zeta} \right\} \right)$$

$$\leq \Phi(h, H; [u;1][u;1]^T),$$

$$(3.82)$$

see (3.78). Moreover, in the case of (3.79), we have also

$$\forall (f = (h, H) \in \mathcal{H}) : \ln \left(\mathbf{E}_{\omega \sim P} \left\{ e^{\langle f, \omega \rangle} \right\} \right) = \ln \left(\mathbf{E}_{\zeta \sim \mathcal{S}G(A[u;1],M(v))} \left\{ e^{h^T \zeta + \frac{1}{2} \zeta^T H \zeta} \right\} \right)$$

$$\leq \Phi_{\delta}(h, H; M(v), [u;1][u;1]^T),$$

$$(3.83)$$

see (3.80);

• we set $v(x = (v, Z)) = \varrho(Z)$,

1

• we define affine functional G(x) on \mathcal{E}_X by the relation

$$G(x := (v, Z)) = q^T v + \operatorname{Tr}(QZ),$$

see (3.76). As a result, for $x \in X$, that is, for $x = (v, [u; 1][u; 1]^T)$ with $v \in V$ and $u \in U$ we have

$$F(u, v) = G(x).$$
 (3.84)

The result. Applying to the just specified data Corollary 3.9 (which is legitimate, since our Φ clearly satisfies (3.39)), we arrive at the result as follows:

Proposition 3.18. In the situation described in Section 3.4.2.1, let us set

$$\begin{split} \widehat{\Psi}_{+}(h,H) &:= \inf_{\alpha} \left\{ \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ (v,Z) \in V \times \mathcal{Z}}} \left[\alpha \Phi(h/\alpha,H/\alpha;Z) - G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] : \\ & \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right] \right\} \\ &= \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1}}} \left[\alpha \Phi(h/\alpha,H/\alpha;Z) - G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] , \\ \widehat{\Psi}_{-}(h,H) &:= \inf_{\alpha} \left\{ \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ (v,Z) \in V \times \mathcal{Z}}} \left[\alpha \Phi(-h/\alpha,-H/\alpha;Z) + G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] : \\ & \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right] \right\} \\ &= \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ -\widehat{\gamma}_{d} \alpha \Theta_{*}^{-1} \preceq H \preceq \widehat{\gamma}_{d} \alpha \Theta_{*}^{-1}}} \left[\alpha \Phi(-h/\alpha,-H/\alpha;Z) + G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] . \end{split}$$
(3.85)

so that the functions $\widehat{\Psi}_{\pm}(h,H)$: $\mathbf{R}^d \times \mathbf{S}^d \to \mathbf{R}$ are convex. Furthermore, whenever $\overline{h}, \overline{H}, \overline{\rho}, \overline{\varkappa}$ form a feasible solution to the system of convex constraints

$$\widehat{\Psi}_{+}(h,H) \le \rho - \varkappa, \ \widehat{\Psi}_{-}(h,H) \le \rho + \varkappa \tag{3.86}$$

in variables $(h, H) \in \mathbf{R}^d \times \mathbf{S}^d$, $\rho \in \mathbf{R}$, $\varkappa \in \mathbf{R}$, setting

$$\widehat{g}(\zeta^K) = \frac{1}{K} \sum_{i=1}^K \left[h^T \zeta_i + \frac{1}{2} \zeta_i^T H \zeta_i \right] + \overline{\varkappa},$$

we get an estimate of the functional of interest $F(u, v) = [u; 1]^T Q[u; 1] + q^T v$ via *i.i.d.* observations

$$\zeta_i \sim \mathcal{S}G(A[u;1], M(v)), \ 1 \le i \le K,$$

with the following property:

$$\forall (u,v) \in U \times V : \operatorname{Prob}_{\zeta^{K} \sim [\mathcal{S}G(A[u;1],M(v))]^{K}} \left\{ |F(u,v) - \widehat{g}(\zeta^{K})| > \bar{\rho} + \varrho([u;1][u;1]^{T}) \right\} \leq \epsilon,$$

$$(3.87)$$

Proof. Under the premise of Proposition, let us fix $u \in U$, $v \in V$, so that $x := (v, Z) := [u; 1][u; 1]^T \in X$. Denoting by P the distribution of $\omega := (\zeta, \zeta\zeta^T)$ with $\zeta \sim SG(A[u; 1], M(v))$, and invoking (3.82), we see that for just defined (x, P), relation (3.34) takes place. Applying Corollary 3.9, we conclude that

$$\operatorname{Prob}_{\zeta^{K} \sim [\mathcal{N}(A[u;1],M(v))]^{K}} \left\{ |\widehat{g}(\zeta^{K}) - G(x)| > \bar{\rho} + \varrho([u;1][u;1]^{T}) \right\} \leq \epsilon.$$

It remains to note that by construction for x = (v, Z) in question it holds

$$G(x) = q^{T}v + \operatorname{Tr}(QZ) = q^{T}v + \operatorname{Tr}(Q[u;1][u;1]^{T}) = q^{T}v + [u;1]^{T}Q[u,1] = F(u,v).$$

Remark 3.19. In the situation described in section 3.4.2.1, let $\delta \in [0,2]$ be such that

$$\|\Theta^{1/2}\Theta_*^{-1/2} - I\| \le \delta \ \forall \Theta \in \mathcal{V}.$$

Then the conclusion of Proposition 3.18 remains valid when the function Φ in (3.85)
is replaced with the function Φ_{δ} , that is, when Ψ_{\pm} are defined as

$$\begin{split} \widehat{\Psi}_{+}(h,H) &:= \inf_{\alpha} \left\{ \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ (v,Z) \in V \times \mathcal{Z}}} \left[\alpha \Phi_{\delta}(h/\alpha,H/\alpha;M(v),Z) - G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] : \\ & \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right] \right\} \\ &= \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1}}} \inf_{\substack{\alpha > 0, \\ (v,Z) \in V \times \mathcal{Z} \\ (v,Z) \in V \times \mathcal{Z}}} \left[\alpha \Phi_{\delta}(-h/\alpha,-H/\alpha;M(v),Z) - G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] : \\ & \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right] \\ &= \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ -\widehat{\gamma}_{d} \alpha \Theta_{*}^{-1} \preceq H \preceq \widehat{\gamma}_{d} \alpha \Theta_{*}^{-1}}} \inf_{\substack{\alpha > 0, \\ -\widehat{\gamma}_{d} \alpha \Theta_{*}^{-1} \preceq H \preceq \widehat{\gamma}_{d} \alpha \Theta_{*}^{-1}}} \left[\alpha \Phi_{\delta}(-h/\alpha,-H/\alpha;M(v),Z) + G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] : \\ & \alpha > 0, -\gamma \alpha \Theta_{*}^{-1} \preceq H \preceq \gamma \alpha \Theta_{*}^{-1} \right] \\ &= \max_{\substack{(v,Z) \in V \times \mathcal{Z} \\ -\widehat{\gamma}_{d} \alpha \Theta_{*}^{-1} \preceq H \preceq \widehat{\gamma}_{d} \alpha \Theta_{*}^{-1}}} \left[\alpha \Phi_{\delta}(-h/\alpha,-H/\alpha;M(v),Z) + G(v,Z) - \varrho(Z) + \alpha K^{-1} \ln(2/\epsilon) \right] : \\ & (3.88) \\ & (3.$$

To justify Remark 3.19, it suffices to use in the proof of Proposition 3.18 relation (3.83) in the role of (3.82). Note that what is better in terms of the risk of the resulting estimate – Proposition 3.18 "as is" or its modification presented in Remark 3.19 – depends on the situation, so that it makes sense to keep in mind both options.

3.4.2.3 Numerical illustration, direct observations

The problem. Our initial illustration is deliberately selected to be extremely simple: given direct noisy observations

$$\zeta = u + \xi$$

of unknown signal $u \in \mathbf{R}^m$ known to belong to a given set U, we want to recover the "energy" $u^T u$ of u; what we are interested in, is the quadratic in ζ estimate with as small ϵ -risk on U as possible; here $\epsilon \in (0, 1)$ is a given design parameter. The details of our setup are as follows:

• U is the "spherical layer" $U = \{u \in \mathbf{R}^m : r^2 \leq u^T u \leq R^2\}$, where $r, R, 0 \leq r < R < \infty$ are given. As a result, the "main ingredient" in constructions from sections 3.4.1.3, 3.4.2.2 – the convex compact subset \mathcal{Z} of the set $\{Z \in \mathbf{S}^{m+1}_+ : Z_{m+1,m+1} = 1\}$ containing all matrices $[u; 1][u; 1]^T, u \in U$, can be specified as

$$\mathcal{Z} = \{ Z \in \mathbf{S}_{+}^{m+1} : Z_{m+1,m+1} = 1, 1 + r^{2} \le \operatorname{Tr}(Z) \le 1 + R^{2} \};$$

- ξ is either ~ N(0, Θ) (Gaussian case), or ~ SG(0, Θ) (sub-Gaussian case), with matrix Θ known to be diagonal with diagonal entries satisfying θσ² ≤ Θ_{ii} ≤ σ², 1 ≤ i ≤ d = m, with known θ ∈ [0, 1] and σ² > 0;
 the calibrating function ρ(Z) is ρ(Z) = ζ(∑_{i=1}^m Z_{ii}), where ζ is a convex contin-
- the calibrating function $\varrho(Z)$ is $\varrho(Z) = \varsigma(\sum_{i=1}^{m} Z_{ii})$, where ς is a convex continuous real-valued function on \mathbf{R}_+ . Note that with this selection, the claim that ϵ -risk of an estimate $\widehat{g}(\cdot)$ is $\leq \rho$ means that whenever $u \in U$, one has

$$\operatorname{Prob}\{|\widehat{g}(u+\xi) - u^T u| > \rho + \varsigma(u^T u)\} \le \epsilon.$$
(3.89)

Processing the problem. It is easily seen that in the situation in question the machinery of sections 3.4.1, 3.4.2 boils down to the following:

1. We lose nothing when restricting ourselves with estimates of the form

$$\widehat{g}(\zeta) = \frac{1}{2}\eta\zeta^T\zeta + \varkappa, \qquad (3.90)$$

with properly selected scalars η and \varkappa ;

 In Gaussian case, η and κ are yielded by the convex optimization problem with just 3 variables α₊, α₋, η, namely the problem

$$\begin{split} \min_{\substack{\alpha_{\pm},\eta}} \left\{ \widehat{\Psi}(\alpha_{+},\alpha_{-},\eta) &= \frac{1}{2} \left[\widehat{\Psi}_{+}(\alpha_{+},\eta) + \widehat{\Psi}_{-}(\alpha_{-},\eta) \right] : \sigma^{2}|\eta| < \alpha_{\pm} \right\}, \\ \widehat{\Psi}_{+}(\alpha_{+},\eta) &= -\frac{d\alpha_{+}}{2} \ln(1 - \sigma^{2}\eta/\alpha_{+}) + \frac{d}{2}\sigma^{2}(1-\theta) \max[-\eta,0] + \frac{d\delta(2+\delta)\sigma^{4}\eta^{2}}{2(\alpha_{+}-\sigma^{2}|\eta|)} \\ &+ \max_{r^{2} \leq t \leq R^{2}} \left[\left[\frac{\alpha_{+}\eta}{2(\alpha_{+}-\sigma^{2}\eta)} - 1 \right] t - \varsigma(t) \right] + \alpha_{+} \ln(2/\epsilon) \\ \widehat{\Psi}_{-}(\alpha_{+},\eta) &= -\frac{d\alpha_{-}}{2} \ln(1 + \sigma^{2}\eta/\alpha_{-}) + \frac{d}{2}\sigma^{2}(1-\theta) \max[\eta,0] + \frac{d\delta(2+\delta)\sigma^{4}\eta^{2}}{2(\alpha_{-}-\sigma^{2}|\eta|)} \\ &+ \max_{r^{2} \leq t \leq R^{2}} \left[\left[-\frac{\alpha_{-}\eta}{2(\alpha_{-}+\sigma^{2}\eta)} + 1 \right] t - \varsigma(t) \right] + \alpha_{-} \ln(2/\epsilon) \ , \end{split}$$

$$(3.91)$$

where $\delta = 1 - \sqrt{\theta}$. Specifically, the η -component of a feasible solution to (3.91) augmented by the quantity

$$\varkappa = \frac{1}{2} \left[\widehat{\Psi}_{-}(\alpha_{-},\eta) - \widehat{\Psi}_{+}(\alpha_{+},\eta) \right]$$

yields estimate (3.90) with ϵ -risk on U not exceeding $\widehat{\Psi}(\alpha_+, \alpha_-, \eta)$;

3. In sub-Gaussian case, η and \varkappa are yielded by convex optimization problem with just 5 variables, $\alpha_{\pm}, g_{\pm}, \eta$, namely, the problem

$$\min_{\alpha_{\pm},g_{\pm},\eta} \left\{ \widehat{\Psi}(\alpha_{\pm},g_{\pm},\eta) = \frac{1}{2} \left[\widehat{\Psi}_{+}(\alpha_{+},\lambda_{+},g_{+}\eta) + \widehat{\Psi}_{-}(\alpha_{-},\lambda_{-},g_{-},\eta) \right] : \\
0 \le \sigma^{2}g_{\pm} \le \gamma\alpha_{\pm}, \sigma^{2}\eta \ge -\alpha_{+}, \sigma^{2}\eta \le \alpha_{-}, \eta \le g_{+}, -\eta \le g_{-} \right\}, \\
\widehat{\Psi}_{+}(\alpha_{+},g_{+},\eta) = -\frac{d\alpha_{+}}{2}\ln(1-\sigma^{2}g_{+}) \\
+\alpha_{+}\ln(2/\epsilon) + \max_{r^{2} \le t \le R^{2}} \left[\left[\frac{\sigma^{2}\eta^{2}}{2(\alpha_{+}-g_{+})} + \frac{1}{2}\eta - 1 \right] r - \varsigma(t) \right] \\
\widehat{\Psi}_{-}(\alpha_{-},g_{-},\eta) = -\frac{d\alpha_{-}}{2}\ln(1-\sigma^{2}g_{-}) + \\
\alpha_{-}\ln(2/\epsilon) + \max_{r^{2} \le t \le R^{2}} \left[\left[\frac{\sigma^{2}\eta^{2}}{2(\alpha_{-}-g_{-})} - \frac{1}{2}\eta + 1 \right] r - \varsigma(t) \right] \\$$
(3.92)

where $\gamma \in (0, 1)$ is construction's parameter (we used $\gamma = 0.99$). Specifically, the η -component of a feasible solution to (3.92) augmented by the quantity

$$\varkappa = \frac{1}{2} \left[\widehat{\Psi}_{-}(\alpha_{-}, g_{-}, \eta) - \widehat{\Psi}_{+}(\alpha_{+}, g_{+}, \eta) \right]$$

yields estimate (3.90) with ϵ -risk on U not exceeding $\widehat{\Psi}(\alpha_{\pm}, g_{\pm}, \eta)$.

Note that the Gaussian case of our "energy estimation" problem is well studied in the literature, mainly in the case $\xi \sim \mathcal{N}(0, \sigma^2 I_m)$ of white Gaussian noise with exactly known variance σ^2 ; available results investigate analytically the interplay between the dimension m of signal, noise intensity σ^2 and the parameters R, r and

d	r	R	θ	Relative 0.01-risk, Gaussian case	Relative 0.01-risk, sub-Gaussian case	Optimality ratio
64	0	16	1	0.34808	0.44469	1.22
64	0	16	0.5	0.43313	0.44469	1.48
64	0	128	1	0.04962	0.05181	1.28
64	0	128	0.5	0.05064	0.05181	1.34
64	8	80	1	0.07827	0.08376	1.28
64	8	80	0.5	0.08095	0.08376	1.34
256	0	32	1	0.19503	0.30457	1.28
256	0	32	0.5	0.26813	0.30457	1.41
256	0	512	1	0.01264	0.01314	1.28
256	0	512	0.5	0.01289	0.01314	1.34
256	16	160	1	0.03996	0.04501	1.28
256	16	160	0.5	0.04255	0.04501	1.34
1024	0	64	1	0.10272	0.21923	1.28
1024	0	64	0.5	0.17032	0.21923	1.34
1024	0	2048	1	0.00317	0.00330	1.28
1024	0	2048	0.5	0.00324	0.00330	1.34
1024	32	320	1	0.02019	0.02516	1.28
1024	32	320	0.5	0.02273	0.02516	1.41

Table 3.3: Recovering signal's energy from direct observations

offer provably optimal, up to absolute constant factors, estimates. A nice property of the proposed approach is that (3.91) automatically takes care of the parameters and results in estimates with seemingly near-optimal performance, as is witnessed by the numerical results we are about to present.

Numerical results. In the first series of experiments, we used the trivial calibrating function: $\varsigma(\cdot) \equiv 0$.

A typical sample of numerical results is presented in Table 3.3. To avoid large numbers, we display in the table *relative* 0.01-risk achievable with our machinery, that is, the plain risk divided by R^2 ; keeping this in mind, one should not be surprised that when extending the range [r, R] of allow norms of the observed signal, all other components of the setup being fixed, the relative risk can decrease (the actual risk, of course, can only increase). Note that in all our experiments σ was set to 1.

Along with the values of the relative 0.01-risk, we present also the values of "optimality ratios" – the ratios of the relative risks achievable with our machinery in the Gaussian case to the (lower bounds on the) the best possible under circumstances relative 0.01-risks. These lower bounds are obtained as follows. Let us select somehow values $r_1 < r_2$ in the allowed under the circumstances range [r, R] of $||u||_2$, and two values, σ_1 , σ_2 , in the allowed range $[\theta\sigma,\sigma] = [\theta,1]$ of values of diagonal entries in diagonal matrices Θ , and consider two distributions of observations P_1 and P_2 as follows: P_{χ} is the distribution of random vector $x + \zeta$, where x and ξ are independent, x is uniformly distributed on the sphere $||x||_2 = r_{\chi}$ and $\zeta \sim \mathcal{N}(0, \sigma_{\chi}^2 I_d)$. It is immediately seen that whenever the two simple hypotheses $\omega \sim P_1$, and $\omega \sim P_2$, cannot be decided upon via a single observation by a test with total risk $\leq 2\epsilon$ (with the total risk of a test defined as the sum, over the two hypotheses in question, of probabilities for the test to reject the hypothesis when it is true), the quantity $\delta = \frac{r_2^2 - r_1^2}{2}$ is a lower bound on the optimal ϵ -risk, Risk^{*}_e, defined as the infimum, over all estimates recovering $||u||_2^2$ via single observation

 $\omega = u + \zeta$, of the ϵ -risk of the estimate, where the ϵ -risk is taken w.r.t. u running through the spherical layer $U = \{u : r^2 \leq u^T u \leq R^2\}$, and the covariance matrices Θ of Gaussian zero mean noise running through the set of scalar matrices with diagonal entries varying in $[\theta, 1]$. In other words, denoting by $p_{\chi}(\cdot)$ the density of P_{χ} , we have

$$0.02 < \int_{\mathbf{R}^d} \min[p_1(\omega), p_2(\omega)] d\omega \Rightarrow \operatorname{Risk}_{0.01}^* \ge \frac{r_2^2 - r_1^2}{2}$$

Now, the densities p_{χ} are spherically symmetric, whence, denoting by $q_{\chi}(\cdot)$ the univariate density of the energy $\omega^T \omega$ of observation $\omega \sim P_{\chi}$, we have

$$\int_{\mathbf{R}^d} \min[p_1(\omega), p_2(\omega)] d\omega = \int_0^\infty \min[q_1(s), q_2(s)] ds,$$

so that

$$0.02 < \int_0^\infty \min[q_1(s), q_2(s)] ds \Rightarrow \operatorname{Risk}_{0.01}^* \ge \frac{r_2^2 - r_1^2}{2}.$$
 (3.93)

Now, on a closest inspection, q_{χ} is the convolution of two univariate densities representable by explicit computation-friendly formulas, implying that given $r_1, r_2, \sigma_1, \sigma_2$, we can check numerically whether the premise in (3.93) indeed takes place, and whenever the latter is the case, the quantity $\frac{r_2^2 - r_1^2}{2}$ is a lower bound on Risk^{*}_{0.01}. In our experiments, we used a simple search strategy (not described here) aimed at crude maximizing this bound in $r_1, r_2, \sigma_1, \sigma_2$ and used the resulting lower bounds on Risk^{*}_{0.01} to compute the optimality ratios presented in the table⁵¹.

We believe that quite moderate values of the optimality ratios presented in the table (these results are typical for a much larger series of experiments we have conducted) witness quite good performance of our machinery.

Optimizing the relative risk. The "relative risk" displayed in Table 3.3 is the corresponding to the trivial calibrating function 0.01-risk in recovery $u^T u$ divided by the largest value R^2 of this risk allowed by the inclusion $u \in U$. When R is large, low relative risk can correspond to pretty high "actual" risk. For example, with $d := \dim u = 1024$, $\theta = 1$, and $U = \{u \in \mathbb{R}^d : ||u||_2 \leq 1.e6\}$, the 0.01-risk becomes as large as $\rho \approx 6.5e6$; for "relatively small" signals, like $u^T u \approx 10^4$, recovering $u^T u$ within accuracy ρ does not make much sense. In order to allow for "large" domains U, it makes sense to pass from the trivial calibrating function to a nontrivial one, like $\varsigma(t) = \alpha t$, with small positive α . With this calibrating function, (3.89) reads

$$\operatorname{Prob}\left\{|\widehat{g}(u+\xi) - u^T u| > \rho + \alpha u^T u\right\} \le \epsilon.$$

It turns out that (quite reasonable when U is large) "relative" characterization of risk results in much smaller values of ρ as compared to the case $\alpha = 0$ of "plain"

 $^{^{51}}$ The reader should not be surprised by "narrow numerical spectrum" of optimality ratios displayed in Table 3.3: our lower bounding scheme was restricted to identify actual optimality ratios among the candidate values 1.05^i , i = 1, 2, ...

241

ris	k.	Here	is	instructive	numerical	d	lata:

r	R	0.01-Risk, $\alpha = 0$	0.01-Risk, $\alpha = 0.01$	0.01-Risk, $\alpha = 0.1$		
0	1.e7	6.51e7/6.51e7	1.33e3/1.58e3	474/642		
1.e2	1.e7	6.51e7/6.51e7	1.33e3/1.58e3	-123/92.3		
1.1e3	1.e7	6.51e7/6.51e7	$-4.73e_{3}/-4.48e_{3}$	-1.14e5/-1.14e5		
$U = \{ u \in \mathbf{R}^{1024} : r < u _2 < R \}, \ \theta = 1/2$						

Left/Right: risks in Gaussian/sub-Gaussian cases

3.4.2.4 Numerical illustration, indirect observations

The problem. The estimation problem we are about to process numerically is as follows. Our observations are

$$\zeta = Au + \xi, \tag{3.94}$$

where

- A is a given $d \times m$ matrix, with m > d ("under-determined observations"),
- $u \in \mathbf{R}^m$ is a signal known to belong to a compact set U,
- $\xi \sim \mathcal{N}(0,\Theta)$ (Gaussian case) of $\xi \sim \mathcal{S}G(0,\Theta)$ (sub-Gaussian case) is the observation noise; Θ is positive semidefinite $d \times d$ matrix known to belong to a given convex compact set $\mathcal{V} \subset \mathbf{S}^d_+$.

Our goal is to recover the energy

$$F(u) = \frac{1}{m} \|u\|_2^2$$

of the signal from a single observation (3.94).

In our experiment, the data is specified as follows:

- 1. We think of $u \in \mathbf{R}^m$ as of discretization of a smooth function x(t) of continuous argument $t \in [0;1]$: $u_i = x(\frac{i}{m}), 1 \leq i \leq m$. We set $U = \{u : ||Su||_2 \leq 1\}$, where $u \mapsto Su$ is the finite-difference approximation of the mapping $x(\cdot) \mapsto (x(0), x'(0), x''(\cdot))$, so that U is a natural discrete-time analogy of the Sobolev-type ball $\{x : [x(0)]^2 + [x'(0)]^2 + \int_0^1 [x''_{\pi}(t)]^2 dt \leq 1\}$.
- 2. $d \times m$ matrix A is of the form UDV^T , where U and V are randomly selected $d \times d$ and $m \times m$ orthogonal matrices, and the d diagonal entries in diagonal $d \times m$ matrix D are of the form $\theta^{-\frac{i-1}{d-1}}$, $1 \le i \le d$; the "condition number" θ of A is design parameter.
- 3. The set \mathcal{V} of allowed matrices Θ is the set of all diagonal $d \times d$ matrices with diagonal entries varying from 0 to σ^2 , where the "noise intensity" σ is design parameter.

Processing the problem. Our estimating problem clearly is covered by the setups considered in sections 3.4.1 (Gaussian case) and 3.4.2 (sub-Gaussian case); in terms of these setups, it suffices to specify Θ_* as $\sigma^2 I_d$, M(v) as the identity mapping of \mathcal{V} onto itself, the mapping $u \mapsto A[u; 1]$ as the mapping $u \mapsto Pu$, and the set \mathcal{Z} (which should be a convex compact subset of the set $\{Z \in \mathbf{S}^{d+1}_+ : Z_{d+1,d+1} = 0\}$

containing all matrices of the form $[u; 1][u; 1]^T$, $u \in U$) as the set

$$\mathcal{Z} = \{ Z \in \mathbf{S}_{+}^{d+1} : Z_{d+1,d+1} = 1, \text{Tr} \left(Z \text{Diag} \{ S^T S, 0 \} \right) \le 1 \}.$$

As suggested by Propositions 3.13 (Gaussian case) and 3.18 (sub-Gaussian case), the linear in "lifted observation" $\omega = (\zeta, \zeta\zeta^T)$ estimates of $F(u) = \frac{1}{m} ||u||_2^2$ stem from the optimal solution (h_*, H_*) to the convex optimization problem

Opt =
$$\min_{h,H} \frac{1}{2} \left[\widehat{\Psi}_{+}(h,H) + \widehat{\Psi}_{-}(h,H) \right],$$
 (3.95)

with $\widehat{\Psi}_{\pm}(\cdot)$ given by (3.69) in the Gaussian, and by (3.85) in the sub-Gaussian cases, with the number K of observations in (3.69), (3.85) set to 1. The resulting estimate is

$$\zeta \mapsto h_*^T \zeta + \frac{1}{2} \zeta^T H_* \zeta + \varkappa, \ \varkappa = \frac{1}{2} \left[\widehat{\Psi}_-(h_*, H_*) - \widehat{\Psi}_+(h_*, H_*) \right]$$
(3.96)

and the ϵ -risk of the estimate is (upper-bounded by) Opt.

Problem (3.95) is a well-structured convex-concave saddle point problem and as such is beyond the "immediate scope" of the standard Convex Programming software toolbox primarily aimed at solving well-structured convex minimization problems. However, applying conic duality, one can easily eliminate in (3.69), (3.85) the inner maxima over v, Z and end up with reformulation which can be solved numerically by CVX [69], and this is how (3.95) was processed in our experiments.

Numerical results. In the experiments to be reported, we used the trivial calibrating function: $\varrho(\cdot) \equiv 0$.

Table 3.4 displays typical numerical results of our experiments. To give an impression of the performance of our approach, we present, along with the upper risk bounds for the estimates yielded by our machinery, simple lower bounds on ϵ -risk achievable under the circumstances. The origin of the lower bounds is as follows. Assume we are speaking about ϵ -risk and have at our disposal a signal $w \in U$, and let $t(w) = ||Aw||_2$, $\rho = 2\sigma \text{ErfInv}(\epsilon)$, where ErfInv is the inverse error-function:

 $\operatorname{Prob}_{\xi \sim \mathcal{N}(0,1)} \{ \xi > \operatorname{ErfInv}(\epsilon) \} = \epsilon.$

Setting $\theta(w) = \max[1 - \rho/t(w), 0]$, observe that $w' := \theta(w)w \in U$ and $||Aw - Aw'||_2 \leq \rho$, which, due to the origin of ρ , implies that there is no way to decide via observation $Au + \xi$, $\xi \sim \mathcal{N}(0, \sigma^2)$, with risk $< \epsilon$ on the two simple hypotheses u = w and u = w'. As an immediate consequence, the quantity $\phi(w) := \frac{1}{2}[||w||_2^2 - ||w'||_2^2] = ||w||_2^2[1 - \theta^2(w)]/2$ is a lower bound on the ϵ -risk, on U, of a whatever estimate of $||u||_2^2$. We can now try to maximize the resulting lower risk bound over U, thus arriving at the lower bound

LwBnd =
$$\max_{w \in U} \left\{ \frac{1}{2} \|w\|_2^2 (1 - \theta^2(w)) \right\}.$$

On a closest inspection, the latter problem is not a convex one, which does not prevent building a suboptimal solution to this problem, and this is how the lower risk bounds in Table 3.4 were built (we omit the details). We see that the ϵ -risks

d, m	Opt, Gaussian case	Opt, sub-Gaussian case	LwBnd
8,12	0.1362(+65%)	0.1382(+67%)	0.0825
16, 24	0.1614(+53%)	0.1640(+55%)	0.1058
32,48	0.0687(+46%)	0.0692(+48%)	0.0469

Table 3.4: Upper bound (Opt) on the 0.01-risk of estimate (3.96), (3.95) vs. lower bound (LwBnd) on 0.01-risk achievable under the circumstances. In the experiments, $\sigma = 0.025$ and $\theta = 10$. Data in parentheses: excess of Opt over LwBnd.



Figure 3.3: Histograms of recovery errors in experiments, data over 1000 simulations per experiment.

of our estimates are within a moderate factor from the optimal ones.

Figure 3.3 shows empirical error distributions of the estimates built in the three experiments reported in Table 3.4. When simulating the observations and estimates, we used $\mathcal{N}(0, \sigma^2 I_d)$ obse4rvation noise and selected signals in U by maximizing over U randomly selected linear forms. Finally, we note that with our design parameters d, m, θ, σ fixed, we still deal with a family of estimation problems rather than with a single problem, the reason being that our U is ellipsoid with essentially different from each other half-axes, and achievable risks heavily depend on how the right singular vectors of A are oriented with respect to the directions of the half-axes of U, so that the risks of our estimates vary significantly from instance to instance even when the design parameters are fixed. Note also that the "sub-Gaussian experiments" were conducted on exactly the same data as "Gaussian experiments" of the same sizes d, m.

244

LECTURE 3

3.5 EXERCISES FOR LECTURE 3

[†] marks more difficult exercises.

Exercise 3.20. . The goal of what follows is to refine the change detection procedure (let us refer to it as to "basic") developed in Section 2.9.6.1. The idea is pretty simple. With the notation from Section 2.9.6.1, in basic procedure, when testing the null hypothesis H_0 vs. signal hypothesis H_t^{ρ} , we looked at the difference $\zeta_t = \omega_t - \omega_1$ and were trying to decide whether the energy of the deterministic component $x_t - x_1$ of ζ_t is 0, as is the case under H_0 , or is $\geq \rho^2$, as is the case under H_t^{ρ} . Note that if $\sigma \in [\underline{\sigma}, \overline{\sigma}]$ is the actual intensity of the observation noise, then the noise component of ζ_t is $\mathcal{N}(0, 2\sigma^2 I_d)$; other things being equal, the large is the noise in ζ_t , the larger should be ρ to allow for a reliable, with a given reliability level, decision of this sort. Now note that under the hypothesis H_t^{ρ} , we have $x_1 = \ldots = x_{t-1}$, so that the deterministic component of the difference $\zeta_t = \omega_t - \omega_1$ is exactly the same as for the difference $\tilde{\zeta}_t = \omega_t - \frac{1}{t-1} \sum_{s=1}^{t-1} \omega_s$, while the noise component in $\tilde{\zeta}_t$ is $\mathcal{N}(0, \sigma_t^2 I_d)$ with $\sigma_t^2 = \sigma^2 + \frac{1}{t-1}\sigma^2 = \frac{t}{t-1}\sigma^2$; thus, the intensity of noise in $\tilde{\zeta}_t$ is at most the one in ζ_t , and this intensity, in contrast to the one for ζ_t , decreases as t grows. Now goes the exercise:

Let reliability tolerances $\epsilon, \epsilon \in (0, 1)$ be given, and let our goal be to design a system of inferences \mathcal{T}_t , t = 2, 3, ..., K, which, when used in the same fashion as tests \mathcal{T}_t^{κ} were used in Basic procedure, results in false alarm probability at most ϵ and in probability to miss a change of energy $\geq \rho^2$ at most ϵ ; needless to say, we want to achieve this goal with as small ρ as possible. Think how to utilize the above observation to refine Basic procedure by hopefully reducing (and provably not increasing) the required value of ρ . Implement the Basic and the refined change detection procedures and compare their quality (the resulting values of ρ) on, say, the data used in the experiment reported in Section 2.9.6.1.

Exercise 3.21. In the situation of Section 3.3.4, design of a "good" estimate is reduced to solving convex optimization problem (3.50). Note that the objective in this problem is, in a sense, "implicit" – the design variable is f, and the objective is obtained from an explicit convex-concave function of f and (x, y) by maximization over (x, y). While there exist solvers capable to process problems of this type efficiently; however, commonly used of-the-shelf solvers, like cvx, cannot handle problems like (3.50). The goal of the exercise to follow is to reformulate (3.50) as a semidefinite program, thus making it amenable for cvx.

On an immediate inspection, the situation we are interested in is as follows. We are given

- a nonempty convex compact set $X \subset \mathbf{R}^n$ along with affine function M(x) taking values in \mathbf{S}^d and such that $M(x) \succeq 0$ when $x \in X$, and
- affine function $F(f) : \mathbf{R}^d \to \mathbf{R}^n$.

Given $\gamma > 0$, this data gives rise to the convex function

$$\Psi(f) = \max_{x \in X} \left\{ F^T(f)x + \gamma \sqrt{f^T M(x) f} \right\},\,$$

and we want to find a "nice" representation of this function, specifically, want to represent the inequality $\tau \geq \Psi(f)$ by a bunch of LMI's in variables τ , f, and perhaps additional variables.

To achieve our goal, we assume in the sequel that the set

$$X^+ = \{(x, M) : x \in X, M = M(x)\}$$

can be described by a system of linear and semidefinite constraints in variables x, M and additional variables, specifically, the system

(a)
$$s_i - a_i^T x - b_i^T \xi - \operatorname{Tr}(C_i M) \ge 0, i \le I$$

(b) $S - \mathcal{A}(x) - \mathcal{B}(\xi) - \mathcal{C}(M) \succeq 0$
(c) $M \succeq 0$

where $\mathcal{A}(\cdot), \mathcal{B}(\cdot), \mathcal{C}(\cdot)$ are affine functions taking values in \mathbf{S}^N . We assume that this system of constraints is essentially strictly feasible, meaning that there exists a feasible solution at which the semidefinite constraints (b), (c) are satisfied strictly (i.e., the left hand sides of the LMI's are positive definite).

Now goes the exercise:

1. Check that $\Psi(f)$ is the optimal value in a semidefinite program, specifically,

$$\Psi(f) = \max_{x,M,\xi,t} \left\{ F^{T}(f)x + \gamma t : \left\{ \begin{array}{cc} s_{i} - a_{i}^{T}x - b_{i}^{T}\xi - \operatorname{Tr}(C_{i}M) \ge 0, i \le I & (a) \\ S - \mathcal{A}(x) - \mathcal{B}(\xi) - \mathcal{C}(M) \succeq 0 & (b) \\ M \succeq 0 & (c) \\ \left[\frac{f^{T}Mf \mid t}{t \mid 1} \right] \ge 0 & (d) \end{array} \right\}.$$
(P)

2. Passing from (P) to the semidefinite dual of (P), build explicit semidefinite representation of Ψ , that is, an explicit system S of LMI's in variables f, τ and additional variables u such that

$$\{\tau \ge \Psi(f)\} \Leftrightarrow \{\exists u : (\tau, f, u) \text{ satisfies } \mathcal{S}\}.$$

Exercise 3.22. [†] Consider the situation as follows: given an $m \times n$ "sensing matrix" *A* which is stochastic– with columns from the probabilistic simplex $\Delta_m = \{v \in \mathbf{R}^m : v \ge 0, \sum_i v_i = 1\}$ and a nonempty closed subset *U* of Δ_n , we observe *M*-element, M > 1, i.i.d. sample $\zeta^M = (\zeta_1, ..., \zeta_M)$ with ζ_k drawn from the discrete distribution Au_* , where u_* is an unknown probabilistic vector ("signal") known to belong to *U*. We treat the discrete distribution $Au, u \in \Delta_n$, as a distribution on the vertices $e_1, ..., e_m$ of Δ_m , so that possible values of ζ_k are basic orths $e_1, ..., e_m$ in \mathbf{R}^m . Our goal is to recover the value at u_* of a given quadratic form

$$F(u) = u^T Q u + 2q^T u.$$

Observe that for $u \in \Delta_n$, we have $u = [uu^T]\mathbf{1}_n$, where $\mathbf{1}_k$ is the all-ones vector in \mathbf{R}^k . This observation allows to rewrite F(u) as a homogeneous quadratic form:

$$F(u) = u^T \bar{Q}u, \ \bar{Q} = Q + [q \mathbf{1}_n^T + \mathbf{1}_n q^T].$$
(3.97)

The goal of Exercise is to follow the approach developed in Section 3.4.1 for the Gaussian case in order to build an estimate $\hat{g}(\zeta^M)$ of F(u), specifically, estimate as follows.

Let

$$\mathcal{J}_M = \{(i, j) : 1 \le i < j \le M\}, J_M = \operatorname{Card}(\mathcal{J}_M).$$

For $\zeta^M = (\zeta_1, ..., \zeta_M)$ with $\zeta_k \in \{e_1, ..., e_m\}, 1 \le k \le M$, let $\omega_{ij}[\zeta^M] = \frac{1}{2}[\zeta_i \zeta_j^T + \zeta_j \zeta_i^T], (i, j) \in \mathcal{J}_M.$

The estimates we are interested in are of the form

$$\widehat{g}(\zeta^{M}) = \operatorname{Tr}\left(h\underbrace{\left[\frac{1}{J_{M}}\sum_{(i,j)\in\mathcal{J}_{M}}\omega_{ij}[\zeta^{M}]\right]}_{\omega[\zeta^{M}]}\right) + \kappa$$
(3.98)

where $h \in \mathbf{S}^m$ and $\kappa \in \mathbf{R}$ are the parameters of the estimate.

Now goes the exercise:

1. Verify that when ζ_k 's stem from signal $u \in U$, the expectation of $\omega[\zeta^M]$ is a linear image $Az[u]A^T$ of the matrix $z[u] = uu^T \in \mathbf{S}^n$: denoting by P_u^M the distribution of ζ^M , we have

$$\mathbf{E}_{\zeta^M \sim P_u^M} \{ \omega[\zeta^M] \} = A z[u] A^T.$$
(3.99)

Check that when setting

$$\mathcal{Z}_k = \{ \omega \in \mathbf{S}^k : \omega \succeq 0, \omega \ge 0, \mathbf{1}_k^T \omega \mathbf{1}_k = 1 \},\$$

where $x \ge 0$ for a matrix x means that x is entrywise nonnegative, the image of \mathcal{Z}_n under the mapping $z \mapsto AzA^T$ is contained in \mathcal{Z}_m .

2. Let $\Delta^k = \{z \in \mathbf{S}^k : z \ge 0, \mathbf{1}_n^T z \mathbf{1}_n = 1\}$, so that \mathcal{Z}_k is the set of all positive semidefinite matrices form Δ^k . For $\mu \in \Delta^m$, let P_{μ} be the distribution of the random matrix w taking values in \mathbf{S}^m , namely, as follows: the possible values of w are matrices of the form $e^{ij} = \frac{1}{2}[e_i e_j^T + e_j e_i^T], 1 \le i \le j \le m$; for every $i \le m$, w takes value e^{ii} with probability μ_{ii} , and for every i, j with i < j, w takes value e^{ij} with probability $2\mu_{ij}$. Further, let us set

$$\Phi_1(h;\mu) = \ln\left(\sum_{i,j=1}^m \mu_{ij} \exp\{h_{ij}\}\right) : \mathbf{S}^m \times \Delta^m \to \mathbf{R},$$
(3.100)

so that Φ_1 is a continuous convex-concave function on $\mathbf{S}^m \times \Delta^m$.

2.1. Prove that

$$\forall (h \in \mathbf{S}^m, \mu \in \mathcal{Z}_m) : \ln\left(\mathbf{E}_{w \sim P_{\mu}}\left\{\exp\{\operatorname{Tr}(hw)\}\right\}\right) = \Phi_1(h;\mu).$$
(3.101)

2.2. Derive from 2.1 that setting

$$K = K(M) = \lfloor M/2 \rfloor, \ \Phi_M(h;\mu) = K\Phi_1(h/K;\mu) : \mathbf{S}^m \times \Delta^m \to \mathbf{R}$$

 Φ_M is a continuous convex-concave function on $\mathbf{S}^m \times \Delta^m$ such $\Phi_K(0; \mu) = 0$ for all $\mu \in \mathcal{Z}_m$, and whenever $u \in U$, the following holds true:

Let P_u^M be the distribution of $\zeta^M = (\zeta_1, ..., \zeta_M)$ with independent blocks $\zeta_k \sim Au$, and let $P_{u,M}$ is the distribution of $\omega = \omega[\zeta^M], \, \zeta^M \sim P_u^M$. Then

$$\forall (u \in U, h \in \mathbf{S}^m) : \\ \ln\left(\mathbf{E}_{\omega \sim P_{u,M}} \left\{ \exp\{\operatorname{Tr}(h\omega)\} \right\} \right) \leq \Phi_M(h; Az[u]A^T), \ z[u] = uu^T.$$
 (3.102)

3. Combine the above observations with Corollary 3.7 to arrive at the following result:

Proposition 3.23. In the situation in question, let \mathcal{Z} be a convex compact subset of \mathcal{Z}_n such that $uu^T \in \mathcal{Z}$ for all $u \in U$. Given $\epsilon \in (0, 1)$, let

$$\begin{split} \Psi_{+}(h,\alpha) &= \max_{z\in\mathcal{Z}} \left[\alpha \Phi_{M}(h/\alpha,AzA^{T}) - \operatorname{Tr}(\bar{Q}z) \right] : \mathbf{S}^{m} \times \{\alpha > 0\} \to \mathbf{R}, \\ \Psi_{-}(h,\alpha) &= \max_{z\in\mathcal{Z}} \left[\alpha \Phi_{M}(-h/\alpha,AzA^{T}) + \operatorname{Tr}(\bar{Q}z) \right] : \mathbf{S}^{m} \times \{\alpha > 0\} \to \mathbf{R} \\ \widehat{\Psi}_{+}(h) &:= \inf_{\alpha > 0} \left\{ \Psi_{+}(h,\alpha) + \alpha \ln(2/\epsilon) \right\} \\ &= \max_{z\in\mathcal{Z}} \inf_{\alpha > 0} \left[\alpha \Phi_{M}(h/\alpha,AzA^{T}) - \operatorname{Tr}(\bar{Q}z) + \alpha \ln(2/\epsilon) \right] \\ &= \max_{z\in\mathcal{Z}} \inf_{\beta > 0} \left[\beta \Phi_{1}(h/\beta,AzA^{T}) - \operatorname{Tr}(\bar{Q}z) + \frac{\beta}{K} \ln(2/\epsilon) \right] \quad [\beta = K\alpha], \\ \widehat{\Psi}_{-}(h) &:= \inf_{\alpha > 0} \left\{ \Psi_{-}(h,\alpha) + \alpha \ln(2/\epsilon) \right\} \\ &= \max_{z\in\mathcal{Z}} \inf_{\alpha > 0} \left[\alpha \Phi_{M}(-h/\alpha,AzA^{T}) + \operatorname{Tr}(\bar{Q}z) + \alpha \ln(2/\epsilon) \right] \\ &= \max_{z\in\mathcal{Z}} \inf_{\beta > 0} \left[\beta \Phi_{1}(-h/\beta,AzA^{T}) + \operatorname{Tr}(\bar{Q}z) + \frac{\beta}{K} \ln(2/\epsilon) \right] \quad [\beta = K\alpha]. \end{split}$$

$$(3.103)$$

The functions $\widehat{\Psi}_{\pm}$ are real valued and convex on \mathbf{S}^m , and every candidate solution h to the convex optimization problem

$$Opt = \min_{h} \left\{ \widehat{\Psi}(h) := \frac{1}{2} \left[\widehat{\Psi}_{+}(h) + \widehat{\Psi}_{-}(h) \right] \right\}, \qquad (3.104)$$

induces the estimate

$$\widehat{g}_h(\zeta^M) = \operatorname{Tr}(h\omega[\zeta^M]) + \kappa(h), \ \kappa(h) = \frac{\widehat{\Psi}_-(h) - \widehat{\Psi}_+(h)}{2}$$

of the functional of interest (3.97) via observation ζ^M with ϵ -risk on U not exceeding $\rho = \widehat{\Psi}(h)$:

$$\forall (u \in U) : \operatorname{Prob}_{\zeta^M \sim P_u^M} \{ |F(u) - \widehat{g}_h(\zeta^M)| > \rho \} \leq \epsilon.$$

4. Consider an alternative way to estimate F(u), namely, as follows. Let $u \in U$. Given a pair of independent observations ζ_1, ζ_2 drawn from distribution Au, let us convert them into the symmetric matrix $\omega_{1,2}[\zeta^2] = \frac{1}{2}[\zeta_1\zeta_2^T + \zeta_2\zeta_1^T]$. The distribution $P_{u,2}$ of this matrix is exactly the distribution $P_{\mu(z[u])}$, see item B, where $\mu(z) = AzA^T : \Delta^n \to \Delta^m$. Now, given M = 2K observations $\zeta^{2K} = (\zeta_1, ..., \zeta_{2K})$ stemming from signal u, we can split them into K consecutive pairs giving rise to K observations $\omega^K = (\omega_1, ..., \omega_K), \ \omega_k = \omega[[\zeta_{2k-1}; \zeta_{2k}]]$ drawn independently of each other from probability distribution $P_{\mu(z[u])}$, and the functional of interest (3.97) is a linear function $\operatorname{Tr}(\bar{Q}z[u])$ of z[u]. Assume that we are given a set \mathcal{Z} as in the premise of Proposition 3.23. Observe that we are in the situation as follows:

Given K independent identically distributed observations $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ with $\omega_{k} \sim P_{\mu(z)}$, where z is unknown signal known to belong to \mathcal{Z} , we

want to recover the value at z of linear function $G(v) = \text{Tr}(\bar{Q}v)$ of $v \in \mathbf{S}^n$. Besides this, we know that P_{μ} , for every $\mu \in \Delta^m$, satisfies the relation

$$\forall (h \in \mathbf{S}^m) : \ln\left(\mathbf{E}_{\omega \sim P_u} \{ \exp\{\operatorname{Tr}(h\omega)\} \} \right) \leq \Phi_1(h;\mu).$$

This situation is the one of Section 3.3.3, with the data specified as

$$\mathcal{H} = \mathcal{E}_H = \mathbf{S}^m, \mathcal{M} = \Delta^m \subset \mathcal{E}_M = \mathbf{S}^m, \Phi = \Phi_1, X := \{z[u] : u \in U\} \subset \mathcal{X} := \mathcal{Z} \subset \mathcal{E}_X = \mathbf{S}^n, \mathcal{A}(z) = AzA^T,$$

and we can use the machinery developed in this Section on order to upper-bound $\epsilon\text{-risk}$ of affine estimate

$$\operatorname{Tr}\left(h\frac{1}{K}\sum_{k=1}^{K}\omega_{k}\right)+\kappa$$

of G(z[u]) and to build the best, in terms of the upper risk bound, estimate, see Corollary 3.9. On a closest inspection (carry it out!), the associated with the above data functions $\widehat{\Psi}_{\pm}$ arising in (3.49) are exactly the functions $\widehat{\Psi}_{\pm}$ specified in Proposition 3.23 for M = 2K. Thus, the just outlined approach to estimating F(u) via stemming from $u \in U$ observations ζ^{2K} results in a family of estimates

$$\widetilde{g}_h(\zeta^{2K}) = \operatorname{Tr}\left(h\frac{1}{K}\sum_{k=1}^K \omega[[\zeta_{2k-1};\zeta_{2k}]]\right) + \kappa(h), \ h \in \mathbf{S}^m$$

and the upper bound on ϵ -risk of estimate \tilde{g}_h is $\Psi(h)$, where $\Psi(\cdot)$ is associated with M = 2K according to Proposition 3.23, that is, is exactly the same as the offered by Proposition upper bound on the ϵ -risk of the estimate \hat{g}_h . Note, however, that the estimates \tilde{g}_h and \hat{g}_h are not identical:

$$\widetilde{g}_h(\zeta^{2K}) = \operatorname{Tr} \left(h \frac{1}{K} \sum_{k=1}^K \omega_{2k-1,2k}[\zeta^{2K}] \right) + \kappa(h), \widehat{g}_h(\zeta^{2K}) = \operatorname{Tr} \left(h \frac{1}{K(2K-1)} \sum_{1 \le i < j \le 2K} \omega_{ij}[\zeta^{2K}] \right) + \kappa(h).$$

Now goes the question:

• Which one of the estimates \tilde{g}_h , \hat{g}_h would you prefer, that is, which one of these estimates, in your opinion, exhibits better practical performance? To check your intuition, test performances of the estimates by simulation. Here

is the story underlying the recommended simulation model:

"Tomorrow, tomorrow not today, all the lazy people say." Does it make sense to be lazy? Imagine you are supposed to do some job, and should decide whether to do it today, or tomorrow. The reward for the job is drawn by nature at random, with unknown to you time-invariant distribution u on n-element set $\{r_1, ..., r_n\}$, with $r_1 \leq r_2 \leq ... \leq r_n$. Given 2K historical observations of the rewards, what is better – to do the job today or tomorrow, that is, is the probability of tomorrow reward to be at least the today one greater than 0.5? What is this probability? How to estimate it from historical data?

Pose the above problem as the one of estimating a quadratic functional $u^T \bar{Q} u$ of distribution u from direct observations $(m = n, A = I_n)$. Pick $u \in \Delta_n$ at random and run simulations to check which one of the estimates \hat{g}_h , \tilde{g}_h works better. To

avoid the necessity to solve optimization problem (3.104), you can use $h = \bar{Q}$, resulting in unbiased estimate of $u^T \bar{Q} u$.

Exercise 3.24. [†] What follows is a variation of Exercise 3.22. Consider the situation as follows: We observe K realizations η_k , $k \leq K$, of discrete random variable with ppossible values, and $L \geq K$ realizations ζ_ℓ , $\ell \leq L$, of discrete random variable with q possible values. All realizations are independent of each other; η_k 's are drawn from distribution Pu, and ζ_ℓ – from distribution Qv, where $P \in \mathbf{R}^{p \times r}$, $Q \in \mathbf{R}^{q \times s}$ are given stochastic "sensing matrices," and u, v are unknown "signals" known to belong to given subsets U, resp., V of probabilistic simplexes Δ_r , resp., Δ_s . Our goal is to recover from observations $\{\eta_k, \zeta_\ell\}$ the value at u, v of a given linear function

$$F(u, v) = u^T F v = \text{Tr}(F[uv^T]^T).$$
 (3.105)

The "covering story" could be as follows. Imagine that there are two possible actions, say, administering to a patient drug A and drug B. Let u is the probability distribution of a (somehow quantified) outcome of the first action, and v be similar distribution for the second action. Observing what happens when the first action is utilized K, and the second – L times, we could ask ourselves what is the probability of an outcome of the first action to be better than an outcome of the second action. This amounts to computing the probability p of the event " $\eta > \zeta$," where η , ζ are independent of each other discrete real-valued random variables with distributions u, resp., v, and p is a linear function of the "joint distribution" uv^T of η , ζ . This story gives rise to the aforementioned estimation problem with the unit sensing matrices P, Q. Assuming that there are "measurement errors" – instead of observing action's outcome "as is," we observe a realization of random variable with distribution depending, in a prescribed fashion, on the outcome.

As always, we encode the p possible values of η_k by the basic orths $e_1, ..., e_p$ in \mathbf{R}^p , and the q possible values of ζ – by the basic orths $f_1, ..., f_q$ in \mathbf{R}^q .

We intend to focus on estimates of the form

$$\widehat{g}_{h,\kappa}(\eta^{K},\zeta^{L}) = \left[\frac{1}{K}\sum_{k}\eta_{k}\right]^{T}h\left[\frac{1}{L}\sum_{\ell}\zeta_{\ell}\right] + \kappa \qquad [h \in \mathbf{R}^{p \times q}, \kappa \in \mathbf{R}]$$

This is what you are supposed to do:

1. (cf. item B in Exercise 3.22) Denoting by Δ_{mn} the set of nonnegative $m \times n$ matrices with unit sum of all entries (i.e., the set of all probability distributions on $\{1, ..., m\} \times \{1, ..., n\}$) and assuming $L \geq K$, let us set

$$\mathcal{A}(z) = P z Q^T : \mathbf{R}^{r \times s} \to \mathbf{R}^{p \times q}$$

and

$$\Phi(h;\mu) = \ln\left(\sum_{i=1}^{p} \sum_{j=1}^{q} \mu_{ij} \exp\{h_{ij}\}\right) : \mathbf{R}^{p \times q} \times \Delta_{pq} \to \mathbf{R},
\Phi_K(h;\mu) = K\Phi(h/K;\mu) : \mathbf{R}^{p \times q} \times \Delta_{pq} \to \mathbf{R}.$$

Verify that \mathcal{A} maps Δ_{rs} into Δ_{pq} , Φ and Φ_K are continuous convex-concave functions on their domains, and that for every $u \in \Delta_r$, $v \in \Delta_s$, the following

holds true:

(!) When $\eta^K = (\eta_1, ..., \eta_K, \zeta^L = (\zeta_1, ..., \zeta_K)$ with mutually independent $\eta_1, ..., \zeta_L$ such that $\eta_k \sim Pu$, $\eta_\ell \sim Qv$ for all k, ℓ , we have

$$\ln\left(\mathbf{E}_{\eta,\zeta}\left\{\exp\left\{\left[\frac{1}{K}\sum_{k}\eta_{k}\right]^{T}h\left[\frac{1}{L}\sum_{\ell}\zeta_{\ell}\right]\right\}\right\}\right) \leq \Phi_{K}(h;\mathcal{A}(uv^{T})).$$
(3.106)

2. Combine (!) with Corollary 3.7 to arrive at the following analogy of Proposition 3.23:

Proposition 3.25. In the situation in question, let \mathcal{Z} be a convex compact subset of Δ_{rs} such that $uv^T \in \mathcal{Z}$ for all $u \in U$, $v \in V$. Given $\epsilon \in (0, 1)$, let

$$\begin{split} \Psi_{+}(h,\alpha) &= \max_{z\in\mathcal{Z}} \left[\alpha \Phi_{K}(h/\alpha,PzQ^{T}) - \operatorname{Tr}(Fz^{T}) \right] : \mathbf{R}^{p\times q} \times \{\alpha > 0\} \to \mathbf{R}, \\ \Psi_{-}(h,\alpha) &= \max_{z\in\mathcal{Z}} \left[\alpha \Phi_{K}(-h/\alpha,PzQ^{T}) + \operatorname{Tr}(Fz^{T}) \right] : \mathbf{R}^{p\times q} \times \{\alpha > 0\} \to \mathbf{R} \\ \widehat{\Psi}_{+}(h) &:= \inf_{\alpha > 0} \left\{ \Psi_{+}(h,\alpha) + \alpha \ln(2/\epsilon) \right\} \\ &= \max_{z\in\mathcal{Z}} \inf_{\alpha > 0} \left[\alpha \Phi_{K}(h/\alpha,PzQ^{T}) - \operatorname{Tr}(Fz^{T}) + \alpha \ln(2/\epsilon) \right] \\ &= \max_{z\in\mathcal{Z}} \inf_{\beta > 0} \left[\beta \Phi(h/\beta,PzQ^{T}) - \operatorname{Tr}(Fz^{T}) + \frac{\beta}{K} \ln(2/\epsilon) \right] \quad [\beta = K\alpha], \\ \widehat{\Psi}_{-}(h) &:= \inf_{\alpha > 0} \left\{ \Psi_{-}(h,\alpha) + \alpha \ln(2/\epsilon) \right\} \\ &= \max_{z\in\mathcal{Z}} \inf_{\alpha > 0} \left[\alpha \Phi_{K}(-h/\alpha,PzQ^{T}) + \operatorname{Tr}(Fz^{T}) + \alpha \ln(2/\epsilon) \right] \\ &= \max_{z\in\mathcal{Z}} \inf_{\beta > 0} \left[\beta \Phi(-h/\beta,PzQ^{T}) + \operatorname{Tr}(Fz^{T}) + \frac{\beta}{K} \ln(2/\epsilon) \right] \quad [\beta = K\alpha]. \end{split}$$

$$(3.107)$$

The functions $\widehat{\Psi}_{\pm}$ are real valued and convex on $\mathbf{R}^{p \times q}$, and every candidate solution h to the convex optimization problem

$$Opt = \min_{h} \left\{ \widehat{\Psi}(h) := \frac{1}{2} \left[\widehat{\Psi}_{+}(h) + \widehat{\Psi}_{-}(h) \right] \right\}, \qquad (3.108)$$

induces the estimate

$$\widehat{g}_{h}(\eta^{K},\zeta^{L}) = \operatorname{Tr}\left(h\left[\left[\frac{1}{K}\sum_{k}\eta_{k}\right]\left[\frac{1}{L}\sum_{\ell}\zeta-\ell\right]^{T}\right]^{T}\right) + \kappa(h),$$
$$\kappa(h) = \frac{\widehat{\Psi}_{-}(h) - \widehat{\Psi}_{+}(h)}{2}$$

of the functional of interest (3.105) via observation η^K , ζ^L with ϵ -risk on $U \times V$ not exceeding $\rho = \widehat{\Psi}(h)$:

$$\forall (u \in U, v \in V) : \operatorname{Prob}\{|F(u, v) - \widehat{g}_h(\eta^K, \zeta^L)| > \rho\} \le \epsilon,$$

the probability being taken w.r.t. the distribution of observations η^K, ζ^L stemming from signals u, v.

Exercise 3.26. [recovering mixture weights] The problem to be addressed in this Exercise is as follows. We are given K probability distributions $P_1, ..., P_K$ on observation space Ω , and let these distributions have densities $p_k(\cdot)$ w.r.t. some reference measure Π on Ω ; we assume that $\sum_k p_k(\cdot)$ is positive on Ω . We are given also N independent observations

$$\omega_t \sim P_\mu, t = 1, \dots, N,$$

drawn from distribution

$$P_{\mu} = \sum_{k=1}^{K} \mu_k P_k,$$

where μ is unknown "signal known to belong to the probabilistic simplex $\Delta_K = \{\mu \in \mathbf{R}^K : \mu \ge 0, \sum_k \mu_k = 1\}$. Given $\omega^N = (\omega_1, ..., \omega_N)$, we want to recover the linear image $G\mu$ of μ , where $G \in \mathbf{R}^{\nu \to K}$ is given.

We intend to measure the risk of a candidate estimate $\widehat{G}(\omega^N): \Omega \times ... \times \Omega \to \mathbf{R}^{\nu}$ by the quantity

$$\operatorname{Risk}[\widehat{G}(\cdot)] = \sup_{\mu \in \mathbf{\Delta}} \left[\mathbf{E}_{\omega^N \sim P_{\mu} \times \dots \times P_{\mu}} \left\{ \|\widehat{G}(\omega^N) - G\mu\|_2^2 \right\} \right]^{1/2}$$

3.26.A. Recovering linear form. Let us start with the case when $G = g^T$ is $1 \times K$ matrix.

3.26.A.1. Preliminaries. To motivate the construction to follow, consider the case when Ω is a finite set (obtained, e.g., by "fine discretization" of the "true" observation space). In this situation our problem becomes an estimation problem in Discrete o.s., specifically, as follows: given stationary N-repeated observation stemming from discrete probability distribution P_{μ} affinely parameterized by signal $\mu \in \Delta_K$, we want to recover a linear form of μ . It is shown in [86] that in this case (same as when recovering linear forms of signals observed via other simple o.s.'s), a nearly optimal in terms of its risk estimate (see [86] for details) is of the form

$$\widehat{g}(\omega^N) = \frac{1}{N} \sum_{t=1}^N \Phi(\omega_t), \qquad (3.109)$$

with properly selected Φ ; this "proper selection" is obtained by the techniques of Section 3.3 as applied to regular data specifying discrete distributions, see Section 2.8.1.2 The difficulty with this approach is that as far as computations are concerned, optimal design of Φ requires solving convex optimization problem of design dimension of order of the cardinality of G, and this cardinality could be huge already when d is in the range of tens. By this reason, we intend to simplify the outlined approach: the only thing we intend to inherit from the optimality results of [86] is the simple structure (3.109) of the estimator; taking this structure for granted, we intend to develop an alternative to [86] and the construction from Section 3.3 way to design Φ . With these alternative designs, we have no theoretical guarantees for the resulting estimates to be near-optimal; we sacrifice these guarantees in order to reduce dramatically the computational effort of building the estimates.

3.26.A.2. Generic estimate. Let us select somehow L functions $F_{\ell}(\cdot)$ on Ω such that

$$\int F_{\ell}^{2}(\omega)p_{k}(\omega)\Pi(d\omega) < \infty, \ 1 \le \ell \le L, \ 1 \le k \le K$$
(3.110)

With $\lambda \in \mathbf{R}^L$, consider estimate of the form

$$\widehat{g}_{\lambda}(\omega^{N}) = \frac{1}{N} \sum_{t=1}^{N} \Phi_{\lambda}(\omega_{t}), \ \Phi_{\lambda}(\omega) = \sum_{\ell} \lambda_{\ell} F_{\ell}(\omega).$$
(3.111)

1. Prove that

$$\begin{aligned} \operatorname{Risk}[\widehat{g}_{\lambda}] &\leq \overline{\operatorname{Risk}}(\lambda) \\ &:= \max_{k \leq K} \left[\frac{1}{N} \int \left[\sum_{\ell} \lambda_{\ell} F_{\ell}(\omega) \right]^{2} p_{k}(\omega) \Pi(d\omega) \\ &+ \left[\int \left[\sum_{\ell} \lambda_{\ell} F_{\ell}(\omega) \right] p_{k}(\omega) \Pi(d\omega) - g^{T} e_{k} \right]^{2} \right]^{1/2} \\ &= \max_{k \leq K} \left[\frac{1}{N} \lambda^{T} W_{k} \lambda + \left[e_{k}^{T} [M\lambda - g] \right]^{2} \right]^{1/2}, \end{aligned}$$
(3.112)

where

$$\begin{split} M &= \left[M_{k\ell} := \int F_{\ell}(\omega) p_k(\omega) \Pi(d\omega) \right]_{\substack{k \leq K \\ \ell \leq L}}, \\ W_k &= \left[[W_k]_{\ell\ell'} := \int F_{\ell}(\omega) F_{\ell'}(\omega) p_k(\omega) \Pi(d\omega) \right]_{\substack{\ell \leq L \\ \ell' \leq L}}, \ 1 \leq k \leq K. \end{split}$$

and $e_1, ..., e_K$ are the standard basic orths in \mathbf{R}^K .

Note that $\overline{\text{Risk}}(\lambda)$ is a convex function of λ ; this function is easy to compute, provided the matrices M and W_k , $k \leq K$, are available. Assuming this is the case, we can solve the convex optimization problem

$$Opt = \min_{\lambda \in \mathbf{R}^{K}} \overline{Risk}(\lambda)$$
(3.113)

and use the estimate (3.111) associated with optimal solution to this problem; the risk of this estimate will be upper-bounded by Opt.

3.26.A.3. Implementation. The question we arrive at is the "Measurement Design" question: what is a "presumably good," in terms of the (upper bound Opt on the) risk of the estimate (3.111) yielded by an optimal solution to (3.113), selection of L and of the functions F_{ℓ} , $1 \leq \ell \leq L$? We are about to consider three related options – naive, basic, and Maximum Likelihood (ML).

Naive option is to take $F_{\ell} = p_{\ell}$, $1 \leq \ell \leq L = K$, assuming that this selection meets (3.110). For the sake of definiteness, consider the "Gaussian case," where $\Omega = \mathbf{R}^d$, Π is the Lebesgue measure, and p_k is Gaussian distribution with parameters ν_k, Σ_k :

$$p_k(\omega) = \frac{\exp\{-\frac{1}{2}(\omega - \nu_k)^T \Sigma_k^{-1}(\omega - \nu_k)\}}{\sqrt{(2\pi)^d \operatorname{Det}(\Sigma_k)}}.$$

In this case, the Naive option leads to easily computable matrices M and W_k appearing in (3.112).

2. Check that in the Gaussian case, setting

$$\begin{split} \Sigma_{k\ell} &= [\Sigma_k^{-1} + \Sigma_\ell^{-1}]^{-1}, \\ \Sigma_{k\ell m} &= [\Sigma_k^{-1} + \Sigma_\ell^{-1} + \Sigma_m^{-1}]^{-1}, \\ \chi_k &= \Sigma_k^{-1} \nu_k, \\ \alpha_{k\ell} &= \frac{\sqrt{\operatorname{Det}(\Sigma_k)}}{\sqrt{(2\pi)^d \operatorname{Det}(\Sigma_k)\operatorname{Det}(\Sigma_\ell)}}, \\ \beta_{k\ell m} &= \frac{\sqrt{\operatorname{Det}(\Sigma_k)\operatorname{Det}(\Sigma_\ell)\operatorname{Det}(\Sigma_\ell)}}{(2\pi)^d \sqrt{\operatorname{Det}(\Sigma_k)\operatorname{Det}(\Sigma_\ell)\operatorname{Det}(\Sigma_m)}} \end{split}$$

we have

$$M_{k\ell} := \int p_{\ell}(\omega)p_{k}(\omega)\Pi(d\omega) = \alpha_{k\ell}\exp\left\{\frac{1}{2}\left[[\chi_{k}+\chi_{\ell}]^{T}\Sigma_{k\ell}[\chi_{k}+\chi_{\ell}]-\chi_{k}^{T}\Sigma_{k}\chi_{k}-\chi_{\ell}^{T}\Sigma_{\ell}\chi_{\ell}]\right\}, [W_{k}]_{\ell m} := \int p_{\ell}(\omega)p_{m}(\omega)p_{k}(\omega)\Pi(d\omega) = \beta_{k\ell m}\exp\left\{\frac{1}{2}\left[[\chi_{k}+\chi_{\ell}+\chi_{m}]^{T}\Sigma_{k\ell m}[\chi_{k}+\chi_{\ell}+\chi_{m}]-\chi_{k}^{T}\Sigma_{\ell}\chi_{k}-\chi_{\ell}^{T}\Sigma_{\ell}\chi_{\ell}-\chi_{m}^{T}\Sigma_{m}\chi_{m}]\right\}.$$

Basic option. On a close inspection, Naive option does not make much sense: when replacing the reference measure Π with another measure Π' which has positive density $\theta(\cdot)$ w.r.t. Π , the densities p_k are updated according to $p_k(\cdot) \mapsto p'_k(\cdot) =$ $\theta(\cdot)p(\cdot)$, so that selecting $F'_{\ell} = p'_{\ell}$, the matrices M and W_k become M' and W'_k with

$$M'_{k\ell} = \int \frac{p_k(\omega)p_\ell(\omega)}{\theta^2(\omega)} \Pi'(d\omega) = \int \frac{p_k(\omega)p_\ell(\omega)}{\theta(\omega)} \Pi(d\omega),$$

$$[W'_k]_{\ell m} = \int \frac{p_k(\omega)p_\ell(\omega)p_m(\omega)}{\theta^3(\omega)} \Pi'(d\omega) = \int \frac{p_k(\omega)p_\ell(\omega)}{\theta^2(\omega)} \Pi(d\omega).$$

We see that in general $M \neq M'$ and $W_k \neq W'_k$, which makes the Naive option unnatural. *Basic* option is to take

$$L = K, F_{\ell}(\omega) = \pi(\omega) := \frac{p_{\ell}(\omega)}{\sum_{k} p_{k}(\omega)}.$$

The motivation is that the functions F_{ℓ} remain intact when replacing Π with Π' , so that here M = M' and $W_k = W'_k$, which is natural. Besides this, there are statistical arguments in favor of Basic option, namely, as follows. Let Π_* be the measure with the density $\sum_k p_k(\cdot)$ w.r.t. Π ; taken w.r.t. Π_* , the densities of P_k are exactly the above $\pi_k(\cdot)$, and $\sum_k \pi_k(\omega) \equiv 1$. Now, (3.112) says that the risk of estimate \hat{g}_{λ} can be upper-bounded by the function $\overline{\text{Risk}}(\lambda)$ defined in (3.112), and this function, in turn, can be upper-bounded by the function

$$\operatorname{Risk}^{+}(\lambda) := \left[\frac{1}{N} \sum_{k} \int \left[\sum_{\ell} \lambda_{\ell} F_{\ell}(\omega) \right]^{2} p_{k}(\omega) \Pi(d\omega) + \max_{k} \left[\int \left[\sum_{k} \lambda_{\ell} F_{\ell}(\omega) \right] p_{k}(\omega) \Pi(d\omega) - g^{T} e_{k} \right]^{2} \right]^{1/2} \\ = \left[\frac{1}{N} \int \left[\sum_{\ell} \lambda_{\ell} F_{\ell}(\omega) \right]^{2} \Pi_{*}(d\omega) + \max_{k} \left[\int \left[\sum_{k} \lambda_{\ell} F_{\ell}(\omega) \right] \pi_{k}(\omega) \Pi_{*}(d\omega) - g^{T} e_{k} \right]^{2} \right]^{1/2} \\ \leq \overline{\operatorname{Risk}}(\lambda)$$

(we just have said that the maximum of K nonnegatve quantities is at least their sum, and the latter is at most K times the maximum of the quantities). Consequently, the risk of the estimate (3.111) stemming from an optimal solution to (3.113) can be upper-bounded by the quantity

$$\operatorname{Opt}^+ := \min_{\lambda} \operatorname{Risk}^+(\lambda) \quad [\ge \operatorname{Opt} := \max_{\lambda} \overline{\operatorname{Risk}}(\lambda)].$$

Now goes the punchline:

3.1. Prove that both the quantities Opt defined in (3.113) and the above Opt^+ depend

only on the linear span of the functions F_{ℓ} , $\ell = 1, ..., L$, not on how the functions F_{ℓ} are selected in this span.

3.2. Prove that the selection $F_{\ell} = \pi_{\ell}$, $1 \leq \ell \leq L = K$, minimizes Opt⁺ among all possible selections L, $\{F_{\ell}\}_{\ell=1}^{L}$ satisfying (3.110).

Conclude that the selection $F_{\ell} = \pi_{\ell}$, $1 \leq \ell \leq L = K$, while not necessary optimal in terms of Opt, definitely is meaningful: this selection optimizes the natural upper bound Opt⁺ on Opt. Observe that Opt⁺ $\leq K$ Opt, so that optimizing instead of Opt the upper bound Opt⁺, although crude, is not completely meaningless.

A downside of Basic option is that it seems problematic to get closed form expressions for the associated matrices M and W_k , see (3.112). For example, in the Gaussian case, Naive choice of F_ℓ 's allows to represent M and W_k in an explicit closed form; in contrast to this, when selecting $F_\ell = \pi_\ell$, $\ell \leq L = K$, seemingly the only way to get M and W_k is to use Monte-Carlo simulations. This being said, we indeed can use Monte-Carlo simulations to compute M and W_k , provided we can sample from distributions P_1, \ldots, P_K . In this respect, it should be stressed that with $F_\ell \equiv \pi_\ell$, the entries in M and W_k are expectations, w.r.t. P_1, \ldots, P_K , of bounded in magnitude by 1, and thus well-suited for Monte-Carlo simulation, functions of ω .

Maximum Likelihood option. This choice of $\{F_\ell\}_{\ell \leq L}$ follows the idea of discretization Exercise was started with. Specifically, we split Ω into L cells $\Omega_1, ..., \Omega_L$ in such a way that the intersection of any two different cells is of Π -measure zero, and treat as our observations not the actual observations ω_t , but the indexes of the cells ω_t 's belong to. With our estimation scheme, this is the same as to select F_ℓ as the characteristic function of Ω_ℓ , $\ell \leq L$ Assuming that for distinct k, k' the densities $P_k, p_{k'}$ differ from each other Π -almost surely, the simplest discretization independent of how the reference measure is selected is the Maximum Likelihood discretization

$$\Omega_{\ell} = \{ \omega : \max_{k} p_{k}(\omega) = p_{\ell}(\omega) \}, \ 1 \le \ell \le L = K;$$

with the ML option, we take, as F_{ℓ} 's, the characteristic functions of the just defined sets Ω_{ℓ} , $1 \leq \ell \leq L = K$. Same as with Basic option, the matrices M and W_k associated with ML option can be found by Monte-Carlo simulation.

We have discussed 3 simple options for selecting F_{ℓ} 's. In applications, one can compute the upper risk bounds Opt, see (3.113), associated with every one of these three options, and to use the option with the best – the smallest – risk bound. Alternatively, one can take as $\{F_{\ell}, \ell \leq L\}$ the union of the three collections yielded by the above options (and, perhaps, further extend this union). Note that the larger is the collection of F_{ℓ} 's, the smaller is the associated Opt, so that the only price for combining different selections is in increasing the computational cost of solving (3.113).

3.26.A.4. Illustration. Now goes the experimental part of Exercise:

- 4.1. Run numerical experiments aimed at comparing with each other the estimates yielded by the above three options (Naive, Basic, ML). Recommended setup:
 - d = 8, K = 90;
 - Gaussian case with the covariance matrices Σ_k of P_k selected at random:

$$S_k = rand(d, d), \ \Sigma_k = \frac{S_k S_k^T}{\|S_k\|^2}$$
 [|| · ||: spectral norm]

and the expectations ν_k of P_k selected at random from $\mathcal{N}(0, \sigma^2 I_d)$, with $\sigma = 0.1$;

- values of N: $\{10^s, 1 = 0, 1, ..., 5\};$
- linear form to be recovered: $g^T \mu \equiv \mu_1$.
- 4.2^{\dagger} . Utilize Cramer-Rao lower risk bound (see Proposition 4.77, Exercise 4.75) to upper-bound the level of conservatism $\frac{\text{Opt}}{\text{Risk}_*}$ of the estimates built in item 4.1. Here Risk_{*} is the minimax risk in our estimation problem:

$$\operatorname{Risk}_{*} = \inf_{\widehat{g}(\cdot)} \sup_{\mu \in \mathbf{\Delta}} \left[\mathbf{E}_{\omega^{N} \sim P_{\mu} \times \ldots \times P_{\mu}} \left\{ |\widehat{g}(\omega^{N}) - g^{T} \mu|^{2} \right\} \right]^{1/2},$$

where inf is taken over all estimates.

3.26.B. Recovering linear images. Now consider the case when G is a generaltype $\nu \times K$ matrix. The analogy of the estimate $\hat{g}_{\lambda}(\cdot)$ is now as follows: with somehow chosen F_1, \ldots, F_L satisfying (3.110), we select a $\nu \times L$ matrix $\Lambda = [\lambda_{i\ell}]$, set

$$\Phi_{\Lambda}(\omega) = \left[\sum_{\ell} \lambda_{1\ell} F_{\ell}(\omega); \sum_{\ell} \lambda_{2\ell} F_{\ell}(\omega); ...; \sum_{\ell} \lambda_{\nu\ell} F_{\ell}(\omega)\right]$$

and estimate $G\mu$ by

$$\widehat{G}_{\Lambda}(\omega^N) = \frac{1}{N} \sum_{t=1}^N \Phi_{\lambda}(\omega_t).$$

5. Prove the following analogy of the results of item 3.26.A:

Proposition 3.27. The risk of the proposed estimator can be upper-bounded as follows:

$$\operatorname{Risk}[\widehat{G}_{\Lambda}] := \max_{\mu \in \mathbf{\Delta}_{K}} \left[\mathbf{E}_{\omega^{N} \sim P_{\mu} \times \ldots \times P_{\mu}} \left\{ \|\widehat{G}(\omega^{N}) - G\mu\|_{2}^{2} \right\} \right]^{1/2} \\ \leq \operatorname{Risk}(\Lambda) := \max_{k \leq K} \overline{\Psi}(\Lambda, e_{k}), \\ \overline{\Psi}(\Lambda, \mu) = \left[\frac{1}{N} \sum_{k=1}^{K} \mu_{k} \mathbf{E}_{\omega \sim P_{k}} \left\{ \|\Phi_{\Lambda}(\omega)\|_{2}^{2} \right\} + \|[\psi_{\Lambda} - G]\mu\|_{2}^{2} \right]^{1/2} \\ = \left[\|[\psi_{\Lambda} - G]\mu\|_{2}^{2} + \frac{1}{N} \sum_{k=1}^{K} \mu_{k} \int [\sum_{i \leq \nu} [\sum_{\ell} \lambda_{i\ell} F_{\ell}(\omega)]^{2}] P_{k}(d\omega) \right]^{1/2},$$

$$(3.114)$$

where

$$\operatorname{Col}_{k}[\psi_{\Lambda}] = \mathbf{E}_{\omega \sim P_{k}(\cdot)} \Phi_{\Lambda}(\omega) = \begin{bmatrix} \int [\sum_{\ell} \lambda_{1\ell} F_{\ell}(\omega)] P_{k}(d\omega) \\ \cdots \\ \int [\sum_{\ell} \lambda_{\nu\ell} F_{\ell}(\omega)] P_{k}(d\omega) \end{bmatrix}, \ 1 \le k \le K$$

and $e_1, ..., e_K$ are the standard basic orths in \mathbf{R}^K .

Note that exactly the same reasoning as in the case of $G\mu \equiv g^T\mu$ demonstrates that a reasonable way to select L and F_{ℓ} , $\ell = 1, ..., L$, is to set L = K and $F_{\ell}(\cdot) = \pi_{\ell}(\cdot)$, $1 \leq \ell \leq L$.

3.6 PROOFS

3.6.1 Proof of Proposition 3.5

3.6.1.1 Proof of Proposition 3.5.i

We call step ℓ essential, if at this step rule 2d is invoked.

10. Let $x \in X$ be the true signal underlying our observation $\bar{\omega}^K$, so that $\bar{\omega}_1, ..., \bar{\omega}_K$ are independently of each other drawn from the distribution $p_{A(x)}$. Consider the "ideal" estimate given by exactly the same rules as the procedure above (in the sequel, we call the latter the "true" one), up to the fact that the role of the tests $\mathcal{T}^K_{\Delta_{\ell,\mathrm{If}},\mathrm{I}}(\cdot)$, $\mathcal{T}^K_{\Delta_{\ell,\mathrm{If}},\mathrm{I}}(\cdot)$ in rule 2d is played by the "tests"

$$\widehat{T}_{\Delta_{\ell,\mathrm{rg}},\mathrm{r}} = \widehat{T}_{\Delta_{\ell,\mathrm{lf}},\mathrm{l}} = \begin{cases} \mathrm{red}, & f(x) > c_{\ell} \\ \mathrm{blue}, & f(x) \le c_{\ell} \end{cases}$$

Marking by * the entities produced by the resulting *fully deterministic* procedure, we arrive at nested sequence of segments $\Delta_{\ell}^* = [a_{\ell}^*, b_{\ell}^*], 0 \leq \ell \leq L^* \leq L$, along with subsegments $\Delta_{\ell,rg}^* = [c_{\ell}^*, v_{\ell}^*], \Delta_{\ell,lf}^* = [u_{\ell}^*, c_{\ell}^*]$ of $\Delta_{\ell-1}^*$, defined for all *-essential values of ℓ , and the output segment $\bar{\Delta}^*$ claimed to contain f(x). Note that the ideal procedure cannot terminate due to arriving at a disagreement, and that f(x), as is immediately seen, is contained in all segments $\Delta_{\ell}^*, 0 \leq \ell \leq L^*$, same as $f(x) \in \bar{\Delta}^*$.

Let \mathcal{L}^* be the set of all *-essential values of ℓ . For $\ell \in \mathcal{L}^*$, let the event $\mathcal{E}_{\ell}[x]$ parameterized by x be defined as follows:

$$\mathcal{E}_{\ell}[x] = \begin{cases} \{\omega^{K} : \mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^{K},\mathrm{rf}}^{K}(\omega^{K}) = \mathrm{red} \text{ or } \mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}(\omega^{K}) = \mathrm{red} \}, & f(x) \leq u_{\ell}^{*} \\ \{\omega^{K} : \mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^{K},\mathrm{rf}}^{K}(\omega^{K}) = \mathrm{red} \}, & u_{\ell}^{*} < f(x) \leq c_{\ell}^{*} \\ \{\omega^{K} : \mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}(\omega^{K}) = \mathrm{blue} \}, & c_{\ell}^{*} < f(x) < v_{\ell}^{*} \\ \{\omega^{K} : \mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^{K},\mathrm{rf}}^{K}(\omega^{K}) = \mathrm{blue} \text{ or } \mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}(\omega^{K}) = \mathrm{blue} \}, & f(x) \geq v_{\ell}^{*} \end{cases}$$

$$(3.115)$$

2^{0} . Observe that by construction and in view of Proposition 2.31 we have

$$\forall \ell \in \mathcal{L}^* : \operatorname{Prob}_{\omega^K \sim p_{A(x)} \times \ldots \times p_{A(x)}} \{ \mathcal{E}_{\ell}[x] \} \le 2\delta. \tag{3.116}$$

Indeed, let $\ell \in \mathcal{L}^*$.

- When $f(x) \leq u_{\ell}^*$, we have $x \in X$ and $f(x) \leq u_{\ell}^* \leq c_{\ell}^*$, implying that $\mathcal{E}_{\ell}[x]$ takes place only when either the left-side test $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^*,\mathrm{l}}^K$, or the right side test $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^*,\mathrm{rg}}^K$, or both, accepted wrong red hypotheses from the pairs of red and blue hypotheses the tests were applied to. Since the corresponding intervals $([u_{\ell}^*, c_{\ell}^*]$ for the left side test, $[c_{\ell}^*, v_{\ell}^*]$ for the right side one) are δ -good left/right, respectively, the risks of the tests do not exceed δ , and the $p_{A(x)}$ -probability of the event $\mathcal{E}_{\ell}[x]$ is at most 2δ ;
- when $u_{\ell}^* < f(x) \le c_{\ell}^*$, the event $\mathcal{E}_{\ell}[x]$ takes place only when the right side test $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^r}^K$ accepts wrong red of the hypotheses from the pair it is applied to; similarly to the above, this can happen with $p_{A(x)}$ -probability at most δ ;
- when $c_{\ell} < f(x) \leq v_{\ell}$, the event $\mathcal{E}_{\ell}[x]$ takes place only when the left-side test $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}$ accepted wrong blue hypothesis from the pair it was applied to, which again happens with $p_{A(x)}$ -probability $\leq \delta$;

• finally, when $f(x) > v_{\ell}$, the event $\mathcal{E}_{\ell}[x]$ takes place only when either the leftside test $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}$, or the right side test $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^{K},\mathrm{r}}^{K}$, or both, accepted wrong – blue – hypotheses from the pairs of red and blue hypotheses the tests were applied to; same as above, this can happen with $p_{A(x)}$ -probability at most 2δ .

3⁰. Let $\overline{L} = \overline{L}(\overline{\omega}^K)$ be the last step of the "true" estimating procedure as run on the observation $\overline{\omega}^K$. We claim that the following holds true:

(!) Let $\mathcal{E} := \bigcup_{\ell \in \mathcal{L}^*} \mathcal{E}_{\ell}[x]$, so that the $p_{A(x)}$ -probability of the event \mathcal{E} , the observations stemming from x, is at most

 $2\delta L = \epsilon$

by (3.116). Assume that $\bar{\omega}^K \notin \mathcal{E}$. Then $\bar{L}(\bar{\omega}^K) \leq L^*$, and just two cases are possible:

(!.A) The true estimating procedure was not terminated due to arriving at disagreement. In this case $L^* = \overline{L}(\overline{\omega}^K)$ and the trajectories of the ideal and the true procedures are identical (the same localizers, essential steps, the same output segments, etc.), and in particular $f(x) \in \overline{\Delta}$, or

(1.B) The true estimating procedure was terminated due to arriving at a disagreement. Then $\Delta_{\ell} = \Delta_{\ell}^*$ for $\ell < \bar{L}$, and $f(x) \in \bar{\Delta}$.

In view of \mathbf{A} , \mathbf{B} the $p_{A(x)}$ -probability of the event $f(x) \in \overline{\Delta}$ is at least $1 - \epsilon$, as claimed in Proposition 3.5.

To prove (!), note that the actions at step ℓ in the ideal and the true procedures depend solely on $\Delta_{\ell-1}$ and on the outcome of rule 2d. Taking into account that $\Delta_0 = \Delta_0^*$, all we need to verify is the following claim:

(!!) Let $\bar{\omega}^K \notin \mathcal{E}$, and let $\ell \leq L^*$ be such that $\Delta_{\ell-1} = \Delta_{\ell-1}^*$, whence also $u_\ell = u_\ell^*, c_\ell = c_\ell^*, v_\ell = v_\ell^*$. Assume that ℓ is essential (given that $\Delta_{\ell-1} = \Delta_{\ell-1}^*$, this may happen if and only if ℓ is *-essential as well). Then either

C. At step ℓ the true procedure is terminated due to disagreement, in which case $f(x) \in \overline{\Delta}$, or

D. At step ℓ there was no disagreement, in which case Δ_{ℓ} as given by (3.25) is identical to Δ_{ℓ}^* as given by the ideal counterpart of (3.25) in the case of $\Delta_{\ell-1}^* = \Delta_{\ell-1}$, that is, by the rule

$$\Delta_{\ell}^{*} = \begin{cases} [c_{\ell}, b_{\ell-1}], & f(x) > c_{\ell}, \\ [a_{\ell-1}, c_{\ell}], & f(x) \le c_{\ell} \end{cases}$$
(3.117)

To verify (!!), let $\bar{\omega}^K$ and ℓ satisfy the premise of (!!). Note that due to $\Delta_{\ell-1} = \Delta_{\ell-1}^*$ we have $u_\ell = u_\ell^*$, $c_\ell = c_\ell^*$, and $v_\ell = v_\ell^*$, and thus also $\Delta_{\ell,\mathrm{lf}}^* = \Delta_{\ell,\mathrm{lf}}$, $\Delta_{\ell,\mathrm{rg}}^* = \Delta_{\ell,\mathrm{rg}}$. Consider first the case when at the step ℓ the true estimation procedure is terminated due to disagreement, so that $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^*,\mathrm{l}}^K(\bar{\omega}^K) \neq \mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^*,\mathrm{r}}^K(\bar{\omega}^K)$. Assuming for a moment that $f(x) < u_\ell = u_\ell^*$, the relation $\bar{\omega}^K \notin \mathcal{E}_\ell[x]$ combines with (3.115) to imply that $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^*,\mathrm{r}}^K(\bar{\omega}^K) = \mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^*,\mathrm{l}}^K(\bar{\omega}^K) = \mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^*,\mathrm{r}}^K(\bar{\omega}^K) = \mathrm{red}$, which again is impossible. We conclude that in the case in question $u_\ell \leq f(x) \leq v_\ell$, i.e., $f(x) \in \bar{\Delta}$, as claimed in **C**. **C** is proved.

Now let in the true estimating procedure there was a consensus at the step ℓ . The relation $\bar{\omega}^K \notin \mathcal{E}_{\ell}[x]$ implies that one of the following four options takes place:

1. $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^{K},\mathrm{r}}^{K}(\bar{\omega}^{K}) = \text{blue and } f(x) \leq u_{\ell} = u_{\ell}^{*},$ 2. $\mathcal{T}_{\Delta_{\ell,\mathrm{rg}}^{K},\mathrm{r}}^{K}(\bar{\omega}^{K}) = \text{blue and } u_{\ell} < f(x) \leq c_{\ell} = c_{\ell}^{*},$ 3. $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}(\bar{\omega}^{K}) = \text{red and } c_{\ell} < f(x) < v_{\ell} = v_{\ell}^{*},$ 4. $\mathcal{T}_{\Delta_{\ell,\mathrm{lf}}^{K},\mathrm{l}}^{K}(\bar{\omega}^{K}) = \text{red and } v_{\ell} \leq f(x),$

In situations 1-2 and due to consensus at the step ℓ , (3.25) says that $\Delta_{\ell} = [a_{\ell-1}, c_{\ell}]$, which combines with (3.117) and $v_{\ell} = v_{\ell}^*$ to imply that $\Delta_{\ell} = \Delta_{\ell}^*$. Similarly, in situations 3-4 and due to consensus at the step ℓ , (3.25) says that $\Delta_{\ell} = [c_{\ell}, b_{\ell-1}]$, which combines with $u_{\ell} = u_{\ell}^*$ and (3.117) to imply that $\Delta_{\ell} = \Delta_{\ell}^*$. **B** is proved. \Box

3.6.1.2 Proof of Proposition 3.5.ii

There is nothing to prove when $\frac{b_0-a_0}{2} \leq \hat{\rho}$, since in this case the estimate $\frac{a_0+b_0}{2}$ which does not use observations at all is $(\hat{\rho}, 0)$ -reliable. From now on we assume that $b_0 - a_0 > 2\hat{\rho}$, implying that L is positive integer.

1°. Observe, first, that if a, b are such that a is lower-feasible, b is upper-feasible, and $b - a > 2\rho$, then for every $i \leq I_{b,\geq}$ and $j \leq I_{a,\leq}$ there exists a test, based on \bar{K} observations, which decides upon the hypotheses H_1 , H_2 , stating that the observations are drawn from $p_{A(x)}$ with $x \in Z_b^{i,\geq}$ (H_1) and with $x \in Z_j^{a,\leq}$ (H_2) with risk at most ϵ . Indeed, it suffices to consider the test which accepts H_1 and rejects H_2 when $\hat{f}(\omega^{\bar{K}}) \geq \frac{a+b}{2}$ and accepts H_2 and rejects H_1 otherwise.

 2^{0} . With parameters of Bisection chosen according to (3.27), by already proved Proposition 3.5.i, we have

E.1: For every $x \in X$, the $p_{A(x)}$ -probability of the event $f(x) \in \overline{\Delta}$, $\overline{\Delta}$ being the output segment of our Bisection, is at least $1 - \epsilon$.

3^0 . We claim also that

- F.1. Every segment $\Delta = [a, b]$ with $b a > 2\rho$ and lower-feasible a is δ -good (right),
- F.2. Every segment $\Delta = [a, b]$ with $b a > 2\rho$ and upper-feasible b is δ -good (left),
- F.3. Every \varkappa -maximal δ -good (left or right) segment has length at most $2\rho + \varkappa = \widehat{\rho}$. As a result, for every essential step ℓ , the lengths of the segments $\Delta_{\ell, \text{rg}}$ and $\Delta_{\ell, \text{lf}}$ do not exceed $\widehat{\rho}$.

Let us verify F.1 (verification of F.2 is completely similar, and F.3 is an immediate consequence of the definitions and F.1-2). Let [a, b] satisfy the premise of F.1. It may happen that b is upper-infeasible, whence $\Delta = [a, b]$ is 0-good (right), and we are done. Now let b be upper-feasible. As we have already seen, whenever $i \leq I_{b,\geq}$ and $j \leq I_{a,\leq}$, the hypotheses stating that ω_k are sampled from $p_{A(x)}$ for some $x \in Z_j^{b,\geq}$, resp., from some $x \in Z_j^{a,\leq}$, can be decided upon with risk $\leq \epsilon$, implying, same as in the proof of Proposition 2.29, that

$$\epsilon_{ij\Delta} \le \left[2\sqrt{\epsilon(1-\epsilon)}\right]^{1/\bar{K}}$$

whence, taking into account that the column and the row sizes of $E_{\Delta,r}$ do not exceed NI,

$$\sigma_{\Delta,\mathbf{r}} \le NI \max_{i,j} \epsilon_{ij\Delta}^K \le NI [2\sqrt{\epsilon(1-\epsilon)}]^{K/\bar{K}} \le \frac{\epsilon}{2L} = \delta$$

(we have used (3.27)), that is, Δ indeed is δ -good (right).

4⁰. Let us fix $x \in X$ and consider a trajectory of Bisection, the observation being drawn from $p_{A(x)}$. The output $\overline{\Delta}$ of the procedure is given by one of the following options:

- 1. At some step ℓ of Bisection, the process was terminated according to rules in 2b or 2c. In the first case, the segment $[c_{\ell}, b_{\ell-1}]$ has lower-feasible left endpoint and is not δ -good (right), implying by F.1 that the length of this segment (which is 1/2 of the length of $\overline{\Delta} = \Delta_{\ell-1}$) is $\leq 2\rho$, so that the length $|\overline{\Delta}|$ of $\overline{\Delta}$ is at most $4\rho \leq 2\widehat{\rho}$. The same conclusion, by completely similar argument, holds true when the process was terminated at step ℓ according to rule 2c.
- 2. At some step ℓ of Bisection, the process was terminated due to disagreement. In this case, by F.3, we have $|\bar{\Delta}| \leq 2\hat{\rho}$.
- 3. Bisection was terminated at step L, and $\overline{\Delta} = \Delta_L$. In this case, termination clauses in rules 2b, 2c and 2d were never invoked, clearly implying that $|\Delta_s| \leq \frac{1}{2}|\Delta_{s-1}|, 1 \leq s \leq L$, and thus $|\overline{\Delta}| = |\Delta_L| \leq \frac{1}{2^L} |\Delta_0| \leq 2\widehat{\rho}$ (see (3.27)).

Thus, along with E.1 we have

E.2: It always holds $|\overline{\Delta}| \leq 2\widehat{\rho}$,

implying that whenever the signal $x \in X$ underlying observations and the output segment $\overline{\Delta}$ are such that $f(x) \in \overline{\Delta}$, the error of the Bisection estimate (which is the midpoint of $\overline{\Delta}$) is at most $\widehat{\rho}$. Invoking E.1, we conclude that the Bisection estimate is $(\widehat{\rho}, \epsilon)$ -reliable.

3.6.2 2-convexity of conditional quantile

A. Let \mathcal{Q} be the family of non-vanishing probability distributions on $S = \{s_1 < s_2 < ... < s_M\} \subset \mathbf{R}$. For $q \in \mathcal{Q}$, let

$$\ell_m(q) = \sum_{i=m}^M q_i, \ 1 \le m \le M,$$

so that $1 = \ell_1(q) > \ell_2(q) > \dots > \ell_M(q) > 0.$

Given $\alpha \in [0, 1]$, let us define (regularized) α -quantile of $q \in \mathcal{Q}, \chi_{\alpha}(q)$, as follows:

- if $\ell_M(q) > \alpha$, we set $\chi_\alpha(q) = s_M$;
- otherwise, there exists $m \in \{1, ..., M-1\}$ such that $\ell_m(q) \ge \alpha \ge \ell_{m+1}(q)$. We select an m with this property, set

$$\beta = \frac{\alpha - \ell_{m+1}(q)}{\ell_m(q) - \ell_{m+1}(q)},$$

so that
$$\beta \in [0,1]$$
 and $\beta \ell_m(q) + (1-\beta)\ell_{m+1}(q) = \alpha$, and set

$$\chi_{\alpha}(q) = \beta s_m + (1 - \beta) s_{m+1}.$$

Note that for some q, the above m is not uniquely defined; this happens if and only if $\ell_{\mu}(q) = \alpha$ for some μ , $1 < \mu < M$, in which case there are exactly two options for selecting m, one $m = \mu$, and another $m = \mu - 1$. The first option results in

$$\beta = \frac{\alpha - \ell_{\mu+1}(q)}{\ell_{\mu}(q) - \ell_{\mu+1}(q)} = \frac{\ell_{\mu}(q) - \ell_{\mu+1}(q)}{\ell_{\mu}(q) - \ell_{\mu+1}(q)} = 1 \Rightarrow \beta s_m + (1 - \beta)s_{m+1} = s_{\mu},$$

and the second option results in

$$\beta = \frac{\alpha - \ell_{\mu}(q)}{\ell_{\mu-1}(q) - \ell_{\mu}(q)} = \frac{\ell_{\mu}(q) - \ell_{\mu}(q)}{\ell_{\mu-1}(q) - \ell_{\mu}(q)} = 0 \Rightarrow \beta s_m + (1 - \beta)s_{m+1} = s_{\mu};$$

Thus, in spite of the fact that m above is not always uniquely defined by $\alpha, q, \chi_{\alpha}(q)$ is well defined – the value we assign to $\chi_{\alpha}(q)$ according to the above construction is independent of how the required m is selected.

B. From what was said so far, it is immediately seen that for $s_1 \leq s < s_M$, the relation $\chi_{\alpha}(q) = s$ is equivalent to the relation

(!) For some
$$m \in \{1, ..., M-1\}$$
, we have $\ell_m(q) \ge \alpha \ge \ell_{m+1}(q)$ and

$$\frac{[\alpha - \ell_{m+1}(q)]s_m + [\ell_m(q) - \alpha]s_{m+1}}{\ell_m(q) - \ell_{m+1}(q)} = s.$$

C. Let $\theta_q(s), s_1 \leq s \leq s_M$, be the piecewise linear version to the cumulative distribution of q, that is, the piecewise linear function on $\Delta = [s_1, s_M]$ with breakpoints at $s_1, ..., s_M$ and such that $\theta_q(s_m) = \ell_m(q), 1 \leq m \leq M$; this is a strictly decreasing function mapping Δ onto $\Delta^+ := [\ell_M(q), 1]$. For given $q, \chi_\alpha(q)$, as a function of $\alpha \in [0, 1]$, is obtained from the inverse of $\theta_q(\cdot)$ by extending this inverse from its natural domain $\Delta^+ \subset [0, 1]$ to the entire [0, 1] by the value S_M to the left of the left endpoint, $\ell_M(q)$, of Δ^+ . As a consequence of this representation, it is immediately seen that $\chi_\alpha(q)$ is continuous in $(\alpha, q) \in [0, 1] \times Q$. Note that $\chi_\alpha(q)$ takes all its values in $\Delta = [s_1, s_M]$.

Note that we have demonstrated the equivalence between the definition of $\chi_{\alpha}(\cdot)$ via "spreading masses" used in Section 3.2.2.1 and the definition we started with in this Section.

D. Let us fix $\alpha \in (0, 1)$. Given $s \in \Delta$, let us look at the set $Q_s^- := \{q \in \mathcal{Q} : \chi_\alpha(q) \le s\}$. This set is as follows:

- 1. When $s = s_M$, we have $Q_s^- = \mathcal{Q}$.
- 2. Now let $s \in \Delta$ be $\langle s_M$, so that $s \in [s_1, s_M)$. Then for some $\mu = \mu(s) \in \{1, ..., M-1\}$ we have $s_{\mu} \leq s < s_{\mu+1}$. We claim that now the set Q_s^- is the union of two convex sets:

$$Q_{s}^{-} = A \cup B,$$

$$A = \{q \in \mathcal{Q} : \ell_{\mu}(q) \le \alpha\},$$

$$B = \left\{q \in \mathcal{Q} : \frac{\ell_{\mu}(q) \ge \alpha, \ \ell_{\mu+1}(q) \le \alpha,}{[\alpha - \ell_{\mu+1}(q)]s_{\mu} + [\ell_{\mu}(q) - \alpha]s_{\mu+1} \le s[\ell_{\mu}(q) - \ell_{\mu+1}(q)]}\right\}$$
(3.118)

Indeed, when $q \in A$, then $\mu > 1$, since $\ell_1(q) = 1 > \alpha$; thus, $\ell_1(q) > \alpha$ and $\ell_{\mu}(q) \leq \alpha$, so that we can find $m \in \{1, ..., \mu - 1\}$ such that $\ell_m(q) \geq \alpha \geq \ell_{m+1}(q)$. By **A**, it implies that $\chi_{\alpha}(q)$ is a convex combination of s_m and s_{m+1} , and both these quantities are $\leq s_{\mu} \leq s$, so that $\chi_{\alpha}(q) \leq s$ as well, i.e., $q \in Q_s^-$; thus, $A \subset Q_s^-$. Now let $q \in B$. In this case we have $\ell_{\mu}(q) \geq \alpha$, $\ell_{\mu+1}(q) \leq \alpha$, implying that when specifying $\chi_{\alpha}(q)$ by **A**, we can take $m = \mu$, resulting, by (!), in $\chi_{\alpha}(q) = \frac{[\alpha - \ell_{\mu+1}(q)]s_{\mu} + [\ell_{\mu}(q) - \alpha]s_{\mu+1}}{\ell_{\mu}(q) - \ell_{\mu+1}(q)}$. The latter expression, by the third inequality in the description of B, is $\leq s$, and we end up with $\chi_{\alpha}(q) \leq s$ and thus $q \in Q_s^-$. Thus, $A \cup B \subset Q_s^-$. To verify the inverse inclusion, let $q \in \mathcal{Q}$ be such that $\chi_{\alpha}(q) \leq s$, and let us verify that $q \in A \cup B$. It may happen that $\ell_{\mu}(q) \leq \alpha$; then $q \in A$, and we are done. Now assume that $\ell_{\mu}(q) > \alpha$. Observe that in this case $\ell_{\mu+1}(q) \leq \alpha$, since otherwise, by **A**, we would have $\chi_{\alpha}(q) \geq s_{\mu+1}$, while we are in the case $\chi_{\alpha}(q) \leq s < s_{\mu+1}$. Thus, $\ell_{\mu}(q) \geq \alpha \geq \ell_{\mu+1}(q)$, i.e., q satisfies the first two inequalities from the description of B. As a result, by (!), it holds $\chi_{\alpha}(q) = \frac{[\alpha - \ell_{\mu+1}(q)]s_{\mu} + [\ell_{\mu}(q) - \alpha]s_{\mu+1}}{\ell_{\mu}(q) - \ell_{\mu+1}(q)}$, which combines with $\chi_{\alpha}(q) \leq s$ to imply the validity at q of the third inequality in the description of B, and we are done.

The bottom line is that Q_s^- is the union of two closed in \mathcal{Q} convex sets, A and B.

Now let us look at the set $Q_s^+ = \{q \in \mathcal{Q} : \chi_\alpha(q) \ge s\}$, where $s \in \Delta$. This set is as follows:

- 1. When $s = s_M$, Q_s^+ , by **A**, is exactly the set $\{q \in \mathcal{Q} : \ell_M(q) \ge \alpha\}$.
- 2. Now let $s \in \Delta$ be $\langle s_M$, so that for some $\mu = \mu(s) \in \{1, ..., M-1\}$ we have $s_{\mu} \leq s < s_{\mu+1}$. We claim that now the set Q_s^+ is the union of two convex sets:

$$Q_{s}^{+} = A' \cup B',$$

$$A' = \{q \in \mathcal{Q} : \ell_{\mu+1}(q) \ge \alpha\},$$

$$B' = \begin{cases} \ell_{\mu}(q) \ge \alpha, \\ q \in \mathcal{Q} : \ell_{\mu+1}(q) \le \alpha, \\ [\alpha - \ell_{\mu+1}(q)]s_{\mu} + [\ell_{\mu}(q) - \alpha]s_{\mu+1} \ge s[\ell_{\mu}(q) - \ell_{\mu+1}(q)] \end{cases}$$
(3.110)

Indeed, if $q \in A'$, then, by **A**, we either have $\chi_{\alpha}(q) = s_M > s$, or m in **A** can be chosen to be $\geq \mu + 1$ implying by **A** that $\chi_{\alpha}(q) \geq s_{\mu+1} > s$; thus, $A' \subset Q_s^+$. Now let $q \in B'$. From the first two inequalities in the description of B', by (!), we conclude that $\chi_{\alpha}(q) = \frac{[\alpha - \ell_{\mu+1}(q)]s_{\mu} + [\ell_{\mu}(q) - \alpha]s_{\mu+1}}{\ell_{\mu}(q) - \ell_{\mu+1}(q)}, \text{ and the latter quantity, by the third inequality in$ the description of B', is $\geq s$, implying that $q \in Q_s^+$. Thus, $B' \subset Q_s^+$, and we see that $A' \cup B' \subset Q_s^+$. To verify the inverse inclusion, let $q \in \mathcal{Q}$ be such that $\chi_{\alpha}(q) \geq s$, and let us prove that $\chi \in A' \cup B'$. It may happen that $\ell_{\mu+1}(q) \geq \alpha$, in which case $q \in A'$, and we are done. Now let $\ell_{\mu+1}(q) < \alpha$. We claim that $\ell_{\mu}(q) \geq \alpha$. Indeed, otherwise m in A is $< \mu$, implying by A $\chi_{\alpha}(q) \le s_{m+1} \le s_{\mu}$; the equality $\chi_{\alpha}(q) = s_{\mu}$ would be possible only when $m = \mu - 1$ and β , as defined in **A**, is equal to 0, that is, $\alpha = \ell_{m+1}(q) = \ell_{\mu}(q)$, which is not the case. Thus, under assumption $\ell_{\mu}(q) < \alpha$ it holds $\chi_{\alpha}(q) < s_{\mu}$, which is impossible due to $s_{\mu} \leq s$ and $\chi_{\alpha}(q) \geq s$. Thus, we are in the case when $\ell_{\mu}(q) \geq \alpha > \ell_{\mu+1}(q)$, that is, the first two inequalities in the description of B' hold true, which, by (!), implies that $\chi_{\alpha}(q) = \frac{[\alpha - \ell_{\mu+1}(q)]s_{\mu} + [\ell_{\mu}(q) - \alpha]s_{\mu+1}}{\ell_{\mu}(q) - \ell_{\mu+1}(q)}$; the latter combines with $\chi_{\alpha}(q) \geq s$ to imply that q satisfies the third inequality in the description of B', that is, $q \in B'$, and we are done.

The bottom line is that Q_s^+ is the union of two closed in \mathcal{Q} convex sets, A' and B'. We have arrived at the following

Proposition 3.28. Let $S = \{s_1 < s_2 < ... < s_M\}$ be a finite subset of \mathbf{R} , T be a finite set, and \mathcal{P} be the set of non-vanishing probability distributions on $\Omega = S \times T$. Given $\tau \in T$ and $\alpha \in (0,1)$, let $\zeta_{\tau,\alpha}(p) : \mathcal{P} \to [s_1, s_M]$ be the α -quantile of the conditional distribution on S induced by a distribution $p \in \mathcal{P}$ and the condition $t = \tau$:

$$\zeta_{\tau,\alpha}(p) = \chi_{\alpha}(q_{\tau}[p]), \ (q_{\tau}[p])_m = \frac{p(m,\tau)}{\sum_{\mu=1}^M p(\mu,\tau)}, \ 1 \le m \le M.$$

The function $\zeta_{\tau,\alpha}(\cdot)$ is 2-convex on \mathcal{P} : for every $s \in [s_1, s_M)$, selecting $\mu \in \{1, ..., M-1\}$ in such a way that $s_{\mu} \leq s < s_{\mu+1}$, we have

$$\left\{ p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \leq s \right\} = \left\{ p \in \mathcal{P} : \ell_{\mu}(p) - \alpha \ell(p) \leq 0 \right\} \\ \bigcup \left\{ p \in \mathcal{P} : \begin{array}{c} \ell_{\mu}(p) \geq \alpha \ell(p) \geq \ell_{\mu+1}(p), \\ [\alpha \ell(p) - \ell_{\mu+1}(p)]s_{\mu} + [\ell_{\mu}(p) - \alpha \ell(p)]s_{\mu+1} \leq s[\ell_{\mu}(p) - \ell_{\mu+1}(p)] \end{array} \right\}, \\ \left\{ p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \geq s \right\} = \left\{ p \in \mathcal{P} : \ell_{\mu+1}(p) - \alpha \ell(p) \geq 0 \right\} \\ \bigcup \left\{ p \in \mathcal{P} : \begin{array}{c} \ell_{\mu}(p) \geq \alpha \ell(p) \geq \ell_{\mu+1}(p), \\ [\alpha \ell(p) - \ell_{\mu+1}(p)]s_{\mu} + [\ell_{\mu}(p) - \alpha \ell(p)]s_{\mu+1} \geq s[\ell_{\mu}(p) - \ell_{\mu+1}(p)] \end{array} \right\},$$

where

$$\ell_m(p) = \sum_{i=m}^{M} p(i,\tau), \quad \ell(p) = \sum_{m=1}^{M} p(m,\tau),$$

and

$$\begin{array}{ll} s < s_1 & \Rightarrow \\ s = s_M & \Rightarrow \\ s > s_M & \Rightarrow \\ s > s_M & \Rightarrow \end{array} \begin{cases} \{p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \le s\} = \emptyset, \\ \{p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \ge s\} = \mathcal{P}, \\ \{p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \le s\} = \{p \in P : \ell_M(p) \ge \alpha \ell(p)\}, \\ \{p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \le s\} = \{p \in P : \ell_M(p) \ge \alpha \ell(p)\}, \\ \{p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \le s\} = \mathcal{P}, \\ \{p \in \mathcal{P} : \zeta_{\tau,\alpha}(p) \ge s\} = \emptyset \end{cases}$$

Indeed, it suffices to apply (3.119) and (3.118) to $q = q_{\tau}[p]$.

3.6.3 Proof of Proposition 3.15

Let us fix $\epsilon \in (0, 1)$. Setting

$$\rho_K = \frac{1}{2} \left[\widehat{\Psi}^{K,\epsilon}_+(\bar{h},\bar{H}) + \widehat{\Psi}^{K,\epsilon}_-(\bar{h},\bar{H}) \right]$$

and invoking Corollary 3.14, all we need to prove is that in the case of A.1-3 one has

$$\lim \sup_{K \to \infty} \left[\widehat{\Psi}_{+}^{K,\epsilon}(\bar{h},\bar{H}) + \widehat{\Psi}_{-}^{K,\epsilon}(\bar{h},\bar{H}) \right] \le 0.$$
(3.120)

To this end note that in our current situation, (3.64) and (3.69) simplify to

$$\begin{split} \Phi(h,H;Z) &= -\frac{1}{2} \ln \operatorname{Det}(I - \Theta_*^{1/2} H \Theta_*^{1/2}) \\ &+ \frac{1}{2} \operatorname{Tr} \left(Z \underbrace{ \left(B^T \left[\left[\begin{array}{c} H & h \\ h^T & \end{array} \right] + [H,h]^T \left[\Theta_*^{-1} - H \right]^{-1} [H,h] \right] B \right) } \right) \\ \widehat{\Psi}_+^{K,\epsilon}(h,H) &= \inf_{\alpha} \left\{ \max_{Z \in \mathcal{Z}} \left[\alpha \Phi(h/\alpha,H/\alpha;Z) - \operatorname{Tr}(QZ) + K^{-1}\alpha \ln(2/\epsilon) \right] : \\ &\alpha > 0, -\gamma \alpha \Theta_*^{-1} \preceq H \preceq \gamma \alpha \Theta_*^{-1} \right\}, \\ \widehat{\Psi}_-^{K,\epsilon}(h,H) &= \inf_{\alpha} \left\{ \max_{Z \in \mathcal{Z}} \left[\alpha \Phi(-h/\alpha,-H/\alpha;Z) + \operatorname{Tr}(QZ) + K^{-1}\alpha \ln(2/\epsilon) \right] : \\ &\alpha > 0, -\gamma \alpha \Theta_*^{-1} \preceq H \preceq \gamma \alpha \Theta_*^{-1} \right\}, \end{split}$$

whence

$$\begin{split} \left[\widehat{\Psi}_{+}^{K}(\bar{h},\bar{H}) + \widehat{\Psi}_{-}^{K}(\bar{h},\bar{H}) \right] &\leq \inf_{\alpha} \left\{ \max_{Z_{1},Z_{2}\in\mathcal{Z}} \left[\alpha \Phi(\bar{h}/\alpha,\bar{H}/\alpha;Z_{1}) - \operatorname{Tr}(QZ_{1}) \right. \\ \left. + \Phi(-\bar{h}/\alpha,-\bar{H}/\alpha;Z_{1}) + \operatorname{Tr}(QZ_{2}) + 2K^{-1}\alpha \ln(2/\epsilon) \right] : \\ \alpha &> 0, -\gamma \alpha \Theta_{*}^{-1} \preceq \bar{H} \preceq \gamma \alpha \Theta_{*}^{-1} \right\} \\ &= \inf_{\alpha} \max_{Z_{1},Z_{2}\in\mathcal{Z}} \left\{ -\frac{1}{2}\alpha \ln \operatorname{Det} \left(I - \left[\Theta_{*}^{1/2}\bar{H}\Theta_{*}^{1/2}\right]^{2}/\alpha^{2} \right) + 2K^{-1}\alpha \ln(2/\epsilon) \\ \left. + \operatorname{Tr}(Q[Z_{2}-Z_{1}]) + \frac{1}{2} \left[\alpha \operatorname{Tr}(Z_{1}Q(\bar{h}/\alpha,\bar{H}/\alpha)) + \alpha \operatorname{Tr}(Z_{2}Q(-\bar{h}/\alpha,-\bar{H}/\alpha)) \right] : \\ \alpha &> 0, -\gamma \alpha \Theta_{*}^{-1} \preceq \bar{H} \preceq \gamma \alpha \Theta_{*}^{-1} \right\} \\ &= \inf_{\alpha} \max_{Z_{1},Z_{2}\in\mathcal{Z}} \left\{ -\frac{1}{2}\alpha \ln \operatorname{Det} \left(I - \left[\Theta_{*}^{1/2}\bar{H}\Theta_{*}^{1/2}\right]^{2}/\alpha^{2} \right) + 2K^{-1}\alpha \ln(2/\epsilon) \\ \left. +\frac{1}{2}\operatorname{Tr}(Z_{1}B^{T}[\bar{H},\bar{h}]^{T}[\alpha\Theta_{*}^{-1}-\bar{H}]^{-1}[\bar{H},\bar{h}]B) \\ \left. +\frac{1}{2}\operatorname{Tr}(Z_{2}B^{T}[\bar{H},\bar{h}]^{T}[\alpha\Theta_{*}^{-1}+\bar{H}]^{-1}[\bar{H},\bar{h}]B) \\ \left. + \operatorname{Tr}(Q[Z_{2}-Z_{1}]) + \frac{1}{2}\operatorname{Tr}([Z_{1}-Z_{2}]B^{T}\left[\frac{\bar{H}}{\bar{h}^{T}}\right]B) : \right] \\ \left. \right\} \\ \end{array} \right]$$

$$(3.121)$$

By (3.74) we have $\frac{1}{2}B^T \left[\frac{\bar{H}}{\bar{h}^T} \middle| \frac{\bar{h}}{\bar{h}} \right] B = B^T [C^T Q C + J] B$, where the only nonzero entry, if any, in $(d+1) \times (d+1)$ matrix J is in the cell (d+1, d+1). Due to the structure of B, see (3.64), we conclude that the only nonzero element, if any, in $\bar{J} = B^T J B$ is in the cell (m+1, m+1), and that

$$\frac{1}{2}B^T \left[\frac{\bar{H}}{\bar{h}^T} \middle| \frac{\bar{h}}{\bar{h}} \right] B = (CB)^T Q(CB) + \bar{J} = Q + \bar{J}$$

(recall that $CB = I_{m+1}$). Now, when $Z_1, Z_2 \in \mathbb{Z}$, the entries of Z_1, Z_2 in the cell (m+1, m+1) both are equal to 1, whence

$$\frac{1}{2} \operatorname{Tr}([Z_1 - Z_2]B^T \begin{bmatrix} \bar{H} & \bar{h} \\ \bar{h}^T & \end{bmatrix} B) = \operatorname{Tr}([Z_1 - Z_2]Q) + \operatorname{Tr}([Z_1 - Z_2]\bar{J}) = \operatorname{Tr}([Z_1 - Z_2]Q),$$

implying that the quantity $T(Z_1, Z_2)$ in (3.121) is zero, provided $Z_1, Z_2 \in \mathcal{Z}$. Consequently, (3.121) becomes

$$\left[\widehat{\Psi}_{+}^{K}(\bar{h},\bar{H}) + \widehat{\Psi}_{-}^{K}(\bar{h},\bar{H}) \right] \leq \inf_{\alpha} \max_{Z_{1},Z_{2}\in\mathcal{Z}} \left\{ -\frac{1}{2}\alpha \ln \operatorname{Det} \left(I - [\Theta_{*}^{1/2}\bar{H}\Theta_{*}^{1/2}]^{2}/\alpha^{2} \right) \\ + 2K^{-1}\alpha \ln(2/\epsilon) + \frac{1}{2}\operatorname{Tr} \left(Z_{1}B^{T}[\bar{H},h][\alpha\Theta_{*}^{-1}-\bar{H}]^{-1}[\bar{H},\bar{h}]^{T}B \right) \\ + \frac{1}{2}\operatorname{Tr} \left(Z_{2}B^{T}[\bar{H},\bar{h}]^{T}[\alpha\Theta_{*}^{-1}+\bar{H}]^{-1}[\bar{H},\bar{h}]B \right) : \alpha > 0, -\gamma\alpha\Theta_{*}^{-1} \leq \bar{H} \leq \gamma\alpha\Theta_{*}^{-1} \right\}$$
(3.122)

Now, for appropriately selected independent of K real c we have

$$\begin{aligned} &-\frac{1}{2}\alpha \ln \operatorname{Det} \left(I - [\Theta_*^{1/2} \bar{H} \Theta_*^{1/2}]^2 / \alpha^2 \right) \le c/\alpha, \\ &\frac{1}{2} \operatorname{Tr} \left(Z_1 B^T [\bar{H}, \bar{h}]^T [\alpha \Theta_*^{-1} - \bar{H}]^{-1} [\bar{H}, \bar{h}] B \right) \\ &+ \frac{1}{2} \operatorname{Tr} \left(Z_2 B^T [\bar{H}, \bar{h}]^T [\alpha \Theta_*^{-1} + \bar{H}]^{-1} [\bar{H}, \bar{h}] B \right) \le c/\alpha \; \forall Z_1, Z_2 \in \mathcal{Z} \end{aligned}$$

(recall that \mathcal{Z} is bounded). Consequently, given $\omega > 0$, we can find $\alpha = \alpha_{\omega} > 0$ large enough to ensure that

$$-\gamma \alpha_{\omega} \Theta_*^{-1} \preceq \bar{H} \preceq \gamma \alpha_{\omega} \Theta_*^{-1} \& 2c/\alpha_{\omega} \leq \omega,$$

which combines with (3.122) to imply that

$$\left[\widehat{\Psi}_{+}^{K}(\bar{h},\bar{H}) + \widehat{\Psi}_{-}^{K}(\bar{h},\bar{H})\right] \leq \omega + 2K^{-1}\alpha_{\omega}\ln(2/\epsilon),$$

and (3.120) follows.

Lecture Four

Signal Recovery from Gaussian Observations and Beyond

OVERVIEW

In this lecture we address one of the most basic problems of High-Dimensional Statistics, specifically, as follows: given positive definite $m \times m$ matrix Γ , $m \times n$ matrix A, $\nu \times n$ matrix B, and indirect noisy observation

$$\omega = Ax + \xi$$

[$A : m \times n, \xi \sim \mathcal{N}(0, \Gamma)$] (4.1)

of unknown "signal" x known to belong to a given convex compact subset \mathcal{X} of \mathbf{R}^n , we want to recover the image $Bx \in \mathbf{R}^{\nu}$ of x under a given linear mapping. We focus first on the case where the quality of a candidate recovery $\omega \mapsto \hat{x}(\omega)$ is quantified by its worst-case, over $x \in \mathcal{X}$, expected $\|\cdot\|_2^2$ -error, that is, by the risk

$$\operatorname{Risk}[\widehat{x}(\cdot)|\mathcal{X}] = \sup_{x \in \mathcal{X}} \sqrt{\mathbf{E}_{\xi \sim \mathcal{N}(0,\Gamma)} \left\{ \|\widehat{x}(Ax+\xi) - Bx\|_2^2 \right\}}.$$
(4.2)

The simplest and the most studied type of recovery is affine one: $\hat{x}(\omega) = H^T \omega + h$; assuming \mathcal{X} symmetric w.r.t. the origin, we lose nothing when passing from affine estimates to linear ones – those of the form $\hat{x}_H(\omega) = H^T \omega$. An advantage of linear estimates is that under favorable circumstances (e.g., when \mathcal{X} is an ellipsoid), minimizing risk over linear estimates is an efficiently solvable problem, and there exists huge literature on optimal in terms of their risk linear estimates (see, e.g., [95, 96, 129, 130, 41, 54, 123, 3] and references therein). Moreover, in the case of signal recovery from direct observations in white Gaussian noise (the case of $B = A = I_n$, $\Gamma = \sigma^2 I_n$), there is huge body of results on near-optimality of properly selected linear estimates among *all* possible recovery routines, see, e.g., [81, 141] and references therein; a typical result of this type states that when recovering $x \in \mathcal{X}$ from direct observation $\omega = x + \sigma\xi$, $\xi \sim \mathcal{N}(0, I_m)$ and \mathcal{X} being an ellipsoid of the form

$$\{x \in \mathbf{R}^n : \sum_j j^{2\alpha} x_j^2 \le L^2\},\$$

or the box

$$\{x \in \mathbf{R}^n : j^{\alpha} | x_j | \le L, j \le n\}$$

with fixed $L < \infty$ and $\alpha > 0$, the ratio of the risk of a properly selected linear estimate to the *minimax risk*

$$\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}] := \inf_{\widehat{x}(\cdot)} \operatorname{Risk}[\widehat{x}|\mathcal{X}]$$
(4.3)

(the infimum is taken over all estimates, not necessarily linear) remains bounded, or even tends to 1, as $\sigma \to +0$, and this happens *uniformly in n*, α and *L* being fixed. Similar "near-optimality" results are known for "diagonal" case, where \mathcal{X}

is the above ellipsoid/box and A, B, Γ are diagonal matrices. To the best of our knowledge, the only "general" (that is, not imposing severe restrictions on how the geometries of $\mathcal{X}, A, B, \Gamma$ are linked to each other) result on optimality of linear estimates is due to D. Donoho who proved [44, 50] that when recovering a linear form (i.e., in the case of one-dimensional Bx), the best, over all linear estimates, risk is within the factor 1.2 of the minimax risk.

The goal of this lecture is to establish a rather general result on near-optimality of properly built linear estimates as compared to all possible estimates. A result of this type is bound to impose some restrictions on \mathcal{X} , since there are cases (e.g., the one of a high-dimensional $\|\cdot\|_1$ -ball \mathcal{X}) where linear estimates are by far nonoptimal. Our restrictions on \mathcal{X} reduce to the existence of a special type representation of \mathcal{X} and are satisfied, e.g., when \mathcal{X} is the intersection of $K < \infty$ ellipsoids/elliptic cylinders:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : x^T R_k x \le 1, 1 \le k \le K \}.$$

$$[R_k \ge 0, \sum_k R_k \succ 0]$$

$$(4.4)$$

in particular, \mathcal{X} can be a symmetric w.r.t. the origin compact polytope given by 2K linear inequalities $-1 \leq s_k^T x \leq 1$, $1 \leq k \leq K$, or, equivalently, $\mathcal{X} = \{x : x^T \underbrace{(r_k r_k^T)}_{R_k} x \leq 1, k \leq K\}$. Another instructive example is a set of the form

 $\mathcal{X} = \{x : \|Sx\|_p \leq L\}$, where $p \geq 2$ and S is a matrix with trivial kernel. It should be stressed than while imposing some restrictions on \mathcal{X} , we require nothing from $A, B, and \Gamma$, aside of positive definiteness of the latter matrix. Our main result (Proposition 4.5) states, in particular, that with \mathcal{X} given by (4.4) and arbitrary A, B, the risk of properly selected linear estimate \hat{x}_{H_*} with both H_* and the risk efficiently computable, satisfies the bound

$$\operatorname{Risk}[\widehat{x}_{H_*}|\mathcal{X}] \le O(1)\sqrt{\ln(K+1)}\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}], \qquad (*)$$

Risk_{opt} $[\mathcal{X}]$ is the minimax risk, and O(1) is an absolute constant. Note that the outlined result is an "operational" one – the risk of *provably nearly optimal* estimate and the estimate itself are given by efficient computation. This is in sharp contrast with traditional results of non-parametric statistics, where near-optimal estimates and their risks are given in a "closed analytical form," at the price of severe restrictions on the structure of the "data" \mathcal{X} , A, B, Γ . This being said, it should be stressed that one of the crucial components in our construction is quite classical – this is the idea, going back to M.S. Pinsker [124], to bound from below the minimax risk via Bayesian risk associated with properly selected Gaussian prior⁵².

The main body of the lecture originates from [92, 90] is organized as follows. Section 4.1 presents basic results on Conic Programming and Conic Duality – the main optimization tools utilized in all subsequent constructions and proofs. Section 4.2 contains problem formulation (Section 4.2.1), construction of the linear estimate we deal with (Section 4.2.2) and the central result on near-optimality of this estimate (Section 4.2.3). We discuss also the "expressive abilities" of the family

⁵²[124] addresses the problem of $\|\cdot\|_2$ -recovery of a signal x from direct observations (A = B = I) in the case when \mathcal{X} is a high-dimensional ellipsoid with "regularly decreasing half-axes," like $\mathcal{X} = \{x \in \mathbf{R}^n : \sum_j j^{2\alpha} x_j^2 \leq L^2\}$ with $\alpha > 0$. In this case Pinsker's construction shows that as $\sigma \to +0$, the risk of properly built linear estimate is, uniformly in n, (1 + o(1)) times the minimax risk. This is much stronger than (*), and it seems to be unlikely that a similarly strong result holds true in the general case underlying (*).

SIGNAL RECOVERY FROM GAUSSIAN OBSERVATIONS AND BEYOND

of sets (we call them *ellitopes*) to which our main result applies. In Section 4.2.5, we show that the key argument underlying the proof of our main result can be used beyond the scope of statistics, specifically, when quantifying the approximation ratio of the semidefinite relaxation bound on the maximum of a quadratic form over an ellitope. In Section 4.3 we extend the results of previous sections from ellitopes to their "matrix analogies" – *spectratopes* in the role of signal sets, passing simultaneously from $\|\cdot\|_2$ as the norm in which the recovery error is measured to an arbitrary *spectratopic* norm – one for which the unit ball of the conjugate norm is a spectratope; in addition, we allow for observation noise to have nonzero mean and to be non-Gaussian. Concluding Section 4.5 deals with "uncertain-but-bounded" observation noise, that is, noise selected "by nature," perhaps in an adversarial fashion, from a given bounded set.

4.1 PRELIMINARIES: EXECUTIVE SUMMARY ON CONIC PROGRAMMING

4.1.1 Cones

A cone in Euclidean space E is a nonempty set K which is closed w.r.t. taking conic combinations of its elements, that is, linear combinations with nonnegative coefficients. Equivalently: $K \subset E$ is a cone if K is nonempty, and

- $x, y \in K \Rightarrow x + y \in K;$
- $x \in K, \lambda \ge 0 \Rightarrow \lambda x \in K.$

It is immediately seen (check it!) that a cone is a convex set. We call a cone K regular, if it is closed, *pointed* (that is, does not contain lines, or, equivalently, $K \bigcap [-K] = \{0\}$) and possesses a nonempty interior.

Given a cone $K \subset E$, we can associate with it its *dual cone* K^* defined as

$$K^* = \{ y \in E : \langle y, x \rangle \ge 0 \, \forall x \in K \};$$

it is immediately seen that whatever be K, K^* is a closed cone, and $K \subset (K^*)^*$. It is well known that

- if K is a closed cone, it holds $K = (K^*)^*$;
- K is a regular cone if and only if K^* is so.

Examples of "useful in applications" regular cones are as follows:

- 1. Nonnegative orthants $\mathbf{R}^d_+ = \{x \in \mathbf{R}^d : x \ge 0\}$
- 2. Lorentz cones $\mathbf{L}^{d}_{+} = \{ x \in \mathbf{R}^{d} : x_{n} \ge \sqrt{\sum_{i=1}^{n-1} x_{i}^{2}} \};$
- 3. Semidefinite cones \mathbf{S}^d_+ comprised of positive semidefinite symmetric $d \times d$ matrices; Semidefinite cone \mathbf{S}^d_+ lives in the space \mathbf{S}^d of symmetric matrices equipped with the Frobenius inner product $\langle A, B \rangle = \operatorname{Tr}(AB^T) = \operatorname{Tr}(AB) = \sum_{i,j=1}^d A_{ij}B_{ij}, A, B \in \mathbf{S}^d$.

All listed so far cones are self-dual.

4. Let $\|\cdot\|$ be a norm on \mathbb{R}^n . The set $\{[x;t] \in \mathbb{R}^n \times \mathbb{R} : t \ge \|x\|\}$ is a regular cone,

and the dual cone is $\{[y; \tau] : ||y||_* \le \tau\}$, where

$$\|y\|_* = \max_x \{x^T y : \|x\| \le 1\}$$

is the norm on \mathbf{R}^n conjugate to $\|\cdot\|$.

Another useful for the sequel example of a regular cone is the *conic hull* of a convex compact set defined as follows. Let \mathcal{T} be a convex compact set with a nonempty interior in Euclidean space E. We can associate with \mathcal{T} its *closed conic hull*

$$\mathbf{T} = \operatorname{cl}\underbrace{\left\{ [t;\tau] \in E^+ = E \times \mathbf{R} : \tau > 0, t/\tau \in \mathcal{T} \right\}}_{K^o(\mathcal{T})}.$$

It is immediately seen that **T** is a regular cone (check it!), and that to get this cone, one should add to the convex set $K^o(\mathcal{T})$ the origin in E^+ . It is also clear that one can "see \mathcal{T} in **T**:" – \mathcal{T} is nothing but the cross-section of the cone **T** by the hyperplane $\tau = 1$ in $E^+ = \{[t; \tau]\}$:

$$\mathcal{T} = \{t \in E : [t; 1] \in \mathbf{T}\}\$$

It is easily seen (check it!) that the cone \mathbf{T}_* dual to \mathbf{T} is given by

$$\mathbf{T}_* = \{ [g; s] \in \mathbf{E}^+ : s \ge \phi_{\mathcal{T}}(-g) \}$$

where

$$\phi_{\mathcal{T}}(g) = \max_{t \in \mathcal{T}} \langle g, t \rangle$$

is the support function of \mathcal{T} .

4.1.2 Conic problems and their duals

Given regular cones $K_i \subset E_i$, $1 \leq i \leq m$, consider optimization problem of the form

$$Opt(P) = \min\left\{ \langle c, x \rangle : \begin{array}{c} A_i x - b_i \in K_i, \ i = 1, ..., m \\ Rx = r \end{array} \right\}, \tag{P}$$

where $x \mapsto A_i x - b_i$ are affine mappings acting from some Euclidean space E to the spaces E_i where the cones K_i live. Problem in this form is called a *conic* problem on the cones $K_1, ..., K_m$; the constraints $A_i x - b_i \in K_i$ on x are called *conic constraints*. We call a conic problem (P) strictly feasible, if it admits a strictly feasible solution \bar{x} , meaning that \bar{x} satisfies the equality constraints and satisfies strictly: $A_i \bar{x} - b_i \in \text{int } K_i$ – the conic constraints.

One can associate with conic problem (P) its *dual*, which also is a conic problem. The origin of the dual problem is the desire to obtain in a systematic way -by linear aggregation of conic constraints – lover bounds on the optimal value Opt(P) of the primal problem (P). Linear aggregation of constraints works as follows: let us equip every one of conic constraints $A_ix - b_i \in K_i$ with aggregation weight, called Lagrange multiplier, y_i restricted to reside in the cone K_i^* dual to K_i . Similarly, we equip the system Rx = r of equality constraints in (P) with Lagrange multiplier z– a vector of the same dimension as r. Now let x be a feasible solution to the conic problem, and let $y_i \in K_i^*$, $i \leq m, z$ be Lagrange multipliers. By the definition of

SIGNAL RECOVERY FROM GAUSSIAN OBSERVATIONS AND BEYOND

the dual cone and due to $A_i x - b_i \in K_i, y_i \in K_i^*$ we have

$$\langle y_i, A_i x \rangle \ge \langle y_i, b_i \rangle, 1 \le i \le m$$

and of course

$$z^T R x \ge r^T z$$

Summing all resulting inequalities up, we arrive at the scalar linear inequality

$$\langle R^* z + \sum_i A_i^* y_i, x \rangle \ge r^T z + \sum_i \langle b_i, y_i \rangle \tag{!}$$

where A_i^* are the conjugates to A_i : $\langle y, A_i x \rangle_{E_i} \equiv \langle A_i^* y, x \rangle_E$, and R^* is the conjugate of R. By its origin, (!) is a consequence of the system of constraints in (P) and as such is satisfied everywhere on the feasible domain of the problem. If we are lucky to get, as the linear function of x in the left hand side of (!), the objective of (P), that is, if

$$R^*z + \sum_i A_i^* y_i = c,$$

(!) imposes a lower bound on the objective of the primal conic problem (P) everywhere on the feasible domain of the primal problem, and the *conic dual* of (P) is the problem

$$Opt(D) = \max_{y_i, z} \left\{ r^T z + \sum_i \langle b_i, y_i \rangle : \begin{array}{c} y_i \in K_i^*, \ 1 \le i \le m \\ R^* z + \sum_{i=1}^m A_i^* y_i = c \end{array} \right\}$$
(D)

of maximizing this lower bound on Opt(P).

The relations between the primal and the dual conic problems are the subject of the standard *Conic Duality Theorem* as follows:

Theorem 4.1. [Conic Duality Theorem] Consider conic problem (P) (where all K_i are regular cones) along with its dual problem (D). Then

- 1. Duality is symmetric: the dual problem (D) is conic, and the conic dual of (D) is (equivalent to) (P);
- 2. Weak duality: It always holds $Opt(D) \leq Opt(P)$
- 3. Strong duality: If one of the problems (P), (D) is strictly feasible and bounded⁵³, then the other problem in the pair is solvable, and the optimal values of the problems are equal to each other. In particular, if both (P) and (D) are strictly feasible, then both problems are solvable with equal optimal values.

Remark 4.2. While Conic Duality Theorem in the just presented form meets all our subsequent needs, it makes sense to note that in fact Strong Duality part of the theorem can be strengthened by replacing strict feasibility with "essential strict feasibility" defined as follows: conic problem in the form of (P) (or, which is the same, form of (D)) is called essentially strictly feasible, if it admits a feasible solution \bar{x} which satisfies strictly the *non-polyhedral* conic constraints, that is, $A_i\bar{x}$ -

 $^{^{53}}$ For a minimization problem, boundedness means that the objective is bounded from below on the feasible set, for a maximization problem – that it is bounded from above on the feasible set.

 $b_i \in \text{int } K_i \text{ for all } i \text{ for which the cone } K_i \text{ is } not \text{ polyhedral} - \text{ is } not \text{ given by a finite list of homogeneous linear inequality constraints.}$

The proof of Conic Duality Theorem can be found in numerous sources, e.g., in [116, Section 7.1.3].

4.1.3 Schur Complement Lemma

We will use the following extremely useful fact:

Lemma 4.3. [Schur Complement Lemma] Symmetric block matrix

$$A = \begin{bmatrix} P & Q^T \\ \hline Q & R \end{bmatrix}$$

with $R \succ 0$ is positive (semi)definite if and only if the matrix $P - Q^T R^{-1}Q$ is so.

Proof. With u, v of the same sizes as P, respectively, R, we have

$$\min_{u} [u; v]^T A [u; v] = u^T [P - Q^T R^{-1} Q] u$$

(direct computation utilizing the fact that $R \succ 0$). It follows that the quadratic form associated with A is nonnegative everywhere if and only if the quadratic form with the matrix $[P-Q^T R^{-1}Q]$ is nonnegative everywhere (since the latter quadratic form is obtained from the former one by partial minimization).

4.2 NEAR-OPTIMAL LINEAR ESTIMATION

4.2.1 Situation and goal

Given $m \times n$ matrix A, $\nu \times n$ matrix B, and $m \times m$ matrix $\Gamma \succ 0$, consider the problem of estimating linear image Bx of unknown signal x known to belong to a given set $\mathcal{X} \subset \mathbf{R}^n$ via noisy observation

$$\omega = Ax + \xi, \ \xi \sim \mathcal{N}(0, \Gamma), \tag{4.5}$$

where ξ is the observation noise. A candidate estimate in this case is a (Borel) function $\hat{x}(\cdot) : \mathbf{R}^m \to \mathbf{R}^{\nu}$, and the performance of such an estimate in what follows will be quantified by the *Euclidean risk* Risk $[\hat{x}|\mathcal{X}]$ defined by (4.2).

4.2.1.1 Ellitopes

From now on we assume that $\mathcal{X} \subset \mathbf{R}^n$ is a set given by

$$\mathcal{X} = \left\{ x \in \mathbf{R}^n : \exists (y \in \mathbf{R}^{\bar{n}}, t \in \mathcal{T}) : x = Py, y^T R_k y \le t_k, 1 \le k \le K \right\},$$
(4.6)

where

- P is an $n \times \bar{n}$ matrix,
- $R_k \succeq 0$ are $\bar{n} \times \bar{n}$ matrices with $\sum_k R_k \succ 0$,

SIGNAL RECOVERY FROM GAUSSIAN OBSERVATIONS AND BEYOND

• \mathcal{T} is a nonempty computationally tractable⁵⁴ convex compact subset of \mathbf{R}_{+}^{K} intersecting the interior of \mathbf{R}_{+}^{K} and such that \mathcal{T} is monotone, meaning that the relations $0 \leq \tau \leq t$ and $t \in \mathcal{T}$ imply that $\tau \in \mathcal{T}$.⁵⁵ Note that under our assumptions int $\mathcal{T} \neq \emptyset$.

In the sequel, we refer to a set of the form (4.6) with data $[P, \{R_k, 1 \leq k \leq K\}, \mathcal{T}]$ satisfying just formulated assumptions as to an *ellitope*, and to (4.6) – as to *ellitopic representation* of \mathcal{X} . Here are instructive examples of ellitopes (in all these examples, P is the identity mapping; in the sequel, we call ellitopes of this type *basic* ones):

- when K = 1, $\mathcal{T} = [0, 1]$ and $R_1 \succ 0$, \mathcal{X} is the ellipsoid $\{x : x^T R_1 x \leq 1\};$
- when $K \ge 1$, $\mathcal{T} = \{t \in \mathbf{R}^K : 0 \le t_k \le 1, k \le K\}$, and \mathcal{X} is the intersection

$$\bigcap_{1 \le k \le K} \{ x : x^T R_k x \le 1 \}$$

of centered at the origin ellipsoids/elliptic cylinders. In particular, when U is a $K \times n$ matrix of rank n with rows u_k^T , $1 \le k \le K$, and $R_k = u_k u_k^T$, \mathcal{X} is symmetric w.r.t. the origin polytope $\{x : ||Ux||_{\infty} \le 1\}$;

• when U, u_k and R_k are as in the latter example and $\mathcal{T} = \{t \in \mathbf{R}_+^K : \sum_k t_k^{p/2} \leq 1\}$ for some $p \geq 2$, we get $\mathcal{X} = \{x : ||Ux||_p \leq 1\}$.

It should be added that the family of ellitope-representable sets is quite rich: this family admits a "calculus", so that more ellitopes can be constructed by taking intersections, direct products, linear images (direct and inverse) or arithmetic sums of ellitopes given by the above examples. In fact, the property to be an ellitope is preserved by all basic operations with sets preserving convexity and symmetry w.r.t. the origin, see Section 4.8.

As another instructive, in the context of non-parametric statistics, example of an ellitope, consider the situation where our signals x are discretizations of functions of continuous argument running through a compact d-dimensional domain D, and the functions f we are interested in are those satisfying a Sobolev-type smoothness constraint – an upper bound on the $L_p(D)$ -norm of $\mathcal{L}f$, where \mathcal{L} is a linear differential operator with constant coefficients. After discretization, this restriction can be modeled as $||Lx||_p \leq 1$, with properly selected matrix L. As we already know from the above example, when $p \geq 2$, the set $\mathcal{X} = \{x : ||Lx||_p \leq 1\}$ is an ellitope, and as such is captured by our machinery. Note also that by the outlined calculus, imposing on the functions f in question several Sobolev-type smoothness constraints with parameters $p \geq 2$, still results in a set of signals which is an ellitope.

$$\mathcal{T} = \{t : \exists w : A(t, w) \succeq 0\},\$$

where A(t, w) is a symmetric and affine in t, w matrix.

 $^{^{54} {\}rm for}$ all practical purposes, it suffices to assume that ${\cal T}$ is given by an explicit semidefinite representation

⁵⁵The latter relation is "for free" – given a nonempty convex compact set $\mathcal{T} \subset \mathbf{R}_{+}^{K}$, the right hand side of (4.6) remains intact when passing from \mathcal{T} to its "monotone hull" { $\tau \in \mathbf{R}_{+}^{K} : \exists t \in \mathcal{T} : \tau \leq t$ } which already is a monotone convex compact set.

4.2.1.2 Estimates and their risks

In the outlined situation, a candidate estimate is a Borel function $\hat{x}(\cdot) : \mathbf{R}^m \to \mathbf{R}^{\nu}$; given observation (4.5), we recover w = Bx as $\hat{x}(\omega)$. In the sequel, we quantify the quality of an estimate by its worst-case, over $x \in \mathcal{X}$, expected $\|\cdot\|_2^2$ recovery error:

$$\operatorname{Risk}[\widehat{x}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \left[\mathbf{E}_{\xi \sim \mathcal{N}(0,\Gamma)} \left\{ \|\widehat{x}(Ax+\xi) - Bx\|_2^2 \right\} \right]^{1/2}$$
(4.7)

and define the optimal, or the *minimax*, risk as

$$\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}] = \inf_{\widehat{x}(\cdot)} \operatorname{Risk}[\widehat{x}|\mathcal{X}], \qquad (4.8)$$

where inf is taken over all Borel candidate estimates.

4.2.1.3 Main goal

Main goal of what follows is to demonstrate that a *linear in* ω estimate

$$\widehat{x}_H(\omega) = H^T \omega \tag{4.9}$$

with properly selected efficiently computable matrix H is near-optimal in terms of its risk.

Our initial observation is that when replacing matrices A and B with AP and BP, respectively, we pass from the initial estimation problem of interest – one where the signal set \mathcal{X} is given by (4.6), and we want to recover $Bx, x \in \mathcal{X}$, via observation (4.5), to the *transformed problem*, where the signal set is

$$\bar{X} = \{ y \in \mathbf{R}^{\bar{n}} : \exists t \in \mathcal{T} : y^T R_k y \le t_k, \ 1 \le k \le K \},\$$

and we want to recover $[BP]y, y \in \overline{X}$, via observation

$$\omega = [AP]y + \xi.$$

It is obvious that the considered families of estimates (the family of all linear and the family of all estimates), same as the risks of the estimates, remain intact under this transformation; in particular,

$$\operatorname{Risk}[\widehat{x}|\mathcal{X}] = \sup_{y \in \widehat{X}} \left[\mathbf{E}_{\xi} \{ \|\widehat{x}([AP] \, y + \xi) - [BP] \, y\|_2^2 \} \right]^{1/2}.$$

Therefore, to save notation, from now on, unless explicitly stated otherwise, we assume that matrix P is identity, so that \mathcal{X} is the basic ellitope

$$\mathcal{X} = \left\{ x \in \mathbf{R}^n : \exists t \in \mathcal{T}, \ x^T R_k x \le t_k, \ 1 \le k \le K \right\}.$$
(4.10)

We assume in the sequel that $B \neq 0$, since otherwise one has Bx = 0 for all $x \in \mathcal{X}$, and the estimation problem is trivial.

4.2.2 Building linear estimate

We start with building a "presumably good" linear estimate. Restricting ourselves to linear estimates (4.9), we may be interested in the estimate with the smallest
risk, that is, associated with a $\nu \times m$ matrix H which is an optimal solution to the optimization problem

$$\min_{H} \left\{ R(H) := \operatorname{Risk}^2[\widehat{x}_H | \mathcal{X}] \right\}$$

We have

$$R(H) = \max_{x \in \mathcal{X}} \mathbf{E}_{\xi} \{ \| H^{T} \omega - Bx \|_{2}^{2} \} = \mathbf{E}_{\xi} \{ \| H^{T} \xi \|_{2}^{2} \} + \max_{x \in \mathcal{X}} \| H^{T} Ax - Bx \|_{2}^{2}$$

= $\operatorname{Tr}(H^{T} \Gamma H) + \max_{x \in \mathcal{X}} x^{T} (H^{T} A - B)^{T} (H^{T} A - B)x.$

This function, while convex, can be hard to compute. For this reason, we use a linear estimate yielded by minimizing an *efficiently computable convex upper bound* on R(H) which is built as follows. Let $\phi_{\mathcal{T}}$ be the support function of \mathcal{T} :

$$\phi_{\mathcal{T}}(\lambda) = \max_{t \in \mathcal{T}} \lambda^T t : \mathbf{R}^K \to \mathbf{R}$$

Observe that whenever $\lambda \in \mathbf{R}_+^K$ and H are such that

$$(B - H^T A)^T (B - H^T A) \preceq \sum_k \lambda_k R_k, \tag{4.11}$$

for $x \in \mathcal{X}$ it holds

$$||Bx - H^T Ax||_2^2 \le \phi_{\mathcal{T}}(\lambda). \tag{4.12}$$

Indeed, in the case of (4.11) and with $y \in \mathcal{X}$, there exists $t \in \mathcal{T}$ such that $y^T R_k y \leq t_k$ for all t, and consequently the vector \bar{t} with the entries $\bar{t}_k = y^T R_k y$ also belongs to \mathcal{T} , whence

$$||Bx - H^T Ax||_2^2 = ||Bx - H^T Ax||_2^2 \le \sum_k \lambda_k x^T R_k x = \lambda^T \overline{t} \le \phi_{\mathcal{T}}(\lambda),$$

which combines with (4.10) to imply (4.12).

From (4.12) it follows that if H and $\lambda \ge 0$ are linked by (4.11), then

$$\operatorname{Risk}^{2}[\widehat{x}_{H}|\mathcal{X}] = \max_{x \in \mathcal{X}} \mathbf{E} \left\{ \|Bx - H^{T}(Ax + \xi)\|_{2}^{2} \right\}$$
$$= \operatorname{Tr}(H^{T}\Gamma H) + \max_{x \in \mathcal{X}} \|[B - H^{T}A]x\|_{2}^{2}$$
$$\leq \operatorname{Tr}(H^{T}\Gamma H) + \phi_{\mathcal{T}}(\lambda).$$

We see that the efficiently computable convex function

$$\widehat{R}(H) = \inf_{\lambda} \left\{ \operatorname{Tr}(H^T \Gamma H) + \phi_{\mathcal{T}}(\lambda) : (B - H^T A)^T (B - H^T A) \preceq \sum_k \lambda_k R_k, \lambda \ge 0 \right\}$$

(which clearly is well defined due to compactness of \mathcal{T} combined with $\sum_k R_k \succ 0$) is an upper bound on R(H).⁵⁶ We have arrived at the following result:

 $[\]overline{}^{56}$ It is well known that when K = 1 (i.e., \mathcal{X} is n ellipsoid), the above bounding scheme is exact: $R(\cdot) \equiv \widehat{R}(\cdot)$. For more complicated \mathcal{X} 's, $\widehat{R}(\cdot)$ could be larger than $R(\cdot)$, although the ratio

Proposition 4.4. In the situation of this Section, the risk of the "presumably good" linear estimate $\hat{x}_{H_*}(\omega) = H_*^T \omega$ yielded by an optimal solution (H_*, λ_*) to the (clearly solvable) convex optimization problem

$$\begin{array}{lll}
\operatorname{Opt} &=& \min_{H,\lambda} \left\{ \operatorname{Tr}(H^T \Gamma H) + \phi_{\mathcal{T}}(\lambda) : (B - H^T A)^T (B - H^T A) \preceq \sum_k \lambda_k R_k, \lambda \ge 0 \right\} \\
&=& \min_{H,\lambda} \left\{ \operatorname{Tr}(H^T \Gamma H) + \phi_{\mathcal{T}}(\lambda) : \left[\frac{\sum_k \lambda_k R_k}{B - H^T A} \middle| \frac{B^T - A^T H}{I_{\nu}} \right] \succeq 0, \lambda \ge 0 \right\} \\
\end{array} \tag{4.13}$$

is upper-bounded by $\sqrt{\text{Opt}}$.

4.2.2.1 Illustration: Recovering temperature distribution

Situation: A square steel plate was somehow heated at time 0 and left to cool, the temperature along the perimeter of the plate being all the time kept zero. At time t_1 , we measure the temperatures at m points of the plate, and want to recover the distribution of the temperature along the plate at another given time t_0 , $0 < t_0 < t_1$.

Physics, after suitable discretization of spatial variables, offers the following model of our situation. We represent the distribution of temperature at time t as $(2N-1) \times (2N-1)$ matrix $U(t) = [u_{ij}(t)]_{i,j=1}^{2N-1}$, where $u_{ij}(t)$ is the temperature, at time t, at the point

$$P_{ij} = (p_i, p_j), p_k = k/N - 1, \quad 1 \le i, j \le 2N - 1$$

of the plate (in our model, this plate occupies the square $S = \{(p,q) : |p| \le 1, |q| \le 1\}$). Here positive integer N is responsible for spatial discretization.

For $1 \le k \le 2N - 1$, let us specify functions $\phi_k(s)$ on the segment $-1 \le s \le 1$ as follows:

$$\phi_{2\ell-1}(s) = c_{2\ell-1}\cos(\omega_{2\ell-1}s), \ \phi_{2\ell}(s) = c_{2\ell}\sin(\omega_{2\ell}s), \ \omega_{2\ell-1} = (\ell-1/2)\pi, \ \omega_{2\ell} = \ell\pi,$$

where c_k are readily given by the normalization condition $\sum_{i=1}^{2N-1} \phi_k^2(p_i) = 1$; note that $\phi_k(\pm 1) = 0$. It is immediately seen that the matrices

$$\Phi^{k\ell} = [\phi_k(p_i)\phi_\ell(p_j)]_{i,j=1}^{2N-1}, \ 1 \le k, \ell \le 2N-1$$

form an orthonormal basis in the space of $(2N-1) \times (2N-1)$ matrices, so that we can write

$$U(t) = \sum_{k,\ell \le 2N-1} x_{k\ell}(t) \Phi^{k\ell}$$

The advantage of representing temperature fields in the basis $\{\Phi^{k\ell}\}_{k,\ell\leq 2N-1}$ stems from the fact that in this basis the heat equation governing evolution of the temperature distribution in time becomes extremely simple, just

$$\frac{d}{dt}x_{k\ell}(t) = -(\omega_k^2 + \omega_\ell^2)x_{k\ell}(t) \Rightarrow x_{k\ell}(t) = \exp\{-(\omega_k^2 + \omega_\ell^2)t\}x_{k\ell} \quad {}^{57}$$

 $\widehat{R}(\cdot)/R(\cdot)$ is bounded by $O(\log(K))$, see Section 4.2.5.

Now we can convert the situation into the one considered in our general estimation scheme, namely, as follows:

• We select somehow the discretization parameter N and treat $x = \{x_{k\ell}(0), 1 \le k, \ell \le 2N - 1\}$ as the signal underlying our observations.

In every potential application, we can safely upper-bound the magnitudes of the initial temperatures and thus the magnitude of x, say, by a constraint of the form

$$\sum_{k,\ell} x_{k\ell}^2(0) \leq R^2$$

with properly selected R, which allows to specify the domain \mathcal{X} of the signal as the Euclidean ball:

$$\mathcal{X} = \{ x \in \mathbf{R}^{(2N-1) \times (2N-1)} : \|x\|_2^2 \le R^2 \}.$$
(4.14)

• Let the measurements of the temperature at time t_1 be taken along the points $P_{i(\nu),j(\nu)}$, $1 \leq \nu \leq m^{58}$, and let them be affected by $\mathcal{N}(0,\sigma^2 I_m)$ -noise, so that our observation is

$$\omega = A(x) + \xi, \ \xi \sim \mathcal{N}(0, \sigma^2 I_m),$$

where $x \mapsto A(x)$ is the linear mapping from $\mathbf{R}^{(2N-1)\times(2N-1)}$ into \mathbf{R}^m given by

$$[A(x)]_{\nu} = \sum_{k,\ell=1}^{2N-1} e^{-(\omega_k^2 + \omega_\ell^2)t_1} \phi_k(p_{i(\nu)}) \phi_\ell(p_{j(\nu)}) x_{k\ell}(0).$$
(4.15)

• What we want to recover, are the temperatures at time t_0 taken along some grid, say, the square $(2K-1) \times (2K-1)$ grid $\{Q_{ij} = (r_i, r_j), 1 \le i, j \le 2K-1\}$, where $r_i = i/K - 1, 1 \le i \le 2K - 1$. In other words, we want to recover B(x), where the linear mapping $x \mapsto B(x)$ from $\mathbf{R}^{(2N-1)\times(2N-1)}$ into $\mathbf{R}^{(2K-1)\times(2K-1)}$ is given by

$$[B(x)]_{ij} = \sum_{k,\ell=1}^{2N-1} e^{-(\omega_k^2 + \omega_\ell^2)t_0} \phi_k(r_i) \phi_\ell(r_j) x_{k\ell}(0).$$

Ill-posedness. Our problem is a typical example of *ill-posed inverse problem*, where one wants to recover a past state of dynamical system converging exponentially fast to equilibrium and thus "forgetting rapidly" its past. More specifically, in our situation ill-posedness stems from the fact that, as is clearly seen from (4.15), contributions of "high frequency" (i.e., with large $\omega_k^2 + \omega_\ell^2$) components $x_{k\ell}(0)$ of the signal to A(x) decrease exponentially fast, with high decay rate, as t_1 grows. As a result, high frequency components $x_{k\ell}(0)$ are impossible to recover from noisy observations of A(x), unless the corresponding time instant t_1 is very small. As a kind of compensation, contributions of high frequency components $x_{k\ell}(0)$ to B(x) are very small, provided t_0 is not too small, implying that there is no necessity to recover well high frequency components, provided they are not huge. Our linear estimate, roughly speaking, seeks for the best tradeoff between these two opposite

 $^{^{58}}$ the construction can be easily extended to allow for measurement points outside of the grid $\{P_{ij}\}.$



Figure 4.1: True distribution of temperature $U_* = B(x)$ at time $t_0 = 0.01$ (left) along with its recovery \hat{U} via optimal linear estimate (center) and the "naive" recovery \tilde{U} (right).

phenomena, utilizing (4.14) as the source of upper bounds on the magnitudes of high frequency components of the signal.

Numerical results. In the experiment to be reported, we used N = 32, m = 100, K = 6, $t_0 = 0.01$, $t_1 = 0.03$ (i.e., temperature is measured at time 0.03 at 100 points selected at random on 63×63 square grid, and we want to recover the temperatures at time 0.01 along 11×11 square grid). We used R = 15, that is,

$$\mathcal{X} = \{ [x_{k\ell}]_{k,\ell=1}^{63} : \sum_{k,\ell} x_{k\ell}^2 \le 225 \},\$$

and $\sigma = 0.001$.

Under the circumstances, the risk of the best linear estimate turns out to be 0.3968. Figure 4.1 shows a sample temperature distribution $B(x) = U_*(t_0)$ at time t_0 stemming from a randomly selected signal $x \in \mathcal{X}$ along with the recovery $\widehat{U}(t_0)$ of U_* by the optimal linear estimate and the naive "least squares" recovery $\widetilde{U}(t_0)$ of U_* . The latter is defined as $B(x_*)$, where x_* is the least squares recovery of signal underlying observation ω :

$$x = x_*(\omega) := \underset{x}{\operatorname{argmin}} \|A(x) - \omega\|_2.$$

Pay attention to the dramatic difference in performances of the "naive least squares" and the optimal linear estimate

4.2.3 Near-optimality of \hat{x}_{H_*}

Proposition 4.5. The efficiently computable linear estimate $\hat{x}_{H_*}(\omega) = H_*^T \omega$ yielded by an optimal solution to the optimization problem (4.13) is nearly optimal in terms of its risk:

$$\operatorname{Risk}[\widehat{x}_{H_*}|\mathcal{X}] \le \sqrt{\operatorname{Opt}} \le 64\sqrt{(3+18\ln 2)(3\ln K+18\ln 2)\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}]}, \quad (4.16)$$

where the minimax optimal risk $\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}]$ is given by (4.8).

277

For proof, see Section 4.10.5.

4.2.3.1 Relaxing the symmetry requirement

Sets \mathcal{X} of the form (4.6) – we called them ellitopes – are symmetric w.r.t. the origin convex compacts of special structure. This structure is rather flexible, but the symmetry is "built in." We are about to demonstrate that, to some extent, the symmetry requirement can be somehow relaxed. Specifically, assume instead of (4.6) that the convex compact set \mathcal{X} known to contain the signals x underlying observations (4.5) can be "sandwiched" by two known to us and similar to each other, with coefficient $\alpha \geq 1$, ellitopes:

$$\underbrace{\left\{x \in \mathbf{R}^{n} : \exists (y \in \mathbf{R}^{\bar{n}}, t \in \mathcal{T}) : x = Py \& y^{T} S_{k} y \leq t_{k}, 1 \leq k \leq K\right\}}_{\mathcal{X}} \subset \mathcal{X} \subset \alpha \underline{\mathcal{X}},$$

with S_k and \mathcal{T} possessing the properties postulated in Section 4.2.1. Let Opt and H_* be the optimal value and optimal solution of the optimization problem (4.13) associated with the data $S_1, ..., S_K, \mathcal{T}$ and matrices $\overline{A} = AP, \overline{B} = BP$ in the role of A, B, respectively. It is immediately seen that the risk $\operatorname{Risk}[\widehat{x}_{H_*}|\mathcal{X}]$ of the linear estimate $\widehat{x}_{H_*}(\omega)$ is at most $\alpha\sqrt{\operatorname{Opt}}$. On the other hand, we have $\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}] \leq \operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}]$, and by Proposition 4.5 also $\sqrt{\operatorname{Opt}} \leq O(1)\sqrt{\ln(2K)}\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}]$. Taken together, these relations imply that

$$\operatorname{Risk}[\widehat{x}_{H^*}|\mathcal{X}] \le O(1)\alpha \sqrt{\ln(2K)}\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}].$$
(4.17)

In other words, as far as the "level of nonoptimality" of efficiently computable linear estimates is concerned, signal sets \mathcal{X} which can be approximated by ellitopes within a factor α of order of 1 are nearly as good as the ellitopes. To give an example: it is known that whenever the intersection \mathcal{X} of K elliptic cylinders $\{x : (x - c_k)^T S_k(x - c_k) \leq 1\}$, $S_k \succeq 0$, concentric or not, is bounded and has a nonempty interior, \mathcal{X} can be approximated by an ellipsoid within the factor $\alpha = K + 2\sqrt{K}$ ⁵⁹. Assuming w.l.o.g. that the approximating ellipsoid is centered at the origin, the level of nonoptimality of a linear estimate is bounded by (4.17) with O(1)K in the role of α . Note that bound (4.17) rapidly deteriorates when α grows, and this phenomenon to some extent "reflects the reality." For example, a perfect simplex \mathcal{X} inscribed into the unit sphere in \mathbb{R}^n is in-between two centered at the origin Euclidean balls with the ratio of radii equal to n (i.e. $\alpha = n$). It is immediately seen that with A = B = I, $\Gamma = \sigma^2 I$, in the range $\sigma \leq n\sigma^2 \leq 1$ of values of n and σ , we have

$$\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}] \approx \sqrt{\sigma}, \ \operatorname{Risk}_{\operatorname{opt}}[\widehat{x}_{H_*}|\mathcal{X}] = O(1)\sqrt{n\sigma},$$

with \approx meaning "up to logarithmic in n/σ factor." In other words, for large $n\sigma$ linear estimates indeed are significantly (albeit not to the full extent of (4.17))

$$\{x: (x-\bar{x})^T F''(\bar{x})(x-\bar{x}) \le 1\} \subset \mathcal{X} \subset \{x: (x-\bar{x})^T F''(\bar{x})(x-\bar{x}) \le [K+2\sqrt{K}]^2\}$$

 $[\]overline{S^9 \text{specifically, setting } F(x) = -\sum_{k=1}^K \ln(1 - (x - c_k)^T S_k(x - c_k)) : \text{int } \mathcal{X} \to \mathbf{R} \text{ and denoting by } \bar{x} \text{ the analytic center } \operatorname{argmin}_{x \in \text{int } \mathcal{X}} F(x), \text{ one has}$

outperformed by nonlinear ones.

Another "bad for linear estimates" situation suggested by (4.16) is the one where the description (4.6) of \mathcal{X} , albeit possible, requires a huge value of K. Here again (4.16) reflects to some extent the reality: when \mathcal{X} is the unit ℓ_1 ball in \mathbf{R}^n , (4.6) takes place with $K = 2^{n-1}$; consequently, the factor at $\operatorname{Risk}_{opt}[\mathcal{X}]$ in the right hand side of (4.16) becomes at least \sqrt{n} . On the other hand, with A = B = I, $\Gamma = \sigma^2 I$, in the range $\sigma \leq n\sigma^2 \leq 1$ of values of n, σ , the risks $\operatorname{Risk}_{opt}[\mathcal{X}]$, $\operatorname{Risk}_{opt}[\widehat{x}_{H_*}|\mathcal{X}]$ are basically the same as in the case of \mathcal{X} being the perfect simplex inscribed into the unit sphere in \mathbf{R}^n , and linear estimates indeed are "heavily non-optimal" when $n\sigma$ is large.

4.2.4 Numerical illustration

The "non-optimality factor" θ in the upper bound $\sqrt{\text{Opt}} \leq \theta \text{Risk}_{\text{opt}}[\mathcal{X}]$ from Proposition 4.5, while logarithmic, seems to be unpleasantly large. On a closest inspection, one can get less conservative bounds on non-optimality factors. Omitting the details, here are some numerical results. In the six experiments to be reported, we used $n = m = \nu = 32$ and $\Gamma = \sigma^2 I_m$. In the first triple of experiments, \mathcal{X} was the ellipsoid

$$X = \{ x \in \mathbf{R}^{32} : \sum_{j=1}^{32} j^2 x_j^2 \le 1 \},$$

that is, P was the identity, K = 1, $S_1 = \sum_{j=1}^{32} j^2 e_j e_j^T$ (e_j were basic orths), and $\mathcal{T} = [0, 1]$. In the second triple of experiments, \mathcal{X} was the box circumscribed around the above ellipsoid:

$$X = \{x \in \mathbf{R}^{32} : j|x_j| \le 1, \ 1 \le j \le 32\} \\ [P = I, K = 32, S_k = k^2 e_k e_k^T, k \le K, \mathcal{T} = [0, 1]^K]$$

In all six experiments, B was the identity, and A was a common for all experiments randomly rotated matrix with singular values λ_j , $1 \le j \le 32$, forming a geometric progression, with $\lambda_1 = 1$ and $\lambda_{32} = 0.01$. Experiments in a triple differed by the values of σ (0.01,0.001,0.0001).

The results of the experiments are presented in Table 4.1, where, as above, $\sqrt{\text{Opt}}$ is the given by (4.13) upper bound on the risk $\text{Risk}[\hat{x}_{H_*}|X]$ of recovering $Bx = x, x \in X$, by the linear estimate yielded by (4.9), (4.13), LwB is the lower bound on $\text{Risk}_{\text{opt}}[X]$ built as explained above, and the numbers in the last column are (conservative estimates of the) "levels of nonoptimality" of the linear estimates.

4.2.5 Byproduct on semidefinite relaxation

A byproduct of our main observation (Section 4.2.3) we are about to present has nothing to do with statistics; it relates to the quality of the standard semidefinite relaxation. Specifically, given a quadratic from $x^T C x$ and an ellitope \mathcal{X} represented by (4.6), consider the problem

$$Opt_* = \max_{x \in \mathcal{X}} x^T C x = \max_{y \in \bar{X}} y^T P^T C P y.$$
(4.18)

##	X	σ	√Opt	LwB	$\sqrt{\text{Opt}/LwB}$
1	ellipsoid	1.0e-2	0.288	0.153	1.88
2	ellipsoid	1.0e-3	0.103	0.060	1.71
3	ellipsoid	1.0e-4	0.019	0.018	1.06
4	box	1.0e-2	0.698	0.231	3.02
5	box	1.0e-3	0.163	0.082	2.00
6	box	1.0e-4	0.021	0.020	1.06

Table 4.1: Performance of linear estimates (4.9), (4.13), m = n = 32, B = I.

This problem can be NP-hard (this is already so when \mathcal{X} is the unit box and C is positive semidefinite); however, Opt admits an efficiently computable upper bound given by *semidefinite relaxation* as follows: whenever $\lambda \geq 0$ is such that

$$P^T C P \preceq \sum_{k=1}^K \lambda_k S_k,$$

for $y \in \overline{X}$ we clearly have

$$[Py]^T C Py \le \sum_k \lambda_k y^T S_k y \le \phi_{\mathcal{T}}(\lambda)$$

due to the fact that the vector with the entries $y^T S_k y$, $1 \le k \le K$, belongs to \mathcal{T} . As a result, the efficiently computable quantity

$$Opt = \min_{\lambda} \left\{ \phi_{\mathcal{T}}(\lambda) : \lambda \ge 0, P^T C P \preceq \sum_k \lambda_k S_k \right\}$$
(4.19)

is an upper bound on Opt_{*}. We have the following

Proposition 4.6. Let C be a symmetric $n \times n$ matrix and \mathcal{X} be given by ellitopic representation (4.6), and let Opt_* and Opt be given by (4.18) and(4.19). Then

$$\frac{\operatorname{Opt}}{\operatorname{Bln}(\sqrt{3}K)} \le \operatorname{Opt}_* \le \operatorname{Opt}.$$
(4.20)

For proof, see Section 4.10.2.

4.3 FROM ELLITOPES TO SPECTRATOPES

So far, the domains of signals we dealt with were ellitopes. In this section we demonstrate that basically all our constructions and results can be extended onto a much wider family of signal domains, namely, *spectratopes*.

280

LECTURE 4

4.3.1 Spectratopes: definition and examples

We call a set $\mathcal{X} \subset \mathbf{R}^n$ a basic spectratope, if it admits simple spectratopic representation – representation of the form

$$\mathcal{X} = \left\{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \leq t_k I_{d_k}, 1 \leq k \leq K \right\}$$
(4.21)

where

- S.1. $R_k[x] = \sum_{i=1}^n x_i R^{ki}$ are symmetric $d_k \times d_k$ matrices linearly depending on $x \in \mathbf{R}^n$ (i.e., "matrix coefficients" R^{ki} belong to \mathbf{S}^n)
- S.2. $\mathcal{T} \in \mathbf{R}_{+}^{K}$ is the set with the same properties as in the definition of an ellitope, that is, \mathcal{T} is a convex compact subset of \mathbf{R}_{+}^{K} which contains a positive vector and is monotone:

$$0 \le t' \le t \in \mathcal{T} \Rightarrow t' \in \mathcal{T}$$

S.3. Whenever $x \neq 0$, it holds $R_k[x] \neq 0$ for at least one $k \leq K$.

An immediate observation (check it!) is as follows:

Remark 4.7. By Schur Complement Lemma, the set (4.21) given by data satisfying S.1-2 can be represented as

$$\mathcal{X} = \left\{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : \left[\begin{array}{c|c} t_k I_{d_k} & R_k[x] \\ \hline R_k[x] & I_{d_k} \end{array} \right] \succeq 0, \ k \le K \right\}$$

By the latter representation, \mathcal{X} is nonempty, closed, convex, symmetric w.r.t. the origin and contains a neighbourhood of the origin. This set is bounded if and only if the data, in addition to S.1-2, satisfies S.3.

A spectratope $\mathcal{X} \subset \mathbf{R}^{\nu}$ is a set represented as linear image of a basic spectratope:

$$\mathcal{X} = \{ x \in \mathbf{R}^{\nu} : \exists (y \in \mathbf{R}^n, t \in \mathcal{T}) : x = Py, R_k^2[y] \leq t_k I_{d_k}, 1 \leq k \leq K \}, \quad (4.22)$$

where P is a $\nu \times n$ matrix, and $R_k[\cdot]$, \mathcal{T} are as in S.1-3.

We associate with a basic spectratope (4.21), S.1-3 the following entities:

1. The size

$$D = \sum_{k=1}^{K} d_k;$$

2. Linear mappings

$$Q \mapsto \mathcal{R}_k[Q] = \sum_{i,j} Q_{ij} R^{ki} R^{kj} : \mathbf{S}^n \to \mathbf{S}^{d_k}$$
(4.23)

As is immediately seen, we have

$$\mathcal{R}_k[xx^T] \equiv R_k^2[x], \tag{4.24}$$

implying that $\mathcal{R}_k[Q] \succeq 0$ whenever $Q \succeq 0$, whence $\mathcal{R}_k[\cdot]$ is \succeq -monotone:

$$Q' \succeq Q \Rightarrow \mathcal{R}_k[Q'] \succeq \mathcal{R}_k[Q].$$
 (4.25)

Besides this, we have

$$Q \succeq 0 \Rightarrow \mathbf{E}_{\xi \sim \mathcal{N}(0,Q)} \{ R_k^2[\xi] \} = \mathbf{E}_{\xi \sim \mathcal{N}(0,Q)} \{ \mathcal{R}_k[\xi \xi^T] \} = \mathcal{R}_k[Q], \qquad (4.26)$$

where the first equality is given by (4.24).

3. Linear mappings $\Lambda_k \mapsto \mathcal{R}_k^*[\Lambda_k] : \mathbf{S}^{d_k} \to \mathbf{S}^n$ given by

$$[\mathcal{R}_{k}^{*}[\Lambda_{k}]]_{ij} = \frac{1}{2} \operatorname{Tr}(\Lambda_{k}[R^{ki}R^{kj} + R^{kj}R^{ki}]), \ 1 \le i, j \le n.$$
(4.27)

It is immediately seen that $\mathcal{R}_{k}^{*}[\cdot]$ is the conjugate of $\mathcal{R}_{k}[\cdot]$:

$$\langle \Lambda_k, \mathcal{R}_k[Q] \rangle_F = \operatorname{Tr}(\Lambda_k \mathcal{R}_k[Q]) = \operatorname{Tr}(\mathcal{R}_k^*[\Lambda_k]Q) = \langle \mathcal{R}_k^*[\Lambda_k], Q \rangle_F,$$
 (4.28)

where $\langle A, B \rangle_F = \text{Tr}(AB)$ is the Frobenius inner product of symmetric matrices. Besides this, we have

$$\Lambda_k \succeq 0 \Rightarrow \mathcal{R}_k^*[\Lambda_k] \succeq 0. \tag{4.29}$$

Indeed, $\mathcal{R}_k^*[\Lambda_k]$ is linear in Λ_k , so that it suffices to verify (4.29) for dyadic matrices $\Lambda_k = f f^T$; for such a Λ_k , (4.27) reads

$$(\mathcal{R}_k^*[ff^T])_{ij} = [R^{ki}f]^T [R^{kj}f],$$

that is, $\mathcal{R}_k^*[ff^T]$ is a Gram matrix and as such is $\succeq 0$. Another way to arrive at (4.29) is to note that when $\Lambda_k \succeq 0$ and $Q = xx^T$, the first quantity in (4.28) is nonnegative by (4.24), and therefore (4.28) states that $x^T \mathcal{R}_k^*[\Lambda_k] x \ge 0$ for every x, implying $\mathcal{R}_k^*[\Lambda_k] \succeq 0$.

x, implying $\mathcal{R}_{k}^{*}[\Lambda_{k}] \succeq 0$. 4. The linear space $\Lambda^{K} = \mathbf{S}^{d_{1}} \times ... \times \mathbf{S}^{d_{K}}$ of all ordered collections $\Lambda = \{\Lambda_{k} \in \mathbf{S}^{d_{k}}\}_{k \leq K}$ along with the linear mapping

$$\Lambda \mapsto \lambda[\Lambda] := [\operatorname{Tr}(\Lambda_1); ...; \operatorname{Tr}(\Lambda_K)] : \Lambda^K \to \mathbf{R}^K.$$
(4.30)

4.3.1.1 Examples of spectratopes

Example 1: Ellitopes. Every ellitope

$$\mathcal{X} = \{ x \in \mathbf{R}^{\nu} : \exists (y \in \mathbf{R}^n, t \in \mathcal{T}) : x = Py, y^T S_k y \le t_k, k \le K \} \\ [S_k \succeq 0, \sum_k S_k \succ 0]$$

is a spectratope as well. Indeed, let $S_k = \sum_{j=1}^{r_k} s_{kj} s_{kj}^T$, $r_k = \text{Rank}(S_k)$, be a dyadic representation of the positive semidefinite matrix S_k , so that

$$y^T S_k y = \sum_j (s_{kj}^T y)^2 \; \forall y,$$

and let

$$\widehat{\mathcal{T}} = \{\{t_{kj} \ge 0, 1 \le j \le r_k, 1 \le k \le K\} : \exists t \in \mathcal{T} : \sum_j t_{kj} \le t_k\},\\ R_{kj}[y] = s_{kj}^T y \in \mathbf{S}^1 = \mathbf{R}.$$

We clearly have

$$\mathcal{X} = \{ x \in \mathbf{R}^{\nu} : \exists (\{t_{kj}\} \in \widehat{\mathcal{T}}, y) : x = Py, R_{kj}^2[y] \leq t_{kj}I_1 \,\forall k, j \}$$

and the right hand side is a legitimate spectratopic representation of \mathcal{X} .

Example 2: "Matrix box." Let *L* be a positive definite $d \times d$ matrix. Then the "matrix box"

$$\begin{aligned} \mathcal{X} &= \{ X \in \mathbf{S}^d : -L \preceq X \preceq L \} = \{ X \in \mathbf{S}^d : -I_d \preceq L^{-1/2} X L^{-1/2} \preceq I_d \} \\ &= \{ X \in \mathbf{S}^d : R^2[X] := [L^{-1/2} X L^{-1/2}]^2 \preceq I_d \} \end{aligned}$$

is a basic spectratope (augment $R_1[\cdot] := R[\cdot]$ with K = 1, $\mathcal{T} = [0, 1]$). As a result, a *bounded* set $\mathcal{X} \subset \mathbf{R}^{\nu}$ given by a system of "two-sided" Linear Matrix Inequalities, specifically,

$$\mathcal{X} = \{ x \in \mathbf{R}^{\nu} : \exists t \in \mathcal{T} : -\sqrt{t_k} L_k \preceq S_k[x] \preceq \sqrt{t_k} L_k, k \leq K \}$$

where $S_k[x]$ are symmetric $d_k \times d_k$ matrices linearly depending on $x, L_k \succ 0$ and \mathcal{T} satisfies S.2, is a basic spectratope:

$$\mathcal{X} = \{ x \in \mathbf{R}^{\nu} : \exists t \in \mathcal{T} : R_k^2[x] \le t_k I_{d_k}, \, k \le K \} \qquad [R_k[x] = L_k^{-1/2} S_k[x] L_k^{-1/2}]$$

Same as ellitopes, spectratopes admit fully algorithmic calculus, see Section 4.8.

4.3.2 Semidefinite relaxation on spectratopes

Now let us extend to our current situation Proposition 4.6. The extension reads as follows:

Proposition 4.8. Let C be a symmetric $n \times n$ matrix and \mathcal{X} be given by spectratopic representation

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists y \in \mathbf{R}^\mu, t \in \mathcal{T} : x = Py, R_k^2[y] \preceq t_k I_{d_k}, k \leq K \},$$
(4.31)

let

$$Opt = \max_{x \in \mathcal{X}} x^T C x$$

and

$$\begin{aligned}
\operatorname{Opt}_{*} &= \min_{\Lambda = \{\Lambda_{k}\}_{k \leq K}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) : \Lambda_{k} \succeq 0, P^{T}CP \preceq \sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}] \right\} \\
& [\lambda[\Lambda] = [\operatorname{Tr}(\Lambda_{1}); ...; \operatorname{Tr}(\Lambda_{K})]]
\end{aligned} \tag{4.32}$$

Then (4.32) is solvable, and

$$Opt \le Opt_* \le 2\max[\ln(2D), 1]Opt, \ D = \sum_k d_k.$$
(4.33)

Proof. In what follows we restrict ourselves with proving the easy and instructive part of Proposition, namely, the left inequality in (4.33); remaining claims will be proved in Section 4.10.3.

The left inequality in (4.33) is readily given by the following

Lemma 4.9. Let \mathcal{X} be spectratope (4.31) and $Q \in \mathbf{S}^n$. Whenever $\Lambda_k \in \mathbf{S}^{d_k}_+$ and $\tau \geq 0$ satisfy

$$P^T Q P \preceq \sum_k \mathcal{R}_k^*[\Lambda_k],$$

we have

$$x \in \mathcal{X} \Rightarrow x^T Q x \le \phi_{\mathcal{T}}(\lambda[\Lambda]), \ \lambda[\Lambda] = [\operatorname{Tr}(\Lambda_1); ...; \operatorname{Tr}(\Lambda_K)].$$

Proof of Lemma: Let $x \in \mathcal{X}$, so that for some $t \in \mathcal{T}$ and y it holds

$$x = Py, R_k^2[y] \preceq t_k I_{d_k} \ \forall k \le K$$

Consequently,

$$\begin{aligned} x^{T}Qx &= y^{T}P^{T}QPy \leq y^{T}\sum_{k}\mathcal{R}_{k}^{*}[\Lambda_{k}]y = \sum_{k}\operatorname{Tr}(\mathcal{R}_{k}^{*}[\Lambda_{k}][yy^{T}]) \\ &= \sum_{k}\operatorname{Tr}(\Lambda_{k}\mathcal{R}_{k}[yy^{T}]) \text{ [by (4.28)]} \\ &= \sum_{k}\operatorname{Tr}(\Lambda_{k}R_{k}^{2}[y]) \text{ [by (4.24)]} \\ &\leq \sum_{k}t_{k}\operatorname{Tr}(\Lambda_{k}I_{d_{k}}) \text{ [since } \Lambda_{k} \succeq 0 \text{ and } R_{k}^{2}[y] \preceq t_{k}I_{d_{k}}] \\ &\leq \phi_{\mathcal{T}}(\lambda[\Lambda]) \end{aligned}$$

4.3.3 Linear estimates beyond ellitopic signal sets and $\|\cdot\|_2$ -risk

In Section 4.2, we have developed a computationally efficient scheme for building "presumably good" linear estimates of the linear image Bx of unknown signal x known to belong to a given ellitope \mathcal{X} in the case when the (squared) risk is defined as the worst, w.r.t. $x \in \mathcal{X}$, expected squared Euclidean norm $\|\cdot\|_2$. We are about to extend these results to the case when \mathcal{X} is a spectratope, and the norm used to measure the recovery error, while not being completely arbitrarily, is not necessarily $\|\cdot\|_2$, with the ultimate goal to demonstrate that the resulting linear estimates are not just "presumably good," but possess near-optimality properties completely similar to those stated in Propositions 4.5. Besides this, in what follows we somehow relax our assumptions on observation noise.

4.3.3.1 Situation and goal

In what follows we consider the problem of recovering the image $Bx \in \mathbf{R}^{\nu}$ of a signal $x \in \mathbf{R}^n$ known to belong to a given spectratope

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \leq t_k I_{d_k}, 1 \leq k \leq K \}$$

$$(4.34)$$

from noisy observation

$$\omega = Ax + \xi, \tag{4.35}$$

where A is a known $m \times n$ matrix, and ξ is random observation noise.

Observation noise. In typical signal processing applications, the distribution of noise is fixed and is part of the data of the estimation problem. In order to cover some applications (e.g., the one in Section 4.3.3.7), we allow for "ambiguous" noise distributions; all we know is that this distribution belongs to a family \mathcal{P} of Borel probability distributions on \mathbf{R}^m associated with a given convex compact subset Π of the interior of the cone \mathbf{S}^m_+ of positive semidefinite $m \times m$ matrices, "association" meaning that the matrix of second moments of every distribution $P \in \mathcal{P}$ is \succeq -dominated by a matrix from Π :

$$P \in \mathcal{P} \Rightarrow \exists Q \in \Pi : \operatorname{Var}[P] := \mathbf{E}_{\xi \sim P}\{\xi\xi^T\} \preceq Q.$$

$$(4.36)$$

Actual distribution of noise in (4.35) is somehow selected from \mathcal{P} by nature (and

may, e.g., depend on x).

In the sequel, for a probability distribution P on \mathbb{R}^m we write $P \ll \Pi$ to express the fact that the matrix of second moments of P is \succeq -dominated by a matrix from Π :

$$\{P \ll \Pi\} \Leftrightarrow \{\exists \Theta \in \Pi : \operatorname{Var}[P] \preceq \Theta\}$$

Quantifying risk. Given Π and a norm $\|\cdot\|$ on \mathbf{R}^{ν} , we quantify the quality of a candidate estimate $\widehat{x}(\cdot): \mathbf{R}^m \to \mathbf{R}^{\nu}$ by its $(\Pi, \|\cdot\|)$ -risk on \mathcal{X} defined as

$$\operatorname{Risk}_{\Pi, \|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x \in \mathcal{X}, P \ll \Pi} \mathbf{E}_{\xi \sim P} \left\{ \|\widehat{x}(Ax + \xi) - Bx\| \right\}.$$
(4.37)

Goal. As before, our focus is on *linear estimates* – estimates of the form

$$\widehat{x}_H(\omega) = H^T \omega$$

given by $m \times \nu$ matrices H; our ultimate goal is to demonstrate that under some restrictions on the signal domain \mathcal{X} , "presumably good" linear estimate yielded by an optimal solution to an efficiently solvable convex optimization problem is near-optimal in terms of its risk among *all* estimates, linear and nonlinear alike.

4.3.3.2 Assumptions

Preliminaries: conjugate norms. Recall that a norm $\|\cdot\|$ on a Euclidean space \mathcal{E} , e.g., on \mathbf{R}^k , gives rise to its *conjugate* norm

$$||y||_* = \max_{x} \{ \langle y, x \rangle : ||x|| \le 1 \},$$

where $\langle \cdot, \cdot \rangle$ is the inner product in \mathcal{E} . Equivalently, $\|\cdot\|_*$ is the smallest norm such that

$$\langle x, y \rangle \le \|x\| \|y\|_* \ \forall x, y. \tag{4.38}$$

It is well known that taken twice, norm conjugation recovers the initial norm: $(\|\cdot\|_*)_*$ is exactly $\|\cdot\|$; in other words,

$$\|x\| = \max_{y} \{ \langle x, y \rangle : \|y\|_* \le 1 \}.$$

The standard examples are the conjugates to the standard ℓ_p -norms on $\mathcal{E} = \mathbf{R}^k$, $p \in [1, \infty]$; it turns out that

$$(\|\cdot\|_p)_* = \|\cdot\|_{p*},$$

where $p_* \in [1, \infty]$ is linked to $p \in [1, \infty]$ by the symmetric relation

$$\frac{1}{p} + \frac{1}{p_*} = 1,$$

so that $1_* = \infty$, $\infty_* = 1$, $2_* = 2$; the corresponding version of inequality (4.38) is called *Hölder inequality* – an extension of the Cauchy-Schwartz inequality dealing with the case $\|\cdot\| = \|\cdot\|_* = \|\cdot\|_2$.

Assumptions. From now on we make the following assumptions:

Assumption A: The unit ball \mathcal{B}_* of the norm $\|\cdot\|_*$ conjugate to the norm $\|\cdot\|$ participating in the formulation of our estimation problem is a spectratope:

$$\mathcal{B}_* = \{ z \in \mathbf{R}^{\nu} : \exists y \in \mathcal{Y} : z = My \},$$

$$\mathcal{Y} := \{ y \in \mathbf{R}^q : \exists r \in \mathcal{R} : S^2_{\ell}[y] \preceq r_{\ell} I_{f_{\ell}}, 1 \leq \ell \leq L \},$$
(4.39)

where the right hand side data are as required in a spectratopic representation.

Note that Assumption **A** is satisfied when $\|\cdot\| = \|\cdot\|_p$ with $p \in [1, 2]$: in this case,

$$\mathcal{B}_* = \{ u \in \mathbf{R}^{\nu} : ||u||_{p_*} \le 1 \}, \ p_* = \frac{p}{p-1} \in [2, \infty],$$

so that \mathcal{B}_* is an ellitope, see Section 4.2.1.1, and thus is a spectratope. Another potentially useful example of norm $\|\cdot\|$ which obeys Assumption **A** is the *nuclear norm* $\|V\|_{\mathrm{Sh},1}$ on the space $\mathbf{R}^{\nu} = \mathbf{R}^{p \times q}$ of $p \times q$ matrices – the sum of singular values of a matrix V; the conjugate norm is the spectral norm $\|\cdot\| = \|\cdot\|_{\mathrm{Sh},\infty}$ on $\mathbf{R}^{\nu} = \mathbf{R}^{p \times q}$, and the unit ball of the latter norm is a spectratope:

$$\{X \in \mathbf{R}^{p \times q} : \|X\| \le 1\} = \{X : \exists t \in \mathcal{T} = [0, 1] : R^2[X] \le tI_{p+q}\},\ R[X] = \left[\frac{|X^T|}{|X|}\right].$$

Besides Assumption \mathbf{A} , we make

Assumption B: The signal set \mathcal{X} is a basic spectratope:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \leq t_k I_{d_k}, 1 \leq k \leq K \},$$
(4.40)

where the right hand side data are as required in a spectratopic representation.

Note: Similar to ellitopic situation of Section 4.2.1.3, the in our context situation where the s9ignal set is a basic spectratope can be straightforwardly reduces to the one where \mathcal{X} is a *basic* spectratope.

In addition we make the following regularity assumption:

Assumption R: All matrices from Π are positive definite.

(4.41)

4.3.3.3 Building linear estimate

Let $H \in \mathbf{R}^{m \times \nu}$. We clearly have

$$\begin{aligned} \operatorname{Risk}_{\Pi, \|\cdot\|} [\widehat{x}_{H}(\cdot) | \mathcal{X}] &= \sup_{x \in X, P \lll \Pi} \mathbf{E}_{\xi \sim P} \left\{ \| [B - H^{T} A] x - H^{T} \xi \| \right\} \\ &\leq \sup_{inX} \| [B - H^{T} A] x \| + \sup_{P \lll \Pi} \mathbf{E}_{\xi \sim P} \left\{ \| H^{T} \xi \| \right\} \\ &= \| B - H^{T} A \|_{\mathcal{X}, \|\cdot\|} + \Psi_{\Pi}(H), \end{aligned}$$

where

$$\begin{aligned} \|V\|_{\mathcal{X},\|\cdot\|} &= \max_x \left\{ \|Vx\| : x \in \mathcal{X} \right\} : \mathbf{R}^{k \times n} \to \mathbf{R}, \\ \Psi_{\Pi}(H) &= \sup_{P \ll \Pi} \mathbf{E}_{\xi \sim P} \left\{ \|H^T \xi\| \right\}. \end{aligned}$$

Same as in Section 4.2.2, we need to derive efficiently computable convex upper bounds on the norm $\|\cdot\|_{\mathcal{X},\|\cdot\|}$ and the function Ψ_{Π} , which, while being convex, can be difficult to compute.

4.3.3.4 Upper-bounding $\|\cdot\|_{\mathcal{X},\|\cdot\|}$

With Assumptions \mathbf{A} , \mathbf{B} in force, consider the spectratope

with $U_i[\cdot]$ readily given by $R_k[\cdot]$ and $S_\ell[\cdot]$. Given a $\nu \times n$ matrix V and setting

$$W[V] = \frac{1}{2} \left[\begin{array}{c|c} & V^T M \\ \hline M^T V & \end{array} \right]$$

we have

$$\|V\|_{\mathcal{X},\|\cdot\|} = \max_{x\in\mathcal{X}} \|Vx\| = \max_{x\in\mathcal{X},z\in\mathcal{B}_*} z^T Vx = \max_{x\in\mathcal{X},y\in\mathcal{Y}} y^T M^T Vx = \max_{w\in\mathcal{Z}} w^T W[V]w.$$

Applying Proposition 4.8, we arrive at the following

Corollary 4.10. In the just defined situation, the efficiently computable convex function

$$\|V\|_{\mathcal{X},\|\cdot\|}^{+} = \min_{\Lambda,\Upsilon} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) : \\ \Lambda = \left\{ \Lambda_{k} \in \mathbf{S}_{+}^{d_{k}} \right\}_{k \leq K}, \Upsilon = \left\{ \Upsilon_{\ell} \in \mathbf{S}_{+}^{f_{\ell}} \right\}_{\ell \leq L}, \\ \left[\frac{\sum_{k} \mathcal{R}_{k}^{k} [\Lambda_{k}]}{\frac{1}{2} M^{T} V} \frac{1}{\sum_{\ell} \mathcal{S}_{\ell}^{k} [\Upsilon_{\ell}]}{\frac{1}{2} \sum_{\ell} \mathcal{S}_{\ell}^{k} [\Upsilon_{\ell}]} \right] \succeq 0 \right\} \\ \left[\phi_{\mathcal{T}}(\lambda) = \max_{t \in \mathcal{T}} \lambda^{T} t, \ \phi_{\mathcal{R}}(\lambda) = \max_{r \in \mathcal{R}} \lambda^{T} r, \ \lambda[\{\Xi_{1}, ..., \Xi_{N}\}] = [\operatorname{Tr}(\Xi_{1}); ...; \operatorname{Tr}(\Xi_{N})], \\ \left[\mathcal{R}_{k}^{*}[\Lambda_{k}]]_{ij} = \frac{1}{2} \operatorname{Tr}(\Lambda_{k}[\mathcal{R}_{k}^{ki} \mathcal{R}_{k}^{kj} + \mathcal{R}_{k}^{kj} \mathcal{R}_{k}^{ki}]), \ where \ R_{k}[x] = \sum_{i} x_{i} \mathcal{R}^{ki}, \\ \left[\mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]]_{ij} = \frac{1}{2} \operatorname{Tr}(\Upsilon_{\ell}[\mathcal{S}_{\ell}^{\ell i} \mathcal{S}_{\ell}^{\ell j} + \mathcal{S}_{\ell}^{\ell j} \mathcal{S}_{\ell}^{\ell i}]), \ where \ S_{\ell}[y] = \sum_{i} y_{i} \mathcal{S}^{\ell i}. \end{array} \right]$$

$$(4.43)$$

is a norm on $\mathbf{R}^{\nu \times n}$, and this norm is a tight upper bound on $\|\cdot\|_{\mathcal{X}, \|\cdot\|}$, namely,

$$\forall V \in \mathbf{R}^{\nu \times n} : \|V\|_{\mathcal{X}, \|\cdot\|} \le \|V\|^+_{\mathcal{X}, \|\cdot\|} \le 2 \max[\ln(2\mathcal{D}), 1] \|V\|_{\mathcal{X}, \|\cdot\|},$$

$$\mathcal{D} = \sum_k d_k + \sum_\ell f_\ell.$$
(4.44)

4.3.3.5 Upper-bounding $\Psi_{\Pi}(\cdot)$

We are about to derive an efficiently computable convex upper bound on the function Ψ_{Π} stemming from any norm obeying Assumption **B**. The underlying observation is as follows:

Lemma 4.11. Let V be a $m \times \nu$ matrix, $Q \in \mathbf{S}^m_+$, and P be a probability distribution on \mathbf{R}^m with $\operatorname{Var}[P] \preceq Q$. Let, further, $\|\cdot\|$ be a norm on \mathbf{R}^ν with the unit ball \mathcal{B}_* of the conjugate norm $\|\cdot\|_*$ given by (4.39). Finally, let $\Upsilon = {\Upsilon_{\ell} \in \mathbf{S}^{f_{\ell}}_+}_{\ell \leq L}$ and a matrix $\Theta \in \mathbf{S}^m$ satisfy the constraint

$$\begin{bmatrix} \Theta & \frac{1}{2}VM \\ \hline \frac{1}{2}M^TV^T & \sum_{\ell} \mathcal{S}^*_{\ell}[\Upsilon_{\ell}] \end{bmatrix} \succeq 0$$
(4.45)

(for notation, see (4.39), (4.43)). Then

$$\mathbf{E}_{\eta \sim P}\{\|V^T\eta\|\} \le \operatorname{Tr}(Q\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]).$$
(4.46)

Proof is immediate. In the case of (4.45), we have

Taking expectation of both sides of the resulting inequality w.r.t. distribution P of ξ and taking into account that $\operatorname{Tr}(\operatorname{Var}[P]\Theta) \leq \operatorname{Tr}(Q\Theta)$ due to $\Theta \succeq 0$ (by (4.45)) and $\operatorname{Var}[P] \preceq Q$, we get (4.46).

Note that when $P = \mathcal{N}(0, Q)$, the smallest possible upper bound on $\mathbf{E}_{\eta \sim P}\{\|V^T\eta\|\}$ which can be extracted from Lemma 4.11 (this bound is efficiently computable) is tight, see Lemma 4.17 below.

An immediate consequence is

Corollary 4.12. Let

$$\Gamma(\Theta) = \max_{Q \in \Pi} \operatorname{Tr}(Q\Theta) \tag{4.47}$$

and

$$\overline{\Psi}_{\Pi}(H) = \min_{\{\Upsilon_{\ell}\}_{\ell \leq L}, \Theta \in \mathbf{S}^{m}} \left\{ \Gamma(\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) : \Upsilon_{\ell} \succeq 0 \,\forall \ell, \, \left[\frac{\Theta}{\frac{1}{2}M^{T}H^{T}} \left| \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \right| \right] \succeq 0 \right\}$$

$$(4.48)$$

Then $\overline{\Psi}_{\Pi}(\cdot): \mathbf{R}^{m \times \nu} \to \mathbf{R}$ is efficiently computable convex upper bound on $\Psi_{\Pi}(\cdot)$.

Indeed, given Lemma 4.11, the only non-evident part of the corollary is that $\overline{\Psi}_{\Pi}(\cdot)$ is a well-defined real-valued function, which is readily given by Lemma 4.89.

Remark 4.13. When $\Upsilon = {\Upsilon_{\ell}}_{\ell \leq L}$, Θ is a feasible solution to the right hand side problem in (4.48) and s > 0, the pair $\Upsilon' = {s\Upsilon_{\ell}}_{\ell \leq L}$, $\Theta' = s^{-1}\Theta$ also is a feasible solution; since $\phi_{\mathcal{R}}(\cdot)$ and $\Gamma(\cdot)$ are positive homogeneous of degree 1, we conclude that $\overline{\Psi}_{\Pi}$ is in fact the infimum of the function

$$2\sqrt{\Gamma(\Theta)\phi_{\mathcal{R}}(\lambda[\Upsilon])} = \inf_{\theta>0} \left[s^{-1}\Gamma(\Theta) + s\phi_{\mathcal{R}}(\lambda[\Upsilon])\right]$$

over Υ, Θ satisfying the constraints of the problem (4.48).

In addition, for every feasible solution $\Upsilon = {\Upsilon_{\ell}}_{\ell \leq L}$, Θ to the problem (4.48) with $\mathcal{M}[\Upsilon] := \sum_{\ell} \mathcal{S}_{\ell}^*[\Upsilon_{\ell}] \succ 0$, the pair Υ , $\widehat{\Theta} = \frac{1}{4} H M \mathcal{M}^{-1}[\Upsilon] M^T H^T$ is feasible for the problem as well and $0 \preceq \widehat{\Theta} \preceq \Theta$ (Schur Complement Lemma), so that $\Gamma(\widehat{\Theta}) \leq \Gamma(\Theta)$. As a result,

$$\overline{\Psi}_{\Pi}(H) = \inf_{\Upsilon} \left\{ \begin{array}{c} \frac{1}{4} \Gamma(HM\mathcal{M}^{-1}[\Upsilon]M^{T}H^{T}) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) :\\ \Upsilon = \{\Upsilon_{\ell} \in \mathbf{S}_{+}^{f_{\ell}}\}_{\ell \leq L}, \mathcal{M}[\Upsilon] \succ 0 \end{array} \right\}.$$
(4.49)

Illustration. Consider the case when $||u|| = ||u||_p$ with $p \in [1, 2]$, and let us apply the just described scheme for upper-bounding Ψ_{Π} , assuming $\{Q\} \subset \Pi \subset \{S \in \mathbf{S}_{+}^{m} : S \leq Q\}$ for some given $Q \succ 0$, so that $\Gamma(\Theta) = \operatorname{Tr}(Q\Theta), \Theta \succeq 0$. The unit ball of the norm conjugate to $|| \cdot ||$, that is, the norm $|| \cdot ||_q$, $q = \frac{p}{p-1} \in [2, \infty]$, is the basic spectratope (in fact, ellitope)

$$\mathcal{B}_* = \{ y \in \mathbf{R}^{\mu} : \exists r \in \mathcal{R} := \{ \mathbf{R}^{\nu}_+ : \|r\|_{q/2} \le 1 \} : S^2_{\ell}[y] \le r_{\ell}, \ 1 \le \ell \le L = \nu \}, \\ S_{\ell}[y] = y_{\ell}.$$

As a result, Υ 's from Remark 4.13 are collections of ν positive semidefinite 1 × 1 matrices, and we can identify them with ν -dimensional nonnegative vectors v, resulting in $\lambda[\Upsilon] = v$ and $\mathcal{M}[\Upsilon] = \text{Diag}\{v\}$. Besides this, for nonnegative v we clearly have $\phi_{\mathcal{R}}(v) = \|v\|_{p/(2-p)}$. The optimization problem in (4.49) now reads

$$\overline{\Psi}_{\Pi}(H) = \inf_{v \in \mathbf{R}^{\nu}} \left\{ \frac{1}{4} \operatorname{Tr}(V \operatorname{Diag}^{-1}\{v\} V^{T}) + \|v\|_{p/(2-p)} : v > 0 \right\} \qquad [V = Q^{1/2} H]$$

After setting $a_{\ell} = \|\operatorname{Col}_{\ell}[V]\|_2$, (4.49) becomes

$$\overline{\Psi}_{\Pi}(H) = \inf_{\upsilon > 0} \left\{ \frac{1}{4} \sum_{\ell} \frac{a_{\ell}^2}{\upsilon_{\ell}} + \|\upsilon\|_{p/(2-p)} \right\}.$$

This results in $\overline{\Psi}_{\Pi}(H) = \|[a_1; ...; a_{\mu}]\|_p$. Recalling what a_{ℓ} and V are, we end up

with

$$\forall P, \operatorname{Var}[P] \leq Q : \\ \mathbf{E}_{\xi \sim P} \{ \| H^T \xi \| \} \leq \overline{\Psi}_{\Pi}(H) := \left\| \left[\| \operatorname{Row}_1[H^T Q^{1/2}] \|_2; \dots; \| \operatorname{Row}_{\nu}[H^T Q^{1/2}] \|_2 \right] \right\|_p,$$

4.3.3.6 Putting things together

An immediate summary of Corollaries 4.10, 4.12 is the following recipe for building "presumably good" linear estimate:

Proposition 4.14. In the situation of Section 4.3.3.1 and under Assumptions A, B, R (see Section 4.3.3.2) consider the convex optimization problem (for notation, see (4.43) and (4.47))

$$\begin{array}{lll}
\operatorname{Opt} &= & \min_{H,\Lambda,\Upsilon,\Upsilon,\Theta} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \Gamma(\Theta) : \\ & \Lambda = \{\Lambda_k \succeq 0, k \leq K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\}, \Upsilon' = \{\Upsilon'_\ell \succeq 0, \ell \leq L\}, \\ & \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T[B - H^TA]} \middle| \frac{1}{2}[B^T - A^TH]M} \right] \succeq 0, \\ & \left[\frac{\Theta}{\frac{1}{2}M^TH^T} \middle| \sum_\ell \mathcal{S}_\ell^*[\Upsilon'_\ell]} \right] \succeq 0 \end{array} \right\}$$

$$(4.50)$$

The problem is solvable, and the H-component H_* of its optimal solution yields linear estimate $\hat{x}_{H_*}(\omega) = H_*^T \omega$ such that

$$\operatorname{Risk}_{\Pi, \|\cdot\|} [\widehat{x}_{H_*}(\cdot) | \mathcal{X}] \le \operatorname{Opt.}$$

$$(4.51)$$

Note that the only claim in Proposition 4.14 which is not an immediate consequence of Corollaries 4.10, 4.12 is that problem (4.50) is solvable; this claim is readily given by the fact that the objective clearly is coercive on the feasible set (recall that $\Gamma(\Theta)$ is coercive on \mathbf{S}^m_+ due to $\Pi \subset \operatorname{int} \mathbf{S}^m_+$ and that $y \mapsto My$ is an onto mapping, since \mathcal{B}_* is full-dimensional).

4.3.3.7 Illustration: covariance matrix estimation

Suppose that we observe a sample

$$\eta^T = \{\eta_k = A\xi_k\}_{k \le T} \tag{4.52}$$

where A is a given $m \times n$ matrix, and $\xi_1, ..., \xi_T$ are sampled, independently of each other, from zero mean Gaussian distribution with unknown covariance matrix ϑ known to satisfy

$$\gamma \vartheta_* \preceq \vartheta \preceq \vartheta_*, \tag{4.53}$$

where $\gamma \geq 0$ and $\vartheta_* \succ 0$ are given. Our goal is to recover ϑ , and the norm on \mathbf{S}^n in which recovery error is measured satisfies Assumption \mathbf{A}' .

Processing the problem. We can process the just outlined problem as follows.

1. We represent the set $\{\vartheta \in \mathbf{S}^n_+ : \gamma \vartheta_* \preceq \vartheta \preceq \vartheta_*\}$ as the image of the matrix box

$$\mathcal{V} = \{ v \in \mathbf{S}^n : \|v\|_{\mathrm{Sh},\infty} \le 1 \} \qquad \qquad [\|\cdot\|_{\mathrm{Sh},\infty}: \text{ spectral norm}]$$

under affine mapping, specifically, we set

$$\vartheta_0 = \frac{1+\gamma}{2}\vartheta_*, \ \sigma = \frac{1-\gamma}{2}$$

and treat the matrix

$$v = \sigma^{-1} \vartheta_*^{-1/2} (\vartheta - \vartheta_0) \vartheta_*^{-1/2} \quad \left[\Leftrightarrow \vartheta = \vartheta_0 + \sigma \vartheta_*^{1/2} v \vartheta_*^{1/2} \right]$$

as the signal underlying our observations. Note that our a priori information on ϑ reduces to $v \in \mathcal{V}$.

2. We pass from observations η_k to "lifted" observations $\eta_k \eta_k^T \in \mathbf{S}^m$, so that

$$\mathbf{E}\{\eta_k \eta_k^T\} = \mathbf{E}\{A\xi_k \xi_k^T A^T\} = A\vartheta A^T = A\underbrace{(\vartheta_0 + \sigma A\vartheta_*^{1/2} v\vartheta_*^{1/2})}_{\vartheta[v]} A^T,$$

and treat as "actual" observations the matrices

$$\omega_k = \eta_k \eta_k^T - A \vartheta_0 A^T.$$

We have 60

$$\omega_k = \mathcal{A}v + \zeta_k \text{ with } \mathcal{A}v = \sigma A \vartheta_*^{1/2} v \vartheta_*^{1/2} A^T \text{ and } \zeta_k = \eta_k \eta_k^T - A \vartheta[v] A^T.$$
(4.54)

Observe that random matrices $\zeta_1, ..., \zeta_T$ are i.i.d. with zero mean and covariance mapping $\mathcal{Q}[v]$ (that of random matrix-valued variable $\zeta = \eta \eta^T - \mathbf{E}\{\eta \eta^T\}, \eta \sim \mathcal{N}(0, A\vartheta[v]A^T)).$

3. Let us \succeq -upper-bound the covariance mapping of ζ . Observe that $\mathcal{Q}[v]$ is a symmetric linear mapping of \mathbf{S}^m into itself given by

$$\langle h, \mathcal{Q}[v]h \rangle = \mathbf{E}\{\langle h, \zeta \rangle^2\} = \mathbf{E}\{\langle h, \eta \eta^T \rangle^2\} - \langle h, \mathbf{E}\{\eta \eta^T\} \rangle^2, \ h \in \mathbf{S}^m.$$

Given $v \in \mathcal{V}$, let us set $\theta = \vartheta[v]$, so that $0 \leq \theta \leq \theta_*$, and let $\mathcal{H}(h) = \theta^{1/2} A^T h A \theta^{1/2}$. We have

⁶⁰In our current considerations, we need to operate with linear mappings acting from \mathbf{S}^p to \mathbf{S}^q . We treat \mathbf{S}^k as Euclidean space equipped with the Frobenius inner product $\langle u, v \rangle = \operatorname{Tr}(uv)$ and denote linear mappings from \mathbf{S}^p into \mathbf{S}^q by capital calligraphic letters, like \mathcal{A}, \mathcal{Q} , etc. Thus, \mathcal{A} in (4.54) denotes the linear mapping which, on a closest inspection, maps matrix $v \in \mathbf{S}^n$ into the matrix $\mathcal{A}v = A[\vartheta[v] - \vartheta[0]]A^T$.

We have $\mathcal{H}(h) = U \text{Diag}\{\lambda\} U^T$ with orthogonal U, so that

$$\begin{aligned} \mathbf{E}_{\chi \sim \mathcal{N}(0,I_n)} \{ (\chi^T \mathcal{H}(h)\chi)^2 \} &- \operatorname{Tr}^2(\mathcal{H}(h)) \\ &= \mathbf{E}_{\bar{\chi}:=U^T \chi \sim \mathcal{N}(0,I_n)} \{ (\bar{\chi}^T \operatorname{Diag}\{\lambda\}\bar{\chi})^2 \} - (\sum_i \lambda_i)^2 \\ &= \mathbf{E}_{\bar{\chi} \sim \mathcal{N}(0,I_n)} \{ (\sum_i \lambda_i \bar{\chi}_i^2)^2 \} - (\sum_i \lambda_i)^2 = \sum_{i \neq j} \lambda_i \lambda_j + 3 \sum_i \lambda_i^2 - (\sum_i \lambda_i)^2 \\ &= 2 \sum_i \lambda_i^2 = 2 \operatorname{Tr}([\mathcal{H}(h)]^2). \end{aligned}$$

Thus,

$$\begin{aligned} \langle h, \mathcal{Q}[v]h \rangle &= 2\mathrm{Tr}([\mathcal{H}(h)]^2) = 2\mathrm{Tr}(\theta^{1/2}A^T hA\theta A^T hA\theta^{1/2}) \\ &\leq 2\mathrm{Tr}(\theta^{1/2}A^T hA\theta_* A^T hA\theta^{1/2}) \; [\text{since } 0 \preceq \theta \preceq \theta_*] \\ &= 2\mathrm{Tr}(\theta_*^{1/2}A^T hA\theta A^T hA\theta_*^{1/2}) \leq 2\mathrm{Tr}(\theta_*^{1/2}A^T hA\theta_* A^T hA\theta_*^{1/2}) \\ &= 2\mathrm{Tr}(\theta_*A^T hA\theta_* A^T hA). \end{aligned}$$

We conclude that

$$\forall v \in \mathcal{V} : \mathcal{Q}[v] \preceq \mathcal{Q}, \ \langle e, \mathcal{Q}h \rangle = 2 \mathrm{Tr}(\vartheta_* A^T h A \vartheta_* A^T e A), \ e, h \in \mathbf{S}^m.$$
(4.55)

4. To continue, we need to set some additional notation to be used when operating with Euclidean spaces \mathbf{S}^p , p = 1, 2, ...

• We denote $\bar{p} = \frac{p(p+1)}{2} = \dim \mathbf{S}^p$, $\mathcal{I}_p = \{(i,j) : 1 \le i \le j \le p\}$, and for $(i,j) \in \mathcal{I}_p$ set i = j

$$e_{p}^{ij} = \begin{cases} e_{i}e_{i}^{i}, & i = j \\ \frac{1}{\sqrt{2}}[e_{i}e_{j}^{T} + e_{j}e_{i}^{T}], & i < j \end{cases},$$

where e_i are the standard basic orths in \mathbf{R}^p . Note that $\{e_p^{ij} : (i,j) \in \mathcal{I}_p\}$ is the standard orthonormal basis in \mathbf{S}^p . Given $v \in \mathbf{S}^p$, we denote by $X^p(v)$ the vector of coordinates of v in this basis:

$$\mathbf{X}_{ij}^{p}(v) = \mathrm{Tr}(ve_{p}^{ij}) = \begin{cases} v_{ii}, & i = j \\ \sqrt{2}v_{ij}, & i < j \end{cases}, \ (i,j) \in \mathcal{I}_{p}.$$

Similarly, for $x \in \mathbf{R}^{\bar{p}}$, we index the entries in x by pairs ij, $(i, j) \in \mathcal{I}_p$, and set $V^p(x) = \sum_{(i,j)\in\mathcal{I}_p} x_{ij}e_p^{ij}$, so that $v \mapsto X^p(v)$ and $x \mapsto V^p(x)$ are inverse to each other linear norm-preserving maps identifying the Euclidean spaces \mathbf{S}^p and $\mathbf{R}^{\bar{p}}$ (recall that the inner products on these spaces are, respectively, the Frobenius and the standard one).

• Recall that \mathcal{V} is the matrix box $\{v \in \mathbf{S}^n : v^2 \preceq I_n\} = \{v \in \mathbf{S}^n : \exists t \in \mathcal{T} := [0, 1] : v^2 \preceq tI_n\}$. We denote by \mathcal{X} the image of \mathcal{V} under the mapping \mathbf{X}^n :

$$\mathcal{X} = \{ x \in \mathbf{R}^{\bar{n}} : \exists t \in \mathcal{T} : R^2[x] \preceq tI_n \}, \ R[x] = \sum_{(i,j) \in \mathcal{I}_n} x_{ij} e_n^{ij}, \ \bar{n} = \frac{1}{2}n(n+1).$$

Note that \mathcal{X} is a basic spectratope of size n.

Now we can assume that the signal underlying our observations is $x \in \mathcal{X}$, and the observations themselves are

$$w_k = \mathbf{X}^m(\omega_k) = \underbrace{\mathbf{X}^m(\mathcal{A}\mathbf{V}^n(x))}_{=:\overline{A}x} + z_k, \ z_k = \mathbf{X}^m(\zeta_k).$$

Note that $z_k \in \mathbf{R}^{\bar{m}}$, $1 \leq k \leq T$, are zero mean i.i.d. random vectors with covariance matrix Q[x] satisfying, in view of (4.55), the relation

$$Q[x] \preceq Q, \text{ where } Q_{ij,k\ell} = 2\text{Tr}(\vartheta_* A^T e_m^{ij} A \vartheta_* A^T e_m^{k\ell} A), \ (i,j) \in \mathcal{I}_m, (k,\ell) \in \mathcal{I}_m.$$

Our goal is to estimate $\vartheta[v] - \vartheta[0]$, or, what is the same, to recover

$$\overline{B}x := \mathbf{X}^n(\vartheta[\mathbf{V}^n(x)] - \vartheta[0])$$

We assume that the norm in which the estimation error is measured is "transferred" from \mathbf{S}^n to $\mathbf{R}^{\bar{n}}$; we denote the resulting norm on $\mathbf{R}^{\bar{n}}$ by $\|\cdot\|$ and assume that the unit ball \mathcal{B}_* of the conjugate norm $\|\cdot\|_*$ is given by spectratopic representation:

$$\{ u \in \mathbf{R}^{\bar{n}} : \|u\|_* \leq 1 \} = \{ u \in \mathbf{R}^{\bar{n}} : \exists y \in \mathcal{Y} : u = My \}, \\ \mathcal{Y} := \{ y \in \mathbf{R}^q : \exists r \in \mathcal{R} : S_\ell^2[y] \leq r_\ell I_{f_\ell}, 1 \leq \ell \leq L \}.$$

$$(4.56)$$

The formulated description of the estimation problem fit the premises of Proposition 4.14, specifically:

• the signal x underlying our observation $w^{(T)} = [w_1; ...; w_T]$ is known to belong to basic spectratope $\mathcal{X} \in \mathbf{R}^{\bar{n}}$, and the observation itself is of the form

$$w^{(T)} = \overline{A}^{(T)}x + z^{(T)}, \ \overline{A}^{(T)} = [\underline{\overline{A}}; ...; \underline{\overline{A}}], \ z^{(T)} = [z_1; ...; z_T];$$

• the noise $z^{(T)}$ is zero mean, and its covariance matrix is $\leq Q_T := \text{Diag}\{\underbrace{Q, ..., Q}_T\},\$

which allows to set $\Pi = \{Q_T\};$

• our goal is to recover $\overline{B}x$, and the norm $\|\cdot\|$ in which the recovery error is measured satisfies (4.56).

Proposition 4.14 supplies the linear estimate

$$\widehat{x}(w^{(T)}) = \sum_{k=1}^{T} H_{*k}^{T} w_k,$$

of $\overline{B}x$ with $H_* = [H_{*1}; ...; H_{*T}]$ stemming from the optimal solution to the convex optimization problem

$$\begin{array}{lll}
\operatorname{Opt} &= & \min_{H = [H_1; \dots; H_T], \Lambda, \Upsilon} \left\{ \operatorname{Tr}(\Lambda) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \overline{\Psi}_{\{Q_T\}}(H_1, \dots, H_T) : \\ & & \Lambda \in \mathbf{S}_+^n, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\}, \\ & \left[\frac{\mathcal{R}^*[\Lambda]}{\frac{1}{2}M^T[\overline{B} - [\sum_k H_k]^T \overline{A}]} \mid \frac{1}{2}[\overline{B}^T - \overline{A}^T \sum_k H_k]M}{\sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell]} \right] \succeq 0 \right\}, \\
\end{array}$$

$$(4.57)$$

where

$$\mathcal{R}^*[\Lambda] \in \mathbf{S}^{\bar{n}} : \ (\mathcal{R}^*[\Lambda])_{ij,k\ell} = \operatorname{Tr}(\Lambda e_n^{ij} e_n^{k\ell}), \ (i,j) \in \mathcal{I}_n, \ (k,\ell) \in \mathcal{I}_n,$$

and, cf. (4.48),

$$\overline{\Psi}_{\{Q_T\}}(H_1, ..., H_T) = \min_{\Upsilon', \Theta} \left\{ \operatorname{Tr}(Q_T \Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) : \Theta \in \mathbf{S}^{mT}, \ \Upsilon' = \{\Upsilon'_{\ell} \succeq 0, \ell \leq L\}, \\ \left[\frac{\Theta}{\frac{1}{2}[M^T H_1^T, ..., M^T H_T^T]} \mid \frac{1}{\sum_{\ell} \mathcal{S}^*_{\ell}[\Upsilon'_{\ell}]} \right] \succeq 0 \right\},$$

5. Evidently, the function $\overline{\Psi}_{\{Q_T\}}([H_1, ..., H_T])$ remains intact when permuting $H_1, ..., H_T$; with this in mind, it is clear that permuting $H_1, ..., H_T$ and keeping intact Λ and Υ is a symmetry of (4.57) – such a transformation maps feasible set onto itself and preserves the value of the objective. Since (4.57) is convex and solvable, it follows that there exists an optimal solution to the problem with $H_1 = ... = H_T = H$. On the other hand,

(we have used Schur Complement Lemma combined with the fact that $\sum_{\ell} S_{\ell}^*[\Upsilon_{\ell}] \succ 0$ whenever $\Upsilon_{\ell}' \succ 0$ for all ℓ , see Lemma 4.89).

In view of the above observations, when replacing variables H and G with $\overline{H} = TH$ and $\overline{G} = T^2G$, respectively, problem (4.57), (4.58) becomes

$$\begin{array}{lll}
\operatorname{Opt} &= \min_{\overline{H},\overline{G},\Lambda,\Upsilon,\Upsilon'} \left\{ \operatorname{Tr}(\Lambda) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \frac{1}{T}\operatorname{Tr}(Q\overline{G}) : \\ \Lambda \in \mathbf{S}_{+}^{n}, \Upsilon = \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \Upsilon' = \{\Upsilon'_{\ell} \succeq 0, \ell \leq L\}, \\ \left[\frac{\mathcal{R}^{*}[\Lambda]}{\frac{1}{2}M^{T}[\overline{B} - \overline{H}^{T}\overline{A}]} \left| \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \right| \geq 0, \\ \left[\frac{\overline{G}}{\frac{1}{2}M^{T}\overline{H}} \left| \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon'_{\ell}] \right| \geq 0 \end{array} \right\}, \\
\end{array} \right\}, \quad (4.59)$$

294

LECTURE 4

and the estimate

$$\widehat{x}(w^T) = \frac{1}{T}\overline{H}^T \sum_{k=1}^T w_k$$

stemming from an optimal solution to (4.59) satisfies

$$\operatorname{Risk}_{\Pi, \|\cdot\|} [\widehat{x} | \mathcal{X}] \leq \operatorname{Opt},$$

where $\Pi = \{Q_T\}.$

4.3.3.8 Estimation from repeated observations

Consider the special case of the situation from Section 4.3.3.1, the case where observation ω in (4.35) is *T*-element sample: $\omega = [\bar{\omega}_1; ...; \bar{\omega}_T]$ with components

$$\bar{\omega}_t = \bar{A}x + \xi_t, \ t = 1, ..., T$$

and ξ_t are i.i.d. observation noises with zero mean distribution \bar{P} satisfying $\bar{P} \ll \bar{\Pi}$ for some convex compact set $\bar{\Pi} \subset \operatorname{int} \mathbf{S}^{\bar{m}}_+$. In other words, we are in the situation where

$$A = [\underbrace{A; ...; A}_{T}] \in \mathbf{R}^{m \times n} \text{ for some } A \in \mathbf{R}^{m \times n} \text{ and } m = T\bar{m}$$
$$\Pi = \{Q = \text{Diag}\{\underbrace{\bar{Q}, ..., \bar{Q}}_{T}\}, \bar{Q} \in \bar{\Pi}\}$$

The same argument as used in item 5 of Section 4.3.3.7 justifies the following

Proposition 4.15. In the situation in question and under Assumptions A, B, R the linear estimate of Bx yielded by an optimal solution to problem (4.50) can be found as follows. We consider the convex optimization problem

$$\overline{\mathrm{Opt}} = \min_{\bar{H},\Lambda,\Upsilon,\Upsilon',\bar{\Theta}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \frac{1}{T}\overline{\Gamma}(\bar{\Theta}) : \\ \Lambda = \{\Lambda_k \succeq 0, k \le K\}, \ \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \le L\}, \ \Upsilon' = \{\Upsilon'_\ell \succeq 0, \ell \le L\}, \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T[B - \bar{H}^T A]} \right] \frac{1}{2} \frac{[B^T - A^T \bar{H}]M}{\sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell]} \right] \succeq 0, \\ \left[\frac{\Theta}{\frac{1}{2}M^T \bar{H}^T} \left| \sum_\ell \mathcal{S}_\ell^*[\Upsilon'_\ell]} \right] \succeq 0 \right]$$

$$(4.60)$$

where

$$\overline{\Gamma}(\bar{\Theta}) = \max_{\bar{Q} \in \bar{\Pi}} \operatorname{Tr}(\bar{Q}\bar{\Theta}).$$

The problem is solvable, and the estimate in question is yielded by the \overline{H} -component \overline{H}_* of the optimal solution according to

$$\widehat{x}([\overline{\omega}_1;...;\overline{\omega}_T]) = \frac{1}{T}\overline{H}_*^T \sum_{t=1}^T \overline{\omega}_t.$$

The provided by Proposition 4.14 upper bound on the risk $\operatorname{Risk}_{\Pi, \|\cdot\|}[\widehat{x}(\cdot)|\mathcal{X}]$ of this estimate is $\overline{\operatorname{Opt}}$.

The advantage of this result as compared to what is stated under the circumstances

by Proposition 4.14 is that the sizes of optimization problem (4.60) are independent of T.

4.3.3.9 Near-optimality in Gaussian case

The risk of the linear estimate $\hat{x}_{H_*}(\cdot)$ constructed in (4.50), (4.51) can be compared to the minimax optimal risk of recovering $Bx, x \in \mathcal{X}$, from observations corrupted by zero mean Gaussian noise with covariance matrix from Π ; formally, this minimax optimal risk is defined as

$$\operatorname{RiskOpt}_{\Pi, \|\cdot\|} [\mathcal{X}] = \sup_{Q \in \Pi} \inf_{\widehat{x}(\cdot)} \left[\sup_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim \mathcal{N}(0, Q)} \{ \|Bx - \widehat{x}(Ax + \xi)\| \} \right]$$
(4.61)

where the infimum is taken over all estimates.

Proposition 4.16. Under the premise and in the notation of Proposition 4.14, we have

$$\operatorname{RiskOpt}_{\Pi, \|\cdot\|}[\mathcal{X}] \ge \frac{\operatorname{Opt}}{64\sqrt{(2\ln F + 10\ln 2)(2\ln D + 10\ln 2)}}.$$
 (4.62)

where

$$D = \sum_{k} d_k, \ F = \sum_{\ell} f_{\ell}.$$
(4.63)

Thus, the upper bound Opt on the risk $\operatorname{Risk}_{\Pi,\|\cdot\|}[\widehat{x}_{H_*}|\mathcal{X}]$ of the presumably good linear estimate \widehat{x}_{H_*} yielded by an optimal solution to optimization problem (4.50) is within logarithmic in the sizes of spectratopes \mathcal{X} and \mathcal{B}_* factor from the Gaussian minimax risk $\operatorname{RiskOpt}_{\Pi,\|\cdot\|}[\mathcal{X}]$.

For the proof, see Section 4.10.5. The key component of the proof is the following important by its own right fact (for proof, see Section 4.10.4):

Lemma 4.17. Let Y be an $N \times \nu$ matrix, let $\|\cdot\|$ be a norm on \mathbb{R}^{ν} such that the unit ball \mathcal{B}_* of the conjugate norm is the spectratope (4.39), and let $\zeta \sim \mathcal{N}(0, Q)$ for some positive semidefinite $N \times N$ matrix Q. Then the best upper bound on $\psi_Q(Y) := \mathbb{E}\{\|Y^T\zeta\|\}$ yielded by Lemma 4.11, that is, the optimal value $\operatorname{Opt}[Q]$ in the convex optimization problem (cf. (4.48))

$$Opt[Q] = \min_{\Theta, \Upsilon} \left\{ \phi_{\mathcal{R}}(\lambda[\Upsilon]) + Tr(Q\Theta) : \Upsilon = \{\Upsilon_{\ell} \succeq 0, 1 \le \ell \le L\}, \Theta \in \mathbf{S}^{N}, \\ \left[\frac{\Theta}{\frac{1}{2}M^{T}Y^{T}} \middle| \frac{\frac{1}{2}YM}{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]} \right] \succeq 0 \right\}$$

$$(4.64)$$

(for notation, see Lemma 4.11 and (4.43)) satisfies the identity

$$\forall (Q \succeq 0) : \\ \operatorname{Opt}[Q] = \overline{\operatorname{Opt}}[Q] := \min_{G, \Upsilon = \{\Upsilon_{\ell}, \ell \leq L\}} \left\{ \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \operatorname{Tr}(G) : \Upsilon_{\ell} \succeq 0, \\ \left[\frac{G}{\frac{1}{2}M^{T}Y^{T}Q^{1/2}} \left| \frac{1}{2}Q^{1/2}YM \right| \right] \succeq 0 \right\},$$

$$(4.65)$$

and is a tight bound on $\psi_Q(Y)$, namely,

$$\psi_Q(Y) \le \operatorname{Opt}[Q] \le 22\sqrt{2\ln F} + 10\ln 2\psi_Q(Y),$$
(4.66)

where $F = \sum_{\ell} f_{\ell}$ is the size of the spectratope (4.39). Besides this, for all $\varkappa \geq 1$ one has

$$\operatorname{Prob}_{\zeta}\left\{\|Y^{T}\zeta\| \geq \frac{\operatorname{Opt}[Q]}{4\varkappa}\right\} \geq \beta_{\varkappa} := 1 - \frac{e^{3/8}}{2} - 2Fe^{-\varkappa^{2}/2}.$$
(4.67)

In particular, when selecting $\varkappa = \sqrt{2 \ln F + 10 \ln 2}$, we obtain

$$\operatorname{Prob}_{\zeta} \left\{ \|Y^{T}\zeta\| \ge \frac{\operatorname{Opt}[Q]}{4\sqrt{2\ln F + 10\ln 2}} \right\} \ge \beta_{\varkappa} = 0.2100 > \frac{3}{16}$$
(4.68)

4.4 LINEAR ESTIMATES OF STOCHASTIC SIGNALS

In the recovery problem considered so far in this Lecture, the signal x underlying observation $\omega = Ax + \xi$ was "deterministic uncertain but bounded" – all a priori information on x was that $x \in \mathcal{X}$ for a given signal set \mathcal{X} . There is a well-known alternative model, where the signal x has a random component, specifically,

$$x = [\eta; u]$$

where the "stochastic component" η is random with (partly) known probability distribution P_{η} , and the "deterministic component" u is known to belong to a given set \mathcal{X} . As a typical example, consider linear dynamical system given by

$$\begin{array}{rcl} y_{t+1} &=& P_t y_t + \eta_t + u_t \\ \omega_t &=& C_t y_t + \xi_t \end{array}, 1 \le t \le K, \tag{4.69}$$

where y_t , η_t , u_t are, respectively, the state, the random "process noise," and the deterministic "uncertain but bounded" disturbance affecting the system at time t, ω_t is the output – it is what we observe at time t, and ξ_t is the observation noise. We assume that the matrices P_t , C_t are known in advance. Note that the trajectory

$$y = [y_1; \dots; y_K]$$

of the states depends not only on the trajectories of process noises η_t and disturbances u_t , but also on the initial state y_1 , which can be modeled as a realization of either the initial noise η_0 , or the initial disturbance u_0 . When $u_t \equiv 0$, $y_1 = \eta_0$ and the random vectors $\{\eta_t, 0 \leq t \leq K, \xi_t, 1 \leq t \leq K\}$ are independent of each other zero mean Gaussian, (4.69) is the model underlying the famous Kalman filter.

Now, given model (4.69), we can use the equations of the model to represent the trajectory of the states as linear image of the trajectory of noises $\eta = \{\eta_t\}$ and the trajectory of disturbances $u = \{u_t\}$:

$$y = P\eta + Qu$$

(recall that the initial state is either the component η_0 of η , or the component u_0 of u), and our "full observation" becomes

$$\omega = [\omega_1; ...; \omega_K] = A[\eta; u] + \xi, \ \xi = [\xi_1, ..., \xi_K].$$

A typical statistical problem associated with the outlined situation is to estimate the linear image $B[\eta; u]$ of the "signal" $x = [\eta; u]$ underlying our observation. For example, when speaking about (4.69), the goal could be to recover y_{K+1} ("forecast").

We arrive at the following estimation problem:

Given noisy observation

$$\omega = Ax + \xi \in \mathbf{R}^m$$

of signal $x = [\eta; u]$ with random component $\eta \in \mathbf{R}^k$ and deterministic component u known to belong to a given set $\mathcal{X} \subset \mathbf{R}^n$, we want to recover the image $Bx \in \mathbf{R}^{\nu}$ of the signal. Here A and B are given matrices, η is independent of ξ , and we have a priori (perhaps, incomplete) information on the probability distribution P_{η} of η , specifically, know that $P_{\eta} \in \mathcal{P}_{\eta}$ for a given family \mathcal{P}_{η} of probability distributions. Similarly, we assume that what we know about the noise ξ is that its distribution belongs to a given family \mathcal{P}_{ξ} of distributions on the observation space.

Given a norm $\|\cdot\|$ on the image space of B, it makes sense to specify the risk of a candidate estimate $\hat{x}(\omega)$ by taking expectation of the recovery error $\|\hat{x}(A[\eta; u] + \xi) - B[\eta; u]\|$ over both ξ and η and then taking supremum of the result over the allowed distributions of η , ξ and over $u \in \mathcal{X}$:

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}] = \sup_{u \in \mathcal{X}} \sup_{P_{\xi} \sim \mathcal{P}_{\xi}, P_{\eta} \sim \mathcal{P}_{\eta}} \mathbf{E}_{[\xi;\eta] \sim P_{\xi} \times P_{\eta}} \left\{ \|\widehat{x}(A[\eta;u] + \xi) - B[\eta;u]\| \right\}.$$

When $\|\cdot\| = \|\cdot\|_2$ and all distributions from \mathcal{P}_{ξ} and \mathcal{P}_{η} are with zero means and finite covariance matrices, it is technically more convenient to operate with the *Euclidean risk*

$$\operatorname{Risk}_{\operatorname{Eucl}}[\widehat{x}] = \left[\sup_{u \in \mathcal{X}} \sup_{P_{\xi} \sim \mathcal{P}_{\xi}, P_{\eta} \sim \mathcal{P}_{\eta}} \mathbf{E}_{[\xi;\eta] \sim P_{\xi} \times P_{\eta}} \left\{ \|\widehat{x}(A[\eta;u] + \xi) - B[\eta;u]\|_{2}^{2} \right\} \right]^{1/2}.$$

Our goal in this Section is to show that as far as the design of "presumably good" linear estimates $\hat{x}(\omega) = H^T \omega$ is concerned, the techniques developed so far can be straightforwardly extended from the case of signals with no random component to the one where this component is present.

4.4.1 Minimizing Euclidean risk

For the time being, assume that \mathcal{P}_{ξ} is comprised of all probability distributions P on \mathbf{R}^m with zero mean and covariance matrices $\operatorname{Cov}[P] = \mathbf{E}_{\xi \sim P} \{\xi\xi^T\}$ running through a computationally tractable convex compact subset $\mathcal{Q}_{\xi} \subset \operatorname{int} \mathbf{S}^m_+$, and \mathcal{P}_{η} is comprised of all probability distributions P on \mathbf{R}^{μ} with zero mean and covariance matrices running through a computationally tractable convex compact subset $\mathcal{Q}_{\eta} \subset \operatorname{int} \mathbf{S}^m_+$. Let, in addition, \mathcal{X} be a basic spectratope:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \leq t_k I_{d_k}, \, k \leq K \}$$

with our standard restrictions on \mathcal{T} and $R_k[\cdot]$. Let us derive efficiently solvable convex optimization problem "responsible" for presumably good, in terms of its Euclidean risk, linear estimate.

For a linear estimate $H^T \omega$, $u \in \mathcal{X}$, $P_{\xi} \in \mathcal{P}_{\xi}$, $P_{\eta} \sim \mathcal{P}_{\eta}$, denoting by Q_{ξ} and Q_{η} the covariance matrices of P_{ξ} and P_{η} , and partitioning A as $A = [A_{\eta}, A_{u}]$ and $B = [B_{\eta}, B_{u}]$ according to the partition $x = [\eta; u]$, we have

$$\begin{split} & \mathbf{E}_{[\xi;\eta]\sim P_{\xi}\times P_{\eta}}\left\{ \|H^{T}(A[\eta;u]+\xi)-B[\eta;u]\|_{2}^{2} \right\} \\ &= \mathbf{E}_{[\xi;\eta]\sim P_{\xi}\times P_{\eta}}\left\{ \|[H^{T}A_{\eta}-B_{\eta}]\eta+H^{T}\xi+[H^{T}A_{u}-B_{u}]u\|_{2}^{2} \right\} \\ &= u^{T}[B_{u}-H^{T}A_{u}]^{T}[B_{u}-H^{T}A_{u}]u+\mathbf{E}_{\xi\sim P_{\xi}}\left\{ \mathrm{Tr}(H^{T}\xi\xi^{T}H) \right\} \\ &\quad + \mathbf{E}_{\eta\sim P_{\eta}}\left\{ \mathrm{Tr}([B_{\eta}-H^{T}A_{\eta}]\eta\eta^{T}[B_{\eta}-H^{T}A_{\eta}]^{T}) \right\} \\ &= u^{T}[B_{u}-H^{T}A_{u}]^{T}[B_{u}-H^{T}A_{u}]u+\mathrm{Tr}(H^{T}Q_{\xi}H)+\mathrm{Tr}([B_{\eta}-H^{T}A_{\eta}]Q_{\eta}[B_{\eta}-H^{T}A_{\eta}]^{T}), \end{split}$$

whence the squared Euclidean risk of the linear estimate $\hat{x}_H(\omega) = H^T \omega$ is

$$\begin{aligned} \operatorname{Risk}_{\operatorname{Eucl}}^{2}[\widehat{x}_{H}] &= \Phi(H) + \Psi_{\xi}(H) + \Psi_{\eta}(H), \\ \Phi(H) &= \max_{u \in \mathcal{X}} u^{T}[B_{u} - H^{T}A_{u}]^{T}[B_{u} - H^{T}A_{u}]u, \\ \Psi_{\xi}(H) &= \max_{Q \in \mathcal{Q}_{\xi}} \operatorname{Tr}(H^{T}QH), \\ \Psi_{\eta}(H) &= \max_{Q \in \mathcal{Q}_{\eta}} \operatorname{Tr}([B_{\eta} - H^{T}A_{\eta}]Q[B_{\eta} - H^{T}A_{\eta}]^{T}) \end{aligned}$$

Functions Ψ_{ξ} and Ψ_{η} are convex and efficiently computable, function $\Phi(H)$, by Proposition 4.8, admits efficiently computable convex upper bound

$$\overline{\Phi}(H) = \min_{\Lambda} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) : \Lambda = \{\Lambda_k \succeq 0, k \le K\}, \\ [B_u - H^T A_u]^T [B_u - H^T A_u] \preceq \sum_k \mathcal{R}_k^* [\Lambda_k] \right\}$$

tight within the factor $2 \max[\ln(2D), 1]$ (for notation, see Proposition 4.8), so that the efficiently solvable convex problem yielding presumably good linear estimate is

$$Opt = \min_{H} \left[\overline{\Phi}(H) + \Psi_{\xi}(H) + \Psi_{\eta}(H) \right];$$

the Euclidean risk of the linear estimate $H_*^T \omega$ yielded by the optimal solution to the problem is upper-bounded by $\sqrt{\text{Opt}}$ and is within factor $\sqrt{2 \max[\ln(2\sum_k D_k), 1]}$ of the minimal Euclidean risk achievable with linear estimates.

4.4.2 Minimizing $\|\cdot\|$ -risk

Now let \mathcal{P}_{ξ} be comprised of all probability distributions P on \mathbb{R}^m with matrices of second moments $\operatorname{Var}[P] = \mathbb{E}_{\xi \sim P} \{\xi \xi^T\}$ running through a computationally tractable convex compact subset $\mathcal{Q}_{\xi} \subset \operatorname{int} \mathbb{S}^m_+$, and \mathcal{P}_{η} be comprised of all probability distributions P on \mathbb{R}^{μ} with matrices of second moments $\operatorname{Var}[P]$ running through a computationally tractable convex compact subset $\mathcal{Q}_{\eta} \subset \operatorname{int} \mathbb{S}^{\mu}_+$. Let, same as above, \mathcal{X} be a basic spectratope:

$$\mathcal{X} = \{ u \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[u] \leq t_k I_{d_k}, \, k \leq K \},\$$

and let $\|\cdot\|$ be such that the unit ball \mathcal{B}_* of the conjugate norm $\|\cdot\|_*$ is a spectratope:

$$\mathcal{B}_* = \{y : \|y\|_* \le 1\} = \{y \in \mathbf{R}^{\nu} : \exists (r \in \mathcal{R}, z \in \mathbf{R}^N) : y = Mz, S_{\ell}^2[z] \preceq r_{\ell} I_{f_{\ell}}, \ell \le L\},\$$

with our standard restrictions on $\mathcal{T}, \mathcal{R}, R_k[\cdot], S_\ell[\cdot]$. Here the efficiently solvable convex optimization problem "responsible" for presumably good, in terms of its

risk Risk_{||·||}, linear estimate can be built as follows. For a linear estimate $H^T \omega$, $u \in \mathcal{X}$, $P_{\xi} \in \mathcal{P}_{\xi}$, $P_{\eta} \in \mathcal{P}_{\eta}$, denoting by Q_{ξ} and Q_{η} the matrices of second moments of P_{ξ} and P_{η} , and partitioning A as $A = [A_{\eta}, A_u]$ and $B = [B_{\eta}, B_u]$ according to the partition $x = [\eta; u]$, we have

$$\begin{aligned} & \mathbf{E}_{[\xi;\eta]\sim P_{\xi}\times P_{\eta}}\left\{ \|H^{T}(A[\eta;u]+\xi) - B[\eta;u]\| \right\} \\ &= \mathbf{E}_{[\xi;\eta]\sim P_{\xi}\times P_{\eta}}\left\{ \|[H^{T}A_{\eta} - B_{\eta}]\eta + H^{T}\xi + [H^{T}A_{u} - B_{u}]u\| \right\} \\ &\leq \|[B_{u} - H^{T}A_{u}]u\| + \mathbf{E}_{\xi\sim P_{\xi}}\left\{ \|H^{T}\xi\| \right\} + \mathbf{E}_{\eta\sim P_{\eta}}\left\{ \|[B_{\eta} - H^{T}A_{\eta}]\eta\| \right\}. \end{aligned}$$

It follows that for a linear estimate $\hat{x}_H(\omega) = H^T \omega$ one has

$$\begin{aligned} \operatorname{Risk}_{\|\cdot\|}[\widehat{x}_{H}] &\leq & \Phi(H) + \Psi_{\xi}(H) + \Psi_{\eta}(H), \\ & \Phi(H) &= & \max_{u \in \mathcal{X}} \|[B_{u} - H^{T}A_{u}]u\|, \\ & \Psi_{\xi}(H) &= & \sup_{P_{\xi} \in \mathcal{P}_{\xi}} \mathbf{E}_{\xi \sim P_{\xi}} \{\|H^{T}\xi\|\}, \\ & \Psi_{\eta}(H) &= & \sup_{P_{\eta} \in \mathcal{P}_{\eta}} \mathbf{E}_{\xi \sim P_{\xi}} \{\|[B_{\eta} - H^{T}A_{\eta}]\eta\|\} \end{aligned}$$

As was shown in Section 4.3.3.3, the functions Φ , Ψ_{ξ} , Ψ_{η} admit efficiently computable upper bounds as follows (for notation, see Section 4.3.3.3):

$$\begin{split} \Phi(H) &\leq \overline{\Phi}(H) := \min_{\Lambda,\Upsilon} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) : \\ \Lambda &= \{\Lambda_k \succeq 0, k \leq K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\} \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T[B_u - H^TAu]} \middle| \frac{1}{2}[B_u^T - A_u^TH]M} \right] \succeq 0 \end{array} \right\}; \\ \Psi_{\xi}(H) &\leq \overline{\Psi}_{\xi}(H) := \min_{\Upsilon,G} \left\{ \phi_{\mathcal{R}}[\lambda[\Upsilon]) + \Gamma_{\xi}(G) : \begin{bmatrix} G & \left| \frac{1}{2}HM \right| \\ \frac{1}{2}M^TH^T & \left| \sum_\ell \mathcal{S}_\ell[\Upsilon_\ell] \right| \right\} \\ \Gamma_{\xi}(G) &= \max_{Q \in \mathcal{Q}_{\xi}} \operatorname{Tr}(GQ); \\ \Psi_{\eta}(H) &\leq \overline{\Psi}_{\eta}(H) := \min_{\Upsilon,G} \left\{ \phi_{\mathcal{R}}[\lambda[\Upsilon]) + \Gamma_{\eta}(G) : \end{bmatrix} \end{split}$$

$$\begin{split} \Upsilon &= \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \\ & \left[\frac{G}{\frac{1}{2} [B_{\eta}^{T} - A_{\eta}^{T} H] M}{\left[\frac{1}{2} [B_{\eta}^{T} - A_{\eta}^{T} H] M \right] \sum_{\ell} \mathcal{S}_{\ell} [\Upsilon_{\ell}]} \right] \succeq 0, \ \ \right\}, \\ \Gamma_{\eta}(G) &= \max_{Q \in \mathcal{Q}_{\eta}} \operatorname{Tr}(GQ), \end{split}$$

and these bounds are reasonably tight (for details on tightness, see Proposition 4.8 and Lemma 4.17). As a result, to get a presumably good linear estimate, one needs to solve the efficiently solvable convex optimization problem

$$Opt = \min_{H} \left[\overline{\Phi}(H) + \overline{\Psi}_{\xi}(H) + \overline{\Psi}_{\eta}(H) \right];$$

the linear estimate $\hat{x}_{H_*} = H_*^T \omega$ yielded by an optimal solution H_* to this problem admits the risk bound

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}_{H_*}] \leq \operatorname{Opt}$$

Note that the above derivation did not use independence of ξ and η at all.

4.5 LINEAR ESTIMATION UNDER UNCERTAIN-BUT-BOUNDED NOISE

So far, the main subject of our interest was recovering (linear images of) signals via indirect observations of these signals corrupted by random noise. In this section, we focus on alternative observation schemes – those with "uncertain-but-bounded" and with "mixed" noise.

4.5.1 Uncertain-but-bounded noise

Consider recovering problem where one, given observation

$$\omega = Ax + \eta \tag{4.70}$$

of unknown signal x known to belong to a given signal set \mathcal{X} , wants to recover linear image Bx of x. Here A and B are given $m \times n$ and $\nu \times n$ matrices. The situation looks exactly as before; the difference with our previous considerations is that now we do not assume the observation noise to be random; all we assume about η is that it belongs to a given compact set \mathcal{H} ("uncertain-but-bounded observation noise"). In the situation in question, a natural definition of the risk on \mathcal{X} of a candidate estimate $\omega \mapsto \hat{x}(\omega)$ is

$$\operatorname{Risk}_{\mathcal{H},\|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x \in X, \eta \in \mathcal{H}} \|Bx - \widehat{x}(Ax + \eta)\|$$
(4.71)

(" \mathcal{H} -risk").

We are about to prove that when \mathcal{X} , \mathcal{H} and the unit ball \mathcal{B}_* of the norm $\|\cdot\|_*$ conjugate to $\|\cdot\|$ is are spectratopes, which we assume from now on, an efficiently computable linear estimate is near-optimal in terms of its \mathcal{H} -risk.

Our initial observation is that the situation in question reduces straightforwardly to the one where there is no observation noise at all. Indeed, let $\mathcal{Y} = \mathcal{X} \times \mathcal{H}$; then \mathcal{Y} is a spectratope, and we lose nothing when assuming that the signal underlying observation ω is $y = [x; \eta] \in \mathcal{Y}$:

$$\omega = Ax + \eta = \bar{A}y, \ \bar{A} = [A, I_m],$$

while the entity to be recovered is

$$Bx = \bar{B}y, \ \bar{B} = [B, 0_{\nu \times m}].$$

With these conventions, the \mathcal{H} -risk of a candidate estimate $\hat{x}(\cdot) : \mathbf{R}^m \to \mathbf{R}^{\nu}$ becomes the quantity

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X} \times \mathcal{H}] = \sup_{y = [x;\eta] \in \mathcal{X} \times \mathcal{H}} \|\overline{B}y - \widehat{x}(\overline{A}y)\|,$$

that is, we indeed arrive at the situation where the observation noise is identically zero.

To avoid messy notation, let us assume that the outlined reduction has been carried out in advance, so that

(!) The problem of interest is to recover the linear image $Bx \in \mathbf{R}^{\nu}$ of

an unknown signal x known to belong to a given spectratope \mathcal{X} (which, as always, we can assume w.l.o.g. to be basic) from noiseless observation

$$\omega = Ax \in \mathbf{R}^m,$$

and the risk of a candidate estimate is defined as

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \|Bx - \widehat{x}(Ax)\|,$$

where $\|\cdot\|$ is a given norm with a spectratope \mathcal{B}_* , see (4.39), as the unit ball of the conjugate norm:

$$\begin{aligned}
\mathcal{X} &= \{x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \leq t_k I_{d_k}, k \leq K\}, \\
\mathcal{B}_* &= \{z \in \mathbf{R}^\nu : \exists y \in \mathcal{Y} : z = My\}, \\
\mathcal{Y} := \{y \in \mathbf{R}^q : \exists r \in \mathcal{R} : S_\ell^2[y] \leq r_\ell I_{f_\ell}, 1 \leq \ell \leq L\},
\end{aligned}$$
(4.72)

with our standard restrictions on \mathcal{T}, \mathcal{R} and $R_k[\cdot], S_\ell[\cdot]$.

4.5.1.1 Building linear estimate

Let us build a presumably good linear estimate. For a linear estimate $\hat{x}_H(\omega) = H^T \omega$, we have

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}_{H}|\mathcal{X}] = \max_{x \in \mathcal{X}} \|(B - H^{T}A)x\|$$
$$= \max_{[u;x] \in \mathcal{B}_{*} \times \mathcal{X}} [u;x]^{T} \left[\frac{|\frac{1}{2}(B - H^{T}A)|}{|\frac{1}{2}(B - H^{T}A)^{T}|} \right] [u;x].$$

Applying Proposition 4.8, we arrive at the following

Proposition 4.18. In the situation of this section, consider the convex optimization problem

$$Opt_{\#} = \min_{H,\Upsilon = \{\Upsilon_{\ell}\},\Lambda = \{\Lambda_{k}\}} \left\{ \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{T}}(\lambda[\Lambda]) : \Upsilon_{\ell} \succeq 0, \ \Lambda_{k} \succeq 0, \ \forall(\ell,k) \\ \left[\frac{\sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}]}{\frac{1}{2}M^{T}[B - H^{T}A]} \middle| \frac{1}{2} \frac{[B - H^{T}A]^{T}M}{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]} \right] \succeq 0 \right\},$$
(4.73)

where $\mathcal{R}_{k}^{*}[\cdot]$, $\mathcal{S}_{\ell}^{*}[\cdot]$ are induced by $R_{k}[\cdot]$, $S_{\ell}[\cdot]$, respectively, as explained in Section 4.3.1. The problem is solvable, and the risk of the linear estimate $\hat{x}_{H_{*}}(\cdot)$ yielded by the H-component of an optimal solution does not exceed $\operatorname{Opt}_{\#}$.

For proof, see Section 4.10.6.1.

4.5.1.2 Near-optimality

Proposition 4.19. The linear estimate \hat{x}_{H_*} yielded by Proposition 4.18 is nearoptimal in terms of its risk:

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}_{H_*}|\mathcal{X}] \le \operatorname{Opt}_{\#} \le O(1)\ln(D)\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}], \quad D = \sum_k d_k + \sum_{\ell} f_{\ell}, \quad (4.74)$$

302

LECTURE 4

(4.75)

where $\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}]$ is the minimax optimal risk:

$$\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}] = \inf_{\widehat{\alpha}} \operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}],$$

where inf is taken w.r.t. all Borel estimates.

Remark 4.20. When \mathcal{X} and \mathcal{B}_* are ellitopes rather than spectratopes:

$$\begin{aligned} \mathcal{X} &= \{x \in \mathbf{R}^n : \exists t \in \mathcal{T}, y : x = Py, y^T R_k y \leq t_k, k \leq K\}, \\ \mathcal{B}_* &:= \{u \in \mathbf{R}^\nu : \|u\|_* \leq 1\} = \{u \in \mathbf{R}^\nu : \exists r \in \mathcal{R}, z : u = Mz, z^T S_\ell z \leq r_\ell, \ell \leq L\} \\ & [R_k \succeq 0, \sum_k R_k \succ 0, S_\ell \succeq 0, \sum_\ell S_\ell \succ 0] \end{aligned}$$

problem (4.73) becomes

$$\begin{aligned}
\operatorname{Opt}_{\#} &= \min_{H,\lambda,\mu} \left\{ \phi_{\mathcal{R}}(\mu) + \phi_{\mathcal{T}}(\lambda) : \lambda \ge 0, \mu \ge 0, \\
& \left[\frac{\sum_{k} \lambda_{k} R_{k}}{\left| \frac{1}{2} [B - H^{T} A]^{T} M \right|} \sum_{\ell} \mu_{\ell} S_{\ell} \right] \ge 0 \right\},
\end{aligned}$$
(4.76)

and (4.74) can be strengthened to

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}_{H_*}|\mathcal{X}] \le \operatorname{Opt}_{\#} \le O(1)\ln(K+L)\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}].$$
(4.77)

For proofs, see Section 4.10.6.

4.5.1.3 Nonlinear estimation

Uncertain-but-bounded model of observation error makes it easy to point out an efficiently computable near-optimal *nonlinear* estimate. Specifically, in the situation described in the beginning of Section 4.5.1, assume that the range of observation error η is

$$\mathcal{H} = \{\eta \in \mathbf{R}^m : \|\eta\|_{(m)} \le \sigma\},\tag{4.78}$$

where $\|\cdot\|_{(m)}$, $\sigma > 0$ are a given norm on \mathbf{R}^m and a given error bound, and let us measure the recovery error by a given norm $\|\cdot\|_{(\nu)}$ on \mathbf{R}^{ν} . We can immediately point out a (nonlinear) estimate optimal, in terms of its \mathcal{H} -risk, within factor 2, specifically, the estimate \hat{x}_* defined as follows:

Given ω , we solve the feasibility problem

find
$$x \in \mathcal{X} : ||Ax - \omega||_{(m)} \le \sigma$$
 $(F[\omega])$

find a feasible solution x_{ω} to the problem, and set $\hat{x}_*(\omega) = Bx_{\omega}$.

Note that the estimate is well defined, since $(F[\omega])$ clearly is solvable, with one of the feasible solutions being the true signal underlying observation ω . When \mathcal{X} is a computationally tractable convex compact set, and $\|\cdot\|_{(m)}$ is an efficiently computable norm, a feasible solution to $(F[\omega])$ can be found in a computationally efficient fashion. Let us make the following immediate observation:

Proposition 4.21. The estimate \hat{x}_* is optimal within factor 2:

$$\operatorname{Risk}_{\mathcal{H}}[\widehat{x}_*|\mathcal{X}] \leq \operatorname{Opt}_* := \sup_{x,y} \left\{ \|Bx - By\|_{(\nu)} : x, y \in \mathcal{X}, \|A(x - y)\|_{(m)} \leq 2\sigma \right\}$$

$$\leq 2\operatorname{Risk}_{\operatorname{opt},\mathcal{H}},$$

where $\operatorname{Risk}_{\operatorname{opt},\mathcal{H}}$ is the infimum, over all estimates, of \mathcal{H} -risks of the estimate.

The proof of Proposition is the subject of Exercise 4.44.

4.5.1.4 Quantifying risk

Note that Proposition 4.21 does not impose restrictions on \mathcal{X} and the norms $\|\cdot\|_{(m)}$, $\|\cdot\|_{(\nu)}$.

The only - but essential – shortcoming of the estimate \hat{x}_* is that we do not know, in general, what is its \mathcal{H} -risk. From (4.79) it follows that this risk is tightly (namely, within factor 2) upper-bounded by Opt_* , but this quantity, being the maximum of a convex function over some domain, can be difficult to compute. Aside from handful of special cases where this difficulty does not arise, there is a generic situation when Opt_* can be tightly upper-bounded by efficient computation. This is the situation where \mathcal{X} is the spectratope defined in (4.72), $\|\cdot\|_{(m)}$ is such that the unit ball of this norm is a basic spectratope:

$$B_{(m)} := \{ u : \|u\|_{(m)} \le 1 \} = \{ u \in \mathbf{R}^m : \exists p \in \mathcal{P} : Q_j^2[u] \preceq p_j I_{e_j}, 1 \le j \le J \},\$$

and the unit ball of the norm $\|\cdot\|_{(\nu),*}$ conjugate to the norm $\|\cdot\|_{(\nu)}$ is a spectratope:

$$\begin{aligned}
B^*_{(\nu)} &:= \{ v \in \mathbf{R}^{\nu} : \|v\|_{(\nu),*} \le 1 \} \\
&= \{ v : \exists (w \in \mathbf{R}^N, r \in \mathcal{R}) : v = Mw, S^2_{\ell}[w] \preceq r_{\ell} I_{f_{\ell}}, 1 \le \ell \le L \},
\end{aligned}$$

with our usual restrictions on $\mathcal{P}, \mathcal{R}, Q_j[\cdot], S_\ell[\cdot]$.

Proposition 4.22. In the situation in question, consider convex optimization problem

$$Opt = \min_{\substack{\Lambda = \{\Lambda_k, k \leq K\}, \\ \Upsilon = \{\Upsilon_\ell, \ell \leq L\}, \\ \Sigma = \{\Sigma_j, j \leq J\}}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \sigma^2 \phi_{\mathcal{P}}(\lambda[\Sigma]) + \phi_{\mathcal{R}}(\lambda([\Sigma])) :$$

$$\prod_{\substack{X \geq 0, \\ Y_\ell \geq 0, \\ \Sigma_\ell \neq 0, \\$$

where $\mathcal{R}_{k}^{*}[\cdot]$ are associated with the mappings $x \mapsto R_{k}[x]$ according to (4.27), and $\mathcal{S}_{\ell}^{*}[\cdot]$ and $\mathcal{Q}_{j}^{*}[\cdot]$ are associated in the same fashion with the mappings $w \mapsto \mathcal{S}_{\ell}[w]$ and $u \mapsto Q_{j}[u]$, respectively, and $\phi_{\mathcal{T}}, \phi_{\mathcal{R}}, \phi_{\mathcal{P}}$ are the support functions of the corresponding sets $\mathcal{T}, \mathcal{R}, \mathcal{P}$.

The optimal value in (4.80) is an efficiently computable upper bound on the quantity $Opt_{\#}$ defined in (4.79), and this bound is tight within factor

$$2\max[\ln(2D), 1], \ D = \sum_{k} d_k + \sum_{\ell} f_{\ell} + \sum_{j} e_j.$$
(4.81)

Proof of Proposition is the subject of Exercise 4.45.

(4.79)

304

LECTURE 4

4.5.2 Mixed noise

So far, we have considered separately the cases of random and uncertain-butbounded observation noises in (4.35). Note that both these observation schemes are covered by the following "mixed" scheme:

$$\omega = Ax + \xi + \eta, \tag{4.82}$$

where, as above, A is a given $m \times n$ matrix, x us unknown deterministic signal known to belong to a given signal set \mathcal{X} , ξ is random noise with distribution known to belong to a family \mathcal{P} of Borel probability distributions on \mathbf{R}^m satisfying (4.36) for a given convex compact set $\Pi \subset \operatorname{int} \mathbf{S}^m_+$, and η is "uncertain-but-bounded" observation error known to belong to a given set \mathcal{H} . As before, our goal is to recover $Bx \in \mathbf{R}^{\nu}$ via observation ω . In our present situation, given a norm $\|\cdot\|$ on \mathbf{R}^{ν} , we can quantify the performance of a candidate estimate $\omega \mapsto \hat{x}(\omega) : \mathbf{R}^m \to \mathbf{R}^{\nu}$ by its risk

$$\operatorname{Risk}_{\Pi,\mathcal{H},\|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x\in\mathcal{X},P\ll \Pi,\eta\in\mathcal{H}} \mathbf{E}_{\xi\sim P}\{\|Bx - \widehat{x}(Ax + \xi + \eta)\|\}.$$

Observe that the estimation problem associated with "mixed" observation scheme straightforwardly reduces to similar problem for random observation scheme, by the same trick we have used in Section 4.5 to eliminate observation noise at all. Indeed, let us treat $x^+ = [x; \eta] \in \mathcal{X}^+ := \mathcal{X} \times \mathcal{H}$ and \mathcal{X}^+ as the new signal/signal set underlying our observation, and set $\bar{A}x^+ = Ax + \eta$, $\bar{B}x^+ = Bx$, where $x^+ = [x; \eta]$. With these conventions, the "mixed" observation scheme reduces to

$$\omega = \bar{A}x^+ + \xi,$$

and for every candidate estimate $\hat{x}(\cdot)$ it clearly holds

$$\operatorname{Risk}_{\Pi,\mathcal{H},\|\cdot\|}[\widehat{x}|\mathcal{X}] = \operatorname{Risk}_{\Pi,\|\cdot\|}[\widehat{x}|\mathcal{X}^+],$$

and we arrive at the situation of Section 4.3.3.1. Assuming that \mathcal{X} and \mathcal{H} are spectratopes, so is \mathcal{X}^+ , meaning that all results of Section 4.3.3 on building presumably good linear estimates and their near-optimality are applicable to our present setup.

4.6 BEYOND THE SCOPE OF LINEAR ESTIMATION: POLYHEDRAL ESTIMATE

4.6.1 Motivation

So far, in this Lecture we were considering the estimation problem as follows:

We want to recover the image $Bx \in \mathbf{R}^{\nu}$ of unknown signal x known to belong to signal set $\mathcal{X} \subset \mathbf{R}^n$ from a noisy observation

$$\omega = Ax + \xi_x \in \mathbf{R}^m,$$

where ξ_x is observation noise; index $_x$ in ξ_x indicates that the distribution P_x of the observation noise may depend on x. Here \mathcal{X} is a given nonempty

convex compact set, and A and B are given $m \times n$ and $\nu \times n$ matrices; in addition, we are given a norm $\|\cdot\|$ on \mathbf{R}^{ν} in which the recovery error is measured.

We have seen that if \mathcal{X} is an ellitope/spectratope, then, under reasonable assumptions on observation noise and $\|\cdot\|$, an appropriate efficiently computable *linear in* ω estimate is near-optimal. Note that the ellitopic/spectratopic structure of \mathcal{X} is crucial here. What follows is motivated by the desire to build an alternative estimation scheme which works beyond the ellitopic/spectratopic case, where linear estimates can become "heavily nonoptimal."

Motivating example. Consider the simply-looking problem of recovering Bx = xin $\|\cdot\|_2$ -norm from *direct* observations (Ax = x) corrupted by the standard Gaussian noise $\xi \sim \mathcal{N}(0, \sigma^2 I)$, and let \mathcal{X} be the unit $\|\cdot\|_1$ =ball:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \sum_i |x_i| \le 1 \}.$$

In this situation, building the optimal, in terms of the worst-case, over $x \in \mathcal{X}$, expected squared risk linear estimate $\widehat{x}_H(\omega) = H^T \omega$ is extremely simple:

$$\operatorname{Risk}^{2}[\widehat{x}_{H}|\mathcal{X}] := \max_{x \in \mathcal{X}} \mathbf{E} \left\{ \|\widehat{x}_{H}(\omega - Bx\|_{2}^{2} \right\} \\ = \max_{x \in \mathcal{X}} \left\{ \|[I - H^{T}]x\|_{2}^{2} + \sigma^{2}\operatorname{Tr}(HH^{T}) \right\} \\ = \max_{i < n} \|\operatorname{Col}_{i}[I - H^{T}]\|_{2}^{2} + \sigma^{2}\operatorname{Tr}(HH^{T}).$$

Clearly, the optimal H is just a scalar matrix hI, the optimal h is the minimizer of the univariate quadratic function $(1-h)^2 + \sigma^2 nh^2$, and the best achievable with linear estimates squared risk is

$$R^{2} = \min_{h} \left[(1-h)^{2} + \sigma^{2} n h^{2} \right] = \frac{n \sigma^{2}}{1 + n \sigma^{2}}.$$

On the other hand, consider *nonlinear* estimate as follows. Given σ , "safety factor" $\rho \geq 1$ and observation ω , specify the estimate $\hat{x}(\omega)$ as an optimal solution to the optimization problem

$$Opt(\omega) = \min_{y \in \mathcal{X}} \|y - \omega\|_{\infty}.$$

Note that the probability for the true signal to satisfy $||x - \omega||_{\infty} \leq \rho\sigma$ is at least $1 - p, p = 2n \exp\{-\rho^2/2\}$, and if this event \mathcal{E} happens, then both x and \hat{x} belong to the box $\{y : ||y - \omega||_{\infty} \leq \rho\sigma\}$, implying that $||x - \hat{x}||_{\infty} \leq 2\rho\sigma$; in addition, we always have $||x - \hat{x}||_2 \leq ||x - \hat{x}||_1 \leq 2$, since $x \in \mathcal{X}$ and $\hat{x} \in \mathcal{X}$. We therefore have

$$\|x - \widehat{x}\|_{2} \leq \sqrt{\|x - \widehat{x}\|_{\infty} \|x - \widehat{x}\|_{1}} \leq \begin{cases} 2\sqrt{\rho\sigma}, & \omega \in \mathcal{E} \\ 2, & \omega \notin \mathcal{E} \end{cases},$$

whence

$$\mathbf{E}\left\{\|\widehat{x} - x\|_{2}^{2}\right\} \le 4\rho\sigma + 4p \le 4\rho\sigma + 8n \exp\{-\rho^{2}/2\}.$$
 (*)

Assuming $\sigma \leq 2n \exp\{-1/2\}$ and specifying ρ as $\sqrt{2\ln(2n/\sigma)}$, we get $\rho \geq 1$ and $2n \exp\{-\rho^2/2\} \leq \sigma$, implying that the right hand side in (*) is at most $8\rho\sigma$. In

other words, for our nonlinear estimate $\hat{x}(\omega)$ it holds

$$\operatorname{Risk}^{2}[\widehat{x}|\mathcal{X}] \leq 8\sqrt{\ln(2n/\sigma)}\sigma.$$

When $n\sigma^2$ is of order of 1, the latter bound on the squared risk is of order of $\sigma\sqrt{\ln(1/\sigma)}$, while the best squared risk achievable with linear estimates under the circumstances is of order of 1. We conclude that when σ is small and n is large (specifically, is of order of $1/\sigma^2$), the best linear estimate is *by far* inferior as compared to our nonlinear estimate – the ratio of the corresponding squared risks is as large as $\frac{O(1)}{\sigma\sqrt{\ln(1/\sigma)}}$ – the factor which is "by far" worse than the nonoptimality factor in the case of ellitope/spectratope \mathcal{X} .

The construction of the nonlinear estimate \hat{x} which we have built⁶¹ admits a natural extension yielding what we shall call *polyhedral estimate*, and our present goal is to design and to analyse presumably good polyhedral estimate.

4.6.2 Generic polyhedral estimate

A generic polyhedral estimate is as follows:

Given the data $A \in \mathbf{R}^{m \times n}, B \in \mathbf{R}^{\nu \times n}, \mathcal{X} \subset \mathbf{R}^n$ of our recovery problem (where \mathcal{X} is a computationally tractable convex compact set) and a "reliability tolerance" $\epsilon \in (0, 1)$, we specify somehow positive integer N and Nlinear forms $h_{\ell}^T z$ on the space \mathbf{R}^m where observations live. These forms define linear forms $g_{\ell}^T x := h_{\ell}^T A x$ on the space of signals \mathbf{R}^n . Assuming that the observation noise ξ_x is zero mean for every $x \in \mathcal{X}$, the "plug-in" estimates $h_{\ell}^T \omega$ are unbiased estimates of the forms $g_i^T x$. Assume that the vectors h_{ℓ} are selected in such a way that

$$\forall (x \in \mathcal{X}) : \operatorname{Prob}\{|h_{\ell}^{T}\xi_{x}| > 1\} \le \epsilon/N.$$

$$(4.83)$$

In this situation, setting $H = [h_1, ..., h_N]$ (in the sequel, H is referred to as *contrast matrix*), we can be sure that whatever be signal $x \in \mathcal{X}$ underlying our observation $\omega = Ax + \xi_x$, the observable vector $H^T \omega$ satisfies the relation

$$\operatorname{Prob}\left\{\|H^T\omega - H^TAx\|_{\infty} > 1\right\} \le \epsilon.$$

$$(4.84)$$

With the polyhedral estimation scheme, we act as if all information about x contained in our observation ω were represented by $H\omega$, and we estimate Bx by $B\bar{x}$, where $\bar{x} = \bar{x}(\omega)$ is a (whatever) vector from \mathcal{X} compatible with this information, specifically, such that

$$\|H^T \omega - H^T A \bar{x}\|_{\infty} \le 1.$$

Note that while the latter relation, up to probability mass ϵ of "bad noise realizations," is satisfied by the signal x underlying observation ω , in which case the problem of finding $\bar{x} \in \mathcal{X}$ satisfying this relation is feasible, this problem with positive probability could be infeasible. To circumvent this

⁶¹ in fact, this estimate is nearly optimal under the circumstances in a meaningful range of values of n and σ).

difficulty, let us define \bar{x} as

$$\bar{x} \in \underset{u}{\operatorname{Argmin}} \left\{ \| H^T \omega - H^T A u \|_{\infty} : u \in \mathcal{X} \right\}$$

$$(4.85)$$

so that \bar{x} always is well defined and belongs to \mathcal{X} , and estimate Bx by $B\bar{x}$.

A polyhedral estimate is specified by $m \times N$ contrast matrix $H = [h_1, ..., h_N]$ with columns h_ℓ satisfying (4.83) and is as follows: given observation ω , we build $\bar{x} = \bar{x}(\omega) \in \mathcal{X}$ according to (4.85) and estimate Bx by $\hat{x}^H(\omega) = B\bar{x}(\omega)$.

The rationale behind polyhedral estimation scheme is the desire to reduce complex estimating problems to those of estimating linear forms. To the best of our knowledge, this approach was first used in [114, Section 2] in connection with recovering (restrictions on regular grids of) multivariate functions from Sobolev balls from direct observations.

 $(\epsilon, \|\cdot\|)$ -risk. Given a desired "reliability tolerance" $\epsilon \in (0, 1)$, it is convenient to quantify the performance of polyhedral estimate by its $(\epsilon, \|\cdot\|)$ -risk

$$\operatorname{Risk}_{\epsilon,\|\cdot\}}[\widehat{x}(\cdot)|\mathcal{X}] = \inf\left\{\rho: \operatorname{Prob}\left\{\|Bx - \widehat{x}(Ax + \xi_x)\| > \rho\right\} \le \epsilon \,\forall x \in \mathcal{X}\right\}, \quad (4.86)$$

that is, the worst, over $x \in \mathcal{X}$, size of " $(1 - \epsilon)$ -reliable $\|\cdot\|$ -confidence interval" associated with a candidate estimate $\hat{x}(\cdot)$.

An immediate observation is as follows:

Proposition 4.23. In the situation in question, denoting by $\mathcal{X}_s = \frac{1}{2}(\mathcal{X} - \mathcal{X})$ the symmeterization of \mathcal{X} , given a contrast matrix $H = [h_1, ..., h_N]$ with columns satisfying (4.83) the quantity

$$\mathfrak{R}[H] = \max_{z} \left\{ \|Bz\| : \|H^T Az\|_{\infty} \le 2, z \in 2\mathcal{X}_{\mathrm{s}} \right\}$$

$$(4.87)$$

is an upper bound on the $(\epsilon, \|\cdot\|)$ -risk of the polyhedral estimate $\widehat{x}^{H}(\cdot)$:

$$\operatorname{Risk}_{\epsilon, \|\cdot\|} [\widehat{x}^H | \mathcal{X}] \le \mathfrak{R}[H].$$
(4.88)

Proof is immediate. Let us fix $x \in \mathcal{X}$, and let \mathcal{E} be the set of all realizations of ξ_x such that $||H^T\xi_x||_{\infty} \leq 1$, so that $P_x(\mathcal{E}) \geq 1 - \epsilon$ by (4.84). Let us fix a realization $\xi \in \mathcal{E}$ of the observation noise, and let $\omega = Ax + \xi$, $\bar{x} = \bar{x}(Ax + \xi)$. Then u = x is a feasible solution to the optimization problem (4.85) with the value of the objective ≤ 1 , implying that the value of this objective at the optimal solution \bar{x} to the problem is ≤ 1 as well, so that $||H^TA[x-\bar{x}]||_{\infty} \leq 2$. Besides this, $z = x - \bar{x} \in 2\mathcal{X}_s$. We see that z is a feasible solution to (4.87), whence $||B[x-\bar{x}]|| = ||Bx - \hat{x}(\omega)|| \leq \Re[H]$. It remains to note that the latter relation holds true whenever $\omega = Ax + \xi$ with $\xi \in \mathcal{E}$, and the P_x -probability of the latter inclusion is at least $1 - \epsilon$, whatever be $x \in \mathcal{X}$.

What is ahead. The basic questions associated with the design of polyhedral estimates are as follows:

1. Given the data of our estimation problem and a tolerance $\delta \in (0, 1)$, how to find

308

LECTURE 4

a set \mathcal{H}_{δ} of vectors $h \in \mathbf{R}^m$ satisfying the relation

$$\forall (x \in \mathcal{X}) : \operatorname{Prob}\left\{ |h^T \xi_x| > 1 \right\} \le \delta.$$

$$(4.89)$$

With our approach, after the decision on the number N of columns in a contrast matrix has been made, we are free to select these columns as we want from the set \mathcal{H}_{δ} , with $\delta = \epsilon/N$, ϵ being a given reliability tolerance of the estimate we are designing. Thus, the larger is \mathcal{H}_{δ} , the better for us.

- 2. The upper bound $\mathfrak{R}[H]$ on the $(\epsilon, \|\cdot\|)$ -risk of the polyhedral estimate \hat{x}^H is, in general, difficult to compute this is the maximum of a convex function over a computationally tractable convex set. Thus, similarly to the case of linear estimates, we need techniques for computationally efficient upper bounding of $\mathfrak{R}[\cdot]$.
- 3. With "raw materials" (sets \mathcal{H}_{δ}) and efficiently computable upper bounds on the risk of candidate polyhedral estimates at our disposal, how to design the best, in terms of (the upper bound on) its risk, polyhedral estimate?

We are about to consider these questions one by one.

4.6.3 Specifying sets \mathcal{H}_{δ} for basic observation schemes

Seemingly the only way to specify reasonable sets \mathcal{H}_{δ} goes via making assumptions on the distributions of observation noises we want to handle. In the sequel we restrict ourselves with 3 special cases as follows:

- Sub-Gaussian case: For every $x \in \mathcal{X}$, the observation noise ξ_x is sub-Gaussian with parameters $(0, \sigma^2 I_m)$, where $\sigma > 0$.
- Discrete case: \mathcal{X} is a convex compact subset of the probabilistic simplex $\Delta_n = \{x \in \mathbf{R}^n : x \ge 0, \sum_i x_i = 1\}$, A is column-stochastic matrix, and

$$\omega = \frac{1}{K} \sum_{k=1}^{K} \zeta_k$$

with independent across $k \leq K$ random vectors ζ_k , with ζ_k taking value e_i with probability $[Ax]_i$, i = 1, ..., m, e_i being the basic orths in \mathbf{R}^m .

• Poisson case: \mathcal{X} is a convex compact subset of the nonnegative orthant \mathbf{R}^n_+ , A is entrywise nonnegative, and the observation ω stemming from $x \in \mathcal{X}$ is random vector with independent across i entries $\omega_i \sim \text{Poisson}([Ax]_i]$.

The associated sets \mathcal{H}_{δ} can be built as follows.

4.6.3.1 Sub-Gaussian case

When $h \in \mathbf{R}^n$ is deterministic and ξ is sub-Gaussian with parameters $0, \sigma^2 I_n$, we have

$$\operatorname{Prob}\{|h^T\xi| > 1\} \le 2\exp\{-\frac{1}{2\sigma^2 \|h\|_2^2}\}.$$

Indeed, when $h \neq 0$ and $\gamma > 0$, we have

$$\operatorname{Prob}\{h^{T}\xi > 1\} \leq \exp\{-\gamma\} \mathbf{E}\left\{\exp\{\gamma h^{T}\xi\}\right\} \leq \exp\{\frac{1}{2}\sigma^{2}\gamma^{2}\|h\|_{2}^{2} - \gamma\}$$
minimizing the resulting bound in $\gamma > 0$, we get $\operatorname{Prob}\{h^T \xi > 1\} \leq \exp\{-\frac{1}{2\|h\|_2^2 \sigma^2}\}$; the same reasoning as applied to -h in the role of h results in $\operatorname{Prob}\{h^T \xi < -1\} \leq \exp\{-\frac{1}{2\|h\|_2^2 \sigma^2}\}$.

Consequently

$$p_G(h) := \underbrace{\sigma\sqrt{2\ln(2/\delta)}}_{\vartheta_G} \|h\|_2 \le 1 \Rightarrow \operatorname{Prob}\{|h^T\xi| > 1\} \le \delta, \tag{4.90}$$

and we can set

$$\mathcal{H}_{\delta} = \mathcal{H}^G_{\delta} := \{h : p_G(h) \le 1\}.$$
(4.91)

4.6.3.2 Discrete case

Given $x \in \mathcal{X}$, setting $\mu = Ax$ and $\eta_k = \zeta_k - \mu$, we get

$$\omega = Ax + \underbrace{\frac{1}{K} \sum_{k=1}^{K} \eta_k}_{\xi_x}.$$

Given $h \in \mathbf{R}^m$,

$$h^T \xi_x = \frac{1}{K} \sum_k \underbrace{h^T \eta_k}_{\chi_k}.$$

Random variables $\chi_1, ..., \chi_K$ are independent zero mean and clearly satisfy

$$\mathbf{E}\left\{\chi_{k}^{2}\right\} \leq \sum_{i} [Ax]_{i}h_{i}^{2}, \ |\chi_{k}| \leq 2||h||_{\infty}.$$

Applying Bernstein Inequality 62 , we get

$$\operatorname{Prob}\{|h^{T}\xi_{x}| > 1\} = \operatorname{Prob}\{\sum_{k} \chi_{k} > K\} \le \exp\{-\frac{K}{2\sum_{i} [Ax]_{i}h_{i}^{2} + \frac{4}{3}||h||_{\infty}}\}.$$
 (4.92)

Setting

$$p_D(h) = \sqrt{\vartheta_D^2 \max_{x \in \mathcal{X}} \sum_i [Ax]_i h_i^2 + \varrho_D^2 \|h\|_{\infty}^2},$$

$$\vartheta_D = 2\sqrt{\frac{2\ln(2/\delta)}{K}},$$

$$\varrho_D = \frac{8\ln(2/\delta)}{3K},$$
(4.93)

we, after completely straightforward computation, conclude from (4.92) that

$$p_D(h) \le 1 \Rightarrow \operatorname{Prob}\{|h^T \xi_x| > 1\} \le \delta, \ \forall x \in \mathcal{X}.$$
 (4.94)

$$\operatorname{Prob}\{X_1 + \dots + X_k > t\} \le \exp\{-\frac{t^2}{2[\sum_k \sigma_k^2 + \frac{1}{3}Mt]}\}.$$

⁶²Classical Bernstein Inequality states that if $X_1, ..., X_K$ are independent zero mean scalar random variables with finite variances σ_k^2 such that $|X_k| \leq M$ a.s., then for every t > 0 one has

310

LECTURE 4

Thus, in the Discrete case we can set

$$\mathcal{H}_{\delta} = \mathcal{H}_{\delta}^{D} := \{h : p_{D}(h) \le 1\}.$$

$$(4.95)$$

4.6.3.3 Poisson case

In the Poisson case, for $x \in \mathcal{X}$, setting $\mu = Ax$, we have

$$\omega = Ax + \xi_x, \, \xi_x = \omega - \mu.$$

It turns out that for every $h \in \mathbf{R}^m$ one has

$$\forall t \ge 0 : \operatorname{Prob}\left\{ |h^T \xi_x| \ge t \right\} \le 2 \exp\{-\frac{t^2}{3[\sum_i h_i^2 \mu_i + \|h\|_{\infty} t]} \}$$
(4.96)

(for verification, see Section 4.10.7). As a result, we conclude via a straightforward computation that setting

$$p_P(h) = \sqrt{\vartheta_P^2 \max_{x \in \mathcal{X}} \sum_i [Ax]_i h_i^2 + \varrho_P^2 \|h\|_{\infty}^2}, \vartheta_P = \sqrt{6 \ln(2/\delta)}, \varrho_P = 6 \ln(2/\delta),$$
(4.97)

we ensure that

$$p_P(h) \le 1 \Rightarrow \operatorname{Prob}\{|h^T \xi_x| > 1\} \le \delta, \ \forall x \in \mathcal{X}.$$
 (4.98)

Thus, in the Poisson case we can set

$$\mathcal{H}_{\delta} = \mathcal{H}_{\delta}^{P} := \{h : p_{P}(h) \le 1\}.$$

$$(4.99)$$

4.6.4 Efficient upper-bounding of $\Re[H]$ and Contrast Design, I.

The scheme for upper-bounding $\Re[H]$ to be presented in this Section (an alternative, completely different, scheme will be presented in Section 4.6.5) is motivated by what happens in our Motivating example. Namely, there is a special case of (4.87) where $\Re[H]$ is easy to compute – the case when $\|\cdot\|$ is the uniform norm $\|\cdot\|_{\infty}$, whence

$$\Re[H] = \widehat{\Re}[H] := 2 \max_{i \le \nu} \max_{x} \left\{ \operatorname{Row}_{i}^{T}[B]x : x \in \mathcal{X}_{s}, \|H^{T}Ax\|_{\infty} \le 1 \right\}$$

is just the maximum of ν efficiently computable convex functions. It turns out that when $\|\cdot\| = \|\cdot\|_{\infty}$, it is easy not only to compute $\Re[H]$, but to optimize this risk bound in H as well⁶³. These observations underly the forthcoming developments in this Section: under appropriate assumptions, we bound the risk of a polyhedral estimate stemming from a contrast matrix H via the efficiently computable quantity $\widehat{\Re}[H]$ and then show that the resulting risk bounds can be efficiently optimized w.r.t. H. We shall also see that in some "simple for analytical analysis" situations, like the one of Motivating example, the resulting estimates turn out to be nearly

⁶³On a closest inspection, in the situation of Motivating example the $\|\cdot\|_{\infty}$ -optimal contrast matrix H is proportional to the unit matrix, and the quantity $\widehat{\mathfrak{R}}[H]$ can be translated into upper bound on, say, $\|\cdot\|_2$ -risk of the associated polyhedral estimate; these are the estimate and the risk bound we dealt with when discussing Motivating example.

minimax optimal.

4.6.4.1 Assumptions

We continue to stay within the setup introduced in Section 4.6.1 which we now augment with the following assumptions:

- **A.1.** $\|\cdot\| = \|\cdot\|_r$ with $r \in [1, \infty]$.
- **A.2.** We have at our disposal a sequence $\gamma = \{\gamma_i > 0, i \leq \nu\}$ and $\rho \in [1, \infty]$ such that the image of \mathcal{X}_s under the mapping $x \mapsto Bx$ is contained in the "scaled $\|\cdot\|_{\rho}$ -ball"

$$\mathcal{Y} = \{ y \in \mathbf{R}^{\nu} : \| \operatorname{Diag}\{\gamma\} y \|_{\rho} \le 1 \}.$$

$$(4.100)$$

4.6.4.2 Simple observation

Let B_{ℓ}^T be ℓ -th row in $B, 1 \leq \ell \leq \nu$. Let us make the following observation:

Proposition 4.24. In the situation described in Section 4.6.1 and with Assumptions A.1-2 in force, let $\epsilon \in (0,1)$ and a positive real $N \ge \nu$ be given, and let $p(\cdot)$ be a norm on \mathbb{R}^m such that

$$\forall (h: p(h) \le 1, x \in \mathcal{X}) : \operatorname{Prob}\{|h^T \xi_x| > 1\} \le \epsilon/N.$$
(4.101)

Let, next, a matrix $H = [H_1, ..., H_{\nu}]$ with $H_{\ell} \in \mathbb{R}^{m \times m_{\ell}}$, $m_{\ell} \ge 1$, and positive reals $\varsigma_{\ell}, \ell \le \nu$, satisfy the relations

(a)
$$p(\operatorname{Col}_{j}[H]) \leq 1, 1 \leq j \leq N;$$

(b) $\max_{x} \{ B_{\ell}^{T} x : x \in \mathcal{X}_{s}, \| H_{\ell}^{T} A x \|_{\infty} \leq 1 \} \leq \varsigma_{\ell}, 1 \leq \ell \leq \nu.$
(4.102)

Then the quantity $\Re[H]$ as defined in (4.87) can be upper-bounded as follows:

$$\Re[H] \le \Psi(\varsigma) := 2 \max_{w} \left\{ \| [w_1/\gamma_1; ...; w_{\nu}/\gamma_{\nu}] \|_r : \| w \|_{\rho} \le 1, 0 \le w_{\ell} \le \gamma_{\ell} \varsigma_{\ell}, \ell \le \nu \right\},$$
(4.103)

which combines with Proposition 4.23 to imply that

$$\operatorname{Risk}_{\epsilon,\|\cdot\|}[\widehat{x}^H|\mathcal{X}] \le \Psi(\varsigma). \tag{4.104}$$

Function Ψ is a nondecreasing on the nonnegative orthant and is easy to compute.

Proof. Let $z = 2\bar{z}$ be a feasible solution to (4.87), so that $\bar{z} \in \mathcal{X}_{s}$ and $\|H^{T}A\bar{z}\|_{\infty} \leq 1$. Let $y = B\bar{z}$, so that $y \in \mathcal{Y}$ (see (4.100)) due to $\bar{z} \in \mathcal{X}_{s}$ and **A.2**. Thus, $\|\text{Diag}\{\gamma\}y\|_{p} \leq 1$. Besides this, by (4.102.b) relations $\bar{z} \in \mathcal{X}_{s}$ and $\|H^{T}A\bar{z}\|_{\infty} \leq 1$ combine with the symmetry of \mathcal{X}_{s} w.r.t. the origin to imply that

$$|y_{\ell}| = |B_{\ell}^T \bar{z}| \le \varsigma_{\ell}, \ \ell \le \nu.$$

Taking into account that $\|\cdot\| = \|\cdot\|_r$ by **A.1**, we see that

 $\begin{aligned} \Re[H] &= \max_{z} \left\{ \|Bz\|_{r} : z \in 2\mathcal{X}_{s}, \|H^{T}Az\|_{\infty} \leq 2 \right\} \\ &\leq 2 \max_{y} \left\{ \|y\|_{r} : |y_{\ell}| \leq \varsigma_{\ell}, \ell \leq \nu, \ \& \ \|\text{Diag}\{\gamma\}y\|_{\rho} \leq 1 \right\} \\ &= 2 \max_{w} \left\{ \|[w_{1}/\gamma_{1}; ...; w_{\nu}/\gamma_{\nu}]\|_{r} : \|w\|_{\rho} \leq 1, 0 \leq w_{\ell} \leq \gamma_{\ell}\varsigma_{\ell}, \ell \leq \nu \right\}, \end{aligned}$

as stated in (4.103).

The fact that Ψ is nondecreasing on the nonnegative orthant is evident. Computing Ψ can be carried out as follows:

- 1. When $r = \infty$, we need to compute $\max_{\ell \leq \nu} \max_{w} \{w_{\ell}/\gamma_{\ell} : \|w\|_{\rho} \leq 1, 0 \leq w_{j} \leq \gamma_{j}\varsigma_{j}, j \leq \nu\}$, so that computing Ψ reduces to solving ν simple convex optimization problems;
- 2. When $\rho = \infty$, we clearly have $\Psi(\varsigma) = \|[\bar{w}_1/\gamma_1; ...; \bar{w}_\nu/\gamma_\nu]\|_r$, $\bar{w}_\ell = \min[1, \gamma_\ell \varsigma_\ell]$;
- 3. When $1 \leq r, \rho < \infty$, passing from variables w_{ℓ} to variables $u_{\ell} = w_{\ell}^{\rho}$, we get

$$\Psi^{r}(\varsigma) = 2^{r} \max_{u} \left\{ \sum_{\ell} \gamma_{\ell}^{-r} u_{\ell}^{r/\rho} : \sum_{\ell} u_{\ell} \le 1, 0 \le u_{\ell} \le (\gamma_{\ell}\varsigma_{\ell})^{\rho} \right\}.$$

When $r \leq \rho$, the right hand side problem here is the easily solvable problem of maximizing a simple concave function over a simple convex compact set. When $\infty > r > \rho$, the right hand side problem can be solved by Dynamic Programming. \Box

4.6.4.3 Specifying contrasts

Risk bound (4.104) allows for an easy design of contrast matrices. Recalling that Ψ is monotone on the nonnegative orthant, all we need is to select h_{ℓ} 's satisfying (4.102) and resulting in the smallest possible ς_{ℓ} 's, which is what we are about to do now.

Preliminaries. Given a vector $b \in \mathbf{R}^m$ and a norm $s(\cdot)$ on \mathbf{R}^m , consider convexconcave saddle point problem

$$Opt = \inf_{g \in \mathbf{R}^m} \max_{x \in \mathcal{X}_s} \phi(g, x) := [b - A^T g]^T x + s(g)$$
(SP)

along with the induced primal and dual problems

where $q(\cdot)$ is the norm conjugate to $s(\cdot)$ (we have used the evident fact that $\inf_{g\in \mathbf{R}^m}[f^Tg+s(g)]$ is either $-\infty$ or 0 depending on whether q(f) > 1 or $q(f) \leq 1$). Since \mathcal{X}_s is compact, we have $\operatorname{Opt}(P) = \operatorname{Opt}(D) = \operatorname{Opt}$ by Sion-Kakutani Theorem. Besides this, (D) is solvable (evident) and (P) is solvable as well, since $\overline{\phi}(g)$ is continuous due to the compactness of \mathcal{X}_s and $\overline{\phi}(g) \geq s(g)$, so that $\overline{\phi}(\cdot)$ has bounded level sets. Let \overline{g} be an optimal solution to (P), and \overline{x} be an optimal solution to (D), and let \overline{h} be $s(\cdot)$ -unit normalization of \overline{g} , so that $s(\overline{h}) = 1$ and $\overline{g} = s(\overline{g})\overline{h}$. Now let us make the observation as follows:

Observation 4.25. In the situation in question, we have

$$\max_{x} \left\{ |b^T x| : x \in \mathcal{X}_{\mathrm{s}}, |\bar{h}^T A x| \le 1 \right\} \le \mathrm{Opt}.$$
(4.105)

In addition, whatever be a matrix $G = [g^1, ..., g^M] \in \mathbf{R}^{m \times M}$ with $s(g^j) \leq 1, j \leq M$,

 $one \ has$

$$\max_{x} \left\{ |b^T x| : x \in \mathcal{X}_{s}, \|G^T A x\|_{\infty} \le 1 \right\} = \max_{x} \left\{ b^T x : x \in \mathcal{X}_{s}, \|G^T A x\|_{\infty} \le 1 \right\} \ge \text{Opt.}$$
(4.106)

Proof. Let x be a feasible solution to the left hand side problem in (4.105). Replacing, if necessary, x with -x, we can assume that $|b^T x| = b^T x$. We now have

$$|b^T x| = b^T x = [\bar{g}^T A x - s(\bar{g})] + \underbrace{[b - A^T \bar{g}]^T x + s(\bar{g})}_{\leq \bar{\phi}(\bar{g})} \leq \operatorname{Opt}(P) + [s(\bar{g})\bar{h}^T A x - s(\bar{g})]$$
$$\leq \operatorname{Opt}(P) + s(\bar{g})\underbrace{|\bar{h}^T A x|}_{\leq 1} - s(\bar{g}) \leq \operatorname{Opt}(P) = \operatorname{Opt},$$

as claimed in (4.105). Now, the equality in (4.106) is due to the symmetry of \mathcal{X}_{s} w.r.t. the origin. To verify inequality in (4.106), note that \bar{x} satisfies the relations $\bar{x} \in \mathcal{X}_{s}$ and $q(A\bar{x}) \leq 1$, implying, due to the fact that the columns of G are of $s(\cdot)$ -norm ≤ 1 , that \bar{x} is a feasible solution to the optimization problems in (4.106). As a result, the second quantity in (4.106) is at least $b^T \bar{x} = \text{Opt}(D) = \text{Opt}$, and (4.106) follows.

Designing contrast. With upper-bounding the risk of a polyhedral estimate via Proposition 4.24, Observation 4.25 basically resolves the the associated contrast design problem, at least in the case of Sub-Gaussian, Discrete, and Poisson observation schemes. Indeed, in these cases, when designing contrast matrix with Ncolumns, with our approach we are supposed to select its columns in the respective sets $\mathcal{H}_{\epsilon/N}$, see Section 4.6.3. Note that these sets, while shrinking as N grows, are "nearly independent" of N, since the norms p_G , p_D , p_P participating in the description of the respective sets \mathcal{H}^G_{δ} , \mathcal{H}^D_{δ} , \mathcal{H}^P_{δ} depend on $1/\delta$ just via logarithmic in $1/\delta$ factors. It follows that we lose nearly nothing when assuming that $N \geq \nu$. Let us act as follows:

We set $N = \nu$, specify $\bar{p}(\cdot)$ as the norm $(p_G, \text{ or } p_D, \text{ or } p_P)$ associated with the observation scheme (Sub-Gaussian, or Discrete, or Poisson) in question and with $\delta = \epsilon/\nu$, and solve ν convex optimization problems

$$\begin{array}{lll}
\operatorname{Opt}_{\ell} &=& \min_{g \in \mathbf{R}^m} \left[\overline{\phi}_{\ell}(g) := \max_{x \in \mathcal{X}_{\mathrm{s}}} \phi_{\ell}(g, x) \right] & (P_{\ell}) \\
\phi_{\ell}(g, x) &=& [B_{\ell} - A^T g]^T x + \overline{p}(g)
\end{array}$$

Next, we convert optimal solution g_{ℓ} to (P_{ℓ}) into vector $h_{\ell} \in \mathbb{R}^{m}$ by representing $g_{\ell} = \bar{p}(g_{\ell})h_{\ell}$ with $\bar{p}(h_{\ell}) = 1$, and set $H_{\ell} = h_{\ell}$. As a result, we get $m \times \nu$ contrast matrix $H = [h_{1}, ..., h_{\nu}]$ which, taken along with $N = \nu$, the quantities

$$\varsigma_{\ell} = \operatorname{Opt}_{\ell}, \ 1 \le \ell \le \nu, \tag{4.107}$$

and with $p(\cdot) \equiv \bar{p}(\cdot)$, in view of Observation 4.25 as applied with $s(\cdot) \equiv \bar{p}(\cdot)$ satisfies the premise of Proposition 4.24.

Consequently, by Proposition 4.24 we have

$$\operatorname{Risk}_{\epsilon, \|\cdot\|} [\widehat{x}^H | \mathcal{X}] \le \Psi([\operatorname{Opt}_1; ...; \operatorname{Opt}_{\nu}]).$$
(4.108)

Within the framework set by Proposition 4.24, optimality of the outlined contrast

design stems from Observation 4.25 which states that when $N \ge \nu$ and the columns of $m \times N$ contrast matrix $H = [H_1, ..., H_\nu]$ belong to the set $\mathcal{H}_{\delta/N}$ associated with the observation scheme in question (Sub-Gaussian, or Discrete, or Poisson), i.e., the norm $p(\cdot)$ in Proposition is the norm p_G , or p_D , or p_P associated with $\delta = \epsilon/N$, the quantities ς_ℓ participating in (4.102.b) cannot be less than Opt_{ℓ} .

Indeed, we are in the situation when the norm $p(\cdot)$ from Proposition 4.24 is \geq the norm $\bar{p}(\cdot)$ participating in (P_{ℓ}) (since $\epsilon/N \leq \epsilon/\nu$), implying, by (4.102.*a*), that the columns of matrix H obeying the premise of Proposition satisfy the relation $\bar{p}(\operatorname{Col}_{j}[H]) \leq 1$. Invoking the second part of Observation 4.25 with $s(\cdot) \equiv \bar{p}(\cdot)$, $b = B_{\ell}$, and $G = H_{\ell}$, and looking at (4.102.*b*), we conclude that $\varsigma_{\ell} \geq \operatorname{Opt}_{\ell}$ for all ℓ , as claimed.

Since the bound on the risk of a polyhedral estimate offered by Proposition 4.24 is the better the less are ς_{ℓ} 's, we see that as far as this bound is concerned, the outlined design procedure is the best possible, provided $N \ge \nu$.

An attractive feature of the contrast design we have just presented is that it is completely independent of the entities participating in assumptions A.1-2 – these entities affect theoretical risk bounds of the resulting polyhedral estimate, but not the estimate itself.

4.6.4.4 Illustration: Diagonal case

Let us consider the *diagonal case* of our estimation problem, where

- $\mathcal{X} = \{x \in \mathbf{R}^n : ||Dx||_{\rho} \leq 1\}$, where *D* is diagonal matrix with positive diagonal entries $D_{\ell\ell} =: \delta_{\ell}$,
- $m = \nu = n$, and A and B are diagonal matrices with diagonal entries $0 < A_{\ell\ell} =: \alpha_{\ell}, 0 < B_{\ell\ell} =: \beta_{\ell},$
- $\bullet \|\cdot\| = \|\cdot\|_r,$
- We are in Sub-Gaussian case, that is, observation noise ξ_x is $(0, \sigma^2 I_n)$ -sub-Gaussian for every $x \in \mathcal{X}$.

Let us implement the approach developed so far.

1. Given reliability tolerance ϵ , we set

$$\delta = \epsilon/n, \ \mathcal{H} = \mathcal{H}_{\delta}^{G} = \{h \in \mathbf{R}^{n} : p_{G}(h) := \vartheta_{G} \|h\|_{2} \le 1\}, \\ \vartheta_{G} := \sigma \sqrt{2\ln(2/\delta)} = \sigma \sqrt{2\ln(2n/\epsilon)};$$

$$(4.109)$$

2. We solve $\nu = n$ convex optimization problems (P_{ℓ}) associated with $\bar{p}(\cdot) \equiv p_G(\cdot)$, which is immediate: the resulting contrast matrix is just

$$H = \vartheta_G^{-1} I_n,$$

and

$$Opt_{\ell} = \varsigma_{\ell} := \beta_{\ell} \min[\vartheta_G / \alpha_{\ell}, 1/\delta_{\ell}].$$
(4.110)

Risk analysis. The $(\epsilon, \|\cdot\|)$ -risk of the resulting polyhedral estimate $\hat{x}(\cdot)$ can be bounded by Proposition 4.24. Note that setting

$$\gamma_{\ell} = \delta_{\ell} / \beta_{\ell}, \, 1 \le \ell \le n, \tag{4.111}$$

we meet assumptions A.1-2, and the above $H, N = n, \varsigma_{\ell}$ satisfy the premise of Proposition 4.24. By this Proposition,

$$\operatorname{Risk}_{\epsilon, \|\cdot\|_{r}}[\widehat{x}^{H} | \mathcal{X}] \leq \Psi := 2 \max_{w} \left\{ \|[w_{1}/\gamma_{1}; ...; w_{n}/\gamma_{n}]\|_{r} : \|w\|_{\rho} \leq 1, 0 \leq w_{\ell} \leq \gamma_{\ell}\varsigma_{\ell} \right\}.$$
(4.112)

Let us work out what happens in the simple case where

$$1 \le \rho \le r < \infty \qquad (a)$$

 $\alpha_{\ell}/\delta_{\ell} \text{ and } \beta_{\ell}/\alpha_{\ell} \text{ are nonincreasing in } \ell \qquad (b)$

$$(4.113)$$

Proposition 4.26. In the just defined simple case, let n = n when

$$\sum_{i=1}^{n} \left(\vartheta_G \delta_\ell / \alpha_\ell \right)^s \le 1,$$

otherwise let \mathfrak{n} be the smallest integer such that

$$\sum_{i=1}^{n} \left(\vartheta_G \delta_\ell / \alpha_\ell\right)^{\rho} > 1, \tag{4.114}$$

with ϑ_G given by (4.109). Then for the contrast matrix $H = \vartheta_G^{-1} I_n$ we have built one has

$$\operatorname{Risk}_{\epsilon,\|\cdot\|_{r}}[\widehat{x}^{H}|\mathcal{X}] \leq \Psi \leq 2 \left[\sum_{\ell=1}^{n} (\vartheta_{G}\beta_{\ell}/\alpha_{\ell})^{r} \right]^{1/r}$$
(4.115)

Proof. Consider optimization problem specifying Ψ in (4.112). Setting $\theta = r/\rho \ge 1$, let us pass in this problem from variables w_{ℓ} to variables $z_{\ell} = w_{\ell}^{\rho}$, so that

$$\begin{split} \Psi^r &= 2^r \max_z \left\{ \sum_{\ell} z_{\ell}^{\theta} (\beta_{\ell}/\delta_{\ell})^r : \sum_{\ell} z_{\ell} \leq 1, 0 \leq z_{\ell} \leq (\delta_{\ell} \varsigma_{\ell}/\beta_{\ell})^{\rho} \right\} \leq 2^r \Gamma, \\ \Gamma &= \max_z \left\{ \sum_{\ell} z_{\ell}^{\theta} (\beta_{\ell}/\delta_{\ell})^r : \sum_{\ell} z_{\ell} \leq 1, 0 \leq z_{\ell} \leq \chi_{\ell} := (\vartheta_G \delta_{\ell}/\alpha_{\ell})^{\rho} \right\}. \end{split}$$

 Γ is the optimal value in the problem of maximizing a convex (since $\theta \geq 1$) function $\sum_{\ell} z_{\ell}^{\theta} (\beta_{\ell} / \delta_{\ell})^r$ over a bounded polyhedral set, so that the maximum is achieved at an extreme point \bar{z} of the feasible set. The (clearly nonempty) set I of positive entries in \bar{z} , by the standard characterization of extreme points, is as follows: denoting by I' the set of indexes $\ell \in I$ such that \bar{z}_{ℓ} is on its upper bound: $\bar{z}_{\ell} = \chi_{\ell}$, the cardinality |I'| of I' is at least |I| - 1. Since $\sum_{\ell \in I'} \bar{z}_{\ell} = \sum_{\ell \in I'} \chi_{\ell} \leq 1$ and χ_{ℓ} are nondecreasing in ℓ by (4.113.b), we conclude that

$$\sum_{\ell=1}^{|I'|} \chi_\ell \le 1$$

implying that $|I'| < \mathfrak{n}$ provided that $\mathfrak{n} < n$, so that in this case $|I| \leq \mathfrak{n}$; and of course $|I| \leq \mathfrak{n}$ when $\mathfrak{n} = n$. Next, we have

$$\Gamma = \sum_{\ell \in I} \bar{z}_{\ell}^{\theta} (\beta_{\ell}/\delta_{\ell})^r \le \sum_{\ell \in I} \chi_{\ell}^{\theta} (\beta_{\ell}/\delta_{\ell})^r = \sum_{\ell \in I} (\vartheta_G \beta_{\ell}/\alpha_{\ell})^r,$$

and since $\beta_{\ell}/\alpha_{\ell}$ is nonincreasing in ℓ and $|I| \leq \mathfrak{n}$, the latter quantity is at most

316

LECTURE 4

$$\sum_{\ell=1}^{n} (\vartheta_G \beta_\ell / \alpha_\ell)^r.$$

Illustration. Consider the case when

$$0 < \sqrt{\ln(n/\sigma)}\sigma \le 1, \, \alpha_{\ell} = \ell^{-\alpha}, \, \beta_{\ell} = \ell^{-\beta}, \, \delta_{\ell} = \ell^{\delta}$$

$$(4.116)$$

with

$$\beta \ge \alpha \ge 0, \, \delta \ge 0, \, (\beta - \alpha)r < 1. \tag{4.117}$$

In this case for large n, namely,

$$n \ge O(\vartheta_G^{-\frac{1}{\alpha+\delta+1/\rho}}) \qquad \qquad [\vartheta_G = \sqrt{2\ln(2n/\epsilon)}\sigma] \qquad (4.118)$$

(here and in what follows, the factors hidden in $O(\cdot)$ depend solely on $\alpha, \beta, \delta, r, \rho$) we get

$$\mathfrak{n} = O(\vartheta_G^{-\frac{1}{\alpha+\delta+1/\rho}}) = O(\vartheta_G^{-\frac{1}{\bar{\alpha}+1/\rho}}), \ \bar{\alpha} = \alpha + \delta,$$

resulting in

$$\operatorname{Risk}_{\epsilon,\|\cdot\|_{r}}[\widehat{x}|\mathcal{X}] \leq O(\vartheta_{G}^{\frac{\beta+1/\rho-1/r}{\overline{\alpha}+1/\rho}}), \ \bar{\beta} = \beta + \delta.$$
(4.119)

Setting $x = \text{Diag}\{\delta_1^{-1}, ..., \delta_n^{-1}\}y$, $\bar{\alpha} = \alpha + \delta$, $\bar{\beta} = \beta + \delta$ and treating of y, rather than x, as the signal underlying our observation, the problem becomes the similar problem with $\bar{\alpha}, \bar{\beta}, \bar{\delta} = 0, \bar{\mathcal{X}} = \{x : \|x\|_s \leq 1\}$ in the role of $\alpha, \beta, \delta, \mathcal{X}$, respectively, and $\bar{A} = \text{Diag}\{\ell^{-\bar{\alpha}}, \ell \leq n\}, \bar{B} = \text{Diag}\{\ell^{-\bar{\beta}}, \ell \leq n\}$, in the role of A and B. Now, setting

$$\mathfrak{m} = O(\sigma^{-\frac{1}{\bar{\alpha}+1/\rho}}) \approx \mathfrak{n}$$

and strengthening (4.118) to $n \ge O(\sigma^{-\frac{1}{\alpha+\delta+1/\rho}})$, observe that when $\vartheta_G \le O(1)$, $\bar{\mathcal{X}}$ contains the "coordinate box"

$$\widehat{\mathcal{X}} = \{ x : |x_{\ell}| \le \mathfrak{m}^{-1/\rho}, \mathfrak{m}/2 \le \ell \le \mathfrak{m}, x_{\ell} = 0 \text{ otherwise} \}$$

of dimension $\geq \mathfrak{m}/2$ such that $\|\bar{A}y\|_2 \leq O(1)\mathfrak{m}^{-\bar{\alpha}}\|y\|_2$ and $\|\bar{B}y\|_r \geq O(1)\mathfrak{m}^{-\bar{\beta}}\|y\|_r$ when $y \in \hat{X}$. This observation straightforwardly combines with Fano inequality⁶⁴; to imply that when $\epsilon \ll 1$, the minimax optimal, w.r.t. the family of all Borel estimates, the signal set being $\hat{X} \subset \mathcal{X}$, $(\epsilon, \|\cdot\|_r)$ -risk is at least

$$O(\sigma^{rac{areta+1/
ho-1/r}{arlpha+1/
ho}})$$

i.e., is just by a logarithmic in n/ϵ factor better than the upper bound (4.119) on the risk of our polyhedral estimate; thus, in the case under consideration the latter estimate is nearly minimax optimal.

4.6.5 Efficient upper-bounding of $\Re[H]$ and Contrast Design, II.

4.6.5.1 Outline

Below we develop an approach to the design of polyhedral estimates which is an alternative to the one of Section 4.6.4. Our new approach resembles the one we

⁶⁴For the classical Fano inequality see, e.g., [58].

have developed to build presumably good linear estimates; here is a "high level" outline of similarities and differences of what we are about to develop and what was done for linear estimates.

In the simplest case of linear estimation considered in Section 4.2, where the signal set $\mathcal{X} = \mathcal{X}_{s}$ was an ellitope, we represented the (squared) $\|\cdot\|_{2}$ -risk of a candidate linear estimate $\hat{x}_{H}(\omega) = H^{T}\omega$ as the sum of an easy-to-compute convex function $\operatorname{Tr}(H^{T}\Gamma H)$ and the difficult to compute function

$$\Phi(H) = \max_{x \in \mathcal{X}} x^T [B - H^T A]^T [B - H^T A] x.$$

We then used semidefinite relaxation to upper-bound this difficult to compute function by an efficiently computable convex function $\overline{\Phi}(H)$; minimizing the sum of the latter function and $\operatorname{Tr}(H^T\Gamma H)$, we arrived at the desired linear estimate. Thus, the basic technique underlying our design was upper-bounding of the maximum of a quadratic form, depending on H as on parameter, over \mathcal{X} ; what we were seeking for was a bounding scheme allowing to optimize the bound in H efficiently.

Now, the design of a presumably good polyhedral estimate also reduces to minimizing w.r.t. the parameter the maximum of a parametric quadratic form. Specifically, the (upper bound on) the risk $\Re[H]$ of a candidate polyhedral estimate \hat{x}^H given by (4.87) is nothing but

$$\Re[H] = 2 \max_{[u;z]} \left\{ [u;z]^T \underbrace{\left[\frac{1}{2}B^T \right]}_{B_+} [u;z] : \begin{array}{c} u \in \mathcal{B}_*, z \in \mathcal{X}_{\mathrm{s}}, \\ z^T A^T h_\ell h_\ell^T A z \le 1, \ell \le N \end{array} \right\},$$

$$(4.120)$$

where \mathcal{B}_* is the unit ball of the norm conjugate to $\|\cdot\|$ and h_ℓ are the columns of $m \times N$ contrast matrix H. The difference with the design of linear estimates is twofold:

- our "design parameter" the contrast matrix H affects the constraints of the optimization problem specifying $\Re[H]$ rather than the objective of the optimization problem specifying $\Phi(H)$;
- Design of presumably good linear estimate reduces to unconstrained minimization of $\overline{\Phi}(H) + \text{Tr}(H^T\Gamma H)$, while now we need to minimize (efficiently computable upper bound on) $\Re[H]$ under the restriction on H – the columns h_ℓ of this matrix should satisfy (4.83).

The strategy we intend to use in order to handle the above issues can be outlined as follows. Assume we have at our disposal a technique for bounding quadratic forms on the set $\mathcal{B}_* \times \mathcal{X}_s$, so that we have at our disposal an efficiently computable convex function $\mathcal{M}(M)$ on $\mathbf{S}^{n+\nu}$ such that

$$\mathcal{M}(M) \ge \max_{[u;z]\in\mathcal{B}_*\times\mathcal{X}_{\mathrm{s}}} [u;z]^T M[u;z] \ \forall M \in \mathbf{S}^{n+\nu}.$$
(4.121)

When $\lambda \in \mathbf{R}^N_+$, the constraints $z^T A^T h_\ell h_\ell^T A z \leq 1$ in (4.120) can be aggregated to yield the quadratic constraint

$$z^T A^T \Theta_{\lambda} A z \leq \mu_{\lambda}, \ \Theta_{\lambda} = H \text{Diag}\{\lambda\} H^T, \ \mu_{\lambda} = \sum_{\ell} \lambda_{\ell}.$$

Observe that for every $\lambda \geq 0$ we have

$$\Re[H] \le 2\mathcal{M}\left(\underbrace{\left[\frac{\frac{1}{2}B}{\frac{1}{2}B^{T}} - A^{T}\Theta A\right]}_{B_{+}[\Theta]}\right) + 2\mu_{\lambda}.$$
(4.122)

Indeed, let [u; z] be a feasible solution to the optimization problem (4.120) specifying $\Re[H]$. Then

$$[u;z]^{T}B_{+}[u;z] = [u;z]^{T}B_{+}[\Theta_{\lambda}][u;z] + z^{T}A^{T}\Theta_{\lambda}Az;$$

the first term in the right hand side is $\leq \mathcal{M}(B_+[\Theta_{\lambda}])$ since $[u; z] \in \mathcal{B}_* \times \mathcal{X}_s$, and the second term in the right hand side, as we have already seen is $\leq \mu_{\lambda}$, and (4.122) follows.

Now assume that we have at our disposal a computationally tractable cone

$$\mathbf{H} \subset \mathbf{S}_{+}^{N} \times \mathbf{R}_{+}$$

satisfying the following assumption:

C. Whenever $(\Theta, \nu) \in \mathbf{H}$, we can efficiently find an $n \times N$ matrix $H = [h_1, ..., h_N]$ and a nonnegative vector $\lambda \in \mathbf{R}^N_+$ such that

the columns
$$h_{\ell}$$
 of H satisfy (4.83) (a)
 $\Theta = H \text{Diag}\{\lambda\} H^T$ (b) (4.123)
 $\sum_i \lambda_i \leq \mu$ (c)

The following simple observation is crucial for us:

Proposition 4.27. Consider the estimation problem posed in Section 4.6.1, and let efficiently computable convex function \mathcal{M} and computationally tractable closed convex cone **H** satisfy (4.121) (where \mathcal{B}_* is the unit ball of the norm conjugate to the norm $\|\cdot\|$ in which the recovery error is measured) and Assumption C, respectively. Consider the convex optimization problem

$$Opt = \min_{\tau,\Theta,\mu} \left\{ 2\tau + 2\mu : (\Theta,\mu) \in \mathbf{H}, \ \mathcal{M}(B_{+}[\Theta]) \leq \tau \right\}$$
$$\begin{bmatrix} B_{+}[\Theta] = \left[\frac{\frac{1}{2}B}{\frac{1}{2}B^{T} - A^{T}\Theta A} \right] \end{bmatrix}$$
(4.124)

Given a feasible solution (τ, Θ, μ) to this problem, by C we can efficiently convert it to (H, λ) such that $H = [h_1, ..., h_N]$ with h_ℓ satisfying (4.83) and $\lambda \ge 0$ with $\sum_\ell \lambda_\ell \le \mu$. We have

$$\Re[H] \le 2\tau + 2\mu,$$

whence the $(\epsilon, \|\cdot\|)$ -risk of the polyhedral estimate \hat{x}^H satisfies the bound

$$\operatorname{Risk}_{\epsilon, \|\cdot\|} [\widehat{x}^H | \mathcal{X}] \le 2\tau + 2\mu. \tag{4.125}$$

As a result, we can build efficiently polyhedral estimates with $(\epsilon, \|\cdot\|)$ -risk arbitrarily close to Opt (and with risk exactly Opt, provided problem (4.124) is solvable).

Proof is readily given by the reasoning preceding Proposition. Indeed, with $\tau, \Theta, \mu, H, \lambda$ as described in Proposition, the columns h_{ℓ} of H satisfy (4.83) by **C**,

implying, by Proposition 4.23, that $\operatorname{Risk}_{\epsilon,\|\cdot\|}[\widehat{x}^H|\mathcal{X}] \leq \mathfrak{R}[H]$. Besides this, **C** says that for our H, λ it holds $\Theta = \Theta_{\lambda}$ and $\mu_{\lambda} \leq \mu$, so that (4.122) combines with the constraints of (4.124) to imply that $\mathfrak{R}[H] \leq 2\tau + 2\mu$, and (4.125) follows by Proposition 4.23.

With the approach we are developing now, presumably good polyhedral estimates will be given by (nearly) optimal solutions to (4.124). To apply this approach, we need to develop techniques for building cones **H** satisfying **C** and for building efficiently computable functions $\mathcal{M}(\cdot)$ satisfying (4.121). These tasks are the subjects of the sections to follow.

4.6.5.2 Specifying cones H

We are about to specify cones **H** in the case when the number N of columns in the candidate contrast matrices is m and under the following assumption on the given reliability tolerance ϵ and observation scheme in question:

D. The observation scheme in question is such that for properly built computationally tractable convex compact subset $Z \subset \mathbf{R}^m_+$ intersecting int \mathbf{R}^m_+ , the norm

$$p(h) = \sqrt{\max_{z \in Z} \sum_{i} z_i h_i^2}$$

induced by Z satisfies the relation

$$p(h) \leq 1 \Rightarrow \operatorname{Prob}\{|h^T \xi_x| > 1\} \leq \epsilon/m \ \forall x \in \mathcal{X}.$$

Note that condition \mathbf{D} is satisfied for Sub-Gaussian, Discrete, and Poisson observation schemes: according to the results of Section 4.6.3,

• in the Sub-Gaussian case, it suffices to take

$$Z = \{2\sigma^2 \ln(2m/\epsilon)[1;...;1]\};\$$

• in the Discrete case, it suffices to take

$$Z = \frac{8\ln(2m/\epsilon)}{K} A\mathcal{X} + \frac{64\ln^2(2m/\epsilon)}{9K^2} \mathbf{\Delta}_m, A\mathcal{X} = \{Ax : x \in \mathcal{X}\}, \ \mathbf{\Delta}_m = \{y \in \mathbf{R}^m : y \ge 0, \sum_i y_i = 1\}$$

• in the Poisson case, it suffices to take

$$Z = 6\ln(2m/\epsilon)A\mathcal{X} + 36\ln^2(2m/\epsilon)\mathbf{\Delta}_m,$$

with the same $A\mathcal{X}$ and Δ_m as in Discrete case.

Note that in all these cases Z only "marginally" – logarithmically – depends on ϵ and m.

Under Assumption $\mathbf{D},$ the cone \mathbf{H} can be built as follows:

• When \mathcal{Z} is a singleton: $\mathcal{Z} = \{\bar{z}\}$, so that $p(\cdot)$ is scaled Euclidean norm, we set

$$\mathbf{H} = \{ (\Theta, \mu) \in \mathbf{S}^m_+ \times \mathbf{R}_+ : \mu \ge \sum_i \bar{z}_i \Theta_{ii} \}.$$

Given $(\Theta, \mu) \in \mathbf{H}$, the $m \times m$ matrix H and $\lambda \in \mathbf{R}^m_+$ are built as follows: setting $S = \text{Diag}\{\sqrt{\overline{z_1}}, ..., \sqrt{\overline{z_m}}\}$, we compute the eigenvalue decomposition of the matrix $S\Theta S$:

$$S\Theta S = U \text{Diag}\{\lambda\} U^T,$$

where U is orthonormal, and set $H = S^{-1}U$, thus ensuring that $\Theta = H\text{Diag}\{\lambda\}H^T$. Since $\mu \geq \sum_i \bar{z}_i \Theta_{ii}$, we have $\sum_i \lambda_i = \text{Tr}(S\Theta S) \leq \mu$. Finally, a column h of H is of the form $S^{-1}f$ with $\|\cdot\|_2$ -unit vector f, implying that $p(h) = \sqrt{\sum_i \bar{z}_i [S^{-1}f]_i^2} = \sqrt{\sum_i f_i^2} = 1$, so that h satisfies (4.83) by **D**.

• When Z is not a singleton, we set

$$\begin{aligned}
\phi(r) &= \max_{z \in \mathbb{Z}} z^T r, \\
\varkappa &= 6 \ln(2\sqrt{3}m^2), \\
\mathbf{H} &= \{(\Theta, \mu) \in \mathbf{S}^m_+ \times \mathbf{R}_+ : \mu \ge \varkappa \phi(\mathrm{dg}(\Theta))),
\end{aligned}$$
(4.126)

where dg(Q) is the diagonal of a (square) matrix Q. Note that $\phi(r) > 0$ whenever $r \ge 0, r \ne 0$, since Z contains a positive vector.

The justification of this construction and the efficient (randomized) algorithm for converting a pair $(\Theta, \mu) \in \mathbf{H}$ into (H, λ) satisfying, along with Θ, μ), the requirements of \mathbf{C} are given by the following

Lemma 4.28. (i) Whenever H is an $m \times m$ matrix with columns h_{ℓ} satisfying $p(h_{\ell}) \leq 1$ and $\lambda \in \mathbb{R}^{m}_{+}$, we have

$$(\Theta_{\lambda} = H \operatorname{Diag}\{\lambda\} H^T, \mu = \varkappa \sum_i \lambda_i) \in \mathbf{H}.$$

(ii) Given $(\Theta, \mu) \in \mathbf{H}$ with $\Theta \neq 0$, we find decomposition $\Theta = QQ^T$ with $m \times m$ matrix A, fix an orthonormal $m \times m$ matrix V with magnitudes of entries not exceeding $\sqrt{2/m}$ (e.g., the orthonormal scaling of the matrix of the cosine transform). When $\mu > 0$, we set $\lambda = \frac{\mu}{m}[1; ...; 1] \in \mathbf{R}^m$ and consider the random matrix

$$H_{\chi} = \sqrt{rac{m}{\mu}} Q \mathrm{Diag}\{\chi\} V$$

where χ is the m-dimensional Rademacher random vector. We have

$$H_{\chi} \text{Diag}\{\lambda\} H_{\chi}^T \equiv \Theta, \ \lambda \ge 0, \sum_i \lambda_i = \mu.$$
 (4.127)

Moreover, probability of the event

$$p(\operatorname{Col}_{\ell}[H_{\chi}]) \le 1 \,\forall \ell \le m \tag{4.128}$$

is at least 1/2. Thus, generating independent samples of χ and terminating with $H = H_{\chi}$ when the latter matrix satisfies (4.128), we with probability 1 terminate with (H, λ) satisfying C, and the probability for the procedure to terminate in course of the first M = 1, 2, ... steps is at least $1 - 2^{-M}$.

When $\mu = 0$, we have $\Theta = 0$ (since $\mu = 0$ implies $\phi(dg(\Theta)) = 0$, which with $\Theta \succeq 0$ is possible only when $\Theta = 0$); thus, when $\mu = 0$, we set $H = 0_{m \times m}$ and $\lambda = 0_{m \times 1}$.

Note that Lemma states, essentially, that the cone ${\bf H}$ is a tight, up to logarithmic in m factor, inner approximation of the set

$$\left\{ \begin{aligned} \Theta &= H \mathrm{Diag}\{\lambda\} H^T, \\ (\Theta, \mu) : \exists (\lambda \in \mathbf{R}^m_+, H \in \mathbf{R}^{m \times m}) : & p(\mathrm{Col}_{\ell}[H]) \leq 1, \ \ell \leq m, \\ \mu \geq \sum_{\ell} \lambda_{\ell} \end{aligned} \right\}$$

For proof, see Section 4.10.8.

4.6.5.3 Specifying functions \mathcal{M}

In this section we focus on computationally efficient upper-bounding of maxima of quadratic forms over symmetric w.r.t. the origin convex compact sets, with ultimate goal to specify "presumably good" efficiently computable convex function $\mathcal{M}(\cdot)$ satisfying (4.121). What we intend to use to this end, is a kind of semidefinite relaxation.

Cones compatible with convex sets. Given a nonempty convex compact set $\mathcal{Y} \subset \mathbf{R}^N$, we say that a cone **Y** is *compatible* with \mathcal{Y} , if

- Y is a closed convex computationally tractable cone contained in $\mathbf{S}^N_+ \times \mathbf{R}_+$
- one has

$$\forall (V,\tau) \in \mathbf{Y} : \max_{y \in \mathcal{V}} y^T V y \le \tau \tag{4.129}$$

- **Y** contains a pair (V, τ) with $V \succ 0$.
- relations $(V, \tau) \in \mathbf{Y}$ and $\tau' \geq \tau$ imply that $(V, \tau') \in \mathbf{Y}^{65}$.

We call a cone **Y** sharp, if **Y** is a closed convex cone contained in $\mathbf{S}_{+}^{N} \times \mathbf{R}_{+}$ and such that the only pair $(V, \tau) \in \mathbf{Y}$ with $\tau = 0$ is the pair (0, 0), or, equivalently, a sequence $\{(V_i, \tau_i) \in \mathbf{Y}, i \geq 1\}$ is bounded if and only if the sequence $\{\tau_i, i \geq 1\}$ is bounded.

Note that whenever the linear span of \mathcal{Y} is the entire \mathbf{R}^N , every compatible with \mathcal{Y} cone is sharp.

Observe that if $\mathcal{Y} \subset \mathbf{R}^N$ is a nonempty convex compact set and \mathbf{Y} is a cone compatible with a shift $\mathcal{Y} - a$ of \mathcal{Y} , then \mathbf{Y} is compatible with \mathcal{Y}_s .

Indeed, when shifting a set \mathcal{Y} , its symmeterization $\frac{1}{2}[\mathcal{Y} - cY]$ remains intact, so that we can assume that \mathbf{Y} is compatible with \mathcal{Y} . Now let $(V, \tau) \in \mathbf{Y}$ and $y, y' \in \mathcal{Y}$. We have

$$[y - y']^{T} V[y - y'] + \underbrace{[y + y']^{T} V[y + y']}_{\geq 0} = 2[y^{T} V y + [y']^{T} V y'] \leq 4\tau,$$

whence for $z = \frac{1}{2}[y - y']$ it holds $z^T V z \leq \tau$. Since every $z \in \mathcal{Y}_s$ is of the form $\frac{1}{2}[y - y']$ with $y, y' \in \mathcal{Y}$, the claim follows.

⁶⁵the latter requirement is "for free" – passing from a computationally tractable closed convex cone $\mathbf{Y} \subset \mathbf{S}^N_+ \times \mathbf{R}_+$ satisfying (4.129) only to the cone $\mathbf{Y}^+ = \{(V, \tau) : \exists \overline{\tau} \leq \tau : (V, \overline{\tau}) \in \mathbf{Y}\}$, we get a larger than \mathbf{Y} cone compatible with \mathcal{Y} . It will be clear from the sequel that in our context, the larger is a cone compatible with \mathcal{Y} , the better, so that this extension makes no harm.

322

LECTURE 4

Note that the claim can be "nearly inverted:" ix $0 \in \mathcal{Y}$ and \mathbf{Y} compatible with \mathcal{Y}_s , then the "widening" of \mathbf{Y} – the cone

$$\mathbf{Y}^{+} = \{ (V, \tau) : (V, \tau/4) \in \mathbf{Y} \}$$

is compatible with \mathcal{Y} (evident, since when $0 \in \mathcal{Y}$, every vector from \mathcal{Y} is proportional, with coefficient 2, to a vector from \mathcal{Y}_s).

Building functions \mathcal{M} . The role of compatibility in our context becomes clear from the following observation:

Proposition 4.29. In the situation described in Section 4.6.1, assume that we have at our disposal cones X and U compatible, respectively, with \mathcal{X}_s and with the unit ball

$$\mathcal{B}_* = \{ v \in \mathbf{R}^{\nu} : \|u\|_* \le 1 \}$$

of the norm $\|\cdot\|_*$ conjugate to the norm $\|\cdot\|$. Given $M \in \mathbf{S}^{n+\nu}$, let us set

$$\mathcal{M}(M) = \inf_{X,t,U,s} \left\{ t + s : (X,t) \in \mathbf{X}, (U,s) \in \mathbf{U}, \text{Diag}\{U,X\} \succeq M \right\}$$
(4.130)

Then \mathcal{M} is real-valued efficiently computable convex function on $\mathbf{S}^{n+\nu}$ such that (4.121) takes place: for every $M \in \mathbf{S}^{n+\nu}$ it holds

$$\mathcal{M}(M) \ge \max_{[u;z]\in\mathcal{B}_*\times\mathcal{X}_{\mathrm{s}}} [u;z]^T M[u;z].$$

In addition, when \mathbf{X} and \mathbf{U} are sharp, problem (4.130) is solvable.

Proof is immediate. Given that the objective of the optimization problem specifying $\mathcal{M}(M)$ is nonnegative on the feasible set, the fact that \mathcal{M} is real-valued is equivalent to problem's feasibility, and the latter is readily given by that fact that **X** is a cone containing a pair (X, t) with $X \succ 0$ and similar fact for **U**. Convexity of \mathcal{M} is evident. To verify (4.121), let (X, t, U, s) form a feasible solution to the optimization problem in (4.130). When $[u; z] \in \mathcal{B}_* \times \mathcal{X}_s$ we have

$$[u;z]^T M[u;z] \le u^T U u + z^T X z \le s+t,$$

where the first inequality is due to the \succeq -constraint in (4.130), and the second is due to the fact that **U** is compatible with \mathcal{B}_* , and **X** is compatible with \mathcal{X}_s . Since the resulting inequality holds true for all feasible solutions to the optimization problem in (4.130), (4.121) follows. Finally, when **X** and **U** are sharp, (4.130) is a feasible conic problem with bounded level sets of the objective and as such is solvable. \Box

4.6.5.4 Putting things together

Combining Propositions 4.29 and 4.27, we arrive at the following recipe for designing presumably good polyhedral estimates:

Proposition 4.30. In the situation of Section 4.6.1, assume that we have at our disposal cones **X** and **U** compatible, respectively, with \mathcal{X}_s and with the unit ball \mathcal{B}_* of the norm conjugate to $\|\cdot\|$. Given reliability tolerance $\epsilon \in (0,1)$, assume that we have at our disposal a positive integer N and a computationally tractable cone **H** satisfying, along with ϵ , Assumption **C**. Consider (clearly feasible) convex

optimization problem

$$Opt = \min_{\Theta,\mu,X,t,U,s} \left\{ f(t,s,\mu) := 2(t+s+\mu) : \begin{array}{c|c} (\Theta,\mu) \in \mathbf{H}, (X,t) \in \mathbf{X}, (U,s) \in \mathbf{U} \\ \hline U & \frac{1}{2}B \\ \hline \frac{1}{2}B^T & A^T\Theta A + X \end{array} \right\} \succeq 0$$

$$(4.131)$$

Given a feasible solution Θ, μ, X, t, U, s and invoking C, we can convert, in a computationally efficient manner, (Θ, μ) into (H, λ) such that the columns of the $m \times N$ contrast matrix satisfy (4.83), $\Theta = H \text{Diag}\{\lambda\} H^T$, and $\mu \geq \sum_{\ell} \lambda_{\ell}$. The $(\epsilon, \|\cdot\|)$ risk of the polyhedral estimate \hat{x}^H satisfies the bound

$$\operatorname{Risk}_{\epsilon,\|\cdot\|}[\widehat{x}^H | \mathcal{X}] \le f(t, s, \mu). \tag{4.132}$$

In particular, we can build, in a computationally efficient manner, polyhedral estimates with risks arbitrarily close to Opt (and with risk Opt, provided that (4.131) is solvable).

Proof. Let Θ, μ, X, t, U, s form a feasible solution to (4.131). By the semidefinite constraint in (4.131) we have

$$0 \preceq \left[\begin{array}{c|c} U & -\frac{1}{2}B \\ \hline -\frac{1}{2}B^T & A^T\Theta A + X \end{array} \right] = \operatorname{Diag}\{U, X\} - \underbrace{\left[\begin{array}{c|c} & \frac{1}{2}B \\ \hline \frac{1}{2}B^T & -A^T\Theta A \end{array} \right]}_{=:M},$$

whence by Proposition 4.29 for the function \mathcal{M} defined in this proposition one has

$$\mathcal{M}(M) \le t + s.$$

Since \mathcal{M} , by the same Proposition 4.29, satisfies (4.121), invoking Proposition 4.27 we arrive at

$$\Re[H] \le 2(\mu + \mathcal{M}(M)) \le f(t, s, \mu),$$

implying by Proposition 4.23 the target relation (4.132).

4.6.5.5 Compatibility: basic examples and calculus

What is crucial for the design of presumably good polyhedral estimates via the recipe described in Proposition 4.30, is our ability to equip convex "sets of interest" (in our context, these are the symmeterization \mathcal{X}_s of the signal set and the unit ball \mathcal{B}_* of the norm conjugate to the norm $\|\cdot\|$ in which the recovery error is measured) with cones compatible with these sets⁶⁶. We are about to discuss two major sources of these cones, namely (a) spectratopes/ellitopes, and (b) absolute norms. We develop also "compatibility calculus" which allows to build, in a fully algorithmic fashion, cones compatible with the results of basic convexity-preserving operations with convex sets from the cones compatible with the operands.

In view of Proposition 4.30, the larger are the cones **X** and **U** compatible with \mathcal{X}_{s} and \mathcal{B}_{*} , the better – the wider is the optimization domain in (4.131) and, consequently, the less is (the best) risk bound achievable with the recipe presented

 $^{^{66}}$ recall that we already know how to specify the second element of the construction, the cone ${\bf H}.$

in the proposition. Given convex compact set $\mathcal{Y} \in \mathbf{R}^N$, the "ideal" – the largest candidate to the role of the cone compatible with \mathcal{Y} would be

$$\mathbf{Y}^* = \{ (V, \tau) \in \mathbf{S}^N_+ \times \mathbf{R}_+ : \tau \ge \max_{y \in \mathcal{Y}} y^T V y \}.$$

This cone, however, typically is intractable, which enforces us to look for "as large as possible" *tractable* inner approximations of \mathbf{Y}^* .

4.6.5.5.A. Cones compatible with ellitopes/spectratopes are readily given by semidefinite relaxation. Specifically, when

$$\begin{aligned} \mathcal{Y} &= \{ y \in \mathbf{R}^N : \exists (r \in \mathcal{R}, z \in \mathbf{R}^K) : y = Mz, R_{\ell}^2[z] \preceq r_{\ell} I_{d_{\ell}}, \ell \leq L \} \\ & \left[R_{\ell}[z] = \sum_j z_j R^{\ell j}, \ R^{\ell j} \in \mathbf{S}^{d_{\ell}} \right] \end{aligned}$$

with our standard restrictions on \mathcal{R} , invoking Proposition 4.8 it is immediately seen that the set

$$\mathbf{Y} = \{ (V,\tau) \in \mathbf{S}_{+}^{N} \times \mathbf{R}_{+} : \exists \Lambda = \{ \Lambda_{\ell} \in \mathbf{S}_{+}^{d_{\ell}}, \ell \leq L \} : M^{T}VM \preceq \sum_{\ell} \mathcal{R}^{*}[\Lambda_{\ell}], \ \phi_{\mathcal{R}}(\lambda[\Lambda]) \leq \tau \} \\ \left[[\mathcal{R}_{\ell}^{*}[\Lambda_{\ell}]]_{ij} = \operatorname{Tr}(R^{\ell i}\Lambda_{\ell}R^{\ell j}), \ \lambda[\Lambda] = [\operatorname{Tr}(\Lambda_{1}); ...; \operatorname{Tr}(\Lambda_{L})]), \ \phi_{\mathcal{R}}(\lambda) = \max_{r \in \mathcal{R}} r^{T}\lambda \right]$$

$$(4.133)$$

is a closed convex cone which is compatible with \mathcal{Y} .

Similarly, when \mathcal{Y} is an ellitope:

$$\mathcal{Y} = \{ y \in \mathbf{R}^N : \exists (r \in \mathcal{R}, z \in \mathbf{R}^K) : y = Mz, z^T R_\ell z \le r_\ell, \, \ell \le L \}$$

with our standard restrictions on \mathcal{R} and R_{ℓ} , invoking Proposition 4.6, the set

$$\mathbf{Y} = \{ (V,\tau) \in \mathbf{S}^N \times \mathbf{R}^+ : \exists \lambda \in \mathbf{R}_+^L : M^T V M \preceq \sum_{\ell} \lambda_{\ell} R_{\ell}, \phi_{\mathcal{R}}(\lambda) \leq \tau \}$$
(4.134)

is compatible with \mathcal{Y} . In both cases, **Y** is sharp, provided that the image space of M is the entire \mathbf{R}^{N} .

Note that in both these cases \mathbf{Y} is a reasonably tight inner approximation of \mathbf{Y}^* : whenever $(V, \tau) \in \mathbf{Y}^*$, we have $(V, \theta \tau) \in \mathbf{Y}$, with a moderate θ (specifically, $\theta = O(1) \ln(2 \sum_{\ell} d_{\ell})$ in the spectratopic, and $\theta = O(1) \ln(2L)$ in the ellitopic case, see Propositions 4.8, 4.6, respectively).

4.6.5.5.B. Compatibility via absolute norms.

Preliminaries. Recall that a norm $p(\cdot)$ on \mathbf{R}^N is called *absolute*, if p(x) is a function of the vector $\operatorname{abs}[x] := [|x_1|; ...; |x_N|]$ of the magnitudes of entries in x. It is well known that an absolute norm p is monotone on \mathbf{R}^N_+ , so that $\operatorname{abs}[x] \leq \operatorname{abs}[x']$ implies that $p(x) \leq p(x')$, and that the norm

$$p_*(x) = \max_{y:p(y) \le 1} x^T y$$

conjugate to $p(\cdot)$ is absolute along with p.

Let us say that an absolute norm $r(\cdot)$ fits an absolute norm $p(\cdot)$ on \mathbb{R}^N , if for every vector x with $p(x) \leq 1$ the entrywise square $[x]^2 = [x_1^2; ...; x_N^2]$ of x satisfies $r([x]^2) \leq 1$. For example, the largest norm $r(\cdot)$ which fits the absolute norm

$$p(\cdot) = \| \cdot \|_s, s \in [1, \infty]$$
, is

$$r(\cdot) = \begin{cases} \|\cdot\|_1, & 1 \le s \le 2\\ \|\cdot\|_{s/2}, & s \ge 2 \end{cases}$$

An immediate observation is that an absolute norm $p(\cdot)$ on \mathbb{R}^N can be "lifted" to a norm on \mathbb{S}^N , specifically, the norm

$$p^+(Y) = p([p([\operatorname{Col}_1[Y]); ...; p(\operatorname{Col}_N[Y])]) : \mathbf{S}^N \to \mathbf{R}_+,$$
 (4.135)

where $\operatorname{Col}_{j}[Y]$ is *j*th column in *Y*. It is immediately seen that when *p* is an absolute norm, the right hand side in (4.135) indeed is a norm on \mathbf{S}^{N} satisfying the identity

$$p^+(xx^T) = p^2(x), x \in \mathbf{R}^N.$$
 (4.136)

Absolute norms and compatibility. Our interest in absolute norms is motivated by the following immediate

Observation 4.31. Let $p(\cdot)$ be an absolute norm on \mathbb{R}^N , and $r(\cdot)$ be another absolute norm which fits $p(\cdot)$, both norms being computationally tractable. These norms give rise to the computationally tractable and sharp closed convex cone

$$\mathbf{P} = \mathbf{P}_{p(\cdot), r(\cdot)} = \left\{ (V, \tau) \in \mathbf{S}_{+}^{N} \times \mathbf{R}_{+} : \exists (W \in \mathbf{S}^{N}, w \in \mathbf{R}_{+}^{N}) : \begin{array}{c} V \preceq W + \operatorname{Diag}\{w\}\\ \left[p^{+}\right]_{*}(W) + r_{*}(w) \leq \tau \end{array} \right\},$$

$$(4.137)$$

where $[p^+]_*(\cdot)$ is the norm on \mathbf{S}^N conjugate to the norm $p^+(\cdot)$, and $r_*(\cdot)$ is the norm on \mathbf{R}^N conjugate to the norm $r(\cdot)$, and this cone is compatible with the unit ball of the norm $p(\cdot)$ (and thus – with any convex compact subset of the latter ball).

Verification is immediate. The fact that **P** is a computationally tractable and closed convex cone is evident. Now let $(V, \tau) \in \mathbf{P}$, so that $V \succeq 0$ and $V \preceq W + \text{Diag}\{w\}$ with $[p^+]_*(W) + r_*(w) \leq \tau$. For x with $p(x) \leq 1$ we have

$$x^{T}Vx \leq x^{T}[W + \text{Diag}\{w\}]x = \text{Tr}(W[xx^{T}]) + w^{T}[x]^{2}$$

$$\leq p^{+}(xx^{T})[p^{+}]_{*}(W) + r([x]^{2})r_{*}(w) = p^{2}(x)[p^{+}]_{*}(W) + r_{*}(w)$$

$$\leq [p^{+}]_{*}(W) + r_{*}(w) \leq \tau$$

(we have used (4.137)), whence $x^T V x \leq \tau$ for all x with $p(x) \leq 1$.

Let us look what is our construction in the case when $p(\cdot) = \|\cdot\|_s$, $s \in [1, \infty]$, which allows to take $r(\cdot) = \|\cdot\|_{\bar{s}}$, $\bar{s} = \max[s/2, 1]$. Setting $s_* = \frac{s}{s-1}$, $\bar{s}_* = \frac{\bar{s}}{\bar{s}-1}$, we clearly have

$$[p^+]_*(W) = \|W\|_{s_*} := \begin{cases} \left(\sum_{i,j} |W_{ij}|^{s_*}\right)^{1/s_*}, & s_* < \infty \\ \max_{i,j} |W_{ij}|, & s_* = \infty \end{cases}, \ r_*(w) = \|w\|_{\bar{s}_*}, \tag{4.138}$$

resulting in

$$\mathbf{P}^{s} := \mathbf{P}_{\|\cdot\|_{s}, \|\cdot\|_{\bar{s}}} \\
= \left\{ (V, \tau) : V \in \mathbf{S}^{N}_{+}, \exists (W \in \mathbf{S}^{N}, w \in \mathbf{R}^{N}_{+}) : \frac{V \preceq W + \operatorname{Diag}\{w\},}{\|W\|_{s_{*}} + \|w\|_{\bar{s}_{*}} \leq \tau} \right\}, \\
(4.139)$$

and Observation 4.31 says that \mathbf{P}^s is compatible with the unit ball of $\|\cdot\|_s$ -norm

on \mathbf{R}^N (and therefore with every closed convex subset of this ball).

When s = 1, that is, $s_* = \bar{s}_* = \infty$, (4.139) results in

$$\mathbf{P}^{1} = \left\{ (V,\tau) : V \succeq 0, \exists (W \in \mathbf{S}^{N}, w \in \mathbf{R}^{N}_{+}) : \begin{array}{l} V \preceq W + \operatorname{Diag}\{w\}, \\ \|W\|_{\infty} + \|w\|_{\infty} \leq \tau \end{array} \right\} \\
= \{ (V,\tau) : V \succeq 0, \|V\|_{\infty} \leq \tau \}, \\$$
(4.140)

and it is easily seen that the situation is a good as it could be, namely,

$$\mathbf{P}^{1} = \{ (V, \tau) : V \succeq 0, \max_{\|x\|_{1} \le 1} x^{T} V x \le \tau \}.$$

It can be shown (see Section 4.10.9) that when $s \in [2, \infty]$, so that $\bar{s}_* = \frac{s}{s-2}$, (4.139) results in

$$\mathbf{P}^{s} = \{ (V,\tau) : V \succeq 0, \exists (w \in \mathbf{R}^{N}_{+}) : V \preceq \text{Diag}\{w\} \& \|w\|_{\frac{s}{s-2}} \le \tau \}.$$
(4.141)

Note that

$$\mathbf{P}^{2} = \{ (V, \tau) : V \succeq 0, \|V\|_{\mathrm{Sh}, \infty} \le \tau \}$$

and this is *exactly* the largest cone compatible with the unit Euclidean ball.

When $s \ge 2$, the unit ball \mathcal{Y} of the norm $\|\cdot\|_s$ is an ellitope:

$$\{y \in \mathbf{R}^N : \|y\|_s \le 1\} = \{y \in \mathbf{R}^N : \exists (t \ge 0, \|t\|_{\bar{s}} \le 1) : y^T R_\ell y := y_\ell^2 \le t_\ell, \, \ell \le L = N\},\$$

so that one of the cones compatible with \mathcal{Y} is given by (4.134) with the identity matrix in the role of M. It goes without surprise that, as it is immediately seen, the latter cone is nothing but the cone given by (4.141).

4.6.5.5.C. Calculus of compatibility. Cones compatible with convex sets admit a kind of fully algorithmic calculus with the rules as follows (verification of the rules is straightforward and is skipped):

- 1. [passing to a subset] When $\mathcal{Y}' \subset \mathcal{Y}$ are convex compact subsets of \mathbb{R}^N and a cone \mathbf{Y} is compatible with \mathcal{Y} , the cone is compatible with \mathcal{Y}' as well.
- 2. [finite intersection] Let cones \mathbf{Y}^{j} be compatible with convex compact sets $\mathcal{Y}_{j} \subset \mathbf{R}^{N}, j = 1, ..., J$. Then the cone

$$\mathbf{Y} = \operatorname{cl}\{(V,\tau) \in \mathbf{S}_{+}^{N} \times \mathbf{R}_{+} : \exists ((V_{j},\tau_{j}) \in \mathbf{Y}^{j}, j \leq J) : V \preceq \sum_{j} V_{j}, \sum_{j} \tau_{j} \leq \tau \}$$

is compatible with $\mathcal{Y} = \bigcap_{j} \mathcal{Y}_{j}$. The closure operation can be skipped whenever

all cones \mathbf{Y}^{j} are sharp, in which case \mathbf{Y} is sharp as well.

3. [convex hulls of finite union] Let cones \mathbf{Y}^{j} be compatible with convex compact sets $\mathcal{Y}_{j} \subset \mathbf{R}^{N}, j = 1, ..., J$, and let there exist (V, τ) such that $V \succ 0$ and

$$(V,\tau) \in \mathbf{Y} := \bigcap_{j} \mathbf{Y}^{j}.$$

Then **Y** is compatible with $\mathcal{Y} = \operatorname{Conv}\{\bigcup_{j} \mathcal{Y}_{j}\}$ and is sharp, provided that all **Y**^j.

4. [direct product] Let cones \mathbf{Y}^{j} be compatible with convex compact sets $\mathcal{Y}_{j} \subset \mathbf{R}^{N_{j}}$, j = 1, ..., J. Then the cone

$$\mathbf{Y} = \{(V,\tau) \in \mathbf{S}_+^{N_1+\ldots+N_J} \times \mathbf{R}_+ : \exists ((V_j,\tau_j) \in \mathbf{Y}^j : V \preceq \text{Diag}\{V_1, \ldots, V_J\} \& \tau \ge \sum_j \tau_j \}$$

is compatible with $\mathcal{Y} = \mathcal{Y}_1 \times ... \times \mathcal{Y}_J$. This cone is sharp, provided that all \mathbf{Y}^j are so.

5. [linear image] Let cone **Y** be compatible with convex compact set $\mathcal{Y} \subset \mathbf{R}^N$, let A be a $K \times N$ matrix, and let $\mathcal{Z} = A\mathcal{Y}$. The cone

$$\mathbf{Z} = \operatorname{cl}\{(V,\tau) \in \mathbf{S}_{+}^{K} \times \mathbf{R}_{+} : \exists U \succeq A^{T}VA : (U,\tau) \in \mathbf{Y}\}$$

is compatible with \mathcal{Z} . The closure operation can be skipped whenever \mathbf{Y} is either sharp, or *complete*, completeness meaning that $(V, \tau) \in \mathbf{Y}$ and $0 \leq V' \leq V$ imply that $(V', \tau) \in \mathbf{Y}^{67}$. The cone \mathbf{Z} is sharp, provided \mathbf{Y} is so and the rank of A is K.

6. [inverse linear image] Let cone **Y** be compatible with convex compact set $\mathcal{Y} \subset \mathbf{R}^N$, let A be a $N \times K$ matrix with trivial kernel, and let $\mathcal{Z} = A^{-1}\mathcal{Y} := \{z \in \mathbf{R}^K : Az \in \mathcal{Y}\}$. The cone

$$\mathbf{Z} = \operatorname{cl}\{(V,\tau) \in \mathbf{S}_{+}^{K} \times \mathbf{R}_{+} : \exists U : A^{T}UA \succeq V \& (U,\tau) \in \mathbf{Y}\}$$

is compatible with \mathcal{Z} . The closure operations can be skipped whenever **Y** is sharp, in which case **Z** is sharp as well.

7. [arithmetic summation] Let cones \mathbf{Y}^{j} be compatible with convex compact sets $\mathcal{Y}_{j} \subset \mathbf{R}^{N}, j = 1, ..., J$. Then the arithmetic sum $\mathcal{Y} = \mathcal{Y}_{1} + ... + \mathcal{Y}_{J}$ of the sets \mathcal{Y}_{j} can be equipped with compatible cone readily given by the cones \mathbf{Y}^{j} ; this cone is sharp, provided all \mathbf{Y}^{j} are so.

Indeed, the arithmetic sum of \mathcal{Y}_j is the linear image of the direct product of \mathcal{Y}_j 's under the mapping $[y^1; ...; y^J] \mapsto y^1 + ... + y^J$, and it remains to combine rules 4 and 5; note the cone yielded by rule 4 is complete, so that when applying rule 5, the closure operation can be skipped.

4.6.5.6 Spectratopic Sub-Gaussian case

As an instructive application of the approach developed so far in Section 4.6.5, consider the special case of the estimation problem posed in Section 4.6.1, where

1. The signal set \mathcal{X} and the unit ball \mathcal{B}_* of the norm conjugate to $\|\cdot\|$ are spectratopes:

$$\begin{aligned} \mathcal{X} &= \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \preceq t_k I_{d_k}, \ 1 \le k \le K \}, \\ \mathcal{B}_* &= \{ z \in \mathbf{R}^\nu : \exists y \in \mathcal{Y} : z = My \}, \\ \mathcal{Y} := \{ y \in \mathbf{R}^q : \exists r \in \mathcal{R} : S_\ell^2[y] \preceq r_\ell I_{f_\ell}, \ 1 \le \ell \le L \}, \end{aligned}$$

(cf. Assumptions **A**, **B** in Section 4.3.3.2; as always, we lose nothing when assuming spectratope \mathcal{X} to be basic).

2. For every $x \in \mathcal{X}$, observation noise ξ_x is $(0, \sigma^2 I_m)$ -sub-Gaussian.

Given reliability tolerance $\epsilon \in (0, 1)$ and looking for a polyhedral estimate with $m \times m$ contrast matrix H, the design recipe suggested by Proposition 4.30 is as follows:

• When building the cone H meeting Assumption C according to the construction

⁶⁷note that if **Y** is compatible with \mathcal{Y} , the *completion* of **Y** – the set $\overline{\mathbf{Y}} = \{(V, \tau) : \exists U : 0 \leq V \leq U, (U, \tau) \in \mathbf{Y}\}$ is a complete cone compatible with \mathcal{Y} .

from Section 4.6.5.2, we set

$$Z = \{\vartheta^2[1;...;1]\}, \, \vartheta = \sigma \kappa, \, \kappa = \sqrt{2\ln(2m/\epsilon)},$$

thus arriving at

$$\mathbf{H} = \{(\Theta, \mu) \in \mathbf{S}^m_+ \times \mathbf{R}_+ : \sigma^2 \kappa^2 \operatorname{Tr}(\Theta) \le \mu\}$$

• We specify the cones X and U compatible with $\mathcal{X}_{s} = \mathcal{X}, \mathcal{B}_{*}$, respectively, according to (4.133).

The resulting problem (4.131), after immediate straightforward simplifications, reads

$$\begin{array}{lll} \mathrm{Opt} & = & \min_{\Theta, U, \Lambda, \Upsilon} \left\{ 2 \left[\phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{T}}(\lambda[\Lambda]) + \sigma^2 \kappa^2 \mathrm{Tr}(\Theta) \right] : \\ & \Theta \succeq 0, U \succeq 0, \Lambda = \{ \Lambda_k \succeq 0, k \leq K \}, \Upsilon = \{ \Upsilon_\ell \succeq 0, \ell \leq L \}, \\ & \left[\begin{array}{c|c} U & \frac{1}{2}B \\ \hline \frac{1}{2}B^T & A^T \Theta A + \sum_k \mathcal{R}_k^* [\Lambda_k] \end{array} \right] \succeq 0, \ M^T U M \preceq \sum_\ell \mathcal{S}_\ell^* [\Upsilon_\ell] \end{array} \right\} \ , \end{array}$$

where, as always,

$$\begin{aligned} & [\mathcal{R}_k^*[\Lambda_k]]_{ij} = \operatorname{Tr}(R^{ki}\Lambda_k R^{kj}) \quad [R_k[x] = \sum_i x_i R^{ki}] \\ & [\mathcal{S}_\ell^*[\Upsilon_\ell]]_{ij} = \operatorname{Tr}(S^{\ell i}\Upsilon_\ell S^{\ell j}) \quad [S_\ell[u] = \sum_i u_i S^{\ell i}] \end{aligned}$$

$$\lambda[\Lambda] = [\operatorname{Tr}(\Lambda_1); ...; \operatorname{Tr}(\Lambda_K)], \, \lambda[\Upsilon] = [\operatorname{Tr}(\Upsilon_1); ...; \operatorname{Tr}(\Upsilon_L)], \, \phi_W(f) = \max_{w \in W} w^T f.$$

We are about to demonstrate that the polyhedral estimate yielded by the efficiently computable (high accuracy near-) optimal solution to the above problem is nearoptimal in the minimax sense.

Observe that the matrices
$$Q := \begin{bmatrix} U & \frac{1}{2}B \\ \hline \frac{1}{2}B^T & A^T\Theta A + \sum_k \mathcal{R}_k^*[\Lambda_k] \end{bmatrix}$$
 and
$$\begin{bmatrix} M^T U M & \frac{1}{2}M^T B \\ \hline \frac{1}{2}B^T M & A^T\Theta A + \sum_k \mathcal{R}_k^*[\Lambda_k] \end{bmatrix} = \begin{bmatrix} M^T & \\ \hline & I_n \end{bmatrix} Q \begin{bmatrix} M & \\ \hline & I_n \end{bmatrix}$$

simultaneously are/are not positive semidefinite due to the fact that the image space of M contains the full-dimensional set \mathcal{B}_* and thus is the entire \mathbf{R}^{ν} , so that the image space of $\begin{bmatrix} M \\ \hline & I_n \end{bmatrix}$ is the entire $\mathbf{R}^{\nu} \times \mathbf{R}^n$. Therefore

$$\begin{array}{lll}
\operatorname{Opt} &= & \min_{\Theta, U, \Lambda, \Upsilon} \left\{ 2 \left[\phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{T}}(\lambda[\Lambda]) + \sigma^{2} \kappa^{2} \operatorname{Tr}(\Theta) \right] : \\ &\Theta \succeq 0, U \succeq 0, \Lambda = \{ \Lambda_{k} \succeq 0, k \leq K \}, \Upsilon = \{ \Upsilon_{\ell} \succeq 0, \ell \leq L \}, \\ & \left[\frac{M^{T} U M}{\frac{1}{2} B^{T} M} \middle| \frac{\frac{1}{2} M^{T} B}{A^{T} \Theta A + \sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}]} \right] \succeq 0, \ M^{T} U M \preceq \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \end{array} \right\} \\ & (4.142)
\end{array}$$

Next, if a collection $\Theta, U, \{\Lambda_k\}, \{\Upsilon_\ell\}$ is a feasible solution to the latter problem and $\theta > 0$, the scaled collection $\theta\Theta, \theta^{-1}U, \{\theta\Lambda_k\}, \{\theta^{-1}\Upsilon_\ell\}$ also is a feasible solution to

the problem; optimizing with respect to scaling, we get

$$\begin{array}{lll}
\operatorname{Opt} &= \inf_{\substack{\Theta,U,\Lambda,\Upsilon}} \left\{ 4\sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon]) \left[\phi_{\mathcal{T}}(\lambda[\Lambda] + \sigma^{2}\kappa^{2}\operatorname{Tr}(\Theta)\right]} : \\ &\quad \Theta \succeq 0, U \succeq 0, \Lambda = \{\Lambda_{k} \succeq 0, k \leq K\}, \Upsilon = \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \\ &\quad \left[\frac{M^{T}UM}{\frac{1}{2}B^{T}M} \mid \frac{1}{2}M^{T}B}{\frac{1}{2}M^{T}\ThetaA + \sum_{k}\mathcal{R}_{k}^{*}[\Lambda_{k}]}\right] \succeq 0, \ M^{T}UM \preceq \sum_{\ell}\mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \end{array}\right\} \\ &\leq 2\kappa \operatorname{Opt}_{+}, \\ \operatorname{Opt}_{+} &= \inf_{\substack{\Theta,U,\Lambda,\Upsilon}} \left\{ 2\sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon]) \left[\phi_{\mathcal{T}}(\lambda[\Lambda]) + \sigma^{2}\operatorname{Tr}(\Theta)\right]} : \\ &\quad \Theta \succeq 0, U \succeq 0, \Lambda = \{\Lambda_{k} \succeq 0, k \leq K\}, \Upsilon = \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \\ &\quad \left[\frac{M^{T}UM}{\frac{1}{2}B^{T}M} \mid \frac{1}{2}M^{T}B}{\frac{1}{2}M^{T}\ThetaA + \sum_{k}\mathcal{R}_{k}^{*}[\Lambda_{k}]}\right] \succeq 0, \ M^{T}UM \preceq \sum_{\ell}\mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \end{array}\right\} \\ &\quad [\text{note that } \kappa > 1] \end{aligned}$$

$$(4.143)$$

On the other hand, consider the optimization problem which under the circumstances is responsible for building presumably good *linear* estimate, that is, the problem

$$\begin{aligned} \operatorname{Opt}_{*} &= \min_{\Theta, H, \Lambda, \Upsilon', \Upsilon''} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \phi_{\mathcal{R}}(\lambda[\Upsilon'']) + \sigma^{2} \operatorname{Tr}(\Theta) : \\ \Lambda &= \{\Lambda_{k} \succeq 0, k \leq K\}, \Upsilon' = \{\Upsilon'_{\ell} \succeq 0, \ell \leq L\}, \Upsilon'' = \{\Upsilon''_{\ell} \succeq 0, \ell \leq L\}, \\ & \left[\frac{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon'_{\ell}]}{\frac{1}{2} [B - H^{T} A]^{T} M} \left| \sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}] \right| \right] \succeq 0, \\ & \left[\frac{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon''_{\ell}]}{\frac{1}{2} H M} \left| \frac{1}{\Theta} \right| \geq 0 \end{aligned} \right\} \end{aligned}$$

(cf. (4.50)). Clearly, strengthening $\Lambda_k \succeq 0$ to $\Lambda_k \succ 0$, we still have

$$\begin{aligned}
\operatorname{Opt}_{*} &= \inf_{\substack{\Theta, H, \Lambda, \Upsilon', \Upsilon''}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \phi_{\mathcal{R}}(\lambda[\Upsilon'']) + \sigma^{2} \operatorname{Tr}(\Theta) : \\
\Lambda &= \{\Lambda_{k} \succ 0, k \leq K\}, \Upsilon' = \{\Upsilon'_{\ell} \succeq 0, \ell \leq L\}, \Upsilon'' = \{\Upsilon''_{\ell} \succeq 0, \ell \leq L\}, \\
\begin{bmatrix} \underbrace{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon'_{\ell}]}_{\frac{1}{2}[B - H^{T}A]^{T}M} \mid \sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}]} \\ \underbrace{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon''_{\ell}]}_{\frac{1}{2}HM} \mid \frac{1}{2}M^{T}H^{T}}_{\Theta} \end{bmatrix} \succeq 0, \\
\end{aligned}$$

$$(4.144)$$

Now let $\Theta, H, \Lambda, \Upsilon', \Upsilon''$ be a feasible solution to the latter problem. By the second semidefinite constraint in (4.144) we have

$$\begin{bmatrix} \underline{\sum_{\ell} \mathcal{S}_{\ell}^* [\Upsilon_{\ell}'']} & \frac{1}{2} M^T H^T A \\ \hline \frac{1}{2} A^T H M & A^T \Theta A \end{bmatrix} = \begin{bmatrix} I \\ A \end{bmatrix}^T \begin{bmatrix} \underline{\sum_{\ell} \mathcal{S}_{\ell}^* [\Upsilon_{\ell}'']} & \frac{1}{2} M^T H^T \\ \hline \frac{1}{2} H M & \Theta \end{bmatrix} \begin{bmatrix} I \\ A \end{bmatrix}$$

$$\succeq \quad 0,$$

which combines with the first semidefinite constraint in (4.144) to imply that

$$\begin{bmatrix} \frac{\sum_{\ell} \mathcal{S}_{\ell}^* [\Upsilon_{\ell}' + \Upsilon_{\ell}'']}{\frac{1}{2} B^T M} & \frac{1}{2} M^T B \\ \hline \frac{1}{2} B^T M & A^T \Theta A + \sum_k \mathcal{R}_k^* [\Lambda_k] \end{bmatrix} \succeq 0.$$

By the Schur Complement Lemma (applicable due to $A^T \Theta A + \sum_k \mathcal{R}_k^* [\Lambda_k] \succeq \sum_k \mathcal{R}_k^* [\Lambda_k] \succ 0$, where the concluding \succ is due to Lemma 4.89 and $\Lambda_k \succ 0$),

330

this relation implies that setting

$$\Upsilon_{\ell} = \Upsilon'_{\ell} + \Upsilon''_{\ell},$$

we have

$$\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \succeq M^{T} \underbrace{\left\lfloor \frac{1}{4} B[A^{T}\Theta A + \sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}]]^{-1}B^{T} \right\rfloor}_{U} M.$$

By the same Schur Complement Lemma, for the just defined $U \succeq 0$ it holds

$$\begin{bmatrix} M^T U M & \frac{1}{2} M^T B \\ \hline \frac{1}{2} B^T M & A^T \Theta A + \sum_k \mathcal{R}_k^* [\Lambda_k] \end{bmatrix} \succeq 0,$$

and in addition

$$M^T U M \preceq \sum_{\ell} \mathcal{S}^*_{\ell}[\Upsilon_{\ell}]$$

by the origin of U. We conclude that

$$(\Theta, U, \Lambda, \Upsilon := \{\Upsilon_{\ell} = \Upsilon'_{\ell} + \Upsilon''_{\ell}, \ell \le L\})$$

is a feasible solution to optimization problem specifying Opt_+ , see (4.142), and that the value of the objective of the latter problem at this feasible solution is

$$2\sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon'] + \lambda[\Upsilon''])} [\phi_{\mathcal{T}}(\lambda[\Lambda]) + \sigma^{2}\mathrm{Tr}(\Theta)]$$

$$\leq \phi_{\mathcal{R}}(\lambda[\Upsilon'] + \lambda[\Upsilon'']) + \phi_{\mathcal{T}}(\lambda[\Lambda]) + \sigma^{2}\mathrm{Tr}(\Theta)$$

$$\leq \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \phi_{\mathcal{R}}(\lambda[\Upsilon'']) + \phi_{\mathcal{T}}(\lambda[\Lambda]) + \sigma^{2}\mathrm{Tr}(\Theta)$$

the concluding quantity in the chain being the value of the objective of problem (4.144) at the feasible solution $\Theta, H, \Lambda, \Upsilon', \Upsilon''$ to this problem. Since the resulting inequality holds true for every feasible solution to (4.144), we conclude that $Opt_+ \leq Opt_*$, whence

$$Opt \le 2\kappa Opt_* = 2\sqrt{2\ln(2m/\epsilon)Opt_*}.$$
(4.145)

This is a really good news. Indeed, from the proof of Proposition 4.16 it follows that under the circumstances, Opt_* is within logarithmic factor of the minimax optimal $(\frac{1}{8}, \|\cdot\|)$ -risk corresponding to the independent of x Gaussian noise: $\xi_x \sim \mathcal{N}(0, \sigma^2 I_m)$ for all x:

$$\operatorname{Opt}_{\epsilon} \leq O(1) \ln \left(\sum_{k} d_{k} + \sum_{\ell} f_{\ell} \right) \operatorname{RiskOpt}_{1/8},$$

$$\operatorname{RiskOpt}_{\epsilon} = \inf_{\widehat{x}(\cdot)} \sup_{x \in \mathcal{X}} \inf \left\{ \rho : \operatorname{Prob}_{\xi \sim \mathcal{N}(0, \sigma^{2}I)} \{ \|Bx - \widehat{x}(Ax + \xi)\| > \rho \} \leq \epsilon \, \forall x \in \mathcal{X} \}$$

$$(4.146)$$

Since the minimax optimal $(\epsilon, \|\cdot\|)$ -risk clearly can only grow when ϵ decreases, we conclude that

When $\epsilon \leq 1/8$, the presumably good polyhedral estimate yielded by a feasible near optimal, in terms of the objective, solution to problem (4.142) is minimax optimal within the logarithmic factor $O(1) \ln(\sum_k d_k + \sum_{\ell} f_{\ell}) \sqrt{\ln(2m/\epsilon)}$.

4.6.6 Numerical illustration

To illustrate our developments, we are about to present numerical comparison of presumably good linear and polyhedral estimates. Our setup is deliberately simple: the signal set \mathcal{X} is just the unit box $\{x \in \mathbf{R}^n : \|x\|_{\infty} \leq 1\}, B \in \mathbf{R}^{n \times n}$ is "numerical double integration:"

$$B_{ij} = \begin{cases} \delta^2(i-j+1), & j \le i \\ 0, & j > i \end{cases}$$

so that x, modulo boundary effects, is the second order finite difference derivative of y = Bx:

$$x_i = \frac{y_i - 2y_{i-1} + y_{i-2}}{\delta^2}, \ 2 < i \le n,$$

and Ax is comprised of m randomly selected entries of Bx. The observation is

$$\omega = Ax + \xi, \, \xi \sim \mathcal{N}(0, \sigma^2 I_m).$$

and the recovery norm is $\|\cdot\|_2$. In other words, we want to recover the restriction of twice differentiable function of one variable on the *n*-point regular grid on the segment $\Delta = [0, n\delta]$ from noisy observations of this restriction taken along *m* randomly selected points of the grid. A priori information on the function is that the magnitude of its second order derivative does not exceed 1.

Note that both presumably good linear and polyhedral estimates \hat{x}_H , \hat{x}^H yielded under the circumstances by Propositions 4.16, resp., 4.29, are near-optimal in the minimax sense in terms of their $\|\cdot\|_{2^-}$, resp., $(\epsilon, \|\cdot\|_2)$ -risk. The goal of our numerical illustration is to compare the empirical performance of the estimates as exhibited in simulation. Note that alternative comparison – via theoretical upper risk bounds of the estimates established in Propositions 4.16 and 4.29 – does not make much sense, since these are bounds on *different from each other* risks $\text{Risk}_{\{\sigma^2 I_m\}, \|\cdot\|_2}$ and $\text{Risk}_{\epsilon, \|\cdot\|_2}$.

In the experiments to be reported, we used n = 64, m = 32, $\delta = 4/n$ (i.e., $\Delta = [0, 4]$); the reliability parameter for the polyhedral estimate was set to $\epsilon = 0.1$. We looked through the values $\{0.1, 0.01, 0.001, 0.0001\}$ of noise intensity σ ; for every one of these values, we generated at random 20 signals x from \mathcal{X} and recorded the $\|\cdot\|_2$ -recovery errors of the linear and the polyhedral estimates. In addition to testing the nearly optimal polyhedral estimate PolyI yielded by Proposition 4.29, we tested the performance of the polyhedral estimate PolyII yielded under the circumstances by construction from Section 4.6.4. The observed $\|\cdot\|_2$ -recovery errors of our three estimates, sorted in the non-descending order, are represented on Figure 4.2. We see that the empirical performances of all three estimates are rather similar, with a clear tendency for the polyhedral estimate PolyII seems to work better than, or at the very worst similarly to, PolyI, in spite of the fact that in the situation in question the estimate PolyI, in contrast to PolyII, is provably near-optimal.





Figure 4.2: Recovery errors for near-optimal linear estimate (red) and the polyhedral estimates yielded by Proposition 4.29 (PolyI, blue) and by construction from Section 4.6.4 (PolyII, cyan).

4.7 RECOVERING SIGNALS FROM NONLINEAR OBSERVATIONS BY STOCHASTIC OPTIMIZATION

The "common denominator" of all estimation problems considered so far in Lecture 4 is that what we observed was obtained by adding noise to the *linear* image of the unknown signal we want to recover. In this section we intend to consider the signal recovery problem in the case where the observation is obtained by adding noise to a *nonlinear* transformation of the signal.

4.7.1 Problem's setting

Motivating example for what follows is what is called *logistic regression* model, where

• the unknown signal to be recovered is a vector x known to belong to a given signal set $\mathcal{X} \subset \mathbf{R}^n$, which we assume to be a nonempty convex compact set;

• our observation

$$\omega^K = \{\omega_k = (\eta_k, y_k), 1 \le k \le K\}$$

stemming from a signal x is as follows:

- the regressors $\eta_1, ..., \eta_K$ are i.i.d. realizations of *n*-dimensional random vector η with distribution Q independent of x and such that Q possesses finite and positive definite matrix $\mathbf{E}_{\eta \sim Q} \{\eta \eta^T\}$ of second moments;
- the labels y_k are generated as follows: y_k is independent of the "history" $\eta_1, ..., \eta_{k-1}, y_1, ..., y_{k-1}$ random variable taking values 0 and 1, and the probability, η_k given, for y_k to take value 1 is $\phi(\eta_k^T x)$, where

$$\phi(s) = \frac{\exp\{s\}}{1 + \exp\{s\}}.$$

In this model, the standard (and very well studied) way to recover the signal x underlying observations is to use Maximum Likelihood (ML) estimate: the conditional, given observed regressors η_k , $k \leq K$, probability to get the observed labels as a function of a candidate signal z is

$$p(z, \omega^{K}) = \sum_{k=1}^{K} \left[y_{k} \ln \left(\frac{\exp\{\eta_{k}^{T} z\}}{1 + \exp\{\eta_{k}^{T} z\}} \right) + (1 - y_{k}) \ln \left(\frac{1}{1 + \exp\{\eta_{k}^{T} x\}} \right) \right]$$

= $\left[\sum_{k} y_{k} \eta_{k} \right]^{T} z - \sum_{k} \ln \left(1 + \exp\{\eta_{k}^{T} z\} \right),$ (4.147)

and the ML estimate of the "true" signal x underlying our observation ω^{K} is obtained by maximizing the log-likelihood $p(z, \omega^{K})$ over $z \in \mathcal{X}$:

$$\widehat{x}_{\mathrm{ML}}(\omega^{K}) \in \operatorname*{Argmax}_{z \in \mathcal{X}} p(z, \omega^{K}), \tag{4.148}$$

which is a convex optimization problem.

The problem we intend to consider can be viewed as a natural generalization of the just presented logistic regression and is as follows:

Our observation depends on unknown signal x known to belong to a given convex compact set $\mathcal{X} \subset \mathbf{R}^n$ and is

$$\omega^{K} = \{\omega_{k} = (\eta_{k}, y_{k}), 1 \le k \le K\}$$
(4.149)

with ω_k , $1 \leq k \leq K$ which are i.i.d. realizations of a random pair (η, y) with the distribution P_x such that

- the regressor η is a random $n \times m$ matrix with some independent of x probability distribution Q;
- the label y is m-dimensional random vector such that the conditional, η given, distribution of y induced by P_x is with the mean $f(\eta^T x)$:

$$\mathbf{E}_{|\eta}\{y\} = f(\eta^T x), \tag{4.150}$$

where $\mathbf{E}_{|\eta}$ is the conditional, η given, distribution of y stemming from the distribution P_x of $\omega = (\eta, y)$, and $f(\cdot) : \mathbf{R}^m \to \mathbf{R}^m$ is a given mapping.

Note that the logistic regression model corresponds to the case where m = 1,

 $f(s) = \frac{\exp\{s\}}{1+\exp\{s\}}$, and y takes values 0,1, with the conditional, η given, probability to take value 1 equal to $f(\eta^T x)$.

Another meaningful example is the one where

$$y = f(\eta^T x) + \xi,$$

where ξ is independent of η random vector with zero mean, say, $\xi \sim \mathcal{N}(0, \sigma^2 I_m)$. Note that in the latter case the ML estimate of the signal x underlying observations is

$$\widehat{x}_{\mathrm{ML}}(\omega^{K}) \in \operatorname{Argmin}_{z \in \mathcal{X}} \sum_{k} \|y_{k} - f(\eta_{k}^{T} z)\|_{2}^{2}.$$
(4.151)

In contrast to what happens with logistic regression, now the optimization problem – "nonlinear Least Squares" – responsible for the ML estimate typically is nonconvex and can be computationally difficult.

We intend to impose on the data of the estimation problem we have just described (namely, on \mathcal{X} , $f(\cdot)$, and the distributions P_x , $x \in \mathcal{X}$, of the pair (η, y)) assumptions which allow to reduce our estimation problem to a problem with convex structure — strongly monotone variational inequality represented by stochastic oracle, which will allow at the end of the day to get consistent estimate, with explicit "finite sample" accuracy guarantees, of the signal we want to recover.

4.7.2 Assumptions

Preliminaries: monotone vector fields. A monotone vector field on \mathbb{R}^m is a single-valued everywhere defined mapping $g(\cdot) : \mathbb{R}^m \to \mathbb{R}^m$ which possesses the monotonicity property

$$[g(z) - g(z')]^T[z - z'] \ge 0 \ \forall z, z' \in \mathbf{R}^m.$$

We say that such a field is monotone with modulus $\varkappa \geq 0$ on a closed convex set $Z \subset \mathbf{R}^m$, if

$$[g(z) - g(z')]^T [z - z'] \ge \varkappa ||z - z'||_2^2, \forall z \, z' \in Z,$$

and say that g is strongly monotone on Z if the modulus of monotonicity of g on Z is positive. It is immediately seen that for a monotone vector field which is continuously differentiable on a closed convex set Z with a nonempty interior, the necessary and sufficient condition for being monotone with modulus \varkappa on the set is

$$d^T f'(z) d \ge \varkappa d^T d \ \forall (d \in \mathbf{R}^n, z \in Z).$$
(4.152)

Basic examples of monotone vector fields are:

- gradient fields $\nabla \phi(x)$ of continuously differentiable convex functions of m variables or, more generally the vector fields $[\nabla_x \phi(x, y); -\nabla_y \phi(x, y)]$ stemming from continuously differentiable functions $\phi(x, y)$ which are convex in x and concave in y;
- "diagonal" vector fields $f(x) = [f_i(x_1); f_2(x_2); ...; f_m(x_m)]$ with monotonically nondecreasing univariate components $f_i(\cdot)$. If, in addition, $f_i(\cdot)$ are continuously differentiable with positive derivatives, then the associated filed f is strongly monotone on every compact convex subset of \mathbf{R}^m , the monotonicity modulus depending on the subset.

Monotone vector fields on \mathbb{R}^n admit simple calculus which includes, in particular, the following two rules:

I. [affine substitution of argument]: If $f(\cdot)$ is monotone vector field on \mathbb{R}^m and A is an $n \times m$ matrix, the vector field

$$g(x) = Af(A^T x + a)$$

is monotone on \mathbf{R}^n ; if, in addition, f is monotone with modulus $\varkappa \geq 0$ on a closed convex set $Z \subset \mathbf{R}^m$ and $X \subset \mathbf{R}^n$ is closed, convex, and such that $A^T x + a \in Z$ whenever $x \in X$, g is monotone with modulus $\sigma^2 \varkappa$ on X, where σ is the minimal singular value of A^T .

II. [summation]: If S is a Polish space, $f(x, s) : \mathbf{R}^m \times S \to \mathbf{R}^m$ is a Borel vectorvalued function which is monotone in x for every $s \in S$ and $\mu(ds)$ is a Borel probability measure on S such that the vector field

$$F(x) = \int_{S} f(x,s)\mu(ds)$$

is well defined for all x, then $F(\cdot)$ is monotone. If, in addition, X is a closed convex set in \mathbb{R}^m and $f(\cdot, s)$ is monotone on X with Borel in s modulus $\varkappa(s)$ for every $s \in S$, then F is monotone on X with modulus $\int_S \varkappa(s)\mu(ds)$.

Assumptions. In what follows, we make the following assumptions on the data of the estimation problem posed in Section 4.7.1:

• A.1. The vector field $f(\cdot)$ is continuous and monotone, and the vector field

$$F(z) = \mathbf{E}_{\eta \sim Q} \left\{ \eta f(\eta^T z) \right\}$$

is well defined (and therefore is monotone along with f by \mathbf{I} , \mathbf{II});

- A.2. The signal set \mathcal{X} is a nonempty convex compact set, and the vector field F is monotone with positive modulus \varkappa on \mathcal{X} ;
- A.3. For properly selected $M < \infty$ and every $x \in \mathcal{X}$ it holds

$$\mathbf{E}_{(\eta,y)\sim P_x}\left\{\|\eta y\|_2^2\right\} \le M^2. \tag{4.153}$$

A simple *sufficient* condition for the validity of Assumptions A.1-3 with properly selected $M < \infty$ and $\varkappa > 0$ is as follows:

- The distribution Q of η has finite moments of all orders, and $\mathbf{E}_{\eta \sim Q}\{\eta \eta^T\} \succ 0$;
- f is continuously differentiable, and $d^T f'(z) d > 0$ for all $d \neq 0$ and all z. Besides this, f is with polynomial growth: for some constants $C \ge 0$ and $p \ge 0$ and all z one has $||f(z)||_2 \le C(1 + ||z||_2^p)$.

Verification of sufficiency is straightforward.

4.7.3 Main observation

The main observation underlying the construction we are about to build is as follows.

Observation 4.32. With Assumptions **A.1-3** in force, let us associate with pair $(\eta, y) \in \mathbf{R}^{n \times m} \times \mathbf{R}^m$ the vector field

$$G_{(\eta,y)}(z) = \eta f(\eta^T z) - \eta y : \mathbf{R}^n \to \mathbf{R}^n.$$

$$(4.154)$$

For every $x \in \mathcal{X}$, denoting by P_x the common distribution of observations (4.149) stemming from signal $x \in \mathcal{X}$, we have

$$\mathbf{E}_{(\eta,y)\sim P_x} \left\{ \begin{array}{ll} G_{(\eta,y)}(z) \right\} &= F(z) - F(x) \ \forall z \in \mathbf{R}^n \quad (a) \\ \|F(z)\|_2 &\leq M \ \forall z \in \mathcal{X} \qquad (b) \\ \mathbf{E}_{(\eta,y)\sim P_x} \left\{ \|G_{(\eta,y)}(z)\|_2^2 \right\} &\leq 4M^2 \ \forall z \in \mathcal{X} \qquad (c) \end{array}$$
(4.155)

Proof is immediate. Specifically, let $x \in \mathcal{X}$. Then, taking into account that the marginal distribution of η induced by distribution P_x of (η, y) is Q and denoting by $\mathbf{E}_{|\eta}$ the conditional, η given, expectation over y induced by P_x , we have

$$\mathbf{E}_{(\eta,y)\sim P_x}\{\eta y\} = \mathbf{E}_{\eta\sim Q}\left\{\mathbf{E}_{|\eta}\{\eta y\}\right\} = \mathbf{E}_{\eta}\left\{\eta f(\eta^T x)\right\} = F(x)$$

(we have used (4.150) and the definition of F), whence,

$$\mathbf{E}_{(\eta,y)\sim P_x} \left\{ G_{(\eta,y)}(z) \right\} = \mathbf{E}_{(\eta,y)\sim P_x} \left\{ \eta f(\eta^T z) - \eta y \right\} = \mathbf{E}_{(\eta,y)\sim P_x} \left\{ \eta f(\eta^T z) \right\} - F(x)$$

= $\mathbf{E}_{\eta\sim Q} \left\{ \eta f(\eta^T z) \right\} - F(x) = F(z) - F(x),$

as stated in (4.155.*a*). Besides this, for $x, z \in \mathcal{X}$, denoting by $P_{|\eta}^z$ the conditional, η given, distribution of y induced by the distribution P_z of (η, y) , and taking into account that the marginal distribution of η induced by P_z is Q, we have

$$\begin{aligned} \mathbf{E}_{(\eta,y)\sim P_x}\{\|\eta f(\eta^T z)\|_2^2\} &= \mathbf{E}_{\eta\sim Q}\left\{\|\eta f(\eta^T z)\|_2^2\right\} \\ &= \mathbf{E}_{\eta\sim Q}\left\{\|\mathbf{E}_{y\sim P_{|\eta}^z}\{\eta y\}\|_2^2\right\} \text{ [since } \mathbf{E}_{y\sim P_{|\eta}^z}\{y\} = f(\eta^T z)] \\ &\leq \mathbf{E}_{\eta\sim Q}\left\{\mathbf{E}_{y\sim P_{|\eta}^z}\left\{\|\eta y\|_2^2\right\}\right\} \text{ [by Jensen's inequality]} \\ &= \mathbf{E}_{(\bar{\eta},y)\sim P_z}\left\{\|\eta y\|_2^2\right\} \\ &\leq M^2 \text{ [by } \mathbf{A.3} \text{ due to } z \in \mathcal{X}]. \end{aligned}$$

This combines with the relation $\mathbf{E}_{(\eta,y)\sim P_x}\{\|\eta y\|_2^2\} \leq M^2$ given by **A.3** due to $x \in \mathcal{X}$ to imply (4.155.*b*) and (4.155.*c*).

The consequences. Our goal is to recover the signal $x \in \mathcal{X}$ underlying observations (4.149), and under assumptions A.1-3, x is a root of the monotone vector field

$$G(z) = F(z) - F(x), \ F(z) = \mathbf{E}_{\eta \sim Q} \left\{ \eta f(\eta^T z) \right\},$$
(4.156)

and we know that this root belongs to \mathcal{X} ; moreover, since $G(\cdot)$ is strongly monotone on \mathcal{X} along with $F(\cdot)$, this root is unique. Now, finding a root, known to belong to a given convex compact set \mathcal{X} , of a strongly monotone on this set vector field G is known to be a computationally tractable problem, provided we have access to "oracle" which, given on input a point $z \in \mathcal{X}$, returns the value G(z) of the field at the point. In the situation we are interested in the latter is not exactly the case: the field G is the expectation of a random field:

$$G(z) = \mathbf{E}_{(\eta, y) \sim P_x} \left\{ \eta f(\eta^T z) - \eta y \right\},$$

and we do not know a priori what is the distribution over which the expectation is

taken. We, however, can sample from this distribution – the samples are exactly the observations (4.149), and we can use these samples to approximate somehow G and use this approximation to approximate the signal x. The two standard ways to implement this idea are to use *Sample Average Approximation* (SAA) and *Stochastic Approximation* (SA). We are about to consider these two techniques as applied to the situation we are in.

Sample Average Approximation. The idea underlying SAA is quite transparent: given observations (4.149), let us approximate the field of interest G with its empirical counterpart

$$G_{\omega^{K}}(z) = \frac{1}{K} \sum_{k=1}^{K} \left[\eta_{k} f(\eta_{k}^{T} z) - \eta_{k} y_{k} \right].$$

By the Law of Large Numbers, as $K \to \infty$, the empirical field $G_{\omega^{K}}$ converges to the field of interest G, so that under mild regularity assumptions, when K is large, $G_{\omega^{K}}$, with overwhelming probability, will be uniformly on \mathcal{X} close to G, which, due to strong monotonicity of G, would imply that a "near-zero" of $G_{\omega^{K}}$ on \mathcal{X} will be close to the zero x of G, which is nothing but the signal we want to recover. The only question is how to define a "near-zero" of $G_{\omega^{K}}$ on \mathcal{X}^{68} . The most convenient in our context notion of a "near-zero" is the concept of a *weak solution* to a Variational Inequality with monotone operator, defined as follows (we restrict the general definition to the situation we are interested in):

Let $\mathcal{X} \subset \mathbf{R}^n$ be a nonempty convex compact set, and $H(z) : \mathcal{X} \to \mathbf{R}^n$ be a monotone (i.e., $[H(z) - H(z')]^T [z - z'] \ge 0$ for all $z, z' \in \mathcal{X}$) vector field. A vector $z_* \in \mathcal{X}$ is called a *weak solution* to the variational inequality (VI) associated with H, \mathcal{X} when

$$H^T(z)(z-z_*) \ge 0 \,\forall z \in \mathcal{X}.$$

Let $\mathcal{X} \subset \mathbf{R}^n$ be a nonempty convex compact set and H be monotone on \mathcal{X} . It is well known that

- The VI associated with H, \mathcal{X} (let us denote it $VI(H, \mathcal{X})$) always has a weak solution. Besides this, it is clear that if $\overline{z} \in \mathcal{X}$ is a root of H, then \overline{z} is a weak solution to $VI(H, \mathcal{X})^{69}$.
- When H is continuous on X, every weak solution z̄ to VI(H, X) is also a strong solution, meaning that

$$H^{T}(\bar{z})(z-\bar{z}) \ge 0 \ \forall z \in \mathcal{X}.$$

$$(4.157)$$

Indeed, (4.157) clearly holds true when $z = \bar{z}$. Assuming $z \neq \bar{z}$ and setting $z_t = \bar{z} + t(z-\bar{z}), 0 < t \leq 1$, we have $H^T(z_t)(z_t-\bar{z}) \geq 0$ (since \bar{z} is a weak solution), whence $H^T(z_t)(z-\bar{z}) \geq 0$ (since $z-\bar{z}$ is a positive multiple of $z_t - \bar{z}$). Passing to limit as $t \to +0$ and invoking the continuity of H, we get $H^T(\bar{z})(z-\bar{z}) \geq 0$, as claimed.

⁶⁸note that we in general cannot define a "near-zero" of $G_{\omega K}$ on \mathcal{X} as a root of $G_{\omega K}$ on this set – while G does have root belonging to \mathcal{X} , nobody told us that the same holds true for $G_{\omega K}$. ⁶⁹indeed, when $\bar{z} \in \mathcal{X}$ and $H(\bar{z}) = 0$, monotonicity of H implies that $H^T(z)[z - \bar{z}] = [H(z) - H(\bar{z})]^T[z - \bar{z}] \ge 0$ for all $z \in \mathcal{X}$, that is, \bar{z} is a weak solution to the VI.

• When H is the gradient field of a continuously differentiable convex function on \mathcal{X} (such a field indeed is monotone), weak (or, which in the case of continuous H is the same, strong) solutions to $VI(H, \mathcal{X})$ are exactly the minimizers of the function on \mathcal{X} .

In the sequel, we heavily exploit the following simple and well known fact:

Lemma 4.33. Let \mathcal{X} be a convex compact set, H be a monotone vector field on \mathcal{X} with monotonicity modulus $\varkappa > 0$, and let \overline{z} be a weak solution to $VI(H, \mathcal{X})$. Then the weak solution to $VI(H, \mathcal{X})$ is unique. Besides this,

$$H^{T}(z)[z-\bar{z}] \ge \varkappa ||z-\bar{z}||_{2}^{2}.$$
(4.158)

Proof: Under the premise of Lemma, let $z \in \mathcal{X}$ and let \bar{z} be a weak solution to $VI(H, \mathcal{X})$ (recall that it does exist). Setting $z_t = \bar{z} + t(z - \bar{z})$, for $t \in (0, 1)$ we have

$$H^{T}(z)[z-z_{t}] \geq H^{T}(z_{t})[z-z_{t}] + \varkappa ||z-z_{t}||^{2} \geq \varkappa ||z-z_{t}||^{2},$$

where the first \geq is due to strong monotonicity of H, and the second \geq is due to the fact that $H^T(z_t)[z - z_t]$ is proportional, with positive coefficient, to $H^T(z_t)[z_t - \bar{z}]$, and the latter quantity is nonnegative since \bar{z} is a weak solution to the VI in question. We end up with $H^T(z)(z - z_t) \geq \varkappa ||z - z_t||_2^2$; passing to limit as $t \to +0$, we arrive at (4.158). To prove uniqueness of a weak solution, assume that aside of the weak solution \bar{z} there exists a weak solution \tilde{z} distinct form \bar{z} , and let us set $z' = \frac{1}{2}[\bar{z} + \tilde{z}]$. Since both \bar{z} and \tilde{z} are weak solutions, both the quantities $H^T(z')[z'-\bar{z}]$ and $H^T(z')[z'-\tilde{z}]$ should be nonnegative, and since the sum of these quantities is 0, both of them are zero; applying (4.158) to z = z', we get $z' = \bar{z}$, whence $\tilde{z} = \bar{z}$ as well.

Now let us come back to the estimation problem we are considering, and let Assumptions **A.1-3** be satisfied. Note that then the vector fields $G_{(\eta_k, y_k)}(z)$ defined in (4.154), and therefore the vector field $G_{\omega^{\kappa}}(z)$ are continuous and monotone. With the SAA estimation, we compute a weak solution $\hat{x}(\omega^K)$ to VI($G_{\omega^{\kappa}}, \mathcal{X}$) and treat it as the SAA estimate of signal x underlying observations (4.149). Since the vector field $G_{\omega^{\kappa}}(\cdot)$ is monotone and its values at given points are efficiently computable, provided the values of f are so, computing a (whatever high accuracy approximation to) a weak solution to VI($G_{\omega^{\kappa}}, \mathcal{X}$) is a computationally tractable problem (see, e.g., [118]). Moreover, utilizing the techniques from [133, Chapter 5], under mild additional to **A.1-3** regularity assumptions one can get non-asymptotical upper bound on, say, the expected $\|\cdot\|_2^2$ -deviation of the SAA estimate from the true signal x as a function of the sample size K and find out the rate at which this bound converges to 0 as $K \to \infty$; this analysis, however, goes beyond our scope.

Let us look what is the SAA estimate in the logistic regression model. In this case we have

$$\begin{aligned} G_{(\eta_k, y_k)}(z) &= \left[\frac{\exp\{\eta_k^T z\}}{1 + \exp\{\eta_k^T z\}} - y_k \right] \eta_k, \\ G_{\omega^K}(z) &= \frac{1}{K} \sum_{k=1}^K \left[\frac{\exp\{\eta_k^T z\}}{1 + \exp\{\eta_k^T z\}} - y_k \right] \eta_k \\ &= \frac{1}{K} \nabla_z \sum_k \left[\ln\left(1 + \exp\{\eta_k^T z\}\right) - y_k \eta_k^T z \right], \end{aligned}$$

that is, $G_{\omega^K}(z)$ is proportional, with negative coefficient (-1/K), to the gradient field of the log-likelihood $p(z, \omega^K)$, see (4.147). As a result, in the case in question

the weak solutions to VI($G_{\omega^{K}}, \mathcal{X}$) are exactly the maximizers of the log-likelihood $p(z, \omega^{K})$ over $z \in \mathcal{X}$, that is, for the logistic regression the SAA estimate is nothing but the Maximum Likelihood estimate $\hat{x}_{\mathrm{ML}}(\omega^{K})$ as defined in (4.148). Note that this phenomenon is specific for the logistic regression model⁷⁰. Say, in the "nonlinear least squares" example described in Section 4.7.1 with (for the sake of simplicity, scalar) monotone $f(\cdot)$ the vector field $G_{\omega^{K}}(\cdot)$ is given by

$$G_{\omega^{K}}(z) = \frac{1}{K} \sum_{k=1}^{K} \left[f(\eta_{k}^{T} z) - y_{k} \right] \eta_{k}$$

which is quite different (provided that f is nonlinear) from the gradient field

$$\Gamma_{\omega^{K}}(z) = 2\sum_{k=1}^{K} f'(\eta_{k}^{T}z) \left[f(\eta_{k}^{T}z) - y_{k} \right] \eta_{k}$$

of the minus log-likelihood appearing in (4.151). As a result, in this case the ML estimate (4.151) is, in general, different from the SAA estimate, especially when taking into account that the SAA estimate, in contrast to the ML one, is easy to compute.

Stochastic Approximation estimate. The Stochastic Approximation (SA) estimate stems from a simple algorithm – Subgradient Descent – for solving variational inequality $VI(G, \mathcal{X})$. Were the values of the vector field $G(\cdot)$ are available, one could approximate a root $x \in \mathcal{X}$ of this VI by running the recurrence

$$z_{k+1} = \operatorname{Proj}_{\mathcal{X}}[z_k - \gamma_k G(z_k)], \ k = 1, 2, ..., K,$$

where

• $\operatorname{Proj}_{\mathcal{X}}[z]$ is the metric projection of \mathbf{R}^n onto \mathcal{X} :

$$\operatorname{Proj}_{\mathcal{X}}[z] = \operatorname*{argmin}_{u \in \mathcal{X}} \|z - u\|_2;$$

- $\gamma_k > 0$ are properly selected stepsizes;
- the starting point z_1 is an arbitrary point of \mathcal{X} .

$$\Phi(z) = \sum_{k} \left[\frac{\phi(\eta_k^T z)}{1 + \phi(\eta_k^T z)} - y_k \right] \eta_k,$$

while the gradient field of the quantity we want to minimize when computing the ML estimate (this quantity is the minus log-likelihood $-\sum_k \left[y_k \ln(f(\eta_k^T z)) + (1 - y_k) \ln(1 - f(\eta_k^T z))\right]$) is

$$\Psi(z) = \sum_{k} \frac{\phi'(\eta_k^T z)}{\phi(\eta_k^T z)} \left[\frac{\phi(\eta_k^T z)}{1 + \phi(\eta_k^T z)} - y_k \right] \eta_k;$$

when k > 1 and ϕ is not an exponent, Φ and Ψ are "essentially different," so that the SAA estimate typically will differ from the ML one.

⁷⁰The equality between the SAA and the ML estimates in the case of logistic regression is due to the fact that the logistic sigmoid $f(s) = \exp\{s\}/(1 + \exp\{s\})$ "happens" to satisfy the identity f'(s) = f(s)(1 - f(s)). When replacing in the logistic model the above sigmoid with $f(s) = \phi(s)/(1 + \phi(s))$, with differentiable monotonically nondecreasing positive $\phi(\cdot)$, the SAA estimate becomes the weak solution to VI(Φ, \mathcal{X}) with

It is well known that under Assumptions **A.1-3** this recurrence with properly selected stepsizes and started at a point from \mathcal{X} allows to approximate the root of G (in fact, the unique weak solution to $\operatorname{VI}(G, \mathcal{X})$) to a whatever high accuracy, provided K is large enough. We, however, are in the situation when the actual values of G are not available; the standard way to cope with this difficulty is to replace in the above recurrence the "unobservable" values $G(z_k)$ of G with their unbiased random estimates $G_{(\eta_k, y_k)}(z_k)$. This modification gives rise to *Stochastic Approximation* (coming back to [93]) – the recurrence

$$z_{k+1} = \operatorname{Proj}_{\mathcal{X}}[z_k - \gamma_k G_{(\eta_k, y_k)}(z_k)], \ 1 \le k \le K,$$
(4.159)

where z_1 is a once for ever chosen point from \mathcal{X} , and $\gamma_k > 0$ are deterministic.

What is on our agenda now is the (perfectly well known) convergence analysis of SA under assumptions **A.1-3**. To this end observe that z_k are deterministic functions of the initial fragments $\omega^{k-1} = \{\omega_t, 1 \leq t < k\}$ of our sequence of observations $\omega^K = \{\omega_k = (\eta_k, y_k), 1 \leq k \leq K\}$:

$$z_k = Z_k(\omega^{k-1}).$$

Let us set

$$D_k = \frac{1}{2} \|Z_k(\omega^{k-1}) - x\|_2^2, \ d_k = \mathbf{E}_{\omega^{k-1} \sim P_x \times \dots \times P_x} \{D_k(\omega^{k-1})\},\$$

where $x \in \mathcal{X}$ is the signal underlying observations (4.149). Note that, as it is well known, the metric projection onto a closed convex set \mathcal{X} possesses the property as follows:

$$\forall (z \in \mathbf{R}^n, u \in \mathcal{X}) : \|\operatorname{Proj}_{\mathcal{X}}[z] - u\|_2 \le \|z - u\|_2.$$

Consequently, for $1 \le k \le K$ it holds

$$\begin{aligned} D_{k+1}(\omega^{k}) &= \frac{1}{2} \|\operatorname{Proj}_{\mathcal{X}}[Z_{k}(\omega^{k-1}) - \gamma_{k}G_{\omega_{k}}(Z_{k}(\omega^{k-1}))] - x\|_{2}^{2} \\ &\leq \frac{1}{2} \|Z_{k}(\omega^{k-1}) - \gamma_{k}G_{\omega_{k}}(Z_{k}(\omega^{k-1})) - x\|_{2}^{2} \\ &= \frac{1}{2} \left[\|Z_{k}(\omega^{k-1}) - x\|_{2}^{2} - \gamma_{k}G_{\omega_{k}}^{T}(Z_{k}(\omega^{k-1}))[Z_{k}(\omega^{k-1}) - x] + \frac{1}{2}\gamma_{k}^{2} \|G_{\omega_{k}}(Z_{k}(\omega^{k-1}))\|_{2}^{2} \right]. \end{aligned}$$

Taking expectations w.r.t. $\omega^k \sim \underbrace{P_x \times \ldots \times P_x}_{P_x^k}$ of both sides of the resulting inequal-

ity and keeping in mind relations (4.155) along with the fact that Z_k takes values in \mathcal{X} , we get

$$d_{k+1} \le d_k - \gamma_k \mathbf{E}_{\omega^{k-1} \sim P_x^{k-1}} \left\{ G^T(Z_k(\omega^{k-1}))[Z_k(\omega^{k-1}) - x] \right\} + 2\gamma_k^2 M^2.$$
(4.160)

Recalling that we are in the case when G is strongly monotone, with modulus $\varkappa > 0$, on \mathcal{X} , x is the weak solution VI(G, \mathcal{X}), and Z_k takes values in \mathcal{X} , invoking (4.158), the expectation in (4.160) is at least $\varkappa d_k$, and we arrive at the relation

$$d_{k+1} \le (1 - \varkappa \gamma_k) d_k + 2\gamma_k^2 M^2.$$
(4.161)

Let us set

$$S = \frac{8M^2}{\varkappa^2}, \, \gamma_k = \frac{\varkappa S}{4M^2(k+1)} = \frac{2}{\varkappa(k+1)}$$

and verify by induction in k that for k = 1, 2, ..., K + 1 it holds

$$d_k \le \frac{S}{k+1}.\tag{*}{k}$$

Base k = 1. Let D be the $\|\cdot\|_2$ -diameter of \mathcal{X} , and $z_{\pm} \in \mathcal{Z}$ be such that $\|z_{+} - z_{-}\|_2 = D$. By (4.155) we have $\|F(z)\|_2 \leq M$ for all $z \in \mathcal{X}$, and by strong monotonicity of $G(\cdot)$ on \mathcal{X} we have

$$[G(z_{+}) - G(z_{-})]^{T}[z_{+} - z_{-}] = [F(z_{+}) - F(z_{-})][z_{+} - z_{-}] \ge \varkappa ||z_{+} - z_{-}||_{2}^{2} = \varkappa D^{2};$$

By Cauchy inequality, the left hand side in the concluding \geq is at most 2MD, and we get

$$D \le \frac{2M}{\varkappa},$$

whence $S \ge 2D^2$, so that $\frac{S}{2} \ge D^2$. On the other hand, due to the origin d_1 we have $d_1 \le D^2$. Thus, $(*_1)$ holds true.

Inductive step $(*_k) \Rightarrow (*_{k+1})$. Now assume that $(*_k)$ holds true for some k, $1 \le k \le K$, and let us prove that $(*_{k+1})$ holds true as well. Observe that $\varkappa \gamma_k = 2/(k+1) \le 1$, so that

$$\begin{aligned} &d_{k+1} \leq d_k (1 - \varkappa \gamma_k) + 2\gamma_k^2 M^2 \text{ [by (4.161)]} \\ &\leq \frac{S}{k+1} (1 - \varkappa \gamma_k) + 2\gamma_k^2 M^2 \text{ [by (*_k) and due to } \varkappa \gamma_k \leq 1] \\ &= \frac{S}{k+1} \left[1 - \frac{2}{k+1} \right] + \frac{8M^2}{\varkappa^2 (k+1)^2} = \frac{8M^2 (k-1)}{\varkappa^2 (k+1)^2} + \frac{8M^2}{\varkappa^2 (k+1)^2} \\ &= \frac{8kM^2}{\varkappa^2 (k+1)^2} = \frac{S}{k+2} \frac{k(k+2)}{(k+1)^2} = \frac{S}{k+2} \frac{k^2 + 2k}{k^2 + 2k+1} \leq \frac{S}{k+2}, \end{aligned}$$

so that $(*_{k+1})$ hods true. Induction is complete. Recalling what d_k is, we have arrived at the following

Proposition 4.34. Under Assumptions A.1-3 and with the stepsizes

$$\gamma_k = \frac{2}{\varkappa(k+1)}, \ k = 1, 2, \dots$$
 (4.162)

for every signal $x \in \mathcal{X}$ the sequence of estimates $\hat{x}_k(\omega^k) = z_{k+1}$ with $z_k = Z_k(\omega^k)$ given by the SA recurrence (4.159) and $\omega_k = (\eta_k, y_k)$ given by (4.149) for every k one has

$$\mathbf{E}_{\omega^{k} \sim P_{x}^{k}}\left\{\|\widehat{x}_{k}(\omega^{k}) - x\|_{2}^{2}\right\} \leq \frac{16M^{2}}{\varkappa^{2}(k+2)}, \ k = 1, 2, \dots$$
(4.163)

 P_x being the distribution of (η, y) stemming from signal x.

4.7.4 Numerical illustration

To illustrate the above developments, we present here results of some numerical experiments. Our setup is deliberately simplistic and is as follows:

- $\mathcal{X} = \{ x \in \mathbf{R}^n : ||x||_2 \le 1 \};$
- the distribution Q of η is $\mathcal{N}(0, I_n)$;
- f is the monotone vector field on \mathbf{R} given by one of the following four options:

- A. $f(s) = \exp\{s\}/(1 + \exp\{s\});$
- B. f(s) = s;
- C. $f(s) = \max[s, 0];$
- D. $f(s) = \min[1, \max[s, 0]].$
- conditional, η given, distribution of y induced by P_x is
 - Bernoulli distribution with probability $f(\eta^T x)$ of outcome 1 in the case of A (i.e., A results in the logistic model),
 - Gaussian distribution $\mathcal{N}(f(\eta^T x), I_n)$ in cases B D.

The dimension n in all or experiments was set to 100, and the number of observations K took values 1000, 10000, and 50000. The signal x underlying observations (4.149) in a particular experiment was selected at random from the uniform distribution on the unit sphere (which is the boundary of our \mathcal{X}).

In a particular experiment, we computed the SAA and the SA estimates (note that in the cases A, B the SAA estimate is the Maximum Likelihood estimate as well).

In SA, the stepsizes γ_k were selected according to (4.162) with "empirically selected" \varkappa . Specifically, given observations $\omega_k = (\eta_k, y_k), k \leq K$, see (4.149), we used them to build the SA estimate in two stages:

— training stage, where we generate at random "training signal" $x' \in \mathcal{X}$ and then generate labels y'_k as if x' were the actual signal; for example, in the case of A y'_k is assigned value 1 with probability $f(\eta_k^T x')$ and value 0 with complementary probability. After "training signal" and associated labels are generated, we run on the resulting artificial observations SA with different values of \varkappa , look how well the resulting estimate recovers x', and select the value of \varkappa resulting in the best recovery;

— *execution stage*, where we run on the actual data SA with stepsizes (4.162) specified by \varkappa found at the training stage.

case	K	error, SAA	cpu, SAA	error, SA	cpu, SA
A	1000	0.678	7.2	0.677	1.04
A	10000	0.201	93.8	0.213	3.5
A	50000	0.100	520.0	0.102	15.6
B	1000	0.290	4.3	0.325	0.9
B	10000	0.106	69.4	0.106	1.9
B	50000	0.049	397.5	0.0.49	8.8
C	1000	0.611	4.0	0.656	0.3
C	10000	0.216	76.1	0.223	1.5
C	50000	0.095	420.0	0.098	7.0
D	1000	0.719	4.8	0.743	0.3
D	10000	0.268	72.4	0.280	1.5
D	50000	0.111	425.6	0.112	9.3

The results of typical numerical experiments are as follows:

recovery error in $\|\cdot\|_2$, cpu time in sec

Note that the cpu time for SA includes both the training and the execution stages. The conclusion from these (simplistic!) experiments is that as far as estimation quality is concerned, the SAA estimate marginally outperforms the SA one, while

In principle, we could get (lower bounds on) the modules of strong monotonicity of the vectors fields $F(\cdot)$ we are interested in analytically, but this would be boring and conservative.

being essentially more time consuming. Note also that the observed in our experiments dependence of recovery errors on K is consistent with the convergence rate $O(1/\sqrt{K})$ established by Proposition 4.34.

4.7.5 "Single-observation" case

Let us look at the special case of our estimation problem where the sequence $\eta_1, ..., \eta_K$ of regressors in (4.149) is deterministic. At the first glance, this situation goes beyond our setup, where the regressors should be i.i.d. drawn from some distribution Q. We can, however, circumvent this "contradiction" by saying that now we are speaking about *single-observation case* with the regressor being the matrix $[\eta_1, ..., \eta_K]$ and Q being a degenerate distribution – the one supported at a singleton. Specifically, consider the case where our observation is

$$\omega = (\eta, y) \in \mathbf{R}^{n \times mK} \times \mathbf{R}^{mK} \tag{4.164}$$

(m, n, K are given positive integers), and the distribution P_x of observation stemming from a signal $x \in \mathbf{R}^n$ is as follows:

- η is a given independent of x deterministic matrix;
- y is random, and the distribution of y induced by P_x is with mean $\phi(\eta x)$, where $\phi: \mathbf{R}^{mK} \to \mathbf{R}^{mK}$ is a given mapping.

As an instructive example linking our current setup with the previous one, one can consider the case where $\eta = [\eta_1, ..., \eta_K]$ with $n \times m$ deterministic "individual regressors" $\eta_k, y = [y_1; ...; y_K]$ with random "individual labels" $y_k \in \mathbf{R}^m$ independent, xgiven, across k and such that the induced by x expectations of y_k are $f(\eta_k^T x)$ for some $f : \mathbf{R}^m \to \mathbf{R}^m$, and to set $\phi([u_1; ...; u_K]) = [f(u_1); ...; f(u_K)]$. The resulting "single observation" model is a natural analogy of the K-observation model considered so far, the only difference being that the individual regressors now form a fixed deterministic sequence rather than to be i.i.d. samples of some random $n \times m$ matrix.

As everywhere in this section, our goal is to use observation (4.164) to recover the (unknown) signal x underlying, as explained above, the distribution of the observation. Formally, we are now in the case K = 1 of our previous recovery problem where Q is supported on a singleton $\{\eta\}$ and can use the constructions developed so far. Specifically,

• The vector field F(z) associated with our problem (it used to be $\mathbf{E}_{\eta \sim Q} \eta f(\eta^T z)$) is

$$F(z) = \eta \phi(\eta^T z),$$

and the vector field

$$G(z) = F(z) - F(x),$$

x being the signal underlying observation (4.164), is

$$G(z) = \mathbf{E}_{(\eta, y) \sim P_x} \{ F(z) - \eta y \}$$

(cf. (4.156)); as before, the signal we are interested to recover is a zero of the latter field. Note that now the vector field F(z) is observable, while before it was the expectation of an observable vector field, and the vector field G still is

the expectation, over P_x , of an observable vector field:

$$G(z) = \mathbf{E}_{(\eta, y) \sim P_x} \{ \underbrace{\eta \phi(\eta^T z) - \eta y}_{G_y(z)} \},$$

cf. Observation 4.32.

• Assumptions A.1-2 now read

A.1' The vector field $\phi(\cdot) : \mathbf{R}^{mK} \times \mathbf{R}^{mK}$ is continuous and monotone, so that $F(\cdot)$ is continuous and monotone as well,

A.2' \mathcal{X} is a nonempty compact convex set, and F is strongly monotone, with modulus $\varkappa > 0$, on \mathcal{X} ;

as before, a simple sufficient condition for the validity of the above monotonicity assumptions is positive definiteness of the matrix $\eta\eta^T$ plus strong monotonicity of ϕ on every bounded set.

• For our present purposes, it makes sense to reformulate assumption A.3 in the following equivalent form:

A.3' For properly selected $\sigma \geq 0$ and every $x \in \mathcal{X}$ it holds

$$\mathbf{E}_{(\eta,y)\sim P_x}\{\|\eta[y-\phi(\eta^T x)]\|_2^2\} \le \sigma^2.$$

In our present situation the SAA $\hat{x}(y)$ is the unique weak solution to $VI(G_y, \mathcal{X})$, and we can easily quantify the quality of this estimate:

Proposition 4.35. In the situation in question and under Assumptions A.1'-3' for every $x \in \mathcal{X}$ and every realization (η, y) of induced by x observation (4.164) one has

$$\|\widehat{x}(y) - x\|_{2} \le \varkappa^{-1} \|\underbrace{\eta[y - \phi(\eta^{T}x)]}_{\Delta(x,y)}\|_{2}, \qquad (4.165)$$

whence also

$$\mathbf{E}_{(\eta,y)\sim P_x}\{\|\hat{x}(y) - x\|_2^2\} \le \sigma^2/\varkappa^2.$$
(4.166)

Proof. Let $x \in \mathcal{X}$ be the signal underlying observation (4.164), and G(z) = F(z) - F(x) be the associated vector field G. We have

$$G_y(z) = F(z) - \eta y = F(z) - F(x) + [F(x) - \eta y] = G(z) - \eta [y - \phi(\eta^T x)] = G(z) - \Delta(x, y)$$

For y fixed, $\overline{z} = \widehat{x}(y)$ is the weak, and therefore the strong (since $G_y(\cdot)$ is continuous) solution to $VI(G_y, \mathcal{X})$, implying, due to $x \in \mathcal{X}$, that

$$0 \le G_y^T(\bar{z})[x - \bar{z}] = G^T(\bar{z})[x - \bar{z}] - \Delta^T(y)[x - \bar{z}],$$

whence

$$-G^T(\bar{z})[x-\bar{z}] \le -\Delta^T(x,y)[x-\bar{z}].$$

Besides this, G(x) = 0, whence $G^T(x)[x - \overline{z}] = 0$, and we arrive at

$$[G(x) - G(\bar{z})]^T [x - \bar{z}] \le -\Delta^T (x, y) [x - \bar{z}].$$
K	λ	$\operatorname{Err} = \ \widehat{x} - x\ _2$	$\operatorname{Err}/[\lambda\sqrt{n/K}]$	K	λ	$\mathrm{Err} = \ \widehat{x} - x\ _2$	$\operatorname{Err}/[\lambda\sqrt{n/K}]$
100	1.0	0.886	0.89	8100	1.0	0.239	2.15
100	0.1	0.685	6.85	8100	0.1	0.031	2.79
900	1.0	0.572	1.72	72900	1.0	0.080	2.17
900	0.1	0.098	2.95	72900	0.1	0.009	2.44

Table 4.2: Performance of SAA recovery, n = 100, m = 1

whence also

$$\varkappa \|x - \bar{z}\|_2^2 \le -\Delta^T(x, y) [x - \bar{z}]$$

(recall that G, along with F, is strongly monotone with modulus \varkappa on \mathcal{X} and $x, \overline{z} \in \mathcal{X}$). Applying the Cauchy inequality, we arrive at (4.165). \Box **Example.** Consider the case where m = 1, ϕ is strongly monotone, with modulus $\varkappa_{\phi} > 0$, on the entire \mathbf{R}^{K} , and η in (4.164) is drawn from "Gaussian ensemble" – the rows η_{k}^{T} of the $n \times K$ matrix η are independently of each other drawn from $\mathcal{N}(0, I_n)$. Assume also that the observation noise is Gaussian:

$$y = \phi(\eta^T x) + \lambda \xi, \ \xi \sim \mathcal{N}(0, I_K).$$

It is well known that when $K/n \geq 1$, the minimal singular value of the $n \times n$ matrix $\eta \eta^T$ with overwhelming as $K/n \to \infty$ probability is at least O(1)K, implying that when $K/n \gg 1$, the typical modulus of strong monotonicity of $F(\cdot)$ is $\varkappa \geq O(1)K\varkappa_{\phi}$. Besides this, in our situation the Frobenius norm of η with overwhelming as $K/n \to \infty$ probability is at most $O(1)\sqrt{nK}$. In other words, when K/n is large, "typical" recovery problem from the ensemble we have just described satisfies the premise of Proposition 4.35 with $\varkappa = O(1)K\varkappa_{\phi}$ and $\sigma^2 = O(\lambda^2 nK)$. As a result, (4.166) becomes the bound

$$\mathbf{E}_{(\eta,y)\sim P_x}\{\|\widehat{x}(y) - x\|_2^2\} \le O(1)\frac{\lambda^2 n}{\varkappa_{\phi}^2 K}.$$
 [K >> n]

It is well known that in the standard case of linear regression, where $\phi(x) = \varkappa_{\phi} x$, the resulting bound is near-optimal, provided \mathcal{X} is large enough.

Numerical illustration to follow deals with the situation described in the Example above, where we set m = 1, n = 100 and use

$$\phi(u) = \arctan[u] := [\arctan(u_1); \dots; \arctan(u_K)] : \mathbf{R}^K \to \mathbf{R}^K.$$

The set \mathcal{X} is just the unit ball $\{x \in \mathbf{R}^n : ||x||_2 \leq 1\}$; in a particular experiment, η was chosen at random from the Gaussian ensemble as described above, and the signal $x \in \mathcal{X}$ underlying observation (4.164) was drawn at random. The observation noise $y - \phi(\eta^T x)$ was $\mathcal{N}(0, \lambda^2 I_K)$. The typical simulation results are presented in Table 4.2 and on Figure 4.3.

346

LECTURE 4



Figure 4.3: Performance of SAA recovery, m = 100, m = 1

4.8 APPENDIX: CALCULUS OF ELLITOPES/SPECTRATOPES

We present here the rules of the calculus of ellitopes/spectratopes. We formulate these rules for ellitopes; the "spectratopic versions" of the rules are straightforward modifications of the "ellitopic versions."

• Intersection $\mathcal{X} = \bigcap_{i=1}^{I} \mathcal{X}_i$ of ellitopes

$$\mathcal{X}_i = \{ x \in \mathbf{R}^n : \exists (y^i \in \mathbf{R}^{n_i}, t^i \in \mathcal{T}_i) : x = P_i y^i \& [y^i]^T R_{ik} y^i \le t^i_k, 1 \le k \le K_i \}$$

is an ellitope. Indeed, this is evident when $\mathcal{X} = \{0\}$. Assuming $\mathcal{X} \neq \{0\}$, we have

$$\begin{array}{lll} \mathcal{X} &=& \{x \in \mathbf{R}^{n} : \exists (y = [y^{1}; ...; y^{I}] \in \mathcal{Y}, t = (t^{1}, ..., t^{I}) \in \mathcal{T} = \mathcal{T}_{1} \times ... \times \mathcal{T}_{I}) : \\ && x = Py := P_{1}y^{1} \& \underbrace{[y^{i}]^{T}R_{ik}y^{i}}_{y^{T}R_{ik}^{+}y} \leq t_{k}^{i}, 1 \leq k \leq K_{i}, 1 \leq i \leq I \}, \\ \mathcal{Y} &=& \{[y^{1}; ...; y^{I}] \in \mathbf{R}^{n_{1} + ... + n_{I}} : P_{i}y^{i} = P_{1}y^{1}, 2 \leq i \leq I \} \end{array}$$

(note that \mathcal{Y} can be identified with $\mathbf{R}^{\bar{n}}$ with a properly selected $\bar{n} > 0$).

• Direct product $\mathcal{X} = \prod_{i=1}^{I} \mathcal{X}_i$ of ellitopes

$$\begin{aligned} \mathcal{X}_i &= \{ x^i \in \mathbf{R}^{n_i} : \exists (y^i \in \mathbf{R}^{\bar{n}_i}, t^i \in \mathcal{T}_i) : \\ x^i &= P_i y^i, \, 1 \leq i \leq I \, \& \, [y^i]^T R_{ik} y^i \leq t^i_k, 1 \leq k \leq K_i \} \end{aligned}$$

is an ellitope:

$$\mathcal{X} = \{ [x^{1}; ...; x^{I}] \in \mathbf{R}^{n_{1}} \times ... \times \mathbf{R}^{n_{I}} : \exists \begin{pmatrix} y = [y^{1}; ...; y^{I}] \in \mathbf{R}^{\bar{n}_{1} + ...\bar{n}_{I}} \\ t = (t^{1}, ..., t^{I}) \in \mathcal{T} = \mathcal{T}_{1} \times ... \times \mathcal{T}_{I} \end{pmatrix}) \\ x = Py := [P_{1}y^{1}; ...; P_{I}y^{I}], \underbrace{[y^{i}]^{T}R_{ik}y^{i}}_{y^{T}S^{i}_{ik}y} \leq t^{i}_{k}, 1 \leq k \leq K_{i}, 1 \leq i \leq I \}$$

• The linear image $\mathcal{Z} = \{Rx : x \in \mathcal{X}\}, R \in \mathbf{R}^{p \times n}$, of an ellitope $\mathcal{X} = \{x \in \mathbf{R}^n : \exists (y \in \mathbf{R}^{\bar{n}}, t \in \mathcal{T}) : x = P_y \& y^T R_k y \leq t_k, 1 \leq k \leq K\}$ is an ellitope:

$$\mathcal{Z} = \{ z \in \mathbf{R}^p : \exists (y \in \mathbf{R}^{\bar{n}}, t \in \mathcal{T}) : z = [RP]y \& y^T R_k y \le t_k, 1 \le k \le K \}.$$

• The inverse linear image $\mathcal{Z} = \{z \in \mathbf{R}^q : Rz \in \mathcal{X}\}, R \in \mathbf{R}^{n \times q}$, of an ellitope $\mathcal{X} = \{x \in \mathbf{R}^n : \exists (y \in \mathbf{R}^{\bar{n}}, t \in \mathcal{T}) : x = Py \& y^T R_k y \leq t_k, 1 \leq k \leq K\}$ under linear mapping $z \mapsto Rz : \mathbf{R}^q \to \mathbf{R}^n$ is an ellitope, provided that the mapping is an embedding: Ker $R = \{0\}$. Indeed, setting $E = \{y \in \mathbf{R}^{\bar{n}} : Py \in \mathrm{Im}R\}$, we get a linear subspace in $\mathbf{R}^{\bar{n}}$; if $E = \{0\}, \mathcal{Z} = \{0\}$ is a spectratope; if $E \neq \{0\}$, we have

$$\mathcal{Z} = \{ z \in \mathbf{R}^q : \exists (y \in E, t \in \mathcal{T}) : z = \bar{P}y \& y^T R_k y \leq t_k, 1 \leq k \leq K \}, \bar{P} : \bar{P}y = \Pi R, \text{ where } \Pi : \operatorname{Im} R \to \mathbf{R}^q \text{ is the inverse of } z \mapsto Rz : \mathbf{R}^q \to \operatorname{Im} R$$

(*E* can be identified with some \mathbf{R}^k , and Π is well defined since *R* is an embedding).

• The arithmetic sum $\mathcal{X} = \{x = \sum_{i=1}^{I} x^i : x^i \in \mathcal{X}_i, 1 \leq i \leq I\}$, of ellitopes \mathcal{X}_i is an ellitope, with representation readily given by those of $\mathcal{X}_1, ..., \mathcal{X}_I$.

Indeed, \mathcal{X} is the image of $\mathcal{X}_1 \times ... \times \mathcal{X}_I$ under the linear mapping $[x^1; ...; x^I] \mapsto x^1 + ... + x^I$, and taking direct products and images under linear mappings preserve ellitopes.

• "S-product." Let $\mathcal{X}_i = \{x^i \in \mathbf{R}^{n_i} : \exists (y^i \in \mathbf{R}^{\bar{n}_i}, t^i \in \mathcal{T}_i) : x^i = P_i y^i, 1 \leq i \leq I \& [y^i]^T R_{ik} y^i \leq t^i_k, 1 \leq k \leq K_i \}$ be ellitopes, and let \mathcal{S} be a convex compact set in \mathbf{R}^I_+ which intersects the interior of \mathbf{R}^I_+ and is monotone: $0 \leq s' \leq s \in \mathcal{S}$ implies $s' \in \mathcal{S}$. We associate with \mathcal{S} the set

$$\mathcal{S}^{1/2} = \left\{ s \in \mathbf{R}_{+}^{I} : [s_{1}^{2}; ...; s_{I}^{2}] \in \mathcal{S} \right\}$$

of entrywise square roots of points from S; clearly, $S^{1/2}$ is a convex compact set. \mathcal{X}_i and S specify the *S*-product of the sets \mathcal{X}_i , $i \leq I$, defined as the set

$$\mathcal{Z} = \left\{ z = [z^1; ...; z^I] : \exists (s \in \mathcal{S}^{1/2}, x^i \in \mathcal{X}_i, i \le I) : z^i = s_i x^i, 1 \le i \le I \right\},\$$

or, equivalently,

$$\begin{aligned} \mathcal{Z} &= \left\{ z = [z^1; ...; z^I] : \exists (r = [r^1; ...; r^I] \in \mathcal{R}, y^1, ..., y^I) : \\ z_i = P_i y_i \, \forall i \leq I, [y^i]^T R_{ik} y^i \leq r_k^i \, \forall (i \leq I, k \leq K_i) \right\}, \\ \mathcal{R} &= \{ [r^1; ...; r^I] \geq 0 : \exists (s \in \mathcal{S}^{1/2}, t^i \in \mathcal{T}_i) : r^i = s_i^2 t^i \, \forall i \leq I \}. \end{aligned}$$

We claim that \mathcal{Z} is an ellitope. All we need to verify to this end is that the set \mathcal{R} is as it should be in an ellitopic representation, that is, that \mathcal{R} is compact and monotone subset of $\mathbf{R}_{+}^{\bar{n}_{1}+\ldots+\bar{n}_{I}}$ containing a strictly positive vector (all this is evident), and that \mathcal{R} is convex. To verify convexity, let $\mathbf{T}_{i} = \operatorname{cl}\{[t^{i}; \tau_{i}] : \tau_{i} > 0, t^{i}/\tau_{i} \in \mathcal{T}_{i}\}$ be the conic hulls of \mathcal{T}_{i} 's. We clearly have

$$\begin{aligned} \mathcal{R} &= \{ [r^1; ...; r^I] : \exists s \in \mathcal{S}^{1/2} : [r^i; s_i^2] \in \mathbf{T}_i, \, i \leq I \} \\ &= \{ [r^1; ...; r^I] : \exists \sigma \in \mathcal{S} : [r^i; \sigma_i] \in \mathbf{T}_i, \, i \leq I \}, \end{aligned}$$

where the concluding equality is due to the origin of $S^{1/2}$. The concluding set in the above chain clearly is convex, and we are done.

As an example, consider the situation where the ellitopes \mathcal{X}_i posses nonempty interiors and thus can be thought of as the unit balls of norms $\|\cdot\|_{(i)}$ on the respective spaces $\mathbf{R}^{\bar{n}_i}$, and let $\mathcal{S} = \{s \in \mathbf{R}_+^I : \|s\|_{p/2} \leq 1\}$, where $p \geq 2$. In this situation, $\mathcal{S}^{1/2} = \{s \in \mathbf{R}^I U_+ : \|s\|_p \leq 1\}$, whence \mathcal{Z} is the unit ball of the "block *p*-norm"

$$\|[z^1;...;z^I]\| = \| \left[\|z^1\|_{(1)};...;\|z^I\|_{(I)} \right] \|_p$$

Note also that the usual direct product of I ellitopes is their S-product, with $S = [0, 1]^I$.

• "S-weighted sum." Let $\mathcal{X}_i \subset \mathbf{R}^n$ be ellitopes, $1 \leq i \leq I$, and let $\mathcal{S} \subset \mathbf{R}^I + \mathcal{S}^{1/2}$ be the same as in the previous item. Then the *S*-weighted sum of the sets \mathcal{X}_i , defined as

$$\mathcal{X} = \{ x : \exists (s \in \mathcal{S}^{1/2}, x^i \in \mathcal{X}_i, i \le I) : x = \sum_i s_i x^i \}$$

is an ellitope. Indeed, the set in question is the image of the S-product of \mathcal{X}_i under the linear mappings $[z^1; ...; z^I] \mapsto z^1 + ... + z^I$, and taking S-products and

349

linear images preserves the property to be an ellitope.

It should be stressed that the outlined "calculus rules" are fully algorithmic: representation (4.6) of the result of an operation is readily given by the representations (4.6) of the operands.

4.9 EXERCISES FOR LECTURE 4

[†] marks more difficult exercises.

4.9.1 Linear Estimates vs. Maximum Likelihood

 $Exercise\ 4.36.$. Consider the problem posed in the beginning of Lecture 4: Given observation

$$\omega = Ax + \sigma\xi, \ \xi \sim \mathcal{N}(0, I)$$

of unknown signal x known to belong to a given signal set $\mathcal{X} \subset \mathbf{R}^n$, we want to recover Bx.

Let us restrict ourselves with the case where A is square and invertible matrix, B is the identity, and \mathcal{X} is a computationally tractable convex compact set. As far as computational aspects are concerned, the situation is well suited for utilizing the "magic wand" of Statistics – the *Maximum Likelihood* (ML) estimate where the recovery of x is

$$\widehat{r}_{\mathrm{ML}}(\omega) = \underset{y \in \mathcal{X}}{\operatorname{argmin}} \|\omega - Ay\|_2$$
(ML)

– the signal which maximizes, over $y \in \mathcal{X}$, the likelihood (the probability density) to get the observation we actually got. Indeed, with computationally tractable \mathcal{X} , (ML) is an explicit convex, and therefore efficiently solvable, optimization problem. Given the exclusive role played by ML estimate in Statistics, perhaps the first question about our problem of interest is: how good in the situation in question is the ML estimate?

The goal of what follows is to demonstrate that in the situation we are interested in, the ML estimate can be "heavily nonoptimal," and this may happen even when the techniques we develop in Lecture 4 do result in efficiently computable nearoptimal linear estimate.

To justify the claim, investigate the risk (4.2) of the ML estimate in the case when

$$\mathcal{X} = \{ x \in \mathbf{R}^n : x_1^2 + \epsilon^{-2} \sum_{i=2}^n x_i^2 \le 1 \} \ \& \ A = \mathrm{Diag}\{1, 1/\epsilon, ..., 1/\epsilon\},$$

 ϵ and σ are small, and n is large, specifically, $\sigma^2(n-1) \geq 2$. Accompany your theoretical analysis by numerical experiments – compare the empirical risks of the ML estimate with theoretical and empirical risks of the optimal under the circumstances linear estimate.

Recommended setup: *n* runs through $\{256, 1024, 2048\}$ and $\epsilon = \sigma$ run through $\{0.01; 0.05; 0.1\}$, and signal *x* is generated as

$$x = [\cos(\phi); \sin(\phi)\epsilon\zeta],$$

where $\phi \sim \text{Uniform}[0, 2\pi]$ and random vector ζ is independent of ϕ and is distributed uniformly on the unit sphere in \mathbf{R}^{n-1} .

4.9.2 Measurement Design in Signal Recovery

Exercise 4.37. [Measurement Design in Gaussian o.s.] As a preamble to the Exercise, please read the story about possible "physics" of Gaussian o.s. from Section 2.7.3.3. The summary of this story is as follows:

We consider the Measurement Design version of signal recovery in Gaussian o.s., specifically, we are allowed to use observations

$$\omega = A_q x + \sigma \xi \qquad \qquad [\xi \sim \mathcal{N}(0, I_m)]$$

where

$$A_q = \text{Diag}\{\sqrt{q_1}, \sqrt{q_2}, \dots, \sqrt{q_m}\}A,\$$

with a given $A \in \mathbf{R}^{m \times n}$ and vector q which we can select in a given convex compact set $\mathcal{Q} \subset \mathbf{R}^{m}_{+}$. The signal x underlying the observation is known to belong to a given ellitope \mathcal{X} . Your goal is to select $q \in \mathcal{Q}$ and a linear recovery $\omega \mapsto G^{T}\omega$ of the image Bx of $x \in \mathcal{X}$, with B given, resulting in the minimal worst-case, over $x \in \mathcal{X}$, expected $\|\cdot\|_{2}^{2}$ recovery risk. Modify, according to this goal, problem (4.13). Is it possible to end up with a tractable problem? Work out in full details the case when $\mathcal{Q} = \{q \in \mathbf{R}^{m}_{+} : \sum_{i} q_{i} = m\}$.

Exercise 4.38. [follow-up to Exercise 4.37] A translucent bar of length n = 32 is comprised of 32 consecutive segments of length 1 each, with density ρ_i of *i*-th segment known to belong to the interval $[\mu - \delta_i, \mu + \delta_i]$.

Sample translucent bar

The bar is lightened from the left end; when light passes through a segment with density ρ , light's intensity is reduced by factor $e^{-\alpha\rho}$. The intensity of light at the left endpoint of the bar is 1. You can scan the segments one by one from left to right and measure light intensity ℓ_i at the right endpoint of *i*-th segment for time q_i ; the result z_i of the measurement is $\ell_i e^{\sigma\xi_i/\sqrt{q_i}}$, where $\xi_i \sim \mathcal{N}(0, 1)$ are independent across *i*. The total time budget is *n*, and you are interested to recover the m = n/2-dimensional vector of densities of the right *m* segments. Build optimization problem responsible for near-optimal linear recovery with and without Measurement Design (with no Measurement Design, each segment is observed during unit time) and compare the resulting near-optimal risks.

Recommended data:

$$\alpha = 0.01, \, \delta_i = 1.2 + \cos(4\pi (i-1)/n), \, \mu = 1.1 \max \delta_i, \, \sigma = 0.001.$$

Exercise 4.39. Let $X \subset \mathbf{R}^n$ be a convex compact set, let $b \in \mathbf{R}^n$, and let A be an $m \times n$ matrix. Consider the problem of affine recovery $\omega \mapsto h^T \omega + c$ of the linear function $Bx = b^T x$ of $x \in X$ from indirect observation

$$\omega = Ax + \sigma\xi, \, \xi \sim \mathcal{N}(0, I_m).$$

Given tolerance $\epsilon \in (0, 1)$, we are interested to minimize the worst-case, over $x \in X$,

width of $(1 - \epsilon)$ confidence interval, that is, the smallest ρ such that

$$\operatorname{Prob}\{\xi: b^T x - f^T (Ax + \sigma\xi) > \rho\} \le \epsilon/2 \& \operatorname{Prob}\{\xi: b^T x - f^T (Ax + \sigma\xi) < \rho\} \le \epsilon/2 \ \forall x \in X.$$

Pose the problem as a convex optimization problem and consider in details the case where X is the box $\{x \in \mathbf{R}^n : a_j | x_j | \le 1, 1 \le j \le n\}$, where $a_j > 0$ for all j.

Exercise 4.40. Let X be an ellitope in \mathbb{R}^n :

$$X = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : x^T S_k x \le t_k, \ 1 \le k \le K \}$$

with our usual restrictions on S_k and \mathcal{T} . Let, further, m be a given positive integer, and $x \mapsto Bx : \mathbf{R}^n \to \mathbf{R}^{\nu}$ be a given linear mapping. Consider the Measurement Design problem where you are looking for a linear recovery $\omega \mapsto \hat{x}_H(\omega) := H^T \omega$ of $Bx, x \in X$, from observation

$$\omega = Ax + \sigma\xi \qquad \qquad [\sigma > 0 \text{ is given and } \xi \sim \mathcal{N}(0, I_m)]$$

in which the $m \times n$ sensing matrix A is under your control – it is allowed to be a whatever $m \times n$ matrix of spectral norm not exceeding 1. You are interested to select H and A in order to minimize the worst case, over $x \in X$, expected $\|\cdot\|_2^2$ recovery error. Similarly to (4.13), this problem can be posed as

$$\begin{aligned}
\text{Opt} &= \min_{H,\lambda,A} \left\{ \sigma^2 \text{Tr}(H^T H) + \phi_{\mathcal{T}}(\lambda) : \\
\left[\frac{\sum_k \lambda_k S_k}{B - H^T A} \middle| \frac{B^T - A^T H}{I_{\nu}} \right] \succeq 0, \, \|A\| \le 1, \, \lambda \ge 0 \right\},
\end{aligned}$$
(4.167)

where $\|\cdot\|$ stands for the spectral norm. The objective in this problem is the (upper bound on the) squared risk $\operatorname{Risk}^2[\hat{x}_H|X]$, the sensing matrix being A. The problem is nonconvex, since the matrix participating in the semidefinite constraint is bilinear in H and A.

A natural way to handle an optimization problem with bilinear in the decision variables u, v objective and/or constraints is to use "alternating minimization," where one alternates optimization in v for u fixed and optimization in u for v fixed, where the value of the variable fixed in a round is the result of optimization w.r.t. this variable in the previous round. Alternating minimizations are carried out until the value of the objective (which in the outlined process definitely improves from round to round) stops to improve (or nearly so). Since the algorithm not necessarily converges to the globally optimal solution to the problem of interest, it makes sense to run the algorithm several times from different, say, randomly selected, starting points.

Now goes the Exercise.

1. Implement Alternating Minimization as applied to (4.167) and look how it works. You could restrict your experimentation to the case where the sizes m, n, ν are quite moderate, in the range of tens, and X is either the box $\{x : j^{2\gamma}x_j^2 \leq 1, 1 \leq j \leq n\}$, or the ellipsoid $\{x : \sum_{j=1}^n j^{2\gamma}x_j^2 \leq 1\}$, where γ is a nonnegative parameter (you could try $\gamma = 0, 1, 2, 3$). As about B, you could generate it at random, or enforce B to have prescribed singular values, say, $\sigma_j = j^{-\theta}$, $1 \leq j \leq \nu$, and randomly selected system of singular vectors.

2. Identify cases where a globally optimal solution to (4.167) is easy to identify and use this in order to understand how reliable is Alternating minimization in the application in question, reliability meaning the ability to identify near-optimal, in terms of the objective, solutions.

If you are not satisfied with Alternating Minimization "as it is," try to improve it.

- 3. Modify (4.167) and your experimentation to cover the cases where the restriction $||A|| \leq 1$ on the sensing matrix is replaced with one of the following restrictions:
 - $\|\operatorname{Row}_{i}[A]\|_{2} \leq 1, 1 \leq i \leq m$
 - $|A_{ij}| \le 1$ for all i, j

(note that these two types of restrictions mimic what happens if you are interested to recover (linear image of) the vector of parameters in a linear regression model from noisy observations of model's outputs at m points which you are allowed to select in the unit ball, resp., unit box).

4. [Embedded Exercise] Recall that a $\nu \times n$ matrix G admits singular value decomposition $G = UDV^T$ with orthogonal matrices $U \in \mathbf{R}^{\nu \times \nu}$ and $V \in \mathbf{R}^{n \times n}$ and diagonal $\nu \times n$ matrix D with nonnegative and nonincreasing diagonal entries ⁷¹. These entries are uniquely defined by G and are called singular values $\sigma_i(G), 1 \leq i \leq \min[\nu, n]$. These singular values admit characterization similar to variational characterization of eigenvalues of a symmetric matrix, see, e.g., [11, Section A.7.3]:

Theorem 4.41. [VCSV - Variational Characterization of Singular Values] For $\nu \times n$ matrix G it holds

$$\sigma_i(G) = \min_{E \in \mathcal{E}_i} \max_{e \in E, \|e\|_2 = 1} \|GE\|_2, \ 1 \le i \le \min[\nu, n], \tag{4.168}$$

where \mathcal{E}_i is the family of all subspaces in \mathbf{R}^n of codimension i-1.

Corollary 4.42. [SVI - Singular Value Interlacement] Let G and G' be $\nu \times n$ matrices, and let k = Rank(G'G'). Then

$$\sigma_i(G) \ge \sigma_{i+k}(G'), \ 1 \le i \le \min[\nu, n],$$

where, by definition, singular values of a $\nu \times n$ matrix with indexes $> \min[\nu, n]$ are zeros.

We denote by $\sigma(G)$ the vector of singular values of G arranged in the nonincreasing order. The function $||G||_{\operatorname{Sh},p} = ||\sigma(G)||_p$ is called *Shatten p-norm* of matrix G; this indeed is a norm on the space of $\nu \times n$ matrices, and the conjugate norm is $||\cdot||_{\operatorname{Sh},q}$, with $\frac{1}{p} + \frac{1}{q} = 1$. An easy and important consequence of Corollary 4.42 is the following fact:

Corollary 4.43. Given a $\nu \times n$ matrix G, an integer k, $0 \le k \le \min[\nu, n]$, and $p \in [1, \infty]$, (one of) the best approximations of G in the Shatten p-norm among matrices of rank $\le k$ is obtained from G by zeroing our all but k largest singular values, that is, the matrix $G^k = \sum_{i=1}^k \sigma_i(G) \operatorname{Col}_i[U] \operatorname{Col}_i^T[V]$, where $G = UDV^T$

⁷¹By definition, diagonality of a rectangular matrix D means that all entries D_{ij} in D with $i \neq j$ are zeros.

is the singular value decomposition of G.

Now goes the Embedded Exercise:

Prove Theorem 4.41 and Corollaries 4.42 and 4.43.

5. Consider the Measurement Design problem (4.167) in the case when X is an ellipsoid:

$$X = \{x \in \mathbf{R}^n : \sum_{j=1}^n x_i^2 / a_j^2 \le 1\}$$

A is restricted to be $m \times n$ matrix of spectral norm not exceeding 1, and there is no noise in observations: $\sigma = 0$, and find an optimal solution to this problem. Think how this result can be used to get a hopefully good starting point for Alternating Minimization in the case when X is an ellipsoid and σ is small.

Exercise 4.44. Prove Proposition 4.21.

Exercise 4.45. Prove Proposition 4.22.

4.9.3 Around semidefinite relaxation

Exercise 4.46. Let \mathcal{X} be an ellitope:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists (y \in \mathbf{R}^N, t \in \mathcal{T}) : x = Py, y^T S_k y \le t_k, k \le K \}$$

with our standard restrictions on \mathcal{T} and S_k . Representing $S_k = \sum_{j=1}^{r_k} s_{kj} s_{kj}^T$, we can pass from initial ellitopic representation of \mathcal{X} to the spectratopic representation of the same set:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists (y \in \mathbf{R}^N, t^+ \in \mathcal{T}^+) : x = Py, [s_{kj}^T x]^2 \preceq t^+ kjI_1, 1 \le k \le K, 1 \le j \le r_k \} \\ \left[\mathcal{T}^+ = \{ t^+ = \{ t_{kj}^+ \ge 0 \} : \exists t \in \mathcal{T} : \sum_{j=1}^{r_k} t_{kj}^+ \le t_k, 1 \le k \le K \} \right]$$

If now C is a symmetric $n \times n$ matrix and $Opt = \max_{x \in \mathcal{X}} x^T C x$, we have

$$\begin{array}{lll} \operatorname{Opt}_{*} & \leq & \operatorname{Opt}_{e} := \min_{\lambda = \{\lambda_{k} \in \mathbf{R}_{+}\}} \left\{ \phi_{\mathcal{T}}(\lambda) : P^{T}CP \preceq \sum_{k} \lambda_{k}S_{k} \right\} \\ \operatorname{Opt}_{*} & \leq & \operatorname{Opt}_{s} := \min_{\Lambda = \{\Lambda_{kj} \in \mathbf{R}_{+}\}} \left\{ \phi_{\mathcal{T}^{+}}(\Lambda) : P^{T}CP \preceq \sum_{k,j} \Lambda_{kj}s_{kj}s_{kj}^{T} \right\} \end{array}$$

where the first relation is yielded by ellitopic representation of \mathcal{X} and Proposition 4.6, and the second, on a closest inspection (carry this inspection out!) – by the spectratopic representation of \mathcal{X} and Proposition 4.8. Now goes Exercise: Prove that $Opt_e = Opt_s$.

Exercise 4.47. [estimating Kolmogorov widths of sperctratopes/ellitopes]

4.47.A Preliminaries: Kolmogorov and Gelfand widths. Let \mathcal{X} be a convex compact set in \mathbb{R}^n , and let $\|\cdot\|$ be a norm on \mathbb{R}^n . Given a linear subspace E in \mathbb{R}^n , let

$$\operatorname{dist}_{\|\cdot\|}(x,E) = \min_{z \in E} \|x - z\| : \mathbf{R}^n \to \mathbf{R}_+$$

be the $\|\cdot\|$ -distance from x to E. The quantity

$$\operatorname{dist}_{\|\cdot\|}(\mathcal{X}, E) = \max_{x \in \mathcal{X}} \operatorname{dist}_{\|\cdot\|}(x, E)$$

can be viewed as the worst-case $\|\cdot\|$ -accuracy to which vectors from \mathcal{X} can be approximated by vectors from E. Given positive integer $m \leq n$ and denoting by \mathcal{E}_m the family of all linear subspaces in \mathbf{R}^m of dimension m, the quantity

$$\delta_m(\mathcal{X}, \|\cdot\|) = \min_{E \in \mathcal{E}_m} \operatorname{dist}_{\|\cdot\|}(\mathcal{X}, E)$$

can be viewed as the best achievable quality of approximation, in $\|\cdot\|$, of vectors from \mathcal{X} by vectors from an *m*-dimensional linear subspace of \mathbf{R}^n ; this quantity is called *m*-th Kolmogorov width of \mathcal{X} taken w.r.t. $\|\cdot\|$.

Observe that one has

where E^{\perp} is the orthogonal complement to E.

1. Prove (!).

<u>Hint:</u> Represent dist_{$\|\cdot\|$}(x, E) as the optimal value in a conic problem on the cone $\mathbf{K} = \{[x;t] : t \ge ||x||\}$ and use Conic Duality Theorem.

Now consider the case when \mathcal{X} is the unit ball of some norm $\|\cdot\|_{\mathcal{X}}$. In this case (!) combines with the definition of Kolmogorov width to imply that

$$\delta_m(\mathcal{X}, \|\cdot\|) = \min_{E \in \mathcal{E}_m} \operatorname{dist}_{\|\cdot\|}(x, E) = \min_{E \in \mathcal{E}_m} \max_{x \in \mathcal{X}} \max_{y \in E^{\perp}, \|y\|_* \le 1} y^T x$$

$$= \min_{E \in \mathcal{E}_m} \operatorname{dist}_{\|\cdot\|}(x, E) = \min_{E \in \mathcal{E}_m} \max_{y \in E^{\perp}, \|y\|_* \le 1} \max_{x: \|x\|_{\mathcal{X}} \le 1} y^T x \quad (4.169)$$

$$= \min_{F \in \mathcal{E}_{n-m}} \max_{y \in F, \|y\|_* \le 1} \|y\|_{\mathcal{X},*},$$

where $\|\cdot\|_{\mathcal{X},*}$ is the norm conjugate to $\|\cdot\|_{\mathcal{X}}$. Note that when \mathcal{Y} is a convex compact set in \mathbf{R}^n and $|\cdot|$ is a norm on \mathbf{R}^n , the quantity

$$d^{m}(\mathcal{Y}, |\cdot|) = \min_{F \in \mathcal{E}_{n-m}} \max_{y \in \mathcal{Y} \cap F} |y|$$

has a name – it is called the *m*-th *Gelfand width* of \mathcal{Y} taken w.r.t. $|\cdot|$. "Duality relation" (4.169) states that

When \mathcal{X}, \mathcal{Y} are the unit balls of the respective norms $\|\cdot\|_{\mathcal{X}}, \|\cdot\|_{\mathcal{Y}}$, for every m < n m-th Kolmogorov width of \mathcal{X} taken w.r.t. $\|\cdot\|_{\mathcal{Y},*}$ is the same as m-th Gelfand width of \mathcal{Y} taken w.r.t. $\|\cdot\|_{\mathcal{X},*}$.

The goal of the remaining part of Exercise is to use our results on the quality of semidefinite relaxation on ellitopes/spectratopes to infer efficiently computable upper bounds on Kolmogorov widths of a given set $\mathcal{X} \subset \mathbf{R}^n$. In the sequel we assume that

• \mathcal{X} is a spectratope:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists (t \in \mathcal{T}, u) : x = Pu, R_k^2[u] \preceq t_k I_{d_k}, k \leq K \};$$

• The unit ball \mathcal{B}_* of the norm conjugate to $\|\cdot\|$ is a spectratope:

$$\mathcal{B}_* = \{y : \|y\|_* \le 1\} = \{y \in \mathbf{R}^n : \exists (r \in \mathcal{R}, z) : y = Mz, S_\ell^2[z] \le r_\ell I_{f_\ell}, \ell \le L\}.$$

with our usual restrictions on \mathcal{T}, \mathcal{R} and $R_k[\cdot], S_\ell[\cdot]$.

4.47.B Simple case: $\|\cdot\| = \|\cdot\|_2$. We start with the *simple case* where $\|\cdot\| = \|\cdot\|_2$, so that \mathcal{B}_* is the ellitope $\{y : y^T y \leq 1\}$.

Let $D = \sum_{k} d_{k}$ be the size of the spectratope \mathcal{X} , and let

 $\varkappa = 2 \max[\ln(2D), 1].$

Given integer m < n, consider convex optimization problem

$$Opt(m) = \min_{\Lambda = \{\Lambda_k, k \le K\}, Y} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) : \\ \Lambda_k \succeq 0 \forall k, \sum_k \mathcal{S}_k^*[\Lambda_k] \succeq P^T Y P, 0 \le Y \le I_n, Tr(Y) = n - m \right\}_{(P_m)}$$

2. Prove the following

Proposition 4.48. Whenever $1 \le \mu \le m < n$, one has

$$Opt(m) \le \varkappa \delta_m^2(\mathcal{X}, \|\cdot\|_2) \& \delta_m^2(\mathcal{X}, \|\cdot\|_2) \le \frac{m+1}{m+1-\mu} Opt(\mu).$$
(4.170)

Moreover, the above upper bounds on $\delta_m(\mathcal{X}, \|\cdot\|_2)$ are "constructive", meaning that an optimal solution to $(P_{\mu}), \mu \leq m$, can be straightforwardly converted into a linear subspace $E^{m,\mu}$ of dimension m such that

$$\operatorname{dist}_{\|\cdot\|_2}(\mathcal{X}, E^{m,\mu}) \le \sqrt{\frac{m+1}{m+1-\mu}} \operatorname{Opt}(\mu).$$

Finally, $Opt(\mu)$ is nonincreasing in μ .

4.47.C General case. Now consider the case when both \mathcal{X} and the unit ball \mathcal{B}_* of the norm conjugate to $\|\cdot\|$ are spectratopes. As you are about to see, this case is essentially more difficult than the case of $\|\cdot\| = \|\cdot\|_2$, but something still can be done.

3. Prove the following statement:

(!!) Given m < n, let Y be an orthoprojector of \mathbb{R}^n of rank n - m, and let collections $\Lambda = \{\Lambda_k \succeq 0, k \leq K\}$ and $\Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\}$ satisfy the relation

$$\left[\frac{\sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}] \mid \frac{1}{2} P^{T} Y M}{\left| \frac{1}{2} M^{T} Y P \mid \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \right|} \right] \succeq 0.$$
(4.171)

Then

$$\operatorname{dist}_{\|\cdot\|}(\mathcal{X}, \operatorname{Ker} Y) \le \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]).$$

$$(4.172)$$

As a result,

$$\delta_{m}(\mathcal{X}, \|\cdot\|) \leq \operatorname{dist}_{\|\cdot\|}(\mathcal{X}, \operatorname{Ker} Y)$$

$$\leq \operatorname{Opt} := \min_{\Lambda = \{\Lambda_{k}, k \leq K\}, \Upsilon = \{\Upsilon_{\ell}, \ell \leq L\}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) :$$

$$\left\{ \begin{array}{c} \Lambda_{k} \succeq 0 \,\forall k, \,\Upsilon_{\ell} \succeq 0 \,\forall \ell, \\ \left[\frac{\sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}]}{\frac{1}{2}M^{T}YP} \mid \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]} \right] \succeq 0 \end{array} \right\}$$

$$(4.173)$$

4. Prove the following statement:

(!!!) Let m, n, Y be as in (!!). Then

$$\delta_{m}(\mathcal{X}, \|\cdot\|) \leq \operatorname{dist}_{\|\cdot\|}(\mathcal{X}, \operatorname{Ker} Y) \\
\leq \widehat{\operatorname{Opt}} := \min_{\substack{\nu, \Lambda = \{\Lambda_{k}, k \leq K\}, \Upsilon \in \{\Upsilon_{\ell}, \ell \leq L\} \\ \nu \geq 0, \Lambda_{k} \succeq 0 \forall k, \Upsilon_{\ell} \succeq 0 \forall \ell, \\ \left[\frac{\nu \geq 0, \Lambda_{k} \succeq 0 \forall k, \Upsilon_{\ell} \succeq 0 \forall \ell, \\ \left[\frac{\sum_{k} \mathcal{R}_{k}^{*} [\Lambda_{k}]}{\frac{1}{2} M^{T} P} \right] \sum_{\ell} \mathcal{S}_{\ell}^{*} [\Upsilon_{\ell}] + \nu M^{T} (I - Y) M \end{bmatrix}} \geq 0 \right\} \cdot (4.174)$$

and $\widehat{\operatorname{Opt}} \leq \operatorname{Opt}$, with Opt given by (4.173).

Statements (!!), (!!!) suggest the following policy for upper-bounding the Kolmogorov width $\delta_m(\mathcal{X}, \|\cdot\|)$:

A. First, we select an integer μ , $1 \leq \mu < n$, and solve the convex optimization problem

$$\min_{\Lambda,\Upsilon,Y} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) : \left\{ \begin{array}{l} \Lambda = \{\Lambda_k \succeq 0, k \le K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \le L\}, \\ 0 \preceq Y \preceq I, \operatorname{Tr}(Y) = n - \mu \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k] \mid \frac{1}{2} P^T Y M}{\frac{1}{2} M^T Y P \mid \sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell]} \right] \succeq 0 \end{array} \right\}$$

$$(P^{\mu})$$

- B. Next, we take the Y-component Y^{μ} of the optimal solution to (I^{μ}) and "round" it to a orthoprojector Y of rank n m in the same fashion as in the case of $\|\cdot\| = \|\cdot\|_2$, that is, keep the eigenvectors of Y^{μ} intact and replace m smallest eigenvalues with zeros, and all remaining eigenvalues with ones.
- C. Finally, we solve the convex optimization problem

$$\begin{aligned}
\operatorname{Opt}_{m,\mu} &= \min_{\Lambda,\Upsilon,\nu} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) : \\ &\left\{ \begin{array}{l} \nu \geq 0, \Lambda = \{\Lambda_k \succeq 0, k \leq K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\}, \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T P} \frac{\frac{1}{2}P^T M}{\sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell] + \nu M^T (I - Y)M} \right] \succeq 0 \end{array} \right\} \end{aligned}$$

$$(P^{m,\mu})$$

By (!!!), $\operatorname{Opt}_{m,\mu}$ is an upper bound on the Kolmogorov width $\delta_m(\mathcal{X}, \|\cdot\|)$ (and in fact – also on dist_{$\|\cdot\|}(\mathcal{X}, \operatorname{Ker} Y)$).</sub>

Pay attention to complications incurred by passing from the simple case $\|\cdot\| = \|\cdot\|_2$ to the case of general norm $\|\cdot\|$ with spectratope as the unit ball of the conjugate norm. Indeed, Proposition 4.48 gives both a lower bound $\sqrt{\operatorname{Opt}(m)/\varkappa}$ on the *m*-th Kolmogorov width of \mathcal{X} w.r.t. $\|\cdot\|_2$, and a family of upper bounds $\sqrt{\frac{m+1}{m+1-\mu}\operatorname{Opt}(\mu)}$, $1 \leq \mu \leq m$, on this width. As a result, we can approximate \mathcal{X} by *m*-dimensional subspaces in the Euclidean norm in a "nearly optimal" fashion. Indeed, if for some

 ϵ and k it holds $\delta_k(\mathcal{X}, \|\cdot\|_2) \leq \epsilon$, then $\operatorname{Opt}(k) \leq \varkappa \epsilon^2$ by Proposition 4.48 as applied with m = k. On the other hand, assuming k < n/2, the same Proposition when applied with m = 2k and $\mu = k$ says that

$$\operatorname{dist}_{\|\cdot\|_2}(\mathcal{X}, E^{m,k}) \leq \sqrt{\frac{2k+1}{k+1}} \operatorname{Opt}(k) \leq \sqrt{2}\sqrt{\operatorname{Opt}(k)} \leq \sqrt{2\varkappa}\epsilon.$$

Thus, if "in the nature" \mathcal{X} can be approximated by k-dimensional subspace within $\|\cdot\|_2$ -accuracy ϵ , we can efficiently get approximation of "nearly the same quality" ($\sqrt{2\varkappa\epsilon}$ instead of ϵ ; recall that \varkappa is just logarithmic in D) and "nearly the same dimension" (2k instead of k).

Neither one of these options is preserved when passing from the Euclidean norm to a general one: in the latter case, we do not have neither lower bounds on Kolmogorov widths, just upper ones, nor understanding of how tight our upper bounds are.

Now – the concluding questions:

- 5. Why in step A of the above bounding policy we utilize statement (!!) rather than less conservative (since $\widehat{Opt} \leq Opt$) statement (!!!) ?
- 6. Implement the above scheme numerically and run experiments. Recommended setup:
 - Given positive integers n and κ and a real $\sigma > 0$, specify \mathcal{X} as the set of n-dimensional vectors x which can be obtained when restricting a function f of continuous argument $t \in [0, 1]$ onto n-point equidistant grid $\{t_i = i/n\}_{i=1}^n$, and impose on f the smoothness restriction that $|f^{(k)}(t)| \leq \sigma^k, 0 \leq t \leq 1$, $k = 0, 1, 2, ..., \kappa$; translate this description on f into a bunch of two-sided linear constraints on x, specifically, the constraints

$$|d_{(k)}^T[x_i; x_{i+1}; ...; x_{i+k}]| \le \sigma^k, 1 \le i \le n-k, 0 \le k \le \kappa,$$

where $d_{(k)} \in \mathbf{R}^{k+1}$ is the vector of coefficients of finite-difference approximation, with resolution 1/n, of k-th derivative:

$$\begin{aligned} d_{(0)} &= 1, \, d_{(1)} = n[-1;1], \, d_{(2)} = n^2[1;-2;1], \, d_{(3)} = n^3[-1;3;-3;1], \\ d_{(4)} &= n^4[1;-4;6;-4;1], \ldots \end{aligned}$$

- Recommended parameters: $n = 32, m = 8, \kappa = 5, \sigma \in \{0.25, 0.5; 1, 2, 4\}$.
- Run experiments with $\|\cdot\| = \|\cdot\|_1$ and $\|\cdot\| = \|\cdot\|_2$.

Exercise 4.49. Prove Proposition 4.51

Exercise 4.50. † [more on semidefinite relaxation] The goal of this Exercise is to extend SDP relaxation beyond ellitopes/spectratopes.

SDP relaxation is aimed at upper-bounding the quantity

$$\operatorname{Opt}_{\mathcal{X}}(B) = \max_{x \in \mathcal{X}} x^T B x,$$
 $[B \in \mathbf{S}^n]$

where $\mathcal{X} \subset \mathbf{R}^n$ is a given set (which we from now on assume to be nonempty convex compact). To this end we look for a computationally tractable convex compact set $\mathcal{U} \subset \mathbf{S}^n$ such that for every $x \in \mathcal{X}$ it holds $xx^T \in \mathcal{U}$; in this case, we refer to \mathcal{U} as to a set *matching* \mathcal{X} (equivalent wording: " \mathcal{U} matches \mathcal{X} "). Given such a set \mathcal{U} , the

358

LECTURE 4

optimal value in the convex optimization problem

$$\overline{\operatorname{Opt}}_{\mathcal{U}}(B) = \max_{U \in \mathcal{U}} \operatorname{Tr}(BU)$$
(4.175)

is an efficiently computable convex upper bound on $\operatorname{Opt}_{\mathcal{X}}(B)$.

Given \mathcal{U} matching \mathcal{X} , we can pass from \mathcal{U} to the conic hull of \mathcal{U} – to the set

 $\mathbf{U}[\mathcal{U}] = \mathrm{cl}\{(U,\mu) \in \mathbf{S}^n \times \mathbf{R}_+ : \mu > 0, U/\mu \in \mathcal{U}\}$

which, as it is immediately seen, is a closed convex cone contained in $\mathbf{S}^n \times \mathbf{R}_+$; the only point (U, μ) in this cone with $\mu = 0$ has U = 0 (since \mathcal{U} is compact), and

$$\mathcal{U} = \{ U : (U,1) \in \mathbf{U} \} = \{ U : \exists \mu \le 1 : (U,\mu) \in \mathbf{U} \}$$

so that the definition of $\overline{\operatorname{Opt}}_{\mathcal{U}}$ can be rewritten equivalently as

$$\overline{\operatorname{Opt}}_{\mathcal{U}}(B) = \min_{U,\mu} \left\{ \operatorname{Tr}(BU) : (U,\mu) \in \mathbf{U}, \mu \le 1 \right\}.$$

The question, of course, is where to take a set \mathcal{U} matching \mathcal{X} , and the answer depends on what we know about \mathcal{X} . For example, when \mathcal{X} is a basic ellitope:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : x^T S_k x \le t_k, k \le K \}$$

with our usual restrictions on \mathcal{T} and S_k , it is immediately seen that

$$x \in \mathcal{X} \Rightarrow xx^T \in \mathcal{U} = \{ U \in \mathbf{S}^n : U \succeq 0, \exists t \in \mathcal{T} : \operatorname{Tr}(US_k) \le t_k, k \le K \}.$$

Similarly, when \mathcal{X} is a basic spectratope:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : S_k^2[x] \preceq t_k I_{d_k}, k \leq K \}$$

with our usual restrictions on \mathcal{T} and $S_k[\cdot]$, it is immediately seen that

$$x \in \mathcal{X} \Rightarrow xx^T \in \mathcal{U} = \{U \in \mathbf{S}^n : U \succeq 0, \exists t \in \mathcal{T} : \mathcal{S}_k[U] \preceq t_k I_{d_k}, k \leq K\}$$

One can verify that the semidefinite relaxation bounds on the maximum of a quadratic form on an ellitope/spectratope \mathcal{X} derived in Sections 4.2.5 (for ellitopes) and 4.3.2 (for spectratopes) are nothing but the bounds (4.175) associated with the just defined \mathcal{U} .

4.50.A Matching via absolute norms. There are other ways to specify a set matching \mathcal{X} . The seemingly simplest of them is as follows. Let $p(\cdot)$ be an absolute norm on \mathbb{R}^n (recall that it is a norm p(x) which depends solely on abs[x], where abs[x] is the vector comprised of the magnitudes of entries in x). We can convert $p(\cdot)$ into the norm $p^+(\cdot)$ on the space \mathbb{S}^n , namely, a

$$p^+(U) = p([p(\text{Col}_1[U]); ...; p(\text{Col}_n[U])])$$
 $[U \in \mathbf{S}^n]$

- 1.1. Prove that p^+ indeed is a norm on \mathbf{S}^n , and $p^+(xx^T) = p^2(x)$. Denoting by $q(\cdot)$ the norm conjugate to $p(\cdot)$, what is the relation between the norm $(p^+)_*(\cdot)$ conjugate to $p^+(\cdot)$ and the norm $q^+(\cdot)$?
- 1.2. Derive from 1.1 that whenever $p(\cdot)$ is an absolute norm such that \mathcal{X} is contained

in the unit ball $\mathcal{B}_{p(\cdot)} = \{x : p(x) \leq 1\}$ of the norm p, the set

$$\mathcal{U}_{p(\cdot)} = \{ U \in \mathbf{S}^n : U \succeq 0, p^+(U) \le 1 \}$$

is matching \mathcal{X} . If, in addition,

$$\mathcal{X} \subset \{x : p(x) \le 1, Px = 0\},\tag{4.176}$$

then the set

$$\mathcal{U}_{p(\cdot),P} = \{ U \in \mathbf{S}^n : U \succeq 0, p^+(U) \le 1, PU = 0 \}$$

is matching \mathcal{X} .

Assume that in addition to $p(\cdot)$, we have at our disposal a computationally tractable closed convex set \mathcal{D} such that whenever $p(x) \leq 1$, the vector $[x]^2 := [x_1^2, ...; x_n^2]$ belongs to \mathcal{D} ; in the sequel we call such a set \mathcal{D} square-dominating $p(\cdot)$. For example, when $p(\cdot) = \|\cdot\|_r$, we can take

$$\mathcal{D} = \begin{cases} \{y \in \mathbf{R}^n_+ : \sum_i y_1 \le 1\}, & r \le 2\\ \{y \in \mathbf{R}^n_+ : \|y\|_{r/2} \le 1\}, & r > 2 \end{cases}$$

Prove that in this situation the above construction can be refined: whenever \mathcal{X} satisfies (4.176), the set

$$\mathcal{U}_{p(\cdot),P}^{\mathcal{D}} = \{ U \in \mathbf{S}^{n} : U \succeq 0, p^{+}(U) \le 1, PU = 0, \mathrm{dg}(U) \in \mathcal{D} \} \\ [\mathrm{dg}(U) = [U_{11}; U_{22}; ...; U_{nn}]]$$

matches \mathcal{X} .

Note: in the sequel, we suppress P in the notation $\mathcal{U}_{p(\cdot),P}$ and $\mathcal{U}_{p(\cdot),P}^{\mathcal{D}}$ when P = 0; thus, $\mathcal{U}_{p(\cdot)}$ is the same as $\mathcal{U}_{p(\cdot),0}$.

1.3. Check that when $p(\cdot) = \|\cdot\|_r$ with $r \in [1, \infty]$, one has

$$p^{+}(U) = \|U\|_{r} := \begin{cases} (\sum_{i,j} |U_{ij}|^{r})^{1/r}, & 1 \le r < \infty, \\ \max_{i,j} |U_{ij}|, & r = \infty \end{cases}$$

1.4. Let $\mathcal{X} = \{x \in \mathbf{R}^n : ||x||_1 \le 1\}$ and $p(x) = ||x||_1$, so that $\mathcal{X} \subset \{x : p(x) \le 1\}$, and

Conv{
$$[x]^2 : x \in \mathcal{X}$$
} $\subset \mathcal{D} = \{ y \in \mathbf{R}^n_+ : \sum_i y_1 = 1 \}.$ (4.177)

What are the bounds $\overline{\operatorname{Opt}}_{\mathcal{U}_{p(\cdot)}}(B)$ and $\overline{\operatorname{Opt}}_{\mathcal{U}_{p(\cdot)}}^{\mathcal{D}}(B)$? Is it true that the former (the latter) of the bounds is precise? Is it true that the former (the latter) of the bounds is precise when $B \succeq 0$?

- 1.5. Let $\mathcal{X} = \{x \in \mathbf{R}^n : ||x||_2 \leq 1\}$ and $p(x) = ||x||_2$, so that $\mathcal{X} \subset \{x : p(x) \leq 1\}$ and (4.177) holds true. What are the bounds $\operatorname{Opt}_{\mathcal{U}_{p(\cdot)}}(B)$ and $\operatorname{Opt}_{\mathcal{U}_{p(\cdot)}}(B)$? Is it true that the former (the latter) of the bounds is precise?
- 1.6. Let $\mathcal{X} \subset \mathbf{R}^n_+$ be closed, convex, bounded, and with a nonempty interior. Verify that the set

$$\mathcal{X}^+ = \{ x \in \mathbf{R}^n : \exists y \in \mathcal{X} : \operatorname{abs}[x] \le y \}$$

is the unit ball of an absolute norm $p_{\mathcal{X}}$, and this is the largest absolute norm $p(\cdot)$ such that $\mathcal{X} \subset \{x : p(x) \leq 1\}$. Derive from this observation that the norm

360

LECTURE 4

 $p_{\mathcal{X}}(\cdot)$ is the best (i.e., resulting in the least conservative bounding scheme) among absolute norms which allow to upper-bound $\operatorname{Opt}_{\mathcal{X}}(B)$ via the construction from item 1.2.

4.50.B "Calculus of matchings." Observe that matching we have introduced admits a kind of "calculus." Specifically, consider the situation as follows: for $1 \leq \ell \leq L$, we are given

• nonempty convex compact sets $\mathcal{X}_{\ell} \subset \mathbf{R}^{n_{\ell}}, 0 \in \mathcal{X}_{\ell}$, along with matching \mathcal{X}_{ℓ} convex compact sets $\mathcal{U}_{\ell} \subset \mathbf{S}^{n_{\ell}}$ giving rise to the closed convex cones

$$\mathbf{U}_{\ell} = \operatorname{cl}\{(U_{\ell}, \mu_{\ell}) \in \mathbf{S}^{n_{\ell}} \times \mathbf{R}_{+} : \mu_{\ell} > 0, \mu_{\ell}^{-1} U_{\ell} \in \mathcal{U}_{\ell}\}$$

We denote by $\vartheta_{\ell}(\cdot)$ the Minkovski functions of \mathcal{X}_{ℓ} :

$$\vartheta_{\ell}(y^{\ell}) = \inf\{t : t > 0, t^{-1}y^{\ell} \in \mathcal{X}_{\ell}\} : \mathbf{R}^{n_{\ell}} \to \mathbf{R} \cup \{+\infty\};$$

note that $\mathcal{X}_{\ell} = \{ y^{\ell} : \vartheta_{\ell}(y^{\ell}) \leq 1 \};$ • $n_{\ell} \times n$ matrices A_{ℓ} such that $\sum_{\ell} A_{\ell}^{T} A_{\ell} \succ 0.$

On the top of it, we are given a monotone convex set $\mathcal{T} \subset \mathbf{R}^L_+$ intersecting the interior of \mathbf{R}^{L}_{\perp} .

These data specify the convex set

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : \vartheta_\ell^2(A_\ell x) \le t_\ell, \ell \le L \}$$
(*)

2.1. Prove the following

Lemma 4.51. In the situation in question, the set

$$\mathcal{U} = \left\{ U \in \mathbf{S}^n : U \succeq 0 \& \exists t \in \mathcal{T} : (A_\ell U A_\ell^T, t_\ell) \in \mathbf{U}_\ell, \ell \le L \right\}$$

is a closed and bounded convex set which matches \mathcal{X} . As a result, the efficiently computable quantity

$$\overline{\operatorname{Opt}}_{\mathcal{U}}(B) = \max_{U} \left\{ \operatorname{Tr}(BU) : U \in \mathcal{U} \right\}$$

is an upper bound on

$$\operatorname{Opt}_{\mathcal{X}}(B) = \max_{x \in \mathcal{X}} x^T B x.$$

- 2.2. Prove that if $\mathcal{X} \subset \mathbf{R}^n$ is a nonempty convex compact set, U is $m \times n$ matrix of rank m, and \mathcal{U} is matching \mathcal{X} , then the set $\mathcal{V} = \{V \in \mathcal{S}^m : V \succeq 0, PVP^T \in \mathcal{U}\}$ matches $\mathcal{Y} = \{y : \exists x \in \mathcal{X} : y = Px\}.$
- 2.3. Consider the "direct product" case where $\mathcal{X} = \mathcal{X}_1 \times ... \times \mathcal{X}_L$; specifying A_ℓ as the matrix which "cuts of" a block vector $x = [x^1; ...; x^L] \in \mathbf{R}^{n_1} \times ... \times \mathbf{R}^{n_L} \ell$ -th block: $A_{\ell}x = x^{\ell}$ and setting $\mathcal{T} = [0,1]^L$, we cover this situation by the setup under consideration. In the direct product case, the construction from item 2.1 is as follows: given the sets \mathcal{U}_{ℓ} matching \mathcal{X}_{ℓ} , we build the set

$$\mathcal{U} = \{ U = [U^{\ell\ell'} \in \mathbf{R}^{n_\ell \times n_{\ell'}}]_{\ell,\ell' \le L} \in \mathbf{S}^{n_1 + \dots + n_L} : U \succeq 0, U^{\ell\ell} \in \mathcal{U}_\ell, \ell \le L \}$$

and claim that this set matches \mathcal{X} .

Could we be less conservative? While we do not know how to be less conservative in general, we do know how to be less conservative in the special case when \mathcal{U}_{ℓ} are built via the "absolute norm" machinery. Specifically, let $p_{\ell}(\cdot) : \mathbf{R}^{n_{\ell}} \to \mathbf{R}_{+}, \ell \leq L$, be absolute norms, let sets \mathcal{D}_{ℓ} be square-dominating $p_{\ell}(\cdot)$, let

$$\mathcal{X}^{\ell} \subset \widehat{X}_{\ell} = \{ x^{\ell} \in \mathbf{R}^{n_{\ell}} : P_{\ell} x_{\ell} = 0, p_{\ell}(x^{\ell}) \le 1 \},\$$

and let $\mathcal{U}_{\ell} = \{ U \in \mathbf{S}^{n_{\ell}} : U \succeq 0, P_{\ell}U = 0, p_{\ell}^+(U) \leq 1, \mathrm{dg}(U) \in \mathcal{D}_{\ell} \}$. In this case the above construction results in

$$\mathcal{U} = \left\{ U = [U^{\ell\ell'} \in \mathbf{R}^{n_\ell \times n_{\ell'}}]_{\ell,\ell' \le L} \in \mathbf{S}^{n_1 + \dots + n_L} : \\ U \succeq 0, P_\ell U^{\ell\ell} = 0, p_\ell^+(U^{\ell\ell}) \le 1, \operatorname{dg}(U^{\ell\ell}) \in \mathcal{D}_\ell, \ell \le L \right\}.$$

Now let

$$p([x^1; \dots; x^L]) = \max[p_1(x^1), \dots, p_L(x^L)] : \mathbf{R}^{n_1} \times \dots \times \mathbf{R}^{n_L} \to \mathbf{R}^{n_L}$$

so that p is an absolute norm and $\mathcal{X} \subset \{x = [x^1; ...; x^L] : p(x) \le 1, P_\ell x^\ell = 0, \ell \le L\}.$

Prove that in fact the set

$$\overline{\mathcal{U}} = \left\{ U = [U^{\ell\ell'} \in \mathbf{R}^{n_\ell \times n_{\ell'}}]_{\ell,\ell' \le L} \in \mathbf{S}^{n_1 + \dots + n_L} : \\ U \succeq 0, P_\ell U^{\ell\ell} = 0, \operatorname{dg}(U^{\ell\ell}) \in \mathcal{D}_\ell, \ell \le L, p^+(U) \le 1 \right\}$$

matches \mathcal{X} , and that we always have $\overline{\mathcal{U}} \subset \mathcal{U}$. Verify that in general this inclusion is strict.

4.50.C Illustration: Nullspace property revisited. Recall sparsity-oriented signal recovery via ℓ_1 minimization from Lecture 1: Given $m \times n$ sensing matrix A and (noiseless) observation y = Aw of unknown signal w known to have at most s nonzero entries, we recover w as

$$\widehat{w} \in \operatorname{Argmin}_{z} \left\{ \|z\|_{1} : Az = y \right\}.$$

Matrix A is called s-good, if whenever y = Aw with s-sparse w, the only optimal solution to the right hand side optimization problem is w. The (difficult to verify!) necessary and sufficient condition for s-goodness is the Nullspace property:

Opt :=
$$\max_{z} \{ \|z\|_{(s)} : z \in \operatorname{Ker} A, \|z\|_{1} \le 1 \} < 1/2,$$

where $||z||_{(k)}$ is the sum of the k largest entries in the vector abs[z]. A verifiable sufficient condition for s-goodness is

$$\widehat{\operatorname{Opt}} := \min_{H} \max_{j} \|\operatorname{Col}_{j}[I - H^{T}A]\|_{(s)} < 1/2,$$
(!)

the reason being that, as it is immediately seen, $\widehat{\text{Opt}}$ is an upper bound on Opt (see Proposition 1.9 with q = 1).

An immediate observation is that Opt is nothing but the maximum of quadratic form on appropriate convex compact set. Specifically, let

$$\begin{aligned} \mathcal{X} &= \{ [u; v] \in \mathbf{R}^n \times \mathbf{R}^n : Au = 0, \|u\|_1 \le 1, \sum_i |v_i| \le s, \|v\|_* \le 1 \}, \\ B &= \left[\frac{\frac{1}{2}I_n}{\frac{1}{2}I_n} \right]. \end{aligned}$$

Then

where (a) is due to the well known fact (prove it!) that whenever s is a positive integer $\leq n$, the extreme points of the set

$$V = \{ v \in \mathbf{R}^n : \sum_i |v_i| \le s, \|v\|_{\infty} \le 1 \}$$

are exactly the vectors with at most s nonzero entries, the nonzero entries being ± 1 ; as a result

$$\forall (z \in \mathbf{R}^n) : \max_{v \in V} z^T v = \|z\|_{(s)}.$$

Now, V is the unit ball of the absolute norm

$$r(v) = \min \left\{ t : \|v\|_1 \le st, \|v\|_\infty \le t \right\},\$$

so that \mathcal{X} is contained in the unit ball \mathcal{B} of the absolute norm on \mathbf{R}^{2n} specified as

$$p([u;v]) = \max\{\|u\|_1, r(v)\} \qquad [u, v \in \mathbf{R}^n],\$$

specifically,

$$\mathcal{X} = \{ [u; v] : p([u, v]) \le 1, Au = 0 \}.$$

As a result, whenever $x = [u; v] \in \mathcal{X}$, the matrix

$$U = xx^{T} = \begin{bmatrix} U^{11} = uu^{T} & U^{12} = uv^{T} \\ U^{21} = vu^{T} & U^{22} = vv^{T} \end{bmatrix}$$

satisfies the condition $p^+(U) \leq 1$ (see item 1.2 above). In addition, this matrix clearly satisfies the condition

$$A[U^{11}, U^{12}] = 0.$$

It follows that the set

$$\mathcal{U} = \{ U = \left[\begin{array}{c|c} U^{11} & U^{12} \\ \hline U^{21} & U^{22} \end{array} \right] \in \mathbf{S}^{2n} : U \succeq 0, p^+(U) \le 1, AU^{11} = 0, AU^{12} = 0 \}$$

(which clearly is a nonempty convex compact set) matches \mathcal{X} . As a result, the efficiently computable quantity

$$\overline{\text{Opt}} = \max_{U \in \mathcal{U}} \operatorname{Tr}(BU) = \max_{U} \left\{ \operatorname{Tr}(U^{12}) : U = \left[\frac{U^{11} | U^{12}}{U^{21} | U^{22}} \right] \succeq 0, p^+(U) \le 1, AU^{11} = 0, AU^{12} = 0 \right\}$$
(!!)

is an upper bound on Opt, so that the verifiable condition

is sufficient for s-goodness of A.

Now goes the concluding part of Exercise:

3.1. Prove that $\overline{\text{Opt}} \leq \widehat{\text{Opt}}$, so that (!!) is less conservative than (!). <u>Hint:</u> Apply Conic Duality to verify that

$$\widehat{\operatorname{Opt}} = \max_{V} \left\{ \operatorname{Tr}(V) : V \in \mathbf{R}^{n \times n}, AV = 0, \sum_{i=1}^{n} r(\operatorname{Col}_{i}[V^{T}]) \leq \right\}$$
(!!!)

3.2. Run simulations with randomly generated Gaussian matrices A and play with different values of s to compare $\widehat{\text{Opt}}$ and $\overline{\text{Opt}}$. To save time, you can use toy sizes m, n, say, m = 18, n = 24.

4.9.4 Around Propositions 4.4 and 4.14

4.9.4.1 Optimizing linear estimates on convex hulls of unions of spectratopes

Exercise 4.52. [optimizing linear estimates on convex hull of union of spectratopes] Let

• $\mathcal{X}_1, ..., \mathcal{X}_J$ be spectratopes in \mathbf{R}^n :

$$\begin{aligned} \mathcal{X}_{j} &= \{ x \in \mathbf{R}^{n} : \exists (y \in \mathbf{R}^{N_{j}}, t \in \mathcal{T}_{j}) : x = P_{j}y, R_{kj}^{2}[y] \leq t_{k}I_{d_{kj}}, \leq K_{j} \}, \ 1 \leq j \leq J \\ & \left[R_{kj}[y] = \sum_{i=1}^{N_{j}} y_{i}R^{kji} \right] \end{aligned}$$

- $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{\nu \times n}$ be given matrices,
- $\|\cdot\|$ be a norm on \mathbf{R}^{ν} such that the unit ball \mathcal{B}_* of the conjugate norm $\|\cdot\|_*$ is a spectratope:

$$\begin{aligned} \mathcal{B}_{*} &:= \{ u : \|u\|_{*} \leq 1 \} \\ &= \{ u \in \mathbf{R}^{\nu} : \exists (z \in \mathbf{R}^{N}, r \in \mathcal{R}) : u = Mz, S_{\ell}^{2}[z] \preceq r_{\ell} I_{f_{\ell}}, \ell \leq L \} \\ & \left[S_{\ell}[z] = \sum_{i=1}^{N} z_{i} S^{\ell i} \right] \end{aligned}$$

• Π be a convex compact subset of the interior of the positive semidefinite cone \mathbf{S}^m_+ ,

with our standard restrictions on $R_{kj}[\cdot], S_{\ell}[\cdot], \mathcal{T}_j, \mathcal{R}$. Let, further,

$$\mathcal{X} = \operatorname{Conv}\left(\bigcup_{j} \mathcal{X}_{j}\right)$$

be the convex hull of the union of spectratopes \mathcal{X}_j . Consider the situation where we, given observation

$$\omega = Ax + \xi$$

of unknown signal x known to belong to \mathcal{X} , want to recover Bx. We assume that the matrix of second moments of noise is \succeq -dominated by a matrix from Π , and quantify the performance of a candidate estimate $\hat{x}(\cdot)$ by its $\|\cdot\|$ -risk

$$\operatorname{Risk}_{\Pi, \|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \sup_{P: P \ll \Pi} \mathbf{E}_{\xi \sim P} \left\{ \|Bx - \widehat{x}(Ax + \xi)\| \right\}$$

where $P \ll \Pi$ means that the matrix $\operatorname{Var}[P] = \mathbf{E}_{\xi \sim P}\{\xi\xi^T\}$ of second moments of distribution P is \succeq -dominated by a matrix from Π .

Prove the following

Proposition 4.53. In the situation in question, consider convex optimization problem

$$Opt = \min_{H,\Theta,\Lambda^{j},\Upsilon^{j},\Upsilon'} \left\{ \max_{j} \left[\phi_{\mathcal{T}_{j}}(\lambda[\Lambda^{j}]) + \phi_{\mathcal{R}}(\lambda[\Upsilon^{j}]) \right] + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \Gamma_{\Pi}(\Theta) : \\ \Lambda^{j} = \left\{ \Lambda^{j}_{k} \succeq 0, j \leq K_{j} \right\}, j \leq J, \\ \Upsilon^{j} = \left\{ \Upsilon^{j}_{\ell} \succeq 0, \ell \leq L \right\}, j \leq J, \Upsilon' = \left\{ \Upsilon'_{\ell} \succeq 0, \ell \leq L \right\} \\ \left[\frac{\sum_{k} \mathcal{R}^{k}_{kj}[\Lambda^{j}_{k}]}{\frac{1}{2}M^{T}[B - H^{T}A]P_{j}} \middle| \frac{1}{2}P_{j}^{T}[B^{T} - A^{T}H]M}{\sum_{\ell} \mathcal{S}^{k}_{\ell}[\Upsilon^{j}_{\ell}]} \right] \succeq 0, j \leq J, \right\}, \\ \left[\frac{\Theta}{\frac{1}{2}M^{T}H^{T}} \middle| \frac{1}{2}P_{\ell}^{T}[\Upsilon^{j}_{\ell}]} \right] \succeq 0$$

$$(4.178)$$

where, as usual,

$$\phi_{\mathcal{T}_{j}}(\lambda) = \max_{t \in \mathcal{T}_{j}} t^{T} \lambda, \ \phi_{\mathcal{R}}(\lambda) = \max_{r \in \mathcal{R}} r^{T} \lambda,$$

$$\Gamma_{\Pi}(\Theta) = \max_{Q \in \Pi} \operatorname{Tr}(Q\Theta), \ \lambda[U_{1}, ..., U_{s}] = [\operatorname{Tr}(U_{1}); ...; \operatorname{Tr}(U_{S})],$$

$$\mathcal{S}_{\ell}^{*}[\cdot] : \mathbf{S}^{f_{\ell}} \to \mathbf{S}^{N} : \mathcal{S}_{\ell}^{*}[U] = [\operatorname{Tr}(S^{\ell p} U S^{\ell q})]_{p,q \leq N},$$

$$\mathcal{R}_{kj}^{*}[\cdot] : \mathbf{S}^{d_{kj}} \to \mathbf{S}^{N_{j}} : \mathcal{R}_{kj}^{*}[U] = [\operatorname{Tr}(R^{kjp} U R^{kjq})]_{p,q \leq N_{j}}$$

Problem (4.178) is solvable, and H-component H_* of its optimal solution gives rise to linear estimate $\hat{x}_{H_*}(\omega) = H_*^T \omega$ such that

$$\operatorname{Risk}_{\Pi, \|\cdot\|} [\widehat{x}_{H_*} | \mathcal{X}] \le \operatorname{Opt.}$$

$$(4.179)$$

Moreover, the estimate \hat{x}_{H_*} is near-optimal among linear estimates:

$$Opt \le O(1) \ln(D+F) \text{RiskOpt}_{lin} \left[D = \max_j \sum_{k \le K_j} d_{kj}, \ F = \sum_{\ell \le L} f_\ell \right]$$
(4.180)

where

$$\operatorname{RiskOpt}_{lin} = \inf_{H} \sup_{x \in \mathcal{X}, Q \in \Pi} \mathbf{E}_{\xi \sim \mathcal{N}(0,Q)} \left\{ \|Bx - H^{T}(Ax + \xi)\| \right\}$$

is the best risk achievable under the circumstances with linear estimates under zero mean Gaussian noise with covariance matrix restricted to belong to Π .

It should be stressed that convex hull of unions of spectratopes not necessarily is a spectratope, and that Proposition states that the linear estimate stemming from

(4.178) is near-optimal only among linear, and not among all estimates (the latter can indeed be not the case).

4.9.4.2 Recovering nonlinear vector-valued functions

Exercise 4.54. † [estimating nonlinear vector-valued functions] Consider situation as follows: We are given a noisy observation

$$\omega = Ax + \xi_x \qquad [A \in \mathbf{R}^{\nu \times n}]$$

of the linear image Ax of an unknown signal x known to belong to a given spectratope $\mathcal{X} \subset \mathbf{R}^n$; here ξ_x is the observation noise with distribution P_x which can depend on x. Similarly to Section 4.3.3, we assume that we are given a computationally tractable convex compact set $\Pi \subset \operatorname{int} \mathbf{S}^{\nu}_+$ such that for every $x \in \mathcal{X}$, $\operatorname{Var}[P_x] \preceq \Theta$ for some $\Theta \in \Pi$, cf. (4.36). What we want is to recover the value f(x) of a given vector-valued function $f : \mathcal{X} \to \mathbf{R}^{\nu}$, and we measure the recovery error in a given norm $|\cdot|$ on \mathbf{R}^{ν} .

4.54.A Preliminaries and Main observation. Let $\|\cdot\|$ be a norm on \mathbb{R}^n , and $g(\cdot): \mathcal{X} \to \mathbb{R}^{\nu}$ be a function. Recall that the function is called *Lipschitz continuous* on \mathcal{X} w.r.t. the pair of norms $\|\cdot\|$ on the argument and $|\cdot|$ on the image spaces, if there exist $L < \infty$ such that

$$|g(x) - g(y)| \le L ||x - y|| \ \forall (x, y \in \mathcal{X});$$

every L with this property is called Lipschitz constant of g. It is well known that in our finite-dimensional situation, the property of g to be Lipschitz continuous is independent of how the norms $\|\cdot\|$, $|\cdot|$ are selected; this selection affects only the value(s) of Lipschitz constant(s).

Assume from now on that the function of interest f is Lipschitz continuous on \mathcal{X} . Let us call a norm $\|\cdot\|$ on \mathbb{R}^n appropriate for f, is f is Lipschitz continuous with constant 1 on \mathcal{X} w.r.t. $\|\cdot\|$, $|\cdot|$. Our immediate observation is as follows:

Observation 4.55. In the situation in question, let $\|\cdot\|$ be appropriate for f. Then recovering f(x) is not more difficult than recovering x in the norm $\|\cdot\|$: every estimate $\hat{x}(\omega)$ of x via ω which takes all its values in \mathcal{X} induces the "plug-in" estimate

$$f(\omega) = f(\widehat{x}(\omega))$$

of f(x), and the $\|\cdot\|$ -risk

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim P_x} \left\{ \|\widehat{x}(Ax + \xi) - x\| \right\}$$

of estimate \hat{x} upper-bounds the $|\cdot|$ -risk

$$\operatorname{Risk}_{|\cdot|}[\widehat{f}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim P_x} \left\{ \|\widehat{f}(Ax + \xi) - f(x)\| \right\}$$

of the induced by \hat{x} estimate \hat{f} :

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{f}|\mathcal{X}] \leq \operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}].$$

When f is defined and Lipschitz continuous with constant 1 w.r.t. $\|\cdot\|, |\cdot|$ on the entire \mathbf{R}^n , the conclusion remains true without the assumption that \hat{x} takes all its values in \mathcal{X} .

4.54.B Consequences. Observation 4.55 suggests the following simple approach to solving the estimation problem we started with: assuming that we have at our disposal a norm $\|\cdot\|$ on \mathbb{R}^n such that

- $\|\cdot\|$ is appropriate for f, and
- $\|\cdot\|$ is good, goodness meaning that the unit ball \mathcal{B}_* of the norm $\|\cdot\|_*$ conjugate to $\|\cdot\|$ is a spectratope given by explicit spectratopic representation,

we use the machinery of linear estimation developed in Section 4.3.3 to build a near-optimal, in terms of its $\|\cdot\|$ -risk, linear estimate of x via ω , and convert this estimate in an estimate of f(x); by Observation, the $|\cdot|$ - risk of the resulting estimate is upper-bounded by $\|\cdot\|$ -risk of the underlying linear estimate. The just outlined construction needs a small correction: in general, the linear estimate $\tilde{x}(\cdot)$ yielded by Proposition 4.14 (same as any nontrivial – not identically zero – *linear* estimate) is *not* guaranteed to take all its values in \mathcal{X} , which is, in general, required for Observation to be applicable. This correction is easy: it is enough to convert \tilde{x} into the estimate \hat{x} defined by

$$\widehat{x}(\omega) \in \underset{u \in \mathcal{X}}{\operatorname{Argmin}} \|u - \widetilde{x}(\omega)\|.$$

This transformation preserves efficient computability of the estimate, and ensures that the corrected estimate takes all its values in \mathcal{X} , so that Observation is applicable to \hat{x} ; at the same time, "correction" $\tilde{x} \mapsto \hat{x}$ nearly preserves the $\|\cdot\|$ -risk:

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}] \le 2\operatorname{Risk}_{\|\cdot\|}[\widetilde{x}|\mathcal{X}]. \tag{(*)}$$

Note that when $\|\cdot\|$ is a (general-type) Euclidean norm: $\|x\|^2 = x^T Q x$ for some $Q \succ 0$, factor 2 in the right hand side can be discarded.

1. Justify (*).

4.54.C How to select $\|\cdot\|$. When implementing the outlined approach, the major question is how to select a norm $\|\cdot\|$ appropriate for f. An ideal for our purposes choice would be to select the smallest among the norms appropriate for f (such a norm does exist under pretty mild assumptions), since the smaller $\|\cdot\|$, the smaller is the $\|\cdot\|$ -risk of an estimate of x. This ideal can be achieved in rare cases only: first, it could be difficult to identify the smallest among the norms appropriate for f, and second, our machinery requires from $\|\cdot\|$ to have an explicitly given spectratope as the unit ball of the conjugate norm. Let us look at a couple of "favorable cases," where the just outlined difficulties can be (partially) avoided.

Example 1: a norm-induced f. Let us start with the important by its own right case when f is a scalar functional which itself is a norm, and this norm has a spectratope as the unit ball of the conjugate norm, as is the case when $f(\cdot) = \|\cdot\|_r$, $r \in [1,2]$, or when $f(\cdot)$ is the nuclear norm. In this case the smallest of the norms appropriate for f clearly is f itself, and no one of the outlined difficulties arises. As

an extension, when f(x) is obtained from a good norm $\|\cdot\|$ by operations preserving Lipschitz continuity and constant, like $f(x) = \|x - c\|$, or $f(x) = \sum_i a_i \|x - c_i\|$, $\sum_i |a_i| \le 1$, or

$$f(x) = \sup_{c \in C} / \inf \|x - c\|,$$

or even something like

$$f(x) = \sup_{\alpha \in \mathcal{A}} / \inf_{c \in C_{\alpha}} \left\{ \sup_{c \in C_{\alpha}} / \inf_{c \in C_{\alpha}} \|x - c\| \right\}$$

it seems natural to use this norm in our construction, although now this, perhaps, is not the smallest of the norms appropriate for f.

Now let us address the general case. Note that in principle the smallest of the norms appropriate for a given Lipschitz continuous f admits a description. Specifically, assume that \mathcal{X} has a nonempty interior (this is w.l.o.g. – we can always replace \mathbb{R}^n with the linear span of \mathcal{X}). A well-known fact of Analysis (Rademacher Theorem) states that in this situation (more generally, when \mathcal{X} is convex with a nonempty interior), a Lipschitz continuous f is differentiable almost everywhere in $\mathcal{X}^o = \operatorname{int} \mathcal{X}$, and f is Lipschitz continuous with constant 1 w.r.t. a norm $\|\cdot\|$ if and only if

$$\|f'(x)\|_{\|\cdot\|\to|\cdot\|} \le 1$$

whenever $x \in \mathcal{X}^o$ is such that the derivative (a.k.a. Jacobian) of f at x exists; here $\|Q\|_{\|\cdot\|\to\|\cdot\|}$ is the matrix norm of a $\nu \times n$ matrix Q induced by the norms $\|\cdot\|$ on \mathbf{R}^n and $|\cdot|$ on \mathbf{R}^{ν} :

$$\|Q\|_{\|\cdot\|\to|\cdot|} := \max_{\|x\|\leq 1} |Qx| = \max_{\|x\|\leq 1\atop \|y\|_*\leq 1} y^T Qx = \max_{\|y_*|\leq 1\atop \|\|x\|_*\|_*\leq 1} x^T Q^T y = \|Q_{|\cdot|_*\to\|\cdot\|_*}^T,$$

where $\|\cdot\|_*$, $|\cdot|_*$ are the conjugates of $\|\cdot\|$, $|\cdot|$.

2. Prove that a norm $\|\cdot\|$ is appropriate for f if and only if the unit ball of the *conjugate* to $\|\cdot\|$ norm contains the set

$$\mathcal{B}_{f,*} = \operatorname{cl}\operatorname{Conv}\{z : \exists (x \in \mathcal{X}_o, y, |y|_* \le 1) : z = [f'(x)]^T y\},\$$

where \mathcal{X}_o is the set of all $x \in \mathcal{X}^o$ where f'(x) exists. Geometrically: $\mathcal{B}_{f,*}$ is the closed convex hull of the union of all images of the unit ball \mathcal{B}_* of $|\cdot|_*$ under the linear mappings $y \mapsto [f'(x)]^T y$ stemming from $x \in \mathcal{X}_o$. Equivalently: $\|\cdot\|$ is appropriate for f if and only if

$$\|u\| \ge \|u\|_f := \max_{z \in \mathcal{B}_{f,*}} z^T u.$$
(!)

Check that $||u||_f$ is a norm, provided that $\mathcal{B}_{f,*}$ (this set by construction is a symmetric w.r.t. the origin convex compact set) possesses a nonempty interior; whenever this is the case, $||u||_f$ is the smallest of the norms appropriate for f. Derive from the above that the norms $|| \cdot ||$ we can use in our approach are the norms on \mathbf{R}^n for which the unit ball of the conjugate norm is a spectratope containing $\mathcal{B}_{f,*}$.

Example 2. Consider the case of componentwise quadratic f:

$$f(x) = \left[\frac{1}{2}x^{T}Q_{1}x; \frac{1}{2}x^{T}Q_{2}x; ...; \frac{1}{2}x^{T}Q_{\nu}x\right] \qquad [Q_{i} \in \mathbf{S}^{n}]$$

and $|u| = ||u||_q$ with $q \in [1, 2]$ ⁷². In this case

$$\mathcal{B}_* = \{ u \in \mathbf{R}^{\nu} : \|u\|_p \le 1 \}, \, p = \frac{q}{q-1} \in [2, \infty[, \text{ and } f'(x) = \left[x^T Q_1; x^T Q_2; ...; x^T Q_{\nu} \right].$$

Setting $S = \{s \in \mathbf{R}_+^{\nu} : \|s\|_{p/2} \le 1\}$ and $S^{1/2} = \{s \in \mathbf{R}_+^{\nu} : [s_1^2; ...; s_{\nu}^2] \in S\} = \{s \in \mathbf{R}_+^{\nu} : \|s\|_p \le 1\}$, the set

$$\mathcal{Z} = \{ [f'(x)]^T u : x \in \mathcal{X}, u \in \mathcal{B}_* \}$$

is contained in the set

$$\mathcal{Y} = \{ y \in \mathbf{R}^n : \exists (s \in \mathcal{S}^{1/2}, x^i \in \mathcal{X}, i \le \nu) : y = \sum_i s_i Q_i x_i \},\$$

and the set \mathcal{Y} is a spectratope with spectratopic representation readily given by the one of \mathcal{X} ; indeed, \mathcal{Y} is nothing but the S-sum of the spectratopes $Q_i\mathcal{X}$, $i = 1, ..., \nu$, see Section 4.54. As a result, we can use the spectratope \mathcal{Y} (when $\operatorname{int} \mathcal{Y} \neq \emptyset$) or the arithmetic sum of \mathcal{Y} with a small Euclidean ball (when $\operatorname{int} \mathcal{Y} = \emptyset$) to build an estimate of $f(\cdot)$.

3.1. As a simple illustration, work out the problem of recovering the value of a scalar quadratic form

$$\begin{split} f(x) &= x^T M x, \, M = \text{Diag}\{i^\alpha, i=1,...,n\} \\ & [\nu=1,|\cdot| \text{ is the usual absolute value}] \end{split}$$

from noisy observation

$$\omega = Ax + \sigma\eta, A = \text{Diag}\{i^{\beta}, i = 1, ..., n\}, \eta \sim \mathcal{N}(0, I_n)$$

of a signal x known to belong to the ellipsoid

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \|Px\|_2 \le 1 \}, P = \text{Diag}\{i^{\gamma}, i = 1, ..., n \},\$$

where α , β , γ are given reals satisfying

$$2\alpha - \gamma - 2\beta < -1.$$

You could start with the simplest unbiased estimate

$$\widetilde{x}(\omega) = [1^{-\beta}\omega_1; 2^{-\beta}\omega_2; ...; n^{-\beta}\omega_n]$$

of x.

 $^{^{72}}$ to save notation, we assume that the linear parts in the components of f_i are trivial – just zeros. In this respect, note that we always can subtract from f a whatever linear mapping and reduce our estimation problem to those of estimating separately the values at the signal x of the modified f and the linear mapping we have subtracted (we know how to solve the latter problem reasonably well).

3.2. Work out the problem of recovering the norm

$$f(x) = ||Mx||_p, M = \text{Diag}\{i^{\alpha}, i = 1, ..., n\}, p \in [1, 2],$$

from the same observations as in item 3.1 and with

$$\mathcal{X} = \{x : \|Px\|_r \le 1\}, P = \mathrm{dg}\{i^{\gamma}, i = 1, ..., n\}, r \in [2, \infty].$$

4.9.4.3 Suboptimal linear estimation

Exercise 4.56. [†] [recovery of large-scale signals] When building presumably good linear recovery of the image $Bx \in \mathbf{R}^{\nu}$ of signal $x \in \mathcal{X}$ from observation

$$\omega = Ax + \sigma\xi \in \mathbf{R}^m$$

in the simplest case where $\mathcal{X} = \{x \in \mathbf{R}^n : x^T S x \leq 1\}$ is an ellipsoid (so that $S \succ 0$), the recovery error is measured in $\|\cdot\|_2$, and $\xi \sim \mathcal{N}(0, I_m)$, problem (4.13) reduces to

$$Opt = \min_{H,\lambda} \left\{ \lambda + \sigma^2 \|H\|_F^2 : \left[\frac{\lambda S}{|B^T - A^T H|} \right] \succeq 0 \right\}, \qquad (4.181)$$

where $\|\cdot\|_F$ is the Frobenius norm of a matrix. An optimal solution H_* to this problem results in the linear estimate $\hat{x}_{H_*}(\omega) = H_*^T \omega$ satisfying the risk bound

$$\operatorname{Risk}[\widehat{x}_{H_*}|\mathcal{X}] := \max_{x \in \mathcal{X}} \sqrt{\mathbf{E}\{\|Bx - H_*^T(Ax + \sigma\xi)\|_2^2\}} \le \sqrt{\operatorname{Opt}}.$$

Now, (4.181) is an efficiently solvable convex optimization problem. However, when the sizes m, n of the problem are large, solving the problem by the standard optimization techniques could become prohibitively time-consuming. The goal of what follows is to develop relatively computationally cheap technique for finding a hopefully good *sub*optimal solution to (4.181). In the sequel, we assume that $A \neq 0$, otherwise (4.181) is trivial.

- 1. Prove that problem (4.181) can be reduced to similar problem with $S = I_n$ and diagonal positive semidefinite matrix A, the reduction requiring several singular value decompositions and multiplications of matrices of the same sizes as those of A, B, S.
- 2. By item 1, we can assume from the very beginning that S = I and $A = \text{Diag}\{\alpha_1, ..., \alpha_n\}$ with $0 \le \alpha_1 \le \alpha_2 \le ... \le \alpha_n$. Passing in (4.181) from variables λ, H to variables $\tau = \sqrt{\lambda}, G = H^T$, the problem becomes

$$Opt = \min_{G,\tau} \left\{ \tau^2 + \sigma^2 \|G\|_F^2 : \|B - GA\| \le \tau \right\},$$
(4.182)

where $\|\cdot\|$ is the spectral norm. Now consider the construction as follows:

- We build a partition $\{1, ..., n\} = I_0 \cup I_1 \cup ... \cup I_K$ of the index set $\{1, ..., n\}$ into consecutive segments in such a way that
 - (a) I_0 is the set of those *i*, if any, for which $\alpha_i = 0$, and $I_k \neq \emptyset$ when $k \ge 1$,
 - (b) for $k \ge 1$ the ratios α_j/α_i , $i, j \in I_k$, do not exceed a $\theta > 1$ (θ is the parameter of our construction), while

370

LECTURE 4

(c) for $1 \le k < k' \le K$, the ratios α_j / α_i , $i \in I_k$, $j \in I_{k'}$, are $> \theta$. The recipe for building the partition is self-evident, and we clearly have

$$K \le \ln(\overline{\alpha}/\underline{\alpha})/\ln(\theta) + 1,$$

where $\overline{\alpha}$ is the largest of α_i , and $\underline{\alpha}$ is the smallest of those α_i which are positive.

• For $1 \leq k \leq K$, we denote by i_k the first index in I_k , set $\alpha^k = \alpha_{i_k}$, $n_k = Card I_k$, and define A_k as $n_k \times n_k$ diagonal matrix with diagonal entries α_i , $i \in I_k$.

Now, given $\nu \times n$ matrix C, let us specify C_k , $0 \le k \le K$, as $\nu \times n_k$ submatrix of C comprised of columns with indexes from I_k , and consider the following parametric optimization problems:

where $\tau \ge 0$ is the parameter, and $1 \le k \le K$. Justify the following simple observations:

2.1. G_k is feasible for $(P_k[\tau])$ if and only if the matrix

$$G_k^* = \alpha^k G_k A_k^{-1}$$

is feasible for $(P_k^*[\tau])$, and $||G_k^*||_F \leq ||G_k||_F \leq \theta ||G_k^*||_F$, implying that

$$\operatorname{Opt}_{k}^{*}(\tau) \leq \operatorname{Opt}_{k}(\tau) \leq \theta^{2} \operatorname{Opt}_{k}^{*}(\tau);$$

2.2. Problems $(P_k[\tau])$ is easy to solve: if $B_k = U_k D_k V_k^T$ is singular value decomposition of B_k and $\sigma_{k\ell}$, $1 \le \ell \le \nu_k := \min[\nu, n_k]$, are diagonal entries of D_k , then an optimal solution to $(P_k[\tau])$ is

$$\widehat{G}_k[\tau] = [\alpha^k]^{-1} U_k D_k[\tau] V_k^T,$$

where $D_k[\tau]$ is obtained from D_k by truncating $\sigma_{k\ell} \mapsto [\sigma_{k\ell} - \tau]_+$ of diagonal entries and keeping zero the off-diagonal entries (from now on, $a_+ = \max[a, 0]$, $a \in \mathbf{R}$). The optimal value in $(P_k[\tau])$ is

$$\operatorname{Opt}_{k}(\tau) = [\alpha^{k}]^{-2} \sum_{\ell=1}^{\nu_{k}} [\sigma_{k\ell} - \tau]_{+}^{2}$$

2.3. If (τ, G) is feasible solution to (4.182), then $\tau \ge \underline{\tau} := ||B_0||$ and the matrices $G_k, 1 \le k \le K$, are feasible solutions to problems $(P_k^*[\tau])$, implying that

$$\sum_{k} \operatorname{Opt}_{k}^{*}(\tau) \leq \|G\|_{F}^{2}$$

and nearly vice versa: if $\tau \geq \underline{\tau}$, G_k , $1 \leq k \leq K$, are feasible solutions to problems $(P_k^*[\tau])$, and

$$K_{+} = \begin{cases} K, & I_{0} = \emptyset \\ K+1, & I_{0} \neq \emptyset \end{cases},$$

then the matrix $G = [0_{\nu \times n_0}, G_1, ..., G_k]$ taken along with $\tau^+ = \sqrt{K_+ \tau}$ form a feasible solution to (4.182).

Extract from these observations that if τ_* is an optimal solution to the convex optimization problem

$$\min_{\tau} \left\{ \theta^2 \tau^2 + \sigma^2 \sum_{k=1}^{K} \operatorname{Opt}_k(\tau) : \tau \ge \underline{\tau} \right\}$$
(4.183)

and $G_{k,*}$ are optimal solutions to the problems $(P_k[\tau_*])$, then the pair

$$\widehat{\tau} = \sqrt{K_+} \tau_*, \widehat{G} = [0_{\nu \times n_0}, G^*_{1,*}, ..., G^*_{K,*}] \qquad \qquad [G^*_{k,*} = \alpha^k G_{k,*} A_k^{-1}]$$

is a feasible solution to (4.182), and the value of the objective of the latter problem at this feasible solution is within the factor $\max[K_+, \theta^2]$ of the true optimal value Opt of this problem. As a result, \hat{G} gives rise to a linear estimate with risk on \mathcal{X} which is within factor $\max[\sqrt{K_+}, \theta]$ of the risk $\sqrt{\text{Opt}}$ of the "presumably good" linear estimate yielded by an optimal solution to (4.181).

Pay attention to the facts that

- After carrying out singular value decompositions of matrices B_k , $1 \le k \le K$, specifying τ_* and $G_{k,*}$ requires solving univariate convex minimization problem with easy to compute objective, so that the problem can be easily solved, e.g., by bisection;
- The computationally cheap suboptimal solution we end up with is not that bad, since K is "moderate" just logarithmic in the condition number $\overline{\alpha}/\underline{\alpha}$ of A.

Your next task is a follows:

- 3. To get an idea of the performance of the proposed synthesis of "suboptimal" linear estimation, run numerical experiments as follows:
 - select somehow n and generate at random the $n \times n$ data matrices S, A, B
 - for "moderate" values of n compute both the presumably good linear estimate by solving $(4.13)^{73}$ and the suboptimal estimate as yielded by the above construction and compare their risk bounds and the associated CPU times. For "large" n, where solving (4.13) becomes prohibitively time consuming, compute only suboptimal estimate in order to get an impression how the corresponding CPU time grows with n.

Recommended setup:

- range of n: 50, 100 ("moderate" values), 1000, 2000 ("large" values)
- range of σ : {1.0, 0.01, 0.0001}
- generation of S, A, B: generate the matrices at random according to

$$S = U_S \text{Diag}\{1, 2, ..., n\} U_S^T, A = U_A \text{Diag}\{\mu_1, ..., \mu_n\} V_A^T, B = U_B \text{Diag}\{\mu_1, ..., \mu_n\} V_B^T,$$

⁷³When \mathcal{X} is an ellipsoid, semidefinite relaxation bound on the maximum of a quadratic form over $x \in \mathcal{X}$ is exact, so that we are in the case when an optimal solution to (4.13) yields the best, in terms of risk on \mathcal{X} , linear estimate.

where U_S, U_A, V_A, U_B, V_B are random orthogonal $n \times n$ matrices, and μ_i form a geometric progression with $\mu_1 = 0.01$ and $\mu_n = 1$.

You could run the above construction for several values of θ and select the best, in terms of its risk bound, of the resulting suboptimal estimates.

4.56.A Simple case. There is a trivial case where (4.182) is really easy; this is the case in the singular value decompositions of A and B the right orthogonal factors are the same, that is, when

$$B = WFV^T, A = UDV^T$$

with orthogonal $n \times n$ matrices W, U, V and diagonal F, D. This, at the first glance, very special case is in fact of some importance – it covers the *denoising* situation where B = A, so that our goal is to denoise our observation of Ax given a priori information $x \in \mathcal{X}$ on x. In this situation, setting $W^T H^T U = G$, problem (4.182) becomes

$$Opt = \min_{C} \left\{ \|F - GD\|^2 + \sigma^2 \|G\|_F^2 \right\}.$$
(4.184)

Now goes the concluding part of Exercise:

4. Prove that in the situation in question an optimal solution G_* to (4.184) can be selected to be diagonal, with diagonal entries γ_i , $1 \leq i \leq n$, yielded by the optimal solution to the optimization problem

$$Opt = \min_{\gamma} \left\{ f(G) := \max_{i \le n} (\phi_i - \gamma_i \delta_i)^2 + \sigma^2 \sum_{i=1}^n \gamma_i^2 \right\} \qquad [\phi_i = F_{ii}, \delta_i = D_{ii}]$$

Exercise 4.57. [†] [image reconstruction – follow-up to Exercise 4.56] A grayscale image can be represented by $m \times n$ matrix $x = [x_{pq}]_{\substack{0 \le p \le m, \\ 0 \le q \le n}}$ with entries in the range $[-\overline{x}, \overline{x}]$, with $\overline{x} = 255/2$ ⁷⁴. Taking picture can be modeled as observing in noise the 2D convolution $x \star \kappa$ of image x with known blurring kernel $\kappa = [\kappa_{uv}]_{\substack{0 \le u \le 2\mu, \\ 0 \le v \le 2\nu}}$, so that the observation is the random matrix

$$\omega = \left[\omega_{rs} = \underbrace{\sum_{\substack{0 \le u \le 2\mu, 0 \le v \le 2\nu \\ 0 \le p < m, 0 \le q < n: \\ u+p=r, v+q=s}} x_{pq}\kappa_{uv}}_{[x \neq \kappa]_{rs}} + \sigma\xi_{rs}\right]_{\substack{0 \le r < m+2\mu, \\ 0 \le s < n+2\nu}},$$

where independent of each other random variables $\xi_{rs} \sim \mathcal{N}(0, 1)$ form observation noise⁷⁵. Our goal is to build a presumably good linear estimate of x via ω . To apply the machinery developed in Section 4.2.2, we need to cover the set of signals x allowed by our a priori assumptions by an ellitope \mathcal{X} , to decide in which norm we want to recover x, and then solve the associated optimization problem (4.13). The difficulty, however, is that the dimension of this problem formally will be huge – with 256 × 256 images (a rather poor resolution!), matrix H we are looking for

⁷⁴The actual grayscale image is a matrix with entries, representing pixels' light intensities, in the range [0,255]. It is convenient for us to represent this actual image as the shift, by \bar{x} , of a matrix with entries in $[-\bar{x}, \bar{x}]$.

 $^{^{75}\}mathrm{pay}$ attention to the fact that everywhere in this Exercise indexing of elements of 2D arrays starts from 0, and not from 1!

is of the size dim $\omega \times \dim x = ((256 + 2\mu)(256 + 2\nu)) \times 256^2 \ge 4.295 \times 10^9$; it is impossible just to store such a matrix in the memory of a usual computer, not speaking about optimizing w.r.t. such a matrix. By this reason, in what follows we develop a "practically," and not just theoretically, efficiently computable estimate.

4.57.A The construction. Our key observation is that when passing from representations of x and ω "as they are" to their Discrete Fourier Transforms, the situation simplifies dramatically. Specifically, for matrices y, x of the same sizes, let $y \bullet z$ be the entrywise product of y and z: $[y \bullet z]_{pq} = y_{pq} z_{pq}$. Setting

$$\alpha = 2\mu + m, \ \beta = 2\nu + n,$$

let $F_{\alpha,\beta}$ be the 2D discrete Fourier Transform – a linear mapping from the space $\mathbf{C}^{\alpha \times \beta}$ onto itself given by

$$[F_{\alpha,\beta}y]_{rs} = \frac{1}{\sqrt{\alpha\beta}} \sum_{\substack{0 \le p < \alpha, \\ 0 \le q < \beta}} y_{pq} \exp\left\{-2\pi i r/\alpha - 2\pi i s/\beta\right\},$$

where *i* is the imaginary unit. It is well known that it is a unitary transformation which is easy-to-compute (it can be computed in $O(\alpha\beta \ln(\alpha\beta))$) arithmetic operations) which "nearly diagonalizes" the convolution: whenever $x \in \mathbf{R}^{m \times n}$, setting

$$x^{+} = \begin{bmatrix} x & 0_{m \times 2\nu} \\ \hline 0_{2\mu \times n} & 0_{2\mu \times 2\nu} \end{bmatrix} \in \mathbf{R}^{\alpha \times \beta},$$

we have

$$F_{\alpha,\beta}(x\star\kappa) = \chi \bullet [F_{\alpha,\beta}x^+]$$

with easy-to-compute χ^{76} . Now, let δ be another $(2\mu + 1) \times (2\nu + 1)$ kernel, with the only nonzero entry, equal to 1, in the position (μ, ν) (recall that numeration of indexes starts from 0); then

$$F_{\alpha,\beta}(x\star\delta) = \theta \bullet [F_{\alpha,\beta}x^+]$$

with easy-to-compute θ . Now consider the auxiliary estimation problem as follows:

Given R > 0 and noisy observation

$$\widehat{\omega} = \chi \bullet \widehat{x} + \sigma \underbrace{F_{\alpha,\beta} \xi}_{n} \qquad [\xi = [\xi_{rs}] \text{ with independent } \xi_{rs} \sim \mathcal{N}(0,1)],$$

of signal $\hat{x} \in \mathbf{C}^{\alpha \times \beta}$ known to satisfy $\|\hat{x}\|_2 \leq R$, we want to recover, in the Frobenius norm $\|\cdot\|_2$, the matrix $\theta \bullet \hat{x}$.

Treating signals \hat{x} and noises η as long vectors rather than matrices and taking into account that $F_{\alpha,\beta}$ is a unitary transformation, we see that our auxiliary problem is nothing but the problem of recovery, in $\|\cdot\|_2$ -norm, of the image Θz of signal z known to belong to the centered at the origin Euclidean ball \mathcal{Z}_R of radius R in

⁷⁶Specifically, $\chi = \sqrt{\alpha\beta}F_{\alpha,\beta}\kappa^+$, where κ^+ is the $\alpha \times \beta$ matrix with κ as $(2\mu + 1) \times (2\nu + 1)$ North-Western block and zeros outside this block.

 $\mathbf{C}^{\alpha\beta}$, from noisy observation

$$\zeta = Az + \sigma\eta,$$

where Θ and A are *diagonal* matrices with complex entries, and η is random complex-valued noise with zero mean and unit covariance matrix. Exactly the same argument as in the real case demonstrates that as far as linear estimates $\hat{z} = H\zeta$ are concerned, we lose nothing when restricting ourselves with diagonal matrices $H = \text{Diag}\{h\}$, and the best, in terms of its worst-case over $z \in \mathbb{Z}_R$ expected $\|\cdot\|_2^2$ error, estimate corresponds to h solving the optimization problem

$$R^2 \max_{\ell \le \alpha\beta} |\Theta_{\ell\ell} - h_{\ell} A_{\ell\ell}|^2 + \sigma^2 \sum_{\ell \le \alpha\beta} |h_{\ell}|^2.$$

Coming back to the initial setting of our auxiliary estimation problem, we conclude that the best linear recovery of $\theta \bullet \hat{x}$ via $\hat{\omega}$ is given by

$$\widehat{z} = h \bullet \widehat{\omega},$$

where h is an optimal solution to the optimization problem

$$Opt = \min_{h \in \mathbf{C}^{\alpha \times \beta}} \left\{ R^2 \max_{r,s} |\theta_{rs} - h_{rs} \chi_{rs}|^2 + \sigma^2 \sum_{r,s} |h_{rs}|^2 \right\}, \qquad (!)$$

and the $\|\cdot\|_2$ -risk

$$\operatorname{Risk}_{R}[\widehat{z}] = \max_{\|\widehat{x}\|_{2} \leq R} \mathbf{E} \{ \|\theta \bullet \widehat{x} - h \bullet [\chi \bullet \widehat{x} + \sigma \eta] \|_{2} \}$$

of this estimate does not exceed $\sqrt{\text{Opt}}$. Now goes your first task:

1.1. Prove that the above h induces the estimate

$$\widehat{w}(\omega) = F_{\alpha,\beta}^{-1} \left[h \bullet \left[F_{\alpha,\beta} \omega \right] \right]$$

of $x \star \delta$, $x \in \mathcal{X}_R = \{x \in \mathbb{R}^{m \times n} : ||x||_2 \leq R\}$, via observation $\omega = x \star \kappa + \sigma \xi$, with risk

$$\operatorname{Risk}[\widehat{w}|R] = \max_{x \in \mathbf{R}^{m \times n} : \|x\|_2 \le R} \mathbf{E} \left\{ \|x \star \delta - \widehat{w}(x \star \kappa + \sigma \xi)\|_2 \right\}$$

not exceeding $\sqrt{\text{Opt.}}$ Pay attention to the fact that x itself is nothing but a block in $x \star \delta$; note also that in order for \mathcal{X}_R to cover all images we are interested in, it suffices to take $R = \sqrt{mn}\overline{x}$.

- 1.2. Prove that finding optimal solution to (!) is easy the problem is in fact just one-dimensional one!
- 1.3. What are the sources, if any, of the conservatism of the estimate \hat{w} we have built as compared to the linear estimate given by an optimal solution to (4.13)?
- 1.4. Think how to incorporate in the above construction a small number L (say, 5-10) of additional a priori constraints on x of the form

$$\|x \star \kappa_\ell\|_2 \le R_\ell$$

where $\kappa_{\ell} \in \mathbf{R}^{(2\mu+1)\times(2\nu+1)}$, and a priori upper bounds u_{rs} on the magnitudes of

Fourier coefficients of x^+ :

$$|[F_{\alpha\beta}x^+]_{rs}| \le u_{rs}, \ 0 \le r < \alpha, 0 \le s < \beta.$$

4.57.B Mimicking Total Variation constraints. For an $m \times n$ image $x \in \mathbb{R}^{m \times n}$, its (anisotropic) total variation is defined as the ℓ_1 norm of the "discrete gradient field" of x:

$$TV(x) = \underbrace{\sum_{p=0}^{m-1} \sum_{q=0}^{n} |x_{p+1,q} - x_{p,q}|}_{TV_a(x)} + \underbrace{\sum_{p=0}^{m} \sum_{q=0}^{n-1} |x_{p,q+1} - x_{p,q}|}_{TV_b(x)}.$$

A well established experimental fact is that for naturally arising images, their total variation is essentially less than what could be expected given the magnitudes of entries in x and the sizes m, n of the image. As a result, it is tempting to incorporate a priori upper bounds on total variation of the image into an image reconstruction procedure. We are about to explain how this can be done in our context. Unfortunately, while an upper bound on total variation is a convex constraint on the image, incorporating this constraint into our construction would completely destroy its "practical computability." What we can do, is to guess that bounds on $TV_{a,b}(x)$ can be somehow mimicked by bounds on the energy of two convolutions: one with kernel $\kappa_a \in \mathbf{R}^{(2\mu+1)\times(2\nu+1)}$ with the only nonzero entries

$$[\kappa_a]_{\mu,\nu} = -1, [\kappa_a]_{\mu+1,\nu} = 1,$$

and the other one with kernel $\kappa_b \in \mathbf{R}^{(2\mu+1)\times(2\nu+1)}$ with the only nonzero entries

$$[\kappa_b]_{\mu,\nu} = -1, [\kappa_b]_{\mu,\nu+1} = 1$$

(recall that the indexes start from 0, and not from 1). Note that $x \star \kappa_a$ and $x \star \kappa_b$ are "discrete partial derivatives" of $x \star \delta$.

For a small library of grayscale $m \times n$ images x we dealt with, experiment shows that, in addition to the energy constraint $||x||_2 \leq R = \sqrt{mn\overline{x}}$, the images satisfy the constraints

$$\|x \star \kappa_a\|_2 \le \gamma R, \, \|x \star \kappa_b\|_2 \le \gamma_2 R \tag{(*)}$$

with small γ_2 , specifically, $\gamma_2 = 0.25$. In addition, it turns out that the ∞ -norms of the Fourier transforms of $x \star \kappa_a$ and $x \star \kappa_b$ for these images are much less than one could expect looking at the energy of the transform's argument. Specifically, for all images x from the library it holds

$$\begin{aligned} \|F_{\alpha\beta}[x\star\kappa_a]\|_{\infty} &\leq \gamma_{\infty}R, \\ \|F_{\alpha\beta}[x\star\kappa_b]\|_{\infty} &\leq \gamma_{\infty}R, \end{aligned}, \quad \|\{z_{rs}\}_{r,s}\|_{\infty} &= \max_{r,s}|z_{rs}| \end{aligned}$$
(**)

with $\gamma_{\infty} = 0.01^{77}$. Now, relations (**) read

$$\max[|\omega_{rs}^a|, |\omega_{rs}^b|]|F_{\alpha\beta}x^+]_{rs}| \le \gamma_{\infty}R \,\forall r, s$$

⁷⁷ note that from (*) it follows that (**) holds true with $\gamma_{\infty} = \gamma_2$, while with our empirical γ 's, γ_{∞} is 25 times smaller than γ_2 .

with easy-to-compute ω^a and ω^b , and in addition $|[F_{\alpha\beta}x^+]_{rs}| \leq R$ due to $||F_{\alpha\beta}x^+||_2 = ||x^+||_2 \leq R$. We arrive at the bounds

$$[F_{\alpha\beta}x^+]_{rs} \leq \min\left[1, 1/|\omega_{rs}^a|, 1/|\omega_{rs}^b|\right] R \,\forall r, s.$$

on the magnitudes of entries in $F_{\alpha\beta}x^+$, and can utilize item 1.4 to incorporate these bounds, along with relations (*),

Now goes the exercise:

2. Write software implementing the outlined deblurring and denoising image reconstruction routine and run numerical experiments.

Recommended kernel κ : set $\mu = \lfloor m/32 \lfloor, \nu = \lfloor n/32 \rfloor$, start with

$$\kappa_{uv} = \frac{1}{(2\mu+1)(2\nu+1)} + \begin{cases} \Delta, & u = \mu, v = \nu \\ 0, & \text{otherwise} \end{cases}, 0 \le u \le 2\mu, 0 \le v \le 2\nu,$$

and then normalize this kernel to make the sum of entries equal to 1. In this description, $\Delta \geq 0$ is control parameter responsible for well-posedness of the auxiliary estimation problem we end up with: the smaller is Δ , the smaller is $\min_{r,s} |\chi_{rs}|$ (note that when decreasing the magnitudes of χ_{rs} , we increase the optimal value in (!)).

We recommend to compare what happens when $\Delta = 0$ with what happens when $\Delta = 0.25$, same as compare the estimates accounting and not accounting for the constraints (*), (**). On the top of it, you can compare your results with what is given by " ℓ_1 -minimization recovery" described as follows:

As we remember from item 4.57.A, our problem of interest can be equivalently reformulated as recovering the image Θz of a signal $z \in \mathbf{C}^{\alpha\beta}$ from noisy observation $\hat{\omega} = Az + \sigma\eta$, where Θ and A are diagonal matrices, and η is the zero mean complex Gaussian noise with unit covariance matrix. In other words, the entries η_{ℓ} in η are independent of each other real two-dimensional Gaussian vectors with zero mean and the covariance matrix $\frac{1}{2}I_2$. Given a reasonable "reliability tolerance" ϵ , say, $\epsilon = 0.1$, we can easily point out the smallest "confidence radius" ρ such that for $\zeta \sim \mathcal{N}(0, \frac{1}{2}I_2)$ it holds $\operatorname{Prob}\{\|\zeta\|_2 > \rho\} \leq \frac{\epsilon}{\alpha\beta}$, implying that for every ℓ it holds

$$\operatorname{Prob}_{\eta}\left\{\left|\widehat{\omega}_{\ell} - A_{\ell} z_{\ell}\right| > \sigma\rho\right\} \leq \frac{\epsilon}{\alpha\beta},$$

and therefore

$$\operatorname{Prob}_{\eta} \left\{ \|\widehat{\omega} - Az\|_{\infty} > \sigma\rho \right\} \le \epsilon.$$

We now can easily find the smallest, in $\|\cdot\|_1$, vector $\hat{z} = \hat{z}(\omega)$ which is "compatible with our observation," that is, satisfies the constraint

$$\|\widehat{\omega} - A\widehat{z}\|_{\infty} \le \sigma\rho,$$

and take $\Theta \hat{z}$ as the estimate of the "entity of interest" Θz (cf. Regular ℓ_1 recovery from Section 1.2.3).

Note that this recovery needs no a priori information on z.

Exercise 4.58. [classical periodic nonparametric deconvolution] In classical univariate nonparametric regression, one is interested to recover a function f(t) of con-

tinuous argument $t \in [0, 1]$ from noisy observations $\omega_i = f(i/n) + \sigma \eta_i$, $0 \le i \le n$, where $\eta_i \sim \mathcal{N}(0, 1)$ are independent across *i* observation noises. Usually, a priory restrictions on *f* are *smoothness assumptions* – existence of \varkappa continuous derivatives satisfying the a priori upper bounds

$$\left(\int_0^1 |f^{(k)}(t)|^{p_k} dt\right)^{1/p_k} \le L_k, \ 0 \le \varkappa,$$

on their L_{p_k} -norms. The risk of an estimate is defined as the supremum, over f's of given smoothness, expected L_r -norm of the recovery error; the primary emphasis of classical studies here was how the minimax optimal (i.e., the best, over estimates) risk goes to 0 as the number of observations n goes to infinity, what are near-optimal estimates, etc. Many of these studies were dealing with *periodic case* – one where f can be extended on the entire real axis as \varkappa times continuously differentiable function, or, which is the same, when f is treated as a smooth function on the circumference of length 1 rather than on the unit segment [0, 1]. While being slightly simpler for analysis than the general case, the periodic case turned out to be highly instructive: what was established for the latter, usually extended straightforwardly to the former.

What you are about to do in this Exercise, is to apply our machinery of building linear estimates to the outlined recovery of smooth univariate periodic regressing functions.

4.58.A. Setup. What follows is aimed at handling restrictions of smooth functions on the unit (i.e., of unit length) circumference C onto an equidistant n-point grid Γ_n on the circumference. These restrictions form the usual n-dimensional coordinate space \mathbf{R}^n ; it is convenient to index the entries in $f \in \mathbf{R}^n$ starting from 0 rather than from 1. We equip \mathbf{R}^n with two linear operators:

• Cyclic shift (in the sequel – just shift) Δ :

$$\Delta \cdot [f_0; f_1; \dots; f_{n-2}; f_{n-1}] = [f_{n-1}; f_0; f_1; \dots; f_{n-2}],$$

and

• Derivative D:

$$D = n[I - \Delta];$$

Treating $f \in \mathbf{R}^n$ as a restriction of a function F on C onto Γ_n , Df is the finitedifference version of the first order derivative of the function, and the norms

$$|f|_p = n^{-1/p} ||f||_p, \, p \in [1,\infty]$$

are the discrete versions of the L_p -norms of F.

Next, we can associate with $\chi \in \mathbf{R}^n$ the operator $\sum_{i=0}^{n-1} \chi_i \Delta^i$; the image of $f \in \mathbf{R}^n$ under this operator is denoted $\chi \star f$ and is called (cyclic) *convolution* of χ and f.

The problem we intend to focus on is as follows:

Given are:

• smoothness data represented by a nonnegative integer \varkappa and two collections: $\{L_{\iota} > 0 : 0 \leq \iota \leq \varkappa\}, \{p_{\iota} \in [2, \infty], 0 \leq \iota \leq \varkappa\}$. The smoothness

data specify the set

$$\mathcal{F} = \{ f \in \mathbf{R}^n : |f|_{p_\iota} \le L_\iota, 0 \le \iota \le \varkappa \}$$

of signals we are interested in (this is the discrete analogy of *periodic* Sobolev ball – the set of \varkappa times continuously differentiable functions on C with derivatives of orders up to \varkappa bounded, in integral p_{ι} -norms, by given quantities L_{ι} ;

- two vectors $\alpha \in \mathbf{R}^n$ (sensing kernel) and $\beta \in \mathbf{R}^n$ (decoding kernel);
- positive integer σ (noise intensity) and a real $q \in [1, 2]$.

These data define the estimation problem as follows: given noisy observation

$$\omega = \alpha \star f + \sigma \eta$$

of unknown signal f known to belong to \mathcal{F} , where $\eta \in \mathbf{R}^n$ is random observation noise, we want to recover $\beta \star f$ in norm $|\cdot|_q$.

The only assumption on the noise is that

$$\operatorname{Var}[\eta] := \mathbf{E}\left\{\eta\eta^T\right\} \preceq I_n.$$

The risk of a candidate estimate \hat{f} is defined as

$$\operatorname{Risk}_{r}[\widehat{f}|\mathcal{F}] = \sup_{\substack{f \in \mathcal{F}, \\ \eta: \operatorname{Cov}[\eta] \leq I_{n}}} \mathbf{E}_{\eta} \left\{ |\beta \star f - \widehat{f}(\alpha \star f + \sigma \eta)|_{q} \right\}.$$

Now goes the exercise:

- 1. Check that the situation in question fits the framework of Section 4.3.3 and figure out to what, under the circumstances, boils down the optimization problem (4.50) responsible for the presumably good linear estimate $\hat{f}_H(\omega) = H^T \omega$.
- 2. Prove that in the case in question the linear estimate yielded by an appropriate optimal solution to (4.50) is just the cyclic convolution

$$\widehat{f}(\omega) = h \star \omega$$

and work out a computationally cheap way to identify h.

3. Implement your findings in software and run simulations. You could, in particular, consider the denoising problem (that is, the one where $\alpha \star x \equiv \beta \star x \equiv x$) and compare numerically the computed risks of your estimates with the classical result on the limits of performance in recovering smooth univariate regression functions; according to this results, in the situation in question and under the natural assumption that L_{ι} are nondecreasing in ι , the minimax optimal risk, up to a factor depending solely on \varkappa , is $(\sigma^2/n)^{\frac{\varkappa}{2\varkappa+1}}L_{\frac{2\varkappa+1}{2\varkappa+1}}^{\frac{1}{2\varkappa+1}}$.

4.9.4.4 Probabilities of large deviations in linear estimation under sub-Gaussian noise

Exercise 4.59. The goal of Exercise is to derive bounds for probabilities of large deviations for estimates yielded by Proposition 4.14.

1. Prove the following fact:

Lemma 4.60. Let $\Theta, Q \in \mathbf{S}^m_+$, with $Q \succ 0$, and let ξ be sub-Gaussian, with parameters (μ, S) , random vector, where μ and S satisfy $\mu\mu^T + S \preceq Q$. Setting $\rho = \text{Tr}(\Theta Q)$, we have

$$\mathbf{E}_{\xi}\left\{\exp\{\frac{1}{8\rho}\xi^{T}\Theta\xi\}\right\} \leq \sqrt{2}\exp\{1/4\}.$$
(4.185)

As a result, for t > 0 it holds

$$\operatorname{Prob}\{\sqrt{\xi^T \Theta \xi} \ge t\sqrt{\rho}\} \le \sqrt{2} \exp\{1/4\} \exp\{-t^2/8\}, t \ge 0.$$
(4.186)

<u>Hint:</u> You could use the same trick as in the proof of Lemma 2.78.

2. Recall that (proof of) Proposition 4.14 states that in the situation of Section 4.3.3.1 and under Assumptions \mathbf{A}' , \mathbf{R} , for every feasible solution $(H, \Lambda, \Upsilon, \Upsilon', \Theta)$ to the optimization problem⁷⁸

$$Opt = \min_{H,\Lambda,\Upsilon,\Upsilon',\Theta} \left\{ \underbrace{\phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon])}_{\mathcal{A}=\mathcal{A}(\Lambda,\Upsilon)} + \underbrace{\phi_{\mathcal{R}}(\lambda[\Upsilon']) + \Gamma_{\Pi}(\Theta)}_{\mathcal{B}=\mathcal{B}(\Theta,\Upsilon')} : \\ \Lambda = \{\Lambda_k \succeq 0, k \leq K\}, \ \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\}, \ \Upsilon' = \{\Upsilon'_\ell \succeq 0, \ell \leq L\}, \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T[B - H^TA]} \right] \underbrace{\sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell]}_{\mathcal{L}_\ell \mathcal{S}_\ell^*[\Upsilon'_\ell]} \right] \succeq 0, \\ \left[\frac{\Theta}{\frac{1}{2}M^TH^T} \frac{1}{\sum_\ell \mathcal{S}_\ell^*[\Upsilon'_\ell]} \right] \succeq 0 \right\},$$

$$(4.187)$$

one has

$$\max_{x \in \mathcal{X}} \| [B - H^T A] x \| \le \mathcal{A} \quad \& \quad \max_{P: \operatorname{Var}[P] \lll \Pi} \mathbf{E}_{\xi \sim P} \left\{ \| H^T \xi \| \right\} \le \mathcal{B}, \tag{4.188}$$

implying that the linear estimate $\hat{x}_H(\omega) = H^T \omega$ satisfies the risk bound

$$\operatorname{Risk}_{\Pi, \|\cdot\|} [\widehat{x}_H(\cdot) | \mathcal{X}] \le \mathcal{A} + \mathcal{B}.$$
(4.189)

Prove the following

Proposition 4.61. Let $H, \Lambda, \Upsilon, \Upsilon', \Theta$) be a feasible solution to (4.187), and let $\hat{x}_H(\omega) = H^T \omega$. Let, further, P be sub-Gaussian, with parameters (μ, S) satisfying

$$\mu\mu^T + S \lll \Pi$$

probability distribution on \mathbb{R}^m . Finally, let $x \in \mathcal{X}$. Then (i) One has

$$\mathbf{E}_{\xi \sim P} \left\{ \|Bx - \widehat{x}_H(Ax + \xi)\| \right\} \leq \mathcal{A}_* + \mathcal{B}_*, \\ \mathcal{A}_* = \mathcal{A}_*(\Lambda, \Upsilon) := 2\sqrt{\phi_{\mathcal{T}}(\lambda[\Lambda])\phi_{\mathcal{R}}(\lambda[\Upsilon])} \leq \mathcal{A}(\Lambda, \Upsilon) := \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) \\ \mathcal{B}_* = \mathcal{B}_*(\Theta, \Upsilon') := 2\sqrt{\Gamma_{\Pi}(\Theta)\phi_{\mathcal{R}}(\lambda[\Upsilon'])} \leq \mathcal{B}(\Theta, \Upsilon') := \Gamma_{\Pi}(\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon'])$$

⁷⁸for notation, see Section 4.3.3.1, (4.43), and (4.47). For reader's convenience, we recall part of this notation: for a probability distribution P on \mathbf{R}^m , $\operatorname{Var}[P] = \mathbf{E}_{\xi \sim P}\{\xi^T \xi\}$, Π is a convex compact subset of int \mathbf{S}^m_+ , $Q \ll \Pi$ means that $Q \preceq Q'$ for some $Q' \in \Pi$, and $\Gamma_{\Pi}(\Theta) = \max_{Q \in \Pi} \operatorname{Tr}(\Theta Q)$.

(ii) For every $\epsilon \in (0,1)$ one has

 $\operatorname{Prob}_{\xi \sim P}\{\xi : \|Bx - \widehat{x}_H(Ax + \xi)\| > \mathcal{A}_* + \theta_{\epsilon}\mathcal{B}_*\} \le \epsilon, \ \theta_{\epsilon} = 2\sqrt{2\ln(\sqrt{2}e^{1/4}/\epsilon)},$ (4.190)with $\mathcal{A}_*, \mathcal{B}_*$ defined in (i).

3. Assume we are given observation $\omega = Ax + \xi$ of unknown signal x known to belong to a given spectratope $\mathcal{X} \subset \mathbf{R}^n$ and want to recover the signal, quantifying the error of a candidate recovery \hat{x} as $\max_{k \leq K} \|B_k(\hat{x}-x)\|_{(k)}$, where $B_k \in \mathbf{R}^{\nu_k \times n}$ are given matrices, and $\|\cdot\|_{(k)}$ are given norms on \mathbf{R}^{ν_k} (for example, x can represent a discretization of a continuous-time signal, and $B_k x$ can be finite-difference approximations of signal's derivatives). As about observation noise ξ , assume, same as in item 2, that it is independent of signal x and is sub-Gaussian with sub-Gaussianity parameters μ , S satisfying $\mu\mu^T + S \leq Q$, for some given matrix $Q \succ 0$. Finally, assume that the unit balls of the norms conjugate to the norms $\|\cdot\|_{(k)}$ are spectratopes. In this situation, Proposition 4.14 provides us with Kefficiently computable linear estimates $\hat{x}_k(\omega) = H_k^T \omega : \mathbf{R}^{\dim \omega} \to \mathbf{R}^{\nu_k}$ along with upper bounds Opt_k on their risks $\max_{x \in \mathcal{X}} \mathbf{E} \{\|B_k x - \hat{x}_k(Ax + \xi)\|_{(k)}\}$. Think how, given reliability tolerance $\epsilon \in (0, 1)$, assemble these linear estimates into a single estimate $\hat{x}(\omega) : \mathbf{R}^{\dim \omega} \to \mathbf{R}^n$ such that for every $x \in \mathcal{X}$, the probability of the event

$$\|B_k(\widehat{x}(Ax+\xi)-x)\|_{(k)} \le \theta \operatorname{Opt}_k, \ 1 \le k \le K,$$
(!)

is at least $1 - \epsilon$, for some moderate (namely, logarithmic in K and $1/\epsilon$) "assembling price" θ .

Exercise 4.62. [†] Prove that if ξ is uniformly distributed on the unit sphere $\{x : \|x\|_2 = 1\}$ in \mathbf{R}^n , then ξ is sub-Gaussian with parameters $(0, \frac{1}{n}I_n)$.

4.9.4.5 Linear recovery under signal-dependent noise

Exercise 4.63. [signal recovery in signal-dependent noise] Consider the situation as follows: we observe a realization ω of *m*-dimensional random vector

$$\omega = Ax + \xi_x,$$

where

• x is unknown signal belonging to a given signal set, specifically, spectratope (which, as always in these cases, we can assume to be basic)

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \preceq t_k I_{d_k}, k \leq K \}$$

with our usual restrictions on \mathcal{T} and $R_k[\cdot]$;

• ξ_x is observation noise with distribution which can depend on x; all we know is that

$$\operatorname{Var}[\xi_x] := \mathbf{E}\{\xi_x \xi_x^T\} \preceq \mathcal{C}[x],$$

where the entries of symmetric matrix C[x] are quadratic in x. We assume in the sequel that signals x belong to the subset

$$\mathcal{X}_{\mathcal{C}} = \{ x \in \mathcal{X} : \mathcal{C}[x] \succeq 0 \}$$
of \mathcal{X} :

• Our goal is to recover Bx, with given $B \in \mathbf{R}^{\nu \times n}$, in a given norm $\|\cdot\|$ such that the unit ball \mathcal{B}_* of the conjugate norm is a spectratope:

$$\mathcal{B}_* = \{ u : \|u\|_* \le 1 \} = M\mathcal{V}, \mathcal{V} = \{ v : \exists r \in \mathcal{R} : S_{\ell}^2[v] \le r_{\ell} I_{f_{\ell}}, \ell \le L \}.$$

As always, we quantify the performance of a candidate estimate $\hat{x}(\omega): \mathbf{R}^m \to \mathbf{R}^{\nu}$ by the risk

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}_{\mathcal{C}}] = \sup_{x \in \mathcal{X}_{\mathcal{C}}} \sup_{\xi_x: \operatorname{Cov}[\xi_x] \leq \mathcal{C}[x]} \mathbf{E} \left\{ \|Bx - \widehat{x}(Ax + \xi_x)\| \right\}.$$

1. Utilize semidefinite relaxation to build, in a computationally efficient fashion, a "presumably good" linear estimate, specifically, prove the following

Proposition 4.64. In the situation in question, for $G \in \mathbf{S}^m$ let us define $\alpha_0[G] \in \mathbf{R}$, $\alpha_1[G] \in \mathbf{R}^n$, $\alpha_2[G] \in \mathbf{S}^n$ from the identity

$$\operatorname{Tr}(\mathcal{C}[x]G) = \alpha_0[G] + \alpha_1^T[G]x + x^T \alpha_2[G]x \; \forall (x \in \mathbf{R}^n, G \in \mathbf{S}^m),$$

so that $\alpha_{\chi}[G]$ are affine in G. Consider convex optimization problem

Whenever $H, \mu, D, \Lambda, \Upsilon, \Upsilon', G$ is feasible for the problem, one has

$$\operatorname{Risk}_{\|\cdot\|}[H(\cdot)|\mathcal{X}_{\mathcal{C}}] \leq \mu + \phi_{\mathcal{T}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]).$$

2. Work out the following special case of the above situation dealing with Poisson Imaging, see Section 2.4.3.2: your observation is *m*-dimensional random vector with independent Poisson entries, the vector of parameters of the corresponding Poisson distributions being Py; here P is $m \times n$ entrywise nonnegative matrix, and the unknown signal y is known to belong to a given box $Y = \{y \in \mathbf{R}^n : \underline{a} \leq y \leq \overline{a}\}$, where $0 \leq \underline{a} < \overline{a}$. You want to recover y in $\|\cdot\|_p$ -norm with given $p \in [1, 2]$.

4.9.5 Signal recovery in Discrete and Poisson observation schemes

Exercise 4.65. [†] The goal of what follows is to "transfer" the constructions of linear estimates to the case of multiple indirect observations of discrete random variables. Specifically, we are interested in the situation where

- Our observation is a K-element sample $\omega^K = (\omega_1, .., \omega_K)$ with independent identically distributed components ω_k taking values in m-element set; as always, we encode the points from this m-element set by the standard basic orths $e_1, ..., e_m$ in \mathbf{R}^m .
- The (common for all k) probability distribution of ω_k is Ax, where x is unknown "signal" n-dimensional probabilistic vector known to belong to a closed convex subset \mathcal{X} of n-dimensional probabilistic simplex $\Delta_n = \{x \in \mathbf{R}^n : x \ge 0, \sum_i x_i = 1\}$, and A is a given $m \times n$ column-stochastic matrix (i.e., entrywise nonnegative matrix with unit column sums).
- Our goal is to recover Bx, where B is a given $\nu \times n$ matrix, and we quantify a candidate estimate $\hat{x}(\omega^K) : \mathbf{R}^{mK} \to \mathbf{R}^{\nu}$ by its risk

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \mathbf{E}_{\omega^{K} \sim [Ax] \times \ldots \times [Ax]} \left\{ \|Bx - \widehat{x}(\omega^{K})\| \right\},$$

where $\|\cdot\|$ is a given norm on \mathbf{R}^{ν} .

What we intend to use are *linear* estimates – estimates of the form

$$\widehat{x}_{H}(\omega^{K}) = H^{T} \underbrace{\left[\frac{1}{K} \sum_{k=1}^{K} \omega_{k}\right]}_{\widehat{\omega}_{K}[\omega^{K}]}, \qquad (4.191)$$

where $H \in \mathbf{R}^{m \times \nu}$.

1. In the main body of Lecture 4, X always was assumed to be symmetric w.r.t. the origin, which easily implies that we gain nothing when passing from linear estimates to affine ones (sums of linear estimates and constants). Now we are in the case when X can be "heavily asymmetric," which, in general, can make "genuinely affine" estimates more preferable than linear ones. Show that in the case in question, we still lose nothing when restricting ourselves with linear, rather than affine, estimates.

4.65.A Observation scheme revisited. When observation ω^K stems from a signal $x \in \Delta_n$, we have

$$\widehat{\omega}_K[\omega^K] = Ax + \xi_x,$$

where

$$\xi_x = \frac{1}{K} \sum_{k=1}^{K} [\omega_k - Ax]$$

is the average of K independent identically distributed zero mean random vectors with common covariance matrix Q[x].

2. Check that

$$Q[x] = \operatorname{Diag}\{Ax\} - [Ax][Ax]^T,$$

and derive from this fact that the covariance matrix of ξ_x is

$$Q_K[x] = \frac{1}{K}Q[x].$$

Setting

$$\Pi = \Pi_{\mathcal{X}} = \{Q = \frac{1}{K} \text{Diag}\{Ax\} : x \in \mathcal{X}\},\$$

check that $\Pi_{\mathcal{X}}$ is a convex compact subset of the positive semidefinite cone \mathbf{S}^m_+ , and that whenever $x \in \mathcal{X}$, one has $Q[x] \preceq Q$ for some $Q \in \Pi$.

4.65.B Upper-bounding risk of a linear estimate. We can upper-bound the risk of a linear estimate \hat{x}_H as follows:

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{x}_{H}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \mathbf{E}_{\omega^{K} \sim [Ax] \times \ldots \times [Ax]} \{ \|Bx - H^{T}\widehat{\omega}_{K}[\omega^{K}]\| \}$$

$$= \sup_{x \in \mathcal{X}} \mathbf{E}_{\xi_{x}} \{ \|[Bx - H^{T}A]x - H^{T}\xi_{x}\| \}$$

$$\leq \sup_{x \in \mathcal{X}} \|[B - H^{T}A]x\| + \sup_{\xi:\operatorname{Cov}[\xi] \in \Pi_{\mathcal{X}}} \mathbf{E}_{\xi} \{ \|H^{T}\xi\| \}.$$

$$\underbrace{\Psi^{\mathcal{X}}(H)}$$

As in the main body of Lecture 4, we intend to build a "presumably good" linear estimate by minimizing over H the sum of efficiently computable upper bounds $\overline{\Phi}(H)$ on $\Phi(H)$ and $\overline{\Psi}^{\mathcal{X}}(H)$ on $\Psi^{\mathcal{X}}(H)$.

Assuming from now on that the unit ball \mathcal{B}_* of the norm conjugate to $\|\cdot\|$ is a spectratope:

$$\mathcal{B}_* := \{ u : \|u\|_* \le 1 \} = \{ u : \exists r \in \mathcal{R}, y : u = My, S_{\ell}^2[y] \preceq r_{\ell} I_{f_{\ell}}, \, \ell \le L \}$$

with our usual restrictions of \mathcal{R} and S_{ℓ} , we can take, as $\overline{\Psi}^{\mathcal{X}}(\cdot)$, the function (4.48). What we intend to focus on, is efficient upper-bounding of $\Phi(\cdot)$.

To simplify our task, we from now on focus on the case when \mathcal{X} is cut off Δ_n by a bunch of linear inequalities:

$$\mathcal{X} = \{ x \in \mathbf{\Delta}_n : Gx \le g, \, Ex = e \} \qquad [G \in \mathbf{R}^{p \times n}, E \in \mathbf{R}^{q \times n}]$$

Observe that replacing G with $G - \mathbf{1}_p^T g$ and E with $E - \mathbf{1}_q^T e$, we can reduce the situation to the one when all linear constraints in question are homogeneous, that is,

$$\mathcal{X} = \{ x \in \mathbf{\Delta}_n : Gx \le 0, Ex = 0 \}.$$

which is what we assume from now on. Setting

$$F = [G; E; -E] \in \mathbf{R}^{(p+2q) \times n},$$

we have also

$$\mathcal{X} = \{ x \in \mathbf{\Delta}_n : Fx \le 0 \}$$

We assume also that \mathcal{X} is nonempty. Finally, for the sake of some of the constructions to follow, in addition to what was already assumed about the norm $\|\cdot\|$, let us assume that this norm is *absolute*, that is, $\|u\|$ depends only on the vector of *magnitudes* of entries in u. From this assumption it immediately follows that if $0 \le u \le u'$, then $\|u\| \le \|u'\|$ (why?).

4.65.C Bounding Φ , simple case. Defining the *simple case* as the one where there are no linear constraints (formally, G and E are zero matrices), observe that in this case bounding Φ is trivial:

3. Prove that in the simple case Φ is convex and efficiently computable "as is:"

$$\Phi(H) = \max_{i \le n} \| (B - H^T A) g_i \|,$$

where $g_1, ..., g_n$ are the standard basic orths in \mathbb{R}^n .

4.65.D Lagrange upper bound on Φ .

4. Observing that when $\mu \in \mathbf{R}^{p+2q}_+$, the function

$$||(B-H^TA)x|| - \mu^TFx$$

of x is convex in $x \in \Delta_n$ and overestimates $||(B - H^T A)x||$ everywhere on \mathcal{X} , conclude that the efficiently computable convex function

$$\Phi_{L}(H) = \min_{\mu} \max_{i < n} \{ \| (B - H^{T} A) g_{i} \| - \mu^{T} F g_{i} : \mu \ge 0 \}$$

upper-bounds $\Phi(H)$. In the sequel, we call this function the Lagrange upper bound on Φ .

4.65.E Basic upper bound on Φ . For vectors u, v of the same dimension, say, k, let Max[u, v] stand for the entrywise maximum of u, v:

$$[\operatorname{Max}[u, v]]_i = \max[u_i, v_i],$$

and let

$$[u]_+ = \operatorname{Max}[u, 0_k],$$

where 0_k is the k-dimensional zero vector.

5.1. Let $\Lambda_+ \geq 0$ and $\Lambda_- \geq 0$ be $\nu \times (p+2q)$ matrices, $\Lambda \geq 0$ meaning that matrix Λ is entrywise nonnegative. Prove that whenever $x \in \mathcal{X}$, one has

$$\| (B - H^T A) x \| \le \mathcal{B}(x, H, \Lambda_+, \Lambda_-) := \min_{t} \left\{ \| t \| : t \ge \max \left[[(B - H^T A) x - \Lambda_+ F x]_+, [-(B - H^T A) x - \Lambda_- F x]_+ \right] \right\}$$

and that $\mathcal{B}(x, H, \Lambda_+, \Lambda_-)$ is convex in x. 5.2. Derive from 5.1 that whenever Λ_{\pm} are as in 5.1, one has

$$\Phi(H) \leq \mathcal{B}^+(H, \Lambda_+, \Lambda_-) := \max_{i \leq n} \mathcal{B}(g_i, H, \Lambda_+, \Lambda_-),$$

where, as in item 3, $g_1, ..., g_n$ are the standard basic orths in \mathbb{R}^n . Conclude that

$$\Phi(H) \le \Phi_B(H) = \inf_{\Lambda_{\pm}} \left\{ \mathcal{B}^+(H, \Lambda_+, \Lambda_-) : \Lambda_{\pm} \in \mathbf{R}_+^{\nu \times (p+2q)} \right\}$$

and that Φ_B is convex and real-valued. In the sequel we refer to $\Phi_B(\cdot)$ as to the

Basic upper bound on $\Phi(\cdot)$.

4.65.F Sherali-Adams upper bound on Φ . The approach we intend to consider now is the one which we used in Lecture 1, Section 1.3.2, when explaining the origin of the verifiable sufficient condition for *s*-goodness, see p. 26. Specifically, setting

$$W = \left[\begin{array}{c|c} G & I \\ \hline E & \\ \end{array} \right],$$

let us introduce slack variable $z \in \mathbf{R}^p$ and rewrite the description of \mathcal{X} as

$$\mathcal{X} = \{ x \in \mathbf{\Delta}_n : \exists z \ge 0 : W[x; z] = 0 \},\$$

so that \mathcal{X} is the projection of the polyhedral set

$$\mathcal{X}^{+} = \{ [x; z] : x \in \mathbf{\Delta}_{n}, z \ge 0, W[x; z] = 0 \}$$

on the x-space. Projection of \mathcal{X}^+ on the z-space is a nonempty (since \mathcal{X} is so) and clearly bounded subset of the nonnegative orthant \mathbf{R}^p_+ , and we can in many ways cover Z by the simplex

$$\Delta[\alpha] = \{ z \in \mathbf{R}^p : z \ge 0, \sum_i \alpha_i z_i \le 1 \},\$$

where all α_i are positive.

6.1. Let $\alpha > 0$ be such that $Z \subset \Delta[\alpha]$. Prove that

$$\mathcal{X}^{+} = \{ [x; z] : W[x; z] = 0, [x; z] \in \operatorname{Conv}\{ v_{ij} = [g_i; h_j], 1 \le i \le n, 0 \le j \le p \} \},$$

where g_i are the standard basic orths in \mathbf{R}^n , $h_0 = 0 \in \mathbf{R}^p$, and $\alpha_j h_j$, $1 \le j \le p$, are the standard basic orths in \mathbf{R}^p .

6.2. Derive from 5.1 that the efficiently computable convex function

$$\Phi_{SA}(H) = \inf_{C} \max_{i,j} \left\{ \| (B - H^T A)g_i + C^T W v_{ij} \| : C \in \mathbf{R}^{(p+q) \times \nu} \right\}$$

is an upper bound on $\Phi(H)$. In the sequel, we refer to this bound as to the Sherali-Adams one.

4.65.G Combined bound. We can combine the above bounds, specifically, as follows:

7. Prove that the efficiently computable convex function

$$\begin{split} \Phi_{LBS}(H) &= \inf_{(\Lambda_{\pm}, C_{\pm}, \mu, \mu_{+}) \in \mathcal{R}} \max_{i,j} \mathcal{G}_{ij}(H, \Lambda_{\pm}, C_{\pm}, \mu, \mu_{+}), \\ where \\ \mathcal{G}_{ij}(H, \Lambda_{\pm}, C_{\pm}, \mu, \mu_{+}) &:= -\mu^{T} F g_{i} + \mu^{T}_{+} W v_{ij} + \min_{t} \left\{ \|t\| : \\ t \geq \max\left[[(B - H^{T} A - \Lambda_{+} F)g_{i} + C^{T}_{+} W v_{ij}]_{+}, [(-B + H^{T} A - \Lambda_{-} F)g_{i} + C^{T}_{-} W v_{ij}]_{+} \right] \right\}, \\ \mathcal{R} &= \left\{ (\Lambda_{\pm}, C_{\pm}, \mu, \mu_{+}) : \Lambda_{\pm} \in \mathbf{R}^{\nu \times (p+2q)}_{+}, C_{\pm} \in \mathbf{R}^{(p+q) \times \nu}, \mu \in \mathbf{R}^{p+2q}_{+}, \mu_{+} \in \mathbf{R}^{p+q}_{+} \right\}$$

$$(\#)$$

is an upper bound on $\Phi(H)$, and that this Combined bound is at least as good as the Lagrange, the Basic, and the Sherali-Adams ones.

4.65.H How to select α ? A shortcoming of the Sherali-Adams and the combined upper bounds on Φ is the presence of a "degree of freedom" – the positive vector α . Intuitively, we would like to select α to make the simplex $\Delta[\alpha] \supset Z$ to be "as small as possible." It is unclear, however, what "as small as possible" in our context is, not speaking about how to select the required α after we agree how we measure the "size" of $\Delta[\alpha]$. It turns out, however, that we can select efficiently α resulting in the smallest volume $\Delta[\alpha]$.

8. Prove that minimizing the volume of $\Delta[\alpha] \supset Z$ in α reduces to solving the following convex optimization problem:

$$\inf_{\alpha,u,v} \left\{ -\sum_{s=1}^{p} \ln(\alpha_s) : 0 \le \alpha \le -v, E^T u + G^T v \le \mathbf{1}_n \right\}$$
(*)

9. Run numerical experiments to get an impression of the quality of the above bounds. It makes sense to generate problems where we know in advance the actual value of Φ , specifically, to take

$$\mathcal{X} = \{ x \in \mathbf{\Delta}_n : x \ge a \} \tag{a}$$

with $a \ge 0$ such that $\sum_i a_i \le 1$. In this case, we can easily list the extreme point of \mathcal{X} (how?) and thus can easily compute $\Phi(H)$.

In your experiments, you can use the matrices stemming from "presumably good" linear estimates yielded by the optimization problems

$$Opt = \min_{H,\Upsilon,\Theta} \left\{ \overline{\Phi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \Gamma_{\mathcal{X}}(\Theta) : \begin{bmatrix} \Theta & | \frac{1}{2}HM \\ \frac{1}{2}M^{T}H^{T} & | \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \end{bmatrix} \succeq 0 \right\},$$
$$\Gamma_{\mathcal{X}}(\Theta) = \frac{1}{K} \max_{x \in \mathcal{X}} \operatorname{Tr}(\operatorname{Diag}\{Ax\}\Theta),$$
(P)

see Corollary 4.12, with the actual Φ (which with our \mathcal{X} is available), or the upper bounds on Φ (Lagrange, Basic, Sherali-Adams, and Combined) in the role of $\overline{\Phi}$. Note that it makes sense to test 7 bounds rather than 4 of them. Specifically, with additional constraints on the optimization variables in (#), we can get, aside of "pure" Lagrange, Basic, and Sherali-Adams bounds and their "three-component combination" (Combined bound), pairwise combinations of the pure bounds as well. For example, to combine Lagrange and Sherali-Adams bound, it suffices to add to (#) the constraints $\Lambda_{\pm} = 0$.

Exercise 4.66. The exercise to follow deals with recovering discrete probability distributions in *Wasserstein norm*.

Wasserstein distance between probability distributions extremely popular in today Statistics is defined as follows⁷⁹. Consider discrete random variables taking

 $^{^{79}}$ What we intend to consider, stems from the Wasserstein 1-distance between discrete probability distributions; this is a particular case of the general Wasserstein *p*-distance between (not

values in finite observation space $\Omega = \{1, 2, ..., n\}$ which is equipped with metric $\{d_{ij} : 1 \leq i, j \leq n\}$ satisfying the standard axioms of a metric⁸⁰. As always, we identify probability distributions on Ω with *n*-dimensional probabilistic vectors $p = [p_1; ...; p_n]$, where p_i is the probability mass assigned by p to $i \in \Omega$. The Wasserstein distance between probability distributions p and q is defined as

$$W(p,q) = \min_{x=[x_{ij}]} \left\{ \sum_{i} d_{ij} x_{ij} : x_{ij} \ge 0, \sum_{j} x_{ij} = p_i, \sum_{i} x_{ij} = q_j \ \forall 1 \le i, j \le n \right\}$$

In other words, think about p and q as about two distributions of unit mass between the points of Ω , and about the problem of transporting masses assigned to points by distribution p from point to point in order to get the distribution q. Denoting by x_{ij} the mass transported from point i to point j, the constraints $\sum_j x_{ij} = p_i$ say that the total mass taken from point i is exactly p_i , the constraints $\sum_i p_{ij} = q_j$ say that as the result of transportation, the mass at point j will be exactly q_j , and the constraints $x_{ij} \geq 0$ reflect the fact that transport of a negative mass is forbidden. Assuming that the cost of transporting a mass μ from point i to point j is $d_{ij}\mu$, the Wasserstein distance W(p,q) between p and q is the cheapest transportation plan which converts p into q. As compared to other natural distances between distance is that it allows to model the situation (indeed arising in some applications) where the effect, measured in terms of intended application, of changing probability masses of points from Ω is small when the probability mass of a point is redistributed among *close* points.

Now goes the first part of the exercise:

1. Let p, q be two probability distributions. Prove that

$$W(p,q) = \max_{f \in \mathbf{R}^n} \left\{ \sum_{i} f_i(p_i - q_i) : |f_i - f_j| \le d_{ij} \,\forall i, j \right\}$$
(4.192)

Treating a vector $f \in \mathbf{R}^n$ as a function on Ω , the value of the function at a point $i \in \Omega$ being f_i , (4.192) admits a very transparent interpretation: the Wasserstein distance W(p,q) between probability distributions p,q is the maximum of inner products of p - q and Lipschitz continuous, with constant 1 w.r.t. the metric d, functions f on Ω . When shifting f by a constant, the inner product remains intact (since p - q is a vector with zero sum of entries), and therefore, denoting by

$$D = \max_{i,j} d_{ij}$$

the *d*-diameter of Ω , we have

$$W(p,q) = \max_{f} \left\{ f^{T}(p-q) : |f_{i} - f_{j}| \le d_{ij} \,\forall i, j, |f_{i}| \le D/2 \,\forall i \right\}, \tag{4.193}$$

the reason being that every Lipschitz continuous, with constant 1 w.r.t. metric d,

necessarily discrete) probability distributions.

⁸⁰specifically, $d_{ij} = d_{ji} \ge 0$, with $d_{ij} = 0$ if and only if i = j, and the Triangle inequality $d_{ik} \le d_{ij} + d_{jk}$ for all triples i, j, k.

function f on Ω can be shifted by a constant to ensure $||f||_{\infty} \leq D/2$ (look what happens when the shift ensures that $\min_i f_i = -D/2$).

Representation (4.193) shows that the Wasserstein distance is generated by a norm on \mathbb{R}^n : for all probability distributions on Ω one has

$$W(p,q) = \|p-q\|_W$$

where $\|\cdot\|_W$ is the Wasserstein norm on \mathbf{R}^n given by

$$\|x\|_{W} = \max_{f \in \mathcal{B}_{*}} f^{T} x,$$

$$\mathcal{B}_{*} = \left\{ u \in \mathbf{R}^{n} : u^{T} S_{ij} u \leq 1, 1 \leq i \leq j \leq n \right\},$$

$$S_{ij} = \left\{ \begin{array}{l} d_{ij}^{-2} [e_{i} - e_{j}] [e_{i} - e_{j}]^{T}, & 1 \leq i < j \leq n \\ 4D^{-2} e_{i} e_{i}^{T}, & 1 \leq i = j \leq n \end{array} \right.,$$

$$(4.194)$$

where $e_1, ..., e_n$ are the standard basic orths in \mathbb{R}^n .

The next portion of Exercise is as follows:

2. Let us equip *n*-element set $\Omega = \{1, ..., d\}$ with the metric $d_{ij} = \begin{cases} 2, & i \neq j \\ 0, & i = j \end{cases}$. What is the associated Wasserstein norm?

Note that the set \mathcal{B}_* in (4.194) is the unit ball of the norm conjugate to $\|\cdot\|_W$, and as we see, this set is a basic ellitope. As a result, the estimation machinery developed in Lecture 4 is well-suited for recovering discrete probability distributions in Wasserstein norm. This observation motivates the concluding part of Exercise:

3. Consider the situation as follows: Given $m \times n$ column-stochastic matrix A and $\nu \times n$ column-stochastic matrix B, we observe K independent of each other samples ω_k , $1 \leq k \leq K$, drawn from discrete probability distribution $Ax \in \Delta_m^{81}$, $x \in \Delta_n$ being unknown "signal" underlying observations; as always, realizations of ω_k are identified with respective vertices $f_1, ..., f_m$ of Δ_m . Our goal is to use the observations to recover the distribution $Bx \in \Delta_{\nu}$. We are given a metric d on the set $\Omega_{\nu} = \{1, 2, ..., \nu\}$ of indexes of entries in Bx, and measure the recovery error in the associated with d Wasserstein norm $\|\cdot\|_W$.

Build explicit convex optimization problem responsible for "presumably good" linear recovery of the form

$$\widehat{x}_H = \frac{1}{K} H^T \sum_{k=1}^K \omega_k.$$

Exercise 4.67. [follow-up to Exercise 4.65] In Exercise 4.65, we have built a "presumably good" linear estimate $\hat{x}_{H_*}(\cdot)$, see (4.191), yielded by the *H*-component H_* of an optimal solution to problem (*P*), see p. 386; the optimal value Opt in this problem is an upper bound on the risk $\operatorname{Risk}_{\|\cdot\|}[\hat{x}_{H_*}|\mathcal{X}]$ (here and in what follows we use the same notation and impose the same assumptions as in Exercise 4.65). Now, $\operatorname{Risk}_{\|\cdot\|}$ is the worst, w.r.t. signals $x \in \mathcal{X}$ underlying our observations, expected norm of the recovery error. It makes sense also to provide upper bounds on the probabilities of deviations of error's magnitude from its expected value, and this is

⁸¹as always, $\Delta_{\nu} \subset \mathbf{R}^{\nu}$ is the probabilistic simplex in \mathbf{R}^{ν} .

the problem we intend to focus on, cf. Exercise 4.59. Now goes the exercise:

1) Prove the following

Lemma 4.68. Let $Q \in \mathbf{S}_{+}^{m}$, let K be a positive integer, and let $p \in \Delta_{m}$. Let, further, $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ be i.i.d. random vectors, with ω_{k} taking the value e_{j} $(e_{1}, ..., e_{m}$ are the standard basic orths in \mathbf{R}^{m}) with probability p_{j} . Finally, let $\xi_{k} = \omega_{k} - \mathbf{E}\{\omega_{k}\} = \omega_{k} - p$, and $\hat{\xi} = \frac{1}{K} \sum_{k=1}^{K} \xi_{k}$. Then for every $\epsilon \in (0, 1)$ it holds

$$\operatorname{Prob}\left\{\|\widehat{\xi}\|_{2}^{2} \leq \frac{12\ln(2m/\epsilon)}{K}\right\} \leq \epsilon.$$

Hint: use the classical

Bernstein inequality: Let $X_1, ..., X_K$ be independent zero mean random variables taking values in [-M, M], and let $\sigma_k^2 = \mathbf{E}\{X_k^2\}$. Then for every $t \ge 0$ one has

$$\operatorname{Prob}\left\{\sum_{k=1}^{K} X_k \ge t\right\} \le \exp\{-\frac{t^2}{2\left[\sum_k \sigma_k^2 + \frac{1}{3}Mt\right]}\}.$$

- 2) Consider the situation described in Exercise 4.65 with $\mathcal{X} = \Delta_n$, specifically,
 - Our observation is a sample $\omega^K = (\omega_1, ..., \omega_K)$ with i.i.d. components $\omega_k \sim Ax$, where $X \in \Delta_n$ is unknown *n*-dimensional probabilistic vector, A is $m \times n$ stochastic matrix (nonnegative matrix with unit column sums), and $\omega \sim Ax$ means that ω is random vector taking value e_i (e_i are standard basic orths in \mathbf{R}^m) with probability $[Ax]_i, 1 \leq i \leq m$;
 - Our goal is to recover Bx in a given norm $\|\cdot\|$; here B is a given $\nu \times n$ matrix.
 - We assume that the unit ball \mathcal{B}_* of the norm $\|\cdot\|_*$ conjugate to $\|\cdot\|$ is a spectratope:

$$\mathcal{B}_* = \{ u = My, y \in \mathcal{Y} \}, \ \mathcal{Y} = \{ y \in \mathbf{R}^N : \exists r \in \mathcal{R} : S_\ell^2[y] \preceq r_\ell I_{f_\ell}, \ell \leq L \}.$$

Our goal is to build a presumably good linear estimate

$$\widehat{x}_H(\omega^K) = H^T \widehat{\omega}[\omega^K], \ \widehat{\omega}[\omega^K] = \frac{1}{K} \sum_k \omega_k.$$

Prove the following

Proposition 4.69. Let H, Θ, Υ be a feasible solution to the convex optimization problem

$$\min_{H,\Theta,\Upsilon} \left\{ \Phi(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \Gamma(\Theta)/K : \begin{bmatrix} \Upsilon_{\ell} \succeq 0, \ell \leq L \} \\ \begin{bmatrix} \Theta & \frac{1}{2}HM \\ \frac{1}{2}M^{T}H^{T} & \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \end{bmatrix} \succeq 0 \\ \Phi(H) = \max_{j \leq n} \|\operatorname{Col}_{j}[B - H^{T}A]\|, \ \Gamma(\Theta) = \max_{x \in \Delta_{n}} \operatorname{Tr}(\operatorname{Diag}\{Ax\}\Theta). \end{cases}$$

$$(4.195)$$

Then

(i) For every $x \in \Delta_n$ it holds

$$\mathbf{E}_{\omega^{K} \sim Ax \times ... \times Ax} \left\{ \|Bx - \widehat{x}_{H}(\omega^{K})\| \right\} \leq \Phi(H) + 2K^{-1/2} \sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon]) \Gamma(\Theta)} \\
\begin{bmatrix} \leq \Phi(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \Phi(H) + \Gamma(\Theta)/K \end{bmatrix} \\
(4.196)$$

(ii) Let $\epsilon \in (0, 1)$. For every $x \in \Delta_n$ with

$$\gamma = 2\sqrt{3\ln(2m/\epsilon)}$$

one has

$$\operatorname{Prob}_{\omega^{K} \sim Ax \times \ldots \times Ax} \left\{ \|Bx - \widehat{x}_{H}(\omega^{K})\| > \Phi(H) + 2\gamma K^{-1/2} \sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon])} \|\Theta\|_{Sh,\infty} \right\}$$

$$\geq 1 - \epsilon.$$
(4.197)

3) Look what happens when $\nu = m = n$, A and B are the unit matrices, and H = I, i.e., we want to understand how good is recovery of a discrete probability distribution by empirical distribute derive from K-element i.i.d. sample drawn from this distribution. Take, as $\|\cdot\|$, the norm $\|\cdot\|_p$ with $p \in [1, 2]$, and show that for every $x \in \Delta_n$ and every $\epsilon \in (0, 1)$ one has

$$\begin{aligned} \forall (x \in \mathbf{\Delta}_n) : \\ \mathbf{E} \left\{ \|x - \hat{x}_I(\omega^K)\|_p \right\} &\leq n^{\frac{1}{p} - \frac{1}{2}} K^{-\frac{1}{2}} \\ \operatorname{Prob} \left\{ \|x - \hat{x}_I(\omega^K)\|_p > 2\sqrt{3\ln(2n/\epsilon)} n^{\frac{1}{p} - \frac{1}{2}} K^{-\frac{1}{2}} \right\} &\geq 1 - \epsilon \quad (b) \end{aligned}$$

$$(4.198)$$

Exercise 4.70. [follow-up to Exercise 4.65] Consider the situation as follows. A retailer sells n items by offering customers via internet bundles of m < n items, so that an offer is an m-element subset B of the set $S = \{1, ..., n\}$ of the items. A customer has private preferences represented by a subset P of S – customer's preference set. We assume that if an offer B intersects with the preference set P of a customer, the latter buys an item drawn at random from the uniform distribution on $B \cap P$, and if $B \cap P = \emptyset$, the customer declines the offer. In the pilot stage we are interested in, the seller learns the market by selecting, one by one, K customers and making offers to them. Specifically, the seller draws k-th customer, and makes the selected customer an offer drawn at random from the uniform distribution on the set $S_{m,n}$ of all m-item offers. What is observed in k-th experiment, is the item, if any, bought by customer, and what we want is to make statistical inferences from these observations.

The outlined observation scheme can be formalized as follows. Let S be the set of all subsets of the *n*-element set, so that S is of cardinality $N = 2^n$. The population of customers induces a probability distribution p on S: for $P \in S$, p_P is the fraction of customers with the preference set being P; we refer to p as to the *preference distribution*. An outcome of a single experiment can be represented by a pair (ι, B) , where $B \in S_{m,n}$ is the offer used in the experiment, and ι is either 0 ("nothing is bought", $P \cap B = \emptyset$), or a point from $P \cap B$, the item which was bought, when $P \cap B \neq \emptyset$. Note that A_P is a probability distribution on the $(M = (m+1)\binom{n}{m})$ element set $\Omega = \{(\iota, B)\}$ of possible outcomes. As a result, our observation scheme

is fully specified by known to us $M \times N$ column-stochastic matrix A with the columns A_P indexed by $P \in S$. When a customer is drawn at random from the uniform distribution on the population of customers, the distribution of the outcome clearly is Ap, where p is the (unknown) preference distribution. Our inferences should be based on K-element sample $\omega^K = (\omega_1, ..., \omega_K)$, with $\omega_1, ..., \omega_K$ drawn, independently of each other, from the distribution Ap.

Now we can pose various inference problems, e.g., the one of recovering p. We, however, intend to focus on a simpler problem – one of recovering Ap. In terms of our story, this makes sense: when we know Ap, we know, e.g., what is the probability for every offer to be "successful" (something indeed is bought) and/or to result in a specific profit, etc. With this knowledge at hand, the seller can pass from "blind" offering policy (drawing an offer at random from the uniform distribution on the set $S_{m,n}$) to something more rewarding.

Now goes the exercise:

1. Use the results of Exercise 4.65 to build "presumably good" linear estimate

$$\widehat{x}_{H}(\omega^{K}) = H^{T} \left[\frac{1}{K} \sum_{k=1}^{K} \omega_{k} \right]$$

of Ap (as always, we encode observations ω , which are elements of M-element set Ω , by standard basic orths in \mathbf{R}^M). As the norm $\|\cdot\|$ quantifying the recovery error, use $\|\cdot\|_1$ and/or $\|\cdot\|_2$. In order to avoid computational difficulties, use small m and n (e.g., m = 3 and n = 5). Compare your results with those for the straightforward estimate $\frac{1}{K} \sum_{k=1}^{K} \omega_k$ (the empirical distribution of $\omega \sim Ap$).

2. Assuming that the "presumably good" linear estimate outperforms the straightforward one, how could this phenomenon be explained? Note that we have no nontrivial a priori information on p!

Exercise 4.71. [Poisson Imaging] Poisson Imaging Problem is to recover an unknown signal observed via Poisson observation scheme. More specifically, assume that our observation is a realization of random vector $\omega \in \mathbf{R}^m_+$ with independent of each other Poisson entries $\omega_i = \text{Poisson}([Ax]_i)$. Here A is a given entrywise nonnegative $m \times n$ matrix, and x is unknown signal known to belong to a given compact convex subset \mathcal{X} of \mathbf{R}^n_+ . Our goal is to recover in a given norm $\|\cdot\|$ the linear image Bx of x, where B is a given $\nu \times n$ matrix.

We assume in the sequel that \mathcal{X} is a subset cut off the *n*-dimensional probabilistic simplex Δ_n by a bunch of linear equality and inequality constraints. The assumption $\mathcal{X} \subset \Delta_n$ is not too restrictive. Indeed, assume that we know in advance a linear inequality $\sum_i \alpha_i x_i \leq 1$ with positive coefficients which is valid on \mathcal{X}^{82} . Introducing slack variable *s* given by $\sum_i \alpha_i x_i + s = 1$, we can pass from signal *x* to the new signal $[\alpha_1 x_1; ...; \alpha_n x_n; s]$, which, after straightforward modification of matrices *A* and *B*, brings the situation to the one where \mathcal{X} is a subset of the probabilistic simplex.

Our goal in the sequel is to build a presumably good linear estimate $\hat{x}_H(\omega) = H^T \omega$ of Bx. Acting in the same fashion as in Exercise 4.65, we start with upper-

⁸²For example, in PET, see Section 2.4.3.2, where x is the density of radioactive tracer injected to the patient taking the PET procedure, we know in advance the total amount $\sum_i v_i x_i$ of the tracer, v_i being the volumes of voxels.

bounding the risk of a linear estimate. Specifically, representing

$$\omega = Ax + \xi_x,$$

we arrive at zero mean observation noise ξ_x with independent of each other entries $[\xi_x]_i = \omega_i - [Ax]_i$ and covariance matrix $\text{Diag}\{Ax\}$. We now can upper-bound the risk of a linear estimate $\hat{x}_H(\cdot)$ in the same fashion as in Exercise 4.65. Specifically, denoting by $\Pi_{\mathcal{X}}$ the set of all diagonal matrices $\text{Diag}\{Ax\}$, $x \in \mathcal{X}$ and by $P_{i,x}$ the Poisson distribution with parameter $[Ax]_i$, we have

$$\operatorname{Risk}_{\|\cdot\|}[\hat{x}_{H}|\mathcal{X}] = \sup_{x \in \mathcal{X}} \operatorname{\mathbf{E}}_{\omega \sim P_{1,x} \times \dots \times P_{m,x}} \left\{ \|Bx - H^{T} \hat{\omega}_{K}[\omega^{K}]\| \right\}$$

$$= \sup_{x \in \mathcal{X}} \operatorname{\mathbf{E}}_{\xi_{x}} \left\{ \|[Bx - H^{T}A]x - H^{T}\xi_{x}\| \right\}$$

$$\leq \sup_{x \in \mathcal{X}} \|[B - H^{T}A]x\| + \sup_{\xi : \operatorname{Cov}[\xi] \in \Pi_{\mathcal{X}}} \operatorname{\mathbf{E}}_{\xi} \left\{ \|H^{T}\xi\| \right\}.$$

In order to build a presumably good linear estimate, it suffices to build efficiently computable convex in H upper bounds $\overline{\Phi}(H)$ on $\Phi(H)$ and $\overline{\Psi}^{\mathcal{X}}(H)$ on $\Psi^{\mathcal{X}}(H)$. and then take as H an optimal solution to the convex optimization problem

$$Opt = \min_{H} \left[\overline{\Phi}(H) + \overline{\Psi}^{\mathcal{X}}(H) \right]$$

Same as in Exercise 4.65, assume from now on that $\|\cdot\|$ is an absolute norm, and the unit ball \mathcal{B}_* of the conjugate norm is a spectratope:

 $\mathcal{B}_* := \{ u : \|u\|_* \le 1 \} = \{ u : \exists r \in \mathcal{R}, y : u = My, S_{\ell}^2[y] \preceq r_{\ell} I_{f_{\ell}}, \ell \le L \}$

Observe that

- In order to build $\overline{\Phi}$, we can use exactly the same techniques as those developed in Exercise 4.65. Indeed, as far as building $\overline{\Phi}$ is concerned, the only difference between our present situation and the one of Exercise4.65 is that in the latter, Awas column-stochastic matrix, while now A is just entrywise nonnegative matrix. Note, however, that when upper-bounding Φ in Exercise 4.65, we never used the fact that A is column-stochastic.
- In order to upper-bound $\Psi^{\mathcal{X}}$, we can use the same bound (4.48) as in Exercise 4.65.

The bottom line is that in order to build a presumably good linear estimate, we need to solve the convex optimization problem

$$Opt = \min_{H,\Upsilon,\Theta} \left\{ \overline{\Phi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \Gamma_{\mathcal{X}}(\Theta) : \begin{bmatrix} \Upsilon_{\ell} \succeq 0, \ell \leq L \} \\ \begin{bmatrix} \Theta & | \frac{1}{2}HM \\ \frac{1}{2}M^{T}H^{T} & | \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \end{bmatrix} \succeq 0 \end{bmatrix},$$

$$\Gamma_{\mathcal{X}}(\Theta) = \max_{x \in \mathcal{X}} \operatorname{Tr}(\operatorname{Diag}\{Ax\}\Theta),$$
(P)

(cf. problem (P) on p. 386) with $\overline{\Phi}$ yielded by a whatever construction from Exercise 4.65, e.g., the least conservative Combined upper bound on Φ .

What in our present situation differs significantly from the situation of Exercise 4.65, are the bounds on probabilities of large deviations established in Exercise

- 4.67, and the goal of what follows is to establish these bounds for Poisson Imaging. Here is what you are supposed to do:
- 1. Let ω be *m*-dimensional random vector with independent entries $\omega_i \sim \text{Poisson}(\mu_i)$, and let $\mu = [\mu_1; ...; \mu_m]$. Prove that whenever $h \in \mathbf{R}^m$, $\gamma > 0$, and $\delta \ge 0$, one has

$$\ln\left(\operatorname{Prob}\{h^T\omega > h^T\mu + \delta\}\right) \le \sum_i [\exp\{\gamma h_i\} - 1]\mu_i - \gamma h^T\mu - \gamma \delta. \qquad (*)$$

2. Taking for granted that $e^x \le 1 + x + \frac{3}{4}x^2$ when $|x| \le 2/3$, prove that in the situation of item 1 one has

$$0 \le \gamma \le \frac{2}{3\|h\|_{\infty}} \Rightarrow \ln\left(\operatorname{Prob}\{h^T\omega > h^T\mu + \delta\}\right) \le \frac{3}{4}\gamma^2 \sum_i h_i^2\mu_i - \gamma\delta. \qquad (\#)$$

Derive from the latter fact that

$$\operatorname{Prob}\left\{h^{T}\omega > h^{T}\mu + \delta\right\} \leq \exp\{-\frac{\delta^{2}}{3\left[\sum_{i}h_{i}^{2}\mu_{i} + \|h\|_{\infty}\delta\right]}\}.$$
 (##)

and conclude that

$$\operatorname{Prob}\left\{|h^{T}\omega - h^{T}\mu| > \delta\right\} \le 2\exp\{-\frac{\delta^{2}}{3[\sum_{i}h_{i}^{2}\mu_{i} + \|h\|_{\infty}\delta]}\}.$$
 (!)

3. Extract from (!) the following

Proposition 4.72. In the situation and under the assumptions of Exercise 4.71, let Opt be the optimal value, and H, Υ, Θ be a feasible solution to problem (P). Whenever $x \in \mathcal{X}$ and $\epsilon \in (0, 1)$, denoting by P_x the distribution of observations stemming from x (i.e., the distribution of random vector ω with independent entries $\omega_i \sim Poisson([Ax]_i))$, one has

$$\mathbf{E}\left\{\|Bx - \widehat{x}_{H}(\omega)\|\right\} \leq \overline{\Phi}(H) + 2\sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon])\operatorname{Tr}(\operatorname{Diag}(Ax)\Theta)} \leq \overline{\Phi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \Gamma_{\mathcal{X}}(\Theta)$$
(4.199)

and

$$\operatorname{Prob}_{\omega \sim P_{x}}\left\{ \|Bx - \widehat{x}_{H}(\omega)\| \leq \overline{\Phi}(H) + 2\sqrt{2}\sqrt{9\ln^{2}(2m/\epsilon)\operatorname{Tr}(\Theta)} + 3\ln(2m/\epsilon)\operatorname{Tr}(\operatorname{Diag}\{Ax\}\Theta)\sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon])} \right\} \geq 1 - \epsilon.$$
(4.200)

Note that in the case of $[Ax]_i \geq 1$ for all $x \in \mathcal{X}$ and all i we have $\operatorname{Tr}(\Theta) \leq \operatorname{Tr}(\operatorname{Diag}\{Ax\}\Theta)$, so that in this case the P_x -probability of the event

$$\left\{\omega: \|Bx - \widehat{x}_H(\omega)\| \le \overline{\Phi}(H) + O(1)\ln(2m/\epsilon)\sqrt{\phi_{\mathcal{R}}(\lambda[\Upsilon])\Gamma_{\mathcal{X}}(\Theta)}\right\}$$

is at least $1 - \epsilon$.

Exercise 4.73. [rudimentary discrete stochastic optimization] Let us revisit the story of Exercise 4.70 (in the sequel, we refer to this story as to the *learning* one). From the computational viewpoint, a bottleneck in problem's setup is that already moderate, but not quite small, number n of items and cardinality m of offers lead

to huge cardinality M of the set where the observations take their values; for example, with m = 10 and n = 100, we get $M = 190, 413, 404, 020, 840 \approx 1.9 \cdot 10^{14}$; it is problematic even to write down an approximation to a distribution on the set of that huge cardinality, not speaking about inferring an approximation from observations.

The goal of the exercise to follow is to modify problem's setup to allow for a kind of tractable solution. Specifically, consider the situation as follows: we have at our disposal

1) A finite set of actions \mathcal{I} of (perhaps, huge) cardinality M, such that we can sample from the uniform distribution on the set In the learning story, the actions are offers – collections of a given number m

of elements selected from a given *n*-element set; when *m* and *n* are moderate numbers, it is easy to draw an order at random from the uniform distribution (how?), in spite of the fact that the cardinality $M = \binom{n}{m}$ of the set of offers can be astronomically large.

- 2) A ground set \mathcal{G} of moderate cardinality ν along with a profit vector $c \in \mathbf{R}^{\nu}$. In the learning story, the ground set is the set of possible outcomes of an offer made to a customer, that is, the serial number of the bought item, if any, bought by the customer, or 0, if no item is bought; thus, $\nu = n + 1$, and the ground set is $\mathcal{G} = \{0, 1, ..., n\}$. Now, we can specify $c_j, 1 \leq j \leq n$, as the seller's profit when item j is bought, and set $c_0 = 0$ (or, perhaps, to make c_0 the minus cost of making an offer, if any).
- 3) We assume that "in the nature" every action χ is associated with a probability distribution p^{χ} on the ground set supported on known to us given χ subset G_{χ} of \mathcal{G} . We do not know what exactly p^{χ} is, but can sample from this distribution. Specifically, when making an observation, we select at our will $\chi \in \mathcal{I}$ and observe a realization of random variable $\omega \sim p^{\chi}$.

In the learning story, p^{χ} is the distribution of outcomes of offer χ induced by the unknown preference distribution p of customers, and $G_{\chi} = \{0, \chi_1, ..., \chi_m\}$, where $\chi = \{\chi_1, ..., \chi_m\}$ is the offer.

Now, in the learning story we were interested to recover the distribution of pairs (i, χ) $(i \in \{0\} \cup \chi$ is an outcome of offer χ) induced by the uniform distribution on the set of offers χ and the unknown preference distribution of customers; as we have already mentioned, when M is huge, this goal cannot be achieved – we cannot even write down a candidate solution! Instead, let us look at a different goal: *identifying the action* $\chi \in \mathcal{I}$ resulting in the largest, over $\chi \in \mathcal{I}$, expected profit

$$\pi_{\chi} = \sum_{j \in \mathcal{G}_{\chi}} c_j p_j^{\chi}.$$

Note that $\pi_{\chi}, \chi \in \mathcal{I}$, are well defined in our general setup 1) – 3), not only in the learning story. The problem we are interested in now is exactly the one of maximizing π_{χ} over $\chi \in \mathcal{I}$ via K observations $\omega^{K} = (\omega_{1}, ..., \omega_{K})$, with ω_{k} generated as follows: we select action χ_{k} (which should depend solely on the observations $\omega_{1}, ..., \omega_{k-1}$), and then "the nature" shows us ω_{k} drawn at random from the distribution $p^{\chi_{k}}$.

Literally speaking, the problem we have just posed still is intractable. Indeed, with our setup, there is not enough structure to relate to each other distributions

 p^{χ} associate with different actions χ ; consequently, seemingly the only way to select the best action is to implement *all of them* one by one, observing the outcomes of a particular action long enough to be able to estimate π_{χ} reliably; of course, such a brute force optimization is completely impossible when M is huge. What we intend to do is to relax this problem, namely, as follows: instead of looking for the action with the best possible expected profit, let us look for an action with the expected profit belonging to the top δ -fraction of the profits $\{\pi_{\chi} : \chi \in \mathcal{I}\}$. Specifically, let us select "small, but not very small" threshold $\delta \in (0, 1)$, like $\delta = 0.1$ or $\delta = 0.01$, and let $\pi_*(\delta)$ be the $(1 - \delta)$ -quantile of the set $\{\pi_{\chi} : \chi \in \mathcal{I}\}$, that is, the largest *s* such the cardinality of the set $\chi \in \mathcal{I} : \pi_{\chi} \geq s\}$ is at least δM . In other words, setting $M_{\delta} = \text{Floor}((1 - \delta)M)$ and arranging $\pi_{\chi}, \chi \in \mathcal{I}$, in a non-descending order, $\pi_*(\delta)$ is the M_{δ} -th element in this arrangement. The "relaxed" goal we are aiming at is

(!) Given $\delta \in (0,1)$, we need to identify reliably an action $\chi \in \mathcal{I}$ such that $\pi_{\chi} \geq \pi_*(\delta)$.

Achieving (!) reduces to the purely statistical problem of estimating π_{χ} for χ 's belonging to a set of *moderate*, provided δ is not very small, cardinality due to the following immediate

Observation: Let an "optimization threshold" $\delta \in (0,1)$ and a "reliability threshold" $\epsilon \in (0,1)$ be given, and let

$$J = J(\epsilon, \delta) = Ceil\left(\frac{\ln(1/\epsilon)}{\ln(1/(1-\delta))}\right)$$

Let, next, \mathcal{J} be a random subset of \mathcal{I} comprised of the samples $\chi^1, ..., \chi^J$ drawn, independently of each other, from the uniform distribution on \mathcal{I} , so that the cardinality K of \mathcal{J} is at most J, and let

$$\pi^*[\mathcal{J}] = \max_{\chi \in \mathcal{J}} \pi_{\chi}.$$

Then

$$\operatorname{Prob}\{\pi^*[\mathcal{J}] \ge \pi_*(\delta)\} \ge 1 - \epsilon.$$

Consequently, a maximizer $\chi_* = \chi_*[\mathcal{J}]$ of π_{χ} over $\chi \in \mathcal{J}$, up to probability of bad sampling $\leq \epsilon$, satisfies the relation $\pi_{\chi_*} \geq \pi_*(\delta)$.

This is your first task:

1. Justify Observation.

Observation suggests the following approach to achieving (!): given optimality tolerance δ and reliability tolerance ϵ , we generate a random subset $\mathcal{J} = \{\chi_1, ..., \chi_K\}$ of \mathcal{I} as explained in Observation, and implement actions $\chi \in \mathcal{J}$, several times each, in order to estimate the expected profits π_{χ} for all $\chi \in \mathcal{J}$, and then select the "seemingly best," as suggested by the estimated profits, action.

We are about to consider two implementations of the just outlined strategy.

4.73.A Single-stage estimation. Given \mathcal{J} and the total number N_{tot} of observations, we distribute the "observation resource" N_{tot} equally between the K actions from \mathcal{J} , thus arriving at

$$N = \operatorname{Floor}(N_{\text{tot}}/K).$$

observations per action. Then we implement one by one the K actions from \mathcal{I} , N times each, and use these observations to build confidence intervals Δ_{χ} , $\chi \in \mathcal{J}$, for the respective expected profits π_{χ} . When building these intervals, we set the underlying confidence levels to be equal to $1 - \hat{\epsilon}/K$, where $\hat{\epsilon} \in (0, 1)$ is the additional to ϵ "reliability tolerance" we use; with these reliability levels we have

$$\operatorname{Prob}\{\exists \chi \in \mathcal{J} : \pi_{\chi} \notin \Delta_{\chi}\} \leq \widehat{\epsilon}.$$

We can now select, as a candidate to the role of the best action from \mathcal{J} , the action $\hat{\chi}$ corresponding to the largest, over $\chi \in \mathcal{J}$, midpoint of the confidence interval Δ_{χ} . The outlined implementation gives rise to two questions:

The outlined implementation gives rise to two questions:

- A. How good is the resulting action $\hat{\chi}$?
- B. How to build the intervals Δ_{χ} ?

As far as A is concerned, your task is as follows:

2.1. Let σ_{χ} be the length of Δ_{χ} . Prove that

$$\operatorname{Prob}\{\pi_{\widehat{\chi}} < \pi^*[\mathcal{J}] - [\sigma_{\widehat{\chi}} + \max_{\chi \in \mathcal{J}} \sigma_{\chi}]\} \le \widehat{\epsilon}, \tag{\#}$$

where Prob is taken w.r.t. the conditional, \mathcal{J} given, probability distribution of our observations.

As about question B, it reduces to the question of how to estimate a linear function $h^T p$ of unknown *d*-dimensional probabilistic vector p via stationary Krepeated observation $\xi^K = (\xi_1, ..., \xi_K)$ drawn from p. This problem was considered in Section **3.3.3**, and we refer to this Section for the terminology we use now. Note that in our present situation, we need to solve the problem for several h's and d's (specifically, for h's obtained by restricting the profit vector c on the set \mathcal{G}_{χ} of indexes, with χ running through \mathcal{J}). Now, the constructions from Section **3.3.3** yield both a "presumably good" linear estimate of $h^T p$ and an upper bound on the risk of the estimate. To reduce the computational burden, let us restrict ourselves with the simplest estimate $h^T \hat{\xi}[\xi^K]$, where

$$\widehat{\xi}[\xi^K] = \frac{1}{K} \sum_{k=1}^K \xi_k$$

(as always, we encode realizations ξ_k of discrete random variable taking values in d-element set by the standard basic orths in \mathbf{R}^d).

The related exercise is as follows:

2.2 Prove the following version of Corollary 3.9:

Corollary 4.74. Let $h \in \mathbf{R}^d$, positive integer K and tolerance $\bar{\epsilon} \in (0,1)$ be given, and let $\xi_1, ..., \xi_K$ be drawn independently of each other from a discrete probability distribution $p \in \mathbf{\Delta}_d$. For $g \in \mathbf{R}^d$, let

$$\operatorname{Opt}_{\bar{\epsilon},K}(g) = \inf_{\beta>0} \left\{ \frac{\beta}{K} \ln(2/\bar{\epsilon}) + \max_{p \in \mathbf{\Delta}_d} \left[\beta \ln\left(\sum_i p_i \exp\{g_i/\beta\}\right) - g^T p \right] \right\}.$$

Next, for a sequence $\xi^K = (\xi_1, ..., \xi_K)$ of basic orths in \mathbf{R}^d , let

$$\Delta[\xi^K] = \left[h^T \widehat{\xi}[\xi^K] - \operatorname{Opt}_{\overline{\epsilon},K}(-h), h^T \widehat{\xi}[\xi^K] + \operatorname{Opt}_{\overline{\epsilon},K}(h) \right].$$

Then for every $p \in \Delta_d$ it holds

$$\operatorname{Prob}_{\xi^{K} \sim p \times \ldots \times p} \left\{ \xi^{K} : h^{T} p \notin \Delta[\xi^{K}] \right\} \leq \bar{\epsilon}.$$

4.73.B Multi-stage estimation. With this strategy we, given \mathcal{J} and additional to ϵ reliability tolerance $\hat{\epsilon} \in (0, 1)$, select somehow

- a (positive integer) number of stages L,
- positive integers N_{ini} and grow factor γ ,

We split $\hat{\epsilon}$ into L parts, that is, find positive reals $\hat{\epsilon}_{\ell}$, $1 \leq \ell \leq L$ such that

$$\widehat{\epsilon} = \sum_{\ell=1}^{L} \widehat{\epsilon}_{\ell}$$

and run one by one L stages as follows:

- At a beginning of stage ℓ , we have at our disposal a nonempty set $\mathcal{J}_{\ell} \subset \mathcal{J}$ comprised of K_{ℓ} distinct from each other elements $\chi_{\ell j}$, $1 \leq j \leq K_{\ell}$, with $\mathcal{J}_1 = \mathcal{J}$.
- We set $N_{\ell} = \gamma^{\ell-1} N_{\text{ini}}$ and implement, one by one, each one of $K_{\ell-1}$ actions from the set $\mathcal{J}_{\ell-1}$, collecting N_{ℓ} observations per action, and use these observations to build confidence intervals $\Delta_{\chi_{\ell j}}^{\ell}$ for the quantities $\pi_{\chi_{\ell j}}$, $1 \leq j \leq K_{\ell}$, the confidence levels being $1 - \hat{\epsilon}_{\ell}/K_{\ell}$. These intervals are built exactly in the same fashion as in the single-stage procedure, with N_{ℓ} playing the role of N.
- When $\ell < L$, we build $\mathcal{J}_{\ell+1}$ as follows. Let us say that $\chi_{\ell j}$ is *dominated*, if there exists a confidence interval $\Delta_{\chi_{\ell j'}}^{\ell}$ which is strictly to the right of the right endpoint of $\Delta_{\chi_{\ell j}}^{\ell}$. There clearly exist non-dominated $\chi_{\ell j}$'s (e.g., the one with the largest, among all confidence intervals built at the stage, right endpoint); the set $\mathcal{J}_{\ell+1}$ is comprised of all non-dominated $\chi_{\ell j}$'s. After $\mathcal{J}_{\ell+1}$ is built, we pass to the stage $\ell + 1$.
- When $\ell = L$, we select among all intervals $\Delta_{\chi_{Lj}}^L$ one with the largest midpoint, and output the corresponding $\chi_{Lj} =: \hat{\chi}$.

Now goes the exercise:

3. Prove that the resulting $\hat{\chi}$ meets the same quality guarantees as the output of the single-stage procedure, specifically

$$\operatorname{Prob}\{\pi_{\widehat{\chi}} < \pi^*[\mathcal{J}] - [\sigma_{\widehat{\chi}} + \max_{j \le K_L} \sigma_{\chi_{L_j}}]\} \le \widehat{\epsilon}, \qquad (\#\#)$$

where $\sigma_{\chi_{\ell_j}}$ is the width of the confidence interval $\Delta^L_{\chi_{L_j}}$, and Prob is taken w.r.t. the conditional, \mathcal{J} given, probability distribution of our observations.

Remark. The rationale underlying the multi-stage procedure is quite transparent: we hope that a significant part of actions (those "heavily bad") will be eliminated at early stages taking (since N_{ℓ} grow with ℓ) small part of the total number $N_{\text{tot}} =$

 $\sum_{\ell} K_{\ell} N_{\ell}$ of observations. If it indeed will be the case, then the "major part" of the total number of observations will be spent on a "promising part" of actions, thus improving the quality of the results as compared to the case when our "observation resource" N_{tot} is equally distributed between all actions form \mathcal{J} , good and bad alike.

Your final task is as follows:

4. Implement single- and multi-stage procedures and run simulations in order to get an impression what is better.

Here is the recommended setup (mimicking the learning problem):

- the action set is comprised of all (m = 10)-element subsets of (n = 100)-element set;
- the ground set is $\{0, 1, ..., n = 100\}$, and $c_i = i, 0 \le i \le n$;
- the support \mathcal{G}_{χ} of the probability distribution p^{χ} , $\chi = \{\chi_1, \chi_2, ..., \chi_m\}$ with distinct from each other elements $\chi_j \in \{1, ..., n\}$, is $\{0, \chi_1, ..., \chi_m\}$;
- $\delta = 0.05, \ \epsilon = \widehat{\epsilon} = 0.01.$

When simulating the above procedures, you need to associate with every action $\chi \in \mathcal{J}$ a probability distribution supported on \mathcal{G}_{χ} ; the simplest way is to select the required distributions at random.

4.9.6 Numerical lower-bounding minimax risk

Exercise 4.75.[†] [numerical lower bounding minimax risk]

4.75.A. Motivation. From the theoretical viewpoint, the results on near-optimality of presumably good linear estimates stated in Propositions 4.5, 4.16 seem to be pretty strong and general. This being said, for a practically oriented user the "nonoptimality factors" arising in these propositions can be too large to make practical sense. This practical drawback of our theoretical results is not too crucial – what matters in applications, is whether the risk of a proposed estimate is appropriate for the application in question, and not by how much it could be improved were we smart enough to build the "ideal" estimate; results of the latter type from practical viewpoint offer no more than some "moral support." Nevertheless, the "moral support" has its value, and it makes sense to strengthen it by improving the lower risk bounds as compared to those underlying Propositions 4.5, 4.16. In this respect, an appealing idea is to pass from lower risk bounds yielded by theoretical considerations to *computation-based* ones. The goal of this exercise is to develop some methodology yielding computation-based lower risk bounds. We start with the main ingredient of this methodology – the classical *Cramer-Rao* bound.

4.75.B. Cramer-Rao bound. Consider the situation as follows: we are given

- an observation space Ω equipped with reference measure Π , basic examples being (A) $\omega = \mathbf{R}^m$ with Lebesgue measure Π , and (B) (finite our countable) discrete set Ω with counting measure Π ;
- a convex compact set $\Theta \subset \mathbf{R}^k$ and a family $\Pi = \{p(\omega, \theta) : \theta \in \Theta\}$ of probability densities, taken w.r.t. Π .

Our goal is, given an observation $\omega \sim p(\cdot, \theta)$ stemming from unknown θ known to belong to Θ , to recover θ . We quantify the risk of a candidate estimate $\hat{\theta}$ as

$$\operatorname{Risk}[\widehat{\theta}|\Theta] = \sup_{\theta \in \Theta} \left(\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\widehat{\theta}(\omega) - \theta\|_2^2 \right\} \right)^{1/2}, \quad (4.201)$$

and define the "ideal" minimax risk as

$$\operatorname{Risk}_{\operatorname{opt}} = \inf_{\widehat{\theta}} \operatorname{Risk}[\widehat{\theta}],$$

the infimum being taken w.r.t. all estimates, or, which is the same, all bounded estimates (indeed, passing from a candidate estimate $\hat{\theta}$ to the projected estimate $\hat{\theta}_{\Theta}(\omega) = \operatorname{argmin}_{\theta \in \Theta} \|\hat{\theta}(\omega - \theta)\|_2$ we can only reduce the risk of an estimate.

The classical \overline{C} ramer-Rao inequality, which we intend to use, is certain relation between the covariance matrix of a bounded estimate and its bias; this relation is valid under mild regularity assumptions on the family Π , specifically, as follows:

- 1) $p(\omega, \theta) > 0$ for all $\omega \in \Omega, \theta \in U$, and $p(\omega, \theta)$ is differentiable in θ , the with $\nabla_{\theta} p(\omega, \theta)$ continuous in $\theta \in \Theta$;
- 2) The Fisher Information matrix

$$\mathcal{I}(\theta) = \int_{\Omega} \frac{\nabla_{\theta} p(\omega, \theta) [\nabla_{\theta} p(\omega, \theta)]^T}{p(\omega, \theta)} \Pi(d\omega)$$

is well defined for all $\theta \in \Theta$;

3) There exists function $M(\omega) \ge 0$ such that $\int_{\Omega} M(\omega) \Pi(d\omega) < \infty$ and

$$\|\nabla_{\theta} p(\omega, \theta)\|_2 \le M(\omega) \ \forall \omega \in \Omega, \theta \in \Theta.$$

The derivation of the Cramer-Rao bound is as follows. Let $\hat{\theta}(\omega)$ be a bounded estimate, and let

$$\phi(\theta) = [\phi_1(\theta); ...; \phi_k(\theta)] = \int_{\Omega} \widehat{\theta}(\omega) p(\omega, \theta) \Pi(d\omega)$$

be the expected value of the estimate. By item 3, $\phi(\theta)$ is differentiable on Θ , with the Jacobian $\phi'(\theta) = \left[\frac{\partial \phi_i(\theta)}{\partial \theta_j}\right]_{i,j \le k}$ given by

$$\phi'(\theta)h = \int_{\Omega} \widehat{\theta}(\omega)h^T \nabla_{\theta} p(\omega, \theta) \Pi(d\omega), \ h \in \mathbf{R}^k.$$

Besides this, recalling that $\int_{\Omega} p(\omega, \theta) \Pi(d\omega) \equiv 1$ and invoking item 3, we have $\int_{\Omega} h^T \nabla_{\theta} p(\omega, \theta) \Pi(d\omega) = 0$, whence, in view of the previous equality,

$$\phi'(\theta)h = \int_{\Omega} [\widehat{\theta}(\omega) - \phi(\theta)]h^T \nabla_{\theta} p(\omega, \theta) \Pi(d\omega), \ h \in \mathbf{R}^k.$$

Therefore for all $g, h \in \mathbf{R}^k$ we have

$$\begin{split} [g^{T}\phi'(\theta)h]^{2} &= \begin{bmatrix} \int_{\omega} [g^{T}(\widehat{\theta} - \phi(\theta)][h^{T}\nabla_{\theta}p(\omega,\theta)/p(\omega,\theta)]p(\omega,\theta)\Pi(d\omega)]^{2} \\ &\leq \begin{bmatrix} \int_{\Omega} g^{T}[\widehat{\theta} - \phi(\theta)][\widehat{\theta} - \phi(\theta)]^{T}gp(\omega,\theta)\Pi(d\omega)] \\ &\times \left[\int_{\Omega} [h^{T}\nabla_{\theta}p(\omega,\theta)/p(\omega,\theta)]^{2}p(\omega,\theta)\Pi(d\omega) \right] \\ &\quad \left[\text{Cauchy's Inequality} \right] \\ &= \begin{bmatrix} g^{T}\text{Cov}_{\widehat{\theta}}(\theta)g \end{bmatrix} \begin{bmatrix} h^{T}\mathcal{I}(\theta)h \end{bmatrix}, \end{split}$$

where $\operatorname{Cov}_{\widehat{\theta}}(\theta)$ is the covariance matrix $\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ [\widehat{\theta}(\omega) - \phi(\theta)] [\widehat{\theta}(\omega) - \phi(\theta)]^T \right\}$ of $\widehat{\theta}(\omega)$ induced by $\omega \sim p(\cdot,\theta)$. We have arrived at the inequality

$$\left[g^T \operatorname{Cov}_{\widehat{\theta}}(\theta)g\right] \left[h^T \mathcal{I}(\theta)h\right] \ge \left[g^T \phi'(\theta)h\right]^2 \,\,\forall (g,h \in \mathbf{R}^k, \theta \in \Theta). \tag{*}$$

For $\theta \in \Theta$ fixed, let \mathcal{J} be a positive definite matrix such that $\mathcal{J} \succeq \mathcal{I}(\theta)$, whence by (*) it holds

$$\left[g^T \operatorname{Cov}_{\widehat{\theta}}(\theta)g\right] \left[h^T \mathcal{J}h\right] \ge \left[g^T \phi'(\theta)h\right]^2 \,\forall (g,h \in \mathbf{R}^k).$$
(**)

For g fixed, the maximum of the right hand side quantity in (**) over h satisfying $h^T \mathcal{J}h \leq 1$ is $g^T \phi'(\theta) \mathcal{J}^{-1}[\phi'(\theta)^T g$, and we arrive at the *Cramer-Rao inequality*

$$\forall (\theta \in \Theta, \mathcal{J} \succeq \mathcal{I}(\theta), \mathcal{J} \succ 0) : \operatorname{Cov}_{\widehat{\theta}}(\theta) \succeq \phi'(\theta) \mathcal{J}^{-1}[\phi'(\theta]^{T} \\ \left[\operatorname{Cov}_{\widehat{\theta}}(\theta) = \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ [\widehat{\theta} - \phi(\theta)] [\widehat{\theta} - \phi(\theta)]^{T} \right\}, \ \phi(\theta) = \mathbf{E}_{\omega \sim p(\cdot),\theta} \left\{ \widehat{\theta}(\omega) \right\} \right]$$
(CR)

which holds true for every bounded estimate $\widehat{\theta}(\cdot)$. Note also that for every $\theta \in \Theta$ and every bounded estimate x we have

$$\begin{aligned} \operatorname{Risk}^{2}[\widehat{\theta}] &\geq \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\widehat{\theta}(\omega) - \theta\|_{2}^{2} \right\} = \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|[\widehat{\theta}(\omega) - \phi(\theta)] + [\phi(\theta) - \theta]\|_{2}^{2} \right\} \\ &= \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\widehat{\theta}(\omega) - \phi(\theta)\|_{2}^{2} \right\} - 2 \underbrace{\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left[[\widehat{\theta}(\omega) - \phi(\theta)]^{T} [\phi(\theta) - \theta]] \right\}}_{=0} \\ &+ \|\phi(\theta) - \theta\|_{2}^{2} \end{aligned}$$

whence, in view of (CR), for every bounded estimate $\hat{\theta}$ it holds

$$\forall (\mathcal{J} \succ 0 : \mathcal{J} \succeq \mathcal{I}(\theta) \,\forall \theta \in \Theta) : \operatorname{Risk}^{2}[\widehat{\theta}] \ge \sup_{\theta \in \Theta} \left[\operatorname{Tr}(\phi'(\theta) \mathcal{J}^{-1}[\phi'(\theta)]^{T}) + \|\phi(\theta) - \theta\|_{2}^{2} \right]$$

$$\left[\phi(\theta) = \mathbf{E}_{\omega \sim p(\cdot, \theta)} \{\widehat{\theta}(\omega)\} \right]$$

$$(4.202)$$

The fact that we were speaking about estimating "the entire" θ rather than a given vector-valued function $f(\theta) : \Theta \to \mathbf{R}^{\nu}$ plays no special role, and in fact the Cramer-Rao inequality admits the following modification (yielded by a reasoning completely similar to the one we just have carried out):

Proposition 4.76. In the situation described in the beginning of item **4.75.B** and under assumptions 1 – 3) of this item, let $f(\cdot) : \Theta \to \mathbf{R}^{\nu}$ be a bounded Borel function, and let $\hat{f}(\omega)$ be a bounded estimate of $f(\omega)$ via observation $\omega \sim p(\cdot, \theta)$.

401

Then, setting

$$\phi(\theta) = \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \widehat{f}(\theta) \right\}, \operatorname{Cov}_{\widehat{f}}(\theta) = \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ [\widehat{f}(\omega) - \phi(\theta)] [\widehat{f}(\omega) - \phi(\theta)]^T \right\} \\ [\theta \in \Theta]$$

one has

$$\forall (\theta \in \Theta, \mathcal{J} \succeq \mathcal{I}(\theta), \mathcal{J} \succ 0) : \operatorname{Cov}_{\widehat{f}}(\theta) \succeq \phi'(\theta) \mathcal{J}^{-1}[\phi'(\theta)]^T.$$

As a result, setting

$$\operatorname{Risk}[\widehat{f}] = \sup_{\theta \in \Theta} \left[\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\widehat{f}(\omega) - f(\theta)\|_2^2 \right\} \right]^{1/2},$$

it holds

$$\begin{aligned} \forall (\mathcal{J} \succ 0 : \mathcal{J} \succeq \mathcal{I}(\theta) \, \forall \theta \in \Theta) : \\ \operatorname{Risk}^2[\widehat{f}] \geq \sup_{\theta \in \Theta} \left[\operatorname{Tr}(\phi'(\theta) \mathcal{J}^{-1}[\phi'(\theta)]^T) + \|\phi(\theta) - f(\theta)\|_2^2 \right] \end{aligned}$$

Now goes the first part of the exercise:

1. Derive from (4.202) the following

Proposition 4.77. In the situation of item 4.75.B, let

- $\Theta \subset \mathbf{R}^k$ be $\|\cdot\|_2$ -ball of radius r > 0,
- the family \mathcal{P} be such that $\mathcal{I}(\theta) \preceq \mathcal{J}$ for some $\mathcal{J} \succ 0$ and all $\theta \in \Theta$.

Then the minimax optimal risk satisfies the bound

$$\operatorname{Risk}_{\operatorname{opt}} \ge \frac{rk}{r\sqrt{\operatorname{Tr}(\mathcal{J})} + k}.$$
(4.203)

In particular, when $\mathcal{J} = \alpha^{-1} I_k$, we have

$$\operatorname{Risk}_{\operatorname{opt}} \ge \frac{r\sqrt{\alpha k}}{r + \sqrt{\alpha k}}.$$
(4.204)

<u>Hint.</u> Assuming w.l.o.g. that Θ is centered at the origin, and given a bounded estimate $\hat{\theta}$ with risk \Re , let $\phi(\theta)$ be associated with the estimate via (4.202). Select $\gamma \in (0, 1)$ and consider two cases: (a): there exists $\theta \in \partial \Theta$ such that $\|\phi(\theta) - \theta\|_2 > \gamma r$, and (b): $\|\phi(\theta) - \theta\|_2 \le \gamma r$ for all $\theta \in \partial \Theta$. In the case of (a), lower-bound \Re by $\max_{\theta \in \Theta} \|\phi(\theta) - \theta\|_2$, see (4.202). In the case of(b), lower-bound \Re^2 by $\max_{\theta \in \Theta} \operatorname{Tr}(\phi'(\theta)\mathcal{J}^{-1}[\phi'(\theta)]^T)$, see (4.202), and use Divergence theorem to lower-bound the latter quantity in terms of the flux of the vector field $\phi(\cdot)$ over $\partial \Theta$.

When implementing the above strategy, you could find useful the following fact (prove it!)

Lemma 4.78. Let Φ be an $n \times n$ matrix, and \mathcal{J} be a positive semidefinite $n \times n$ matrix. Then

$$\operatorname{Tr}(\Phi \mathcal{J}^{-1} \Phi^T) \ge \operatorname{Tr}^2(\Phi) / \operatorname{Tr}(\mathcal{J}).$$

4.75.C. Application to signal recovery. Proposition 4.77 allows to build computation-based lower risk bounds in the signal recovery problem considered in Section 4.2, specifically, the problem where one wants to recover the linear image Bx of unknown signal x known to belong to a given ellitope

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : x^T S_\ell x \le t_\ell, \ell \le L \}$$

(with our usual restriction on S_{ℓ} and \mathcal{T}) via observation

$$\omega = Ax + \sigma\xi, \ \xi \sim \mathcal{N}(0, I_m),$$

and the risk of a candidate estimate, same as in Section 4.2, is defined according to $(4.201)^{83}$. It is convenient to assume that the matrix *B* (which in our general setup can be an arbitrary $\nu \times n$ matrix) is a *nonsingular* $n \times n$ matrix⁸⁴ Under this assumption, setting

$$\mathcal{Y} = B^{-1}\mathcal{X} = \{ y \in \mathbf{R}^n : \exists t \in \mathcal{T} : y^T [B^{-1}]^T S_\ell B^{-1} y \le t_\ell, \ell \le L \}$$

and $\bar{A} = AB^{-1}$, we lose nothing when replacing the sensing matrix A with \bar{A} and treating as our signal $y \in \mathcal{Y}$ rather than \mathcal{X} ; thus, we have reduced the situation to the one where A is replaced with \bar{A} , \mathcal{X} with \mathcal{Y} , and B with the unit matrix I_n . For the sake of simplicity, we assume from now on that A (and therefore \bar{A}) is with trivial kernel. Finally, let $\tilde{S}_{\ell} \succeq S_{\ell}$ be close to S_k positive definite matrices, e.g., $\tilde{S}_{\ell} = S_{\ell} + 10^{-100}I_n$; setting $\bar{S}_{\ell} = [B^{-1}]^T \tilde{S}_{\ell} B^{-1}$ and

$$\bar{\mathcal{Y}} = \{ y \in \mathbf{R}^n : \exists t \in \mathcal{T} : y^T \bar{S}_\ell y \le t_\ell, \ell \le L \},\$$

observe that $\bar{S}_{\ell} \succ 0$ and $\bar{\mathcal{Y}} \subset \mathcal{Y}$; this, any lower bound on the $\|\cdot\|_2$ -risk of recovery $y \in \bar{\mathcal{Y}}$ via observation $\omega = AB^{-1}y + \sigma\xi$, $\xi \sim \mathcal{N}(0, I_m)$, automatically is a lower bound on the minimax risk Risk_{opt} corresponding to our original problem of interest.

Now assume that we can point out a k-dimensional linear subspace E in \mathbb{R}^n and positive reals r, γ such that

(i) the centered at the origin $\|\cdot\|_2$ -ball $\Theta = \{\theta \in E : \|\theta\|_2 \le r\}$ is contained in \overline{Y} ; (ii) The restriction \overline{A}_E of \overline{A} onto E satisfies the relation

$$\operatorname{Tr}(\bar{A}_E^* \bar{A}_E) \le \gamma$$

 $(\bar{A}_E^*: \mathbf{R}^m \to E \text{ is the conjugate of the linear map } \bar{A}_E: E \to \mathbf{R}^m).$

Consider the auxiliary estimation problem obtained from the (reformulated) problem of interest by replacing the signal set $\bar{\mathcal{Y}}$ with Θ . Since $\Theta \subset \bar{\mathcal{Y}}$, the minimax

⁸³In fact, the approach to be developed can be applied to signal recovery problems involving Discrete/Poisson observation schemes, different from $\|\cdot\|_2$ norms used to measure the recovery error, signal-dependent noises, etc.

⁸⁴This assumption is nonrestrictive. Indeed, when $B \in \mathbf{R}^{\nu \times n}$ with $\nu < n$, we can add to $B \ n - \nu$ zero rows, which keeps our estimation problem intact. When $\nu \ge n$, we can add to B a small perturbation to ensure Ker $B = \{0\}$, which, for small enough perturbation, again keeps our estimation problem basically intact. It remains to note that when Ker $B = \{0\}$. we can replace \mathbf{R}^{ν} with the image space of B, which again does not affect the estimation problem we are interested in.

risk in the auxiliary problem is a lower bound on the minimax risk $\operatorname{Risk}_{opt}$ we are interested in. On the other hand, the auxiliary problem is nothing but the problem of recovering parameter $\theta \in \Theta$ from observation $\omega \sim \mathcal{N}(\overline{A}\theta, \sigma^2 I)$, which is nothing but a special case of the problem considered in item 4.75.B; as is immediately seen, the Fisher Information matrix in this problem is independent of θ and is $\sigma^{-2}\overline{A}_{E}^{*}\overline{A}_{E}$:

$$e^T \mathcal{I}(\theta) e = \sigma^{-2} e^T \bar{A}_E^* \bar{A}_E e, \ e \in E.$$

Invoking Proposition 4.77, we arrive at the lower bound on the minimax risk in the auxiliary problem (and thus – in the problem of interest as well):

$$\operatorname{Risk}_{\operatorname{opt}} \ge \frac{r\sigma k}{r\sqrt{\gamma} + \sigma k}.$$
(4.205)

The resulting risk bound depends on r, k, γ and is the larger the smaller is γ and the larger are k and r.

Lower-bounding Risk_{opt}. In order to extract from the just outlined bounding scheme its best, we need a mechanism which allows to generate k-dimensional "disks" $\Theta \subset \bar{\mathcal{Y}}$ along with associated quantities r, γ . In order to design such a mechanism, it is convenient to represent k-dimensional linear subspaces of \mathbf{R}^n as the image spaces of orthogonal $n \times n$ projectors P of rank k. Such a projector Pgives rise to the contained in $\bar{\mathcal{Y}}$ disk Θ_P of the radius $r = r_P$, where r_P is the largest ρ such that the set $\{y \in \mathbf{R}^n : y^T P y \leq \rho^2\}$ is contained in $\bar{\mathcal{Y}}$ ("condition $\mathcal{C}(r)$ "), and we can equip the disk with γ satisfying (ii) if and only if

$$\operatorname{Tr}(P\bar{A}^T\bar{A}P) \leq \gamma,$$

or, which is the same (recall that P is orthogonal projector)

$$\operatorname{Tr}(\bar{A}P\bar{A}^T) \le \gamma \tag{4.206}$$

("condition $\mathcal{D}(\gamma)$ "). Now, when P is a nonzero orthogonal projector, the simplest sufficient condition for the validity of $\mathcal{C}(r)$ is the existence of $t \in \mathcal{T}$ such that

$$\forall (y \in \mathbf{R}^n, \ell \le L) : y^T P \bar{S}_{\ell} P y \le t_{\ell} r^{-2} y^T P y,$$

or, which is the same,

$$\exists s : r^2 s \in \mathcal{T} \& P\bar{S}_{\ell}P \preceq s_{\ell}P, \, \ell \leq L. \tag{4.207}$$

We are about to rewrite (4.206), (4.207) as a system of *linear* matrix inequalities. This is what you are supposed to do:

2.1. Prove the following simple fact:

Observation 4.79. Let Q be a positive definite and R be a nonzero positive semidefinite matrix, and s be a real. Then

$$RQR \preceq sR$$

if and only if

 $sQ^{-1} \succeq R.$

2.2. Extract from Observation the conclusion as follows. Let **T** be the conic hull of \mathcal{T} :

$$\mathbf{T} = cl\{[s;\tau] : \tau > 0, s/\tau \in \mathcal{T}\} = \{[s;\tau] : \tau > 0, s/\tau \in \mathcal{T}\} \cup \{0\}$$

Consider the system of constraints

$$s_{\ell} S_{\ell}^{-1} \succeq P, \ell \le L \& \operatorname{Tr}(APA^{T}) \le \gamma$$

$$P \text{ is orthogonal projector of rank } k \ge 1 \tag{#}$$

in variables $[s; \tau] \in \mathbf{T}$, k, γ and P. Every feasible solution to this system gives rise to k-dimensional Euclidean subspace $E \subset \mathbf{R}^n$ (the image space of P) such that the centered at the origin Euclidean ball Θ in E of radius

$$r = 1/\sqrt{\tau}$$

taken along with γ satisfy the conditions (i) - (ii). Consequently, this feasible solution yields the lower bound

$$\operatorname{Risk}_{\operatorname{opt}} \geq \psi_{\sigma,k}(\gamma,\tau) := \frac{\sigma k}{\sqrt{\gamma} + \sigma \sqrt{\tau} k}$$

on the minimax risk in the problem of interest.

An "ideal" way to utilize item 2.2 to lower-bound Risk_{opt} would be to look through k = 1, ..., n and for every k to maximize the lower risk bound $\psi_{\sigma,k}(\gamma, \tau)$ under constraints (#), thus arriving at the problem

$$\min_{[s;\tau],\gamma,P} \left\{ \frac{\sigma}{\psi_{\sigma,k}(\gamma,\tau)} = \sqrt{\gamma}/k + \sigma\sqrt{\tau} : \begin{array}{c} s_{\ell}\bar{S}_{\ell}^{-1} \succeq P, \ell \le L \& \operatorname{Tr}(\bar{A}P\bar{A}^{T}) \le \gamma \\ P \text{ is orthogonal projector of rank } k \end{array} \right\}_{\substack{(P_{k}) \in \mathcal{P}_{k}}}$$

This problem seems to be computationally intractable, since the constraints of (P_k) include the nonconvex restriction on P to be an orthogonal projector of rank k. A natural convex relaxation of this restriction is

$$0 \leq P \leq I_n, \operatorname{Tr}(P) = k.$$

The (minor) remaining difficulty is that the objective in (P) is nonconvex. Note, however, that to minimize $\sqrt{\gamma}/k + \sigma\sqrt{\tau}$ is basically the same as to minimize the convex function $\gamma/k^2 + \sigma^2\tau$ which is a tight "proxy" of the squared objective of (P_k) . We arrive at convex "proxy" of (P_k) – the problem

$$\min_{[s;\tau],\gamma,P} \left\{ \gamma/k^2 + \sigma^2 \tau : \begin{array}{c} [s;\tau] \in \mathbf{T}, 0 \leq P \leq I_n, \operatorname{Tr}(P) = k \\ s_\ell \bar{S}_\ell^{-1} \succeq P, \ell \leq L, \operatorname{Tr}(\bar{A}P\bar{A}^T) \leq \gamma \end{array} \right\}$$
(P[k])

k = 1, ..., n. Problem (P[k]) clearly is solvable, and the *P*-component $P^{(k)}$ of its optimal solution gives rise to a bunch of orthogonal projectors $P_{\kappa}^{(k)}$, $\kappa = 1, ..., n$ obtained from $P^{(k)}$ by "rounding" – to get $P_{\kappa}^{(k)}$, we replace the κ leading eigenvalues of $P^{(k)}$ with ones, and the remaining eigenvalues – with zeros, while keeping the eigenvectors intact. We can now for every $\kappa = 1, ..., n$ fix the *P*-variable in (P_k) as $P_{\kappa}^{(k)}$ and solve the resulting problem in the remaining variables $[s; \tau]$ and γ , which is easy – with *P* fixed, the problem clearly reduces to the one of minimizing τ under

the convex constraints

$$s_{\ell}\bar{S}_{\ell}^{-1} \succeq P, \ell \leq L, \, [s;\tau] \in \mathbf{T}$$

on $[s;\tau]$. As a result, for every $k \in \{1, ..., n\}$, we get n lower bounds on Risk_{opt}, that is, total of n^2 lower risk bounds, of which we select the best – the largest.

Now goes the next part of the exercise:

- 3. Implement the outlined methodology numerically and compare the lower bound on the minimax risk with the upper risk bounds of presumably good linear estimates yielded by Proposition 4.4. Recommended setup:
 - Sizes: $m = n = \nu = 16$
 - A, B: $B = I_n$, $A = \text{Diag}\{a_1, ..., a_n\}$ with $a_i = i^{-\alpha}$ and α running through $\{0, 1, 2\}$;
 - $\mathcal{X} = \{x \in \mathbf{R}^n : x^T S_\ell x \le 1, \ell \le L\}$ (i.e., $\mathcal{T} = [0, 1]^L$) with randomly generated S_ℓ .

Range of L: {1,4,16}. For L in this range, you can generate S_{ℓ} , $\ell \leq L$, as $S_{\ell} = R_{\ell}R_{\ell}^{T}$ with $R_{\ell} = randn(n,p)$, where $p = \lfloor n/L \rfloor$.

• Range of σ : {1.0, 0.1, 0.01, 0.001, 0.0001}

4.75.D. More on Cramer-Rao risk bound. Let us fix $\mu \in (1, \infty)$ and a norm $\|\cdot\|$ on \mathbb{R}^k , and let $\|\cdot\|_*$ be the norm conjugate to $\|\cdot\|$, and $\mu_* = \frac{\mu}{\mu-1}$. Assume that we are in the situation of item 4.75.B and under assumptions 1) and 3) from this item; as about assumption 2) we now replace it with the assumption that the quantity

$$\mathcal{I}_{\|\cdot\|_{*},\mu_{*}}(\theta) := \left[\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\nabla_{\theta} p(\omega,\theta)\|_{*}^{\mu_{*}} \right\} \right]^{1/\mu_{*}}$$

is well defined and bounded on Θ ; in the sequel, we set

$$\mathcal{I}_{\|\cdot\|_*,\mu_*} = \sup_{\theta\in\Theta} \mathcal{I}_{\|\cdot\|_*,\mu_*}(\theta).$$

4. Prove the following variant of Cramer-Rao risk hound:

Proposition 4.80. In the situation described in the beginning of item 4.75.D, let $\Theta \subset \mathbf{R}^k$ be a $\|\cdot\|$ -ball of radius r. Then the minimax $\|\cdot\|$ -risk of recovering $\theta \in \Theta$ via observation $\omega \sim p(\cdot, \theta)$ can be lower-bounded as

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|}[\Theta] := \inf_{\widehat{\theta}(\cdot)} \sup_{\theta \in \Theta} \left[\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\widehat{\theta}(\omega) - \theta\|^{\mu} \right\} \right]^{1/\mu} \ge \frac{rk}{r\mathcal{I}_{\|\cdot\|_{*},\mu_{*}} + k}, \quad (4.208)$$
$$\mathcal{I}_{\|\cdot\|_{*},\mu_{*}} = \max_{\theta \in \Theta} \left[\mathcal{I}_{\|\cdot\|_{*},\mu_{*}}(\theta) := \left[\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\nabla_{\theta} \ln(p(\omega,\theta))\|_{*}^{\mu_{*}} \right\} \right]^{1/\mu_{*}} \right]$$

Example I: Gaussian case, estimating shift. Let $\mu = 2$, and let $p(\omega, \theta) =$

 $\mathcal{N}(A\theta, \sigma^2 I_m)$ with $A \in \mathbf{R}^{m \times k}$. Then

$$\begin{split} \nabla_{\theta} \ln(p(\omega,\theta)) &= \sigma^{-2} A^{T}(\omega - A\theta) \Rightarrow \\ \int \|\nabla_{\theta} \ln(p(\omega,\theta))\|_{*}^{2} p(\omega,\theta) d\omega &= \sigma^{-4} \int \|A^{T}(\omega - A\theta)\|_{*}^{2} p(\omega,\theta) d\omega \\ &= \sigma^{-4} \frac{1}{[\sqrt{2\pi}\sigma]^{m}} \int \|A^{T}\omega\|_{*}^{2} \exp\{-\frac{\omega^{T}\omega}{2\sigma^{2}}\} d\omega \\ &= \sigma^{-4} \frac{1}{[2\pi]^{m/2}} \int \|A^{T}\sigma\xi\|_{*}^{2} \exp\{-\xi^{T}\xi/2\} d\xi \\ &= \sigma^{-2} \frac{1}{[2\pi]^{m/2}} \int \|A^{T}\xi\|_{*}^{2} \exp\{-\xi^{T}\xi/2\} d\xi \end{split}$$

whence

$$\mathcal{I}_{\|\cdot\|_{*},2} = \sigma^{-1} \underbrace{\left[\mathbf{E}_{\xi \sim \mathcal{N}(0,I_m)} \left\{ \|A^T \xi\|_{*}^2 \right\} \right]^{1/2}}_{\gamma_{\|\cdot\|}(A)}.$$

Consequently, assuming Θ to be $\|\cdot\|$ -ball of radius r in \mathbf{R}^k , lower bound (4.208) becomes

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|}[\Theta] \ge \frac{rk}{r\mathcal{I}_{\|\cdot\|_*} + k} = \frac{rk}{r\sigma^{-1}\gamma_{\|\cdot\|}(A) + k} = \frac{r\sigma k}{r\gamma_{\|\cdot\|}(A) + \sigma k}.$$
(4.209)

The case of direct observations. Just to see how it works, consider the case m = k, $A = I_k$ of direct observations, and let $\Theta = \{\theta \in \mathbf{R}^k : ||\theta|| \le r\}$. Then

• We have $\gamma_{\|\cdot\|_1}(I_k) \leq O(1)\sqrt{\ln(n)}$, whence the $\|\cdot\|_1$ -risk bound is

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|_1}[\Theta] \ge O(1) \frac{r\sigma k}{r\sqrt{\ln(n)} + \sigma k}; \qquad [\Theta = \{\theta \in \mathbf{R}^k : \|\theta - a\|_1 \le r\}]$$

• We have $\gamma_{\|\cdot\|_2}(I_k) = \sqrt{k}$, whence the $\|\cdot\|_2$ -risk bound is

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|_{2}}[\Theta] \geq \frac{r\sigma\sqrt{k}}{r+\sigma\sqrt{k}}; \qquad \qquad [\Theta = \{\theta \in \mathbf{R}^{k} : \|\theta - a\|_{2} \leq r\}]$$

• We have $\gamma_{\|\cdot\|_{\infty}}(I_k) \leq O(1)k$, whence the $\|\cdot\|$ -risk bound is

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|_{2}}[\Theta] \ge O(1)\frac{r\sigma}{r+\sigma}. \qquad \qquad [\Theta = \{\theta \in \mathbf{R}^{k} : \|\theta - a\|_{\infty} \le r\}]$$

In fact, the above examples are basically covered by the following

Observation 4.81. Let $\|\cdot\|$ be a norm on \mathbf{R}^k , and let

$$\Theta = \{\theta \in \mathbf{R}^k : \|\theta\| \le r\}$$

Consider the problem of recovering signal $\theta \in \Theta$ via observation $\omega \sim \mathcal{N}(\theta, \sigma^2 I_k)$, let

$$\operatorname{Risk}_{\|\cdot\|}[\widehat{\theta}|\Theta] = \sup_{\theta \in \Theta} \left(\mathbf{E}_{\omega \sim \mathcal{N}(\theta, \sigma^2 I)} \left\{ \|\widehat{\theta}(\omega) - \theta\|^2 \right\} \right)^{1/4}$$

be the $\|\cdot\|$ -risk of an estimate $\widehat{\theta}(\cdot)$, and let

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|}[\Theta] = \inf_{\widehat{\theta}(\cdot)} \operatorname{Risk}_{\|\cdot\|}[\widehat{\theta}|\Theta]$$

407

be the associated minimax risk.

Assume that the norm $\|\cdot\|$ is absolute and symmetric w.r.t permutation of coordinates. Then

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|}[\Theta] \ge \frac{r\sigma k}{2\sqrt{\ln(ek)}r\alpha_* + \sigma k}, \quad \alpha_* = \|[1;...;1]\|_*.$$
(4.210)

Here is the concluding part of the exercise:

5. Prove Observation and compare the lower risk bound from Observation with the $\|\cdot\|$ -risk of the "plug-in" estimate $\widehat{\chi}(\omega) \equiv \omega$.

Example II: Gaussian case, estimating covariance. Let $\mu = 2$, let K be a positive integer, and let our observation ω be a collection of K i.i.d. samples $\omega_t \sim \mathcal{N}(0,\theta), 1 \leq t \leq K$, with unknown θ known to belong to a given convex compact subset Θ of the interior of the positive semidefinite cone \mathbf{S}^n_+ . Given $\omega_1, ..., \omega_K$, we want to recover θ in the Shatten norm $\|\cdot\|_{\mathrm{Sh},s}$ with $s \in [1,\infty]$. Our estimation problem is covered by the setup of Exercise 4.75 with \mathcal{P} comprised of the product probability densities $p(\omega, \theta) = \prod_{t=1}^{K} g(\omega_t, \theta), \ \theta \in \Theta$, where $g(\cdot, \theta)$ is the density of $\mathcal{N}(0, \theta)$. We have

$$\nabla_{\theta} \ln(p(\omega,\theta)) = \frac{1}{2} \sum_{t} \nabla_{\theta} \ln(g(\omega_{t},\theta)) = \frac{1}{2} \sum_{t} \left[\theta^{-1} \omega_{t} \omega_{t}^{T} \theta^{-1} - \theta^{-1} \right] \\
= \frac{1}{2} \theta^{-1/2} \left[\sum_{t} \left[[\theta^{-1/2} \omega_{t}] [\theta^{-1/2} \omega_{t}]^{T} - I_{n} \right] \right] \theta^{-1/2}$$
(4.211)

With some effort it can be proved that when

 $K \ge n$,

which we assume from now on, for independent across t random vectors $\xi_1, ..., \xi_K$ sampled from the standard Gaussian distribution $\mathcal{N}(0, I_n)$ for every $u \in [1, \infty]$ one has

$$\left[\mathbf{E}\left\{\|\sum_{t=1}^{K} [\xi_t \xi_t^T - I_n]\|_{\mathrm{Sh},u}^2\right\}\right]^{1/2} \le Cn^{\frac{1}{2} + \frac{1}{u}}\sqrt{K}$$
(4.212)

with appropriate absolute constant C. Consequently, for $\theta \in \Theta$ and all $u \in [1, \infty]$ we have

$$\begin{split} \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \| \nabla_{\theta} \ln(p(\omega,\theta)) \|_{\mathrm{Sh},u}^{2} \right\} \\ &= \frac{1}{4} \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \| \theta^{-1/2} \left[\sum_{t} \left[[\theta^{-1/2} \omega_{t}] [\theta^{-1/2} \omega_{t}]^{T} - I_{n} \right] \right] \theta^{-1/2} \|_{\mathrm{Sh},u}^{2} \right\} \\ &\quad [by (4.211)] \\ &= \frac{1}{4} \mathbf{E}_{\xi \sim p(\cdot,I_{n})} \left\{ \| \theta^{-1/2} \left[\sum_{t} \left[\xi_{t} \xi_{t}^{T} - I_{n} \right] \right] \theta^{-1/2} \|_{\mathrm{Sh},u}^{2} \right\} \quad [setting \ \theta^{-1/2} \omega_{t} = \xi_{t}] \\ &\leq \frac{1}{4} \| \theta^{-1/2} \|_{\mathrm{Sh},\infty}^{4} \mathbf{E}_{\xi \sim p(\cdot,I_{n})} \left\{ \| \sum_{t} \left[\xi_{t} \xi_{t}^{T} - I_{n} \right] \|_{\mathrm{Sh},u}^{2} \right\} \\ &\quad [since \ \| AB \|_{\mathrm{Sh},u} \| \leq \| A \|_{\mathrm{Sh},\infty} \| B \|_{\mathrm{Sh},u}] \\ &\leq \frac{1}{4} \| \theta^{-1/2} \|_{\mathrm{Sh},\infty}^{4} \left[Cn^{\frac{1}{2} + \frac{1}{u}} \sqrt{K} \right]^{2} \quad [by \ (4.212)] \end{split}$$

and we arrive at

$$\left[\mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\nabla_{\theta} \ln(p(\omega,\theta))\|_{\mathrm{Sh},u}^{2} \right\} \right]^{1/2} \leq \frac{C}{2} \|\theta^{-1}\|_{\mathrm{Sh},\infty} n^{\frac{1}{2} + \frac{1}{u}} \sqrt{K}.$$
(4.213)

408

LECTURE 4

Now assume that Θ is $\|\cdot\|_{\mathrm{Sh},s}$ -ball of radius r < 1 centered at I_n :

$$\Theta = \{ \theta \in \mathbf{S}^n : \| \theta - I_n \|_{\mathrm{Sh}, s} \le r \}.$$

$$(4.214)$$

In this case the estimation problem from Example II is the scope of Proposition 4.80, and the quantity $I_{\|\cdot\|_{*},2}$ as defined in (4.208) can be upper-bounded as follows:

$$I_{\|\cdot\|_{*},2} = \max_{\theta\in\Theta} \left[\mathbf{E}_{\omega\sim p(\cdot,\theta)} \left\{ \|\nabla_{\theta} \ln(p(\omega,\theta))\|_{\mathrm{Sh},s_{*}}^{2} \right\} \right]^{1/2} \\ \leq O(1)n^{\frac{1}{2}+\frac{1}{s_{*}}} \sqrt{K} \max_{\theta\in\Theta} \|\theta^{-1}\|_{\mathrm{Sh},\infty} \text{ [see (4.213)]} \\ \leq O(1)^{\frac{n^{\frac{1}{2}+\frac{1}{s_{*}}}\sqrt{K}}{1-r}}.$$

We can now use Proposition 4.80 to lower-bound the minimax $\|\cdot\|_{\mathrm{Sh},s}$ -risk, thus arriving at

$$\operatorname{Risk}_{\operatorname{opt},\|\cdot\|_{\operatorname{Sh},s}}[\Theta] \ge O(1) \frac{n(1-r)r}{\sqrt{K}n^{\frac{1}{2}-\frac{1}{s}}r + n(1-r)}$$
(4.215)

(note that we are in the case of $k = \dim \theta = \frac{n(n+1)}{2}$). Let us compare this lower risk bound with the $\|\cdot\|_{\mathrm{Sh},s}$ -risk of the "plug-in" estimate

$$\widehat{\theta}(\omega) = \frac{1}{K} \sum_{t=1}^{K} \omega_t \omega_t^T.$$

Assuming $\theta \in \Theta$, we have

$$\begin{aligned} \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|K[\widehat{\theta}(\omega) - \theta]\|_{\mathrm{Sh},s}^{2} \right\} &= \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\sum_{t} [\omega_{t}\omega_{t}^{T} - \theta]\|_{\mathrm{Sh},s}^{2} \right\} \\ &= \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\theta^{1/2} \left[\sum_{t} [[\theta^{-1/2}\omega_{t}][\theta^{-1/2}\omega_{t}]^{T} - I_{n}] \right] \theta^{1/2} \|_{\mathrm{Sh},s}^{2} \right\} \\ &= \mathbf{E}_{\xi \sim p(\cdot,I_{n})} \left\{ \|\theta^{1/2} \left[\sum_{t} [\xi_{t}\xi_{t}^{T} - I_{n}] \right] \theta^{1/2} \|_{\mathrm{Sh},s}^{2} \right\} \\ &\leq \|\theta^{1/2}\|_{\mathrm{Sh},\infty}^{4} \mathbf{E}_{\xi \sim p(\cdot,I_{n})} \left\{ \|\sum_{t} [\xi_{t}\xi_{t}^{T} - I_{n}] \|_{\mathrm{Sh},s}^{2} \right\} \\ &\leq \|\theta^{1/2}\|_{\mathrm{Sh},\infty}^{4} \left[Cn^{\frac{1}{2} + \frac{1}{s}} \sqrt{K} \right]^{2}, \text{ [see (4.212)]} \end{aligned}$$

and we arrive at

$$\operatorname{Risk}_{\|\cdot\|_{\operatorname{Sh},s}}[\widehat{\theta}|\Theta] \le O(1) \max_{\theta \in \Theta} \|\theta\|_{\operatorname{Sh},\infty} \frac{n^{\frac{1}{2} + \frac{1}{s}}}{\sqrt{K}}.$$
(4.216)

In the case of (4.214), the latter bound becomes

$$\operatorname{Risk}_{\|\cdot\|_{\operatorname{Sh},s}}[\widehat{\theta}|\Theta] \le O(1) \max_{\theta \in \Theta} \|\theta\|_{\operatorname{Sh},\infty} \frac{n^{\frac{1}{2} + \frac{1}{s}}}{\sqrt{K}}.$$
(4.217)

For the sake of simplicity, assume that r in (4.214) is 1/2 (what actually matters below is that $r \in (0,1)$ is bounded away from 0 and from 1). In this case the lower bound (4.215) on the minimax $\|\cdot\|_{\mathrm{Sh},s}$ -risk reads

$$\operatorname{Risk}_{\operatorname{opt}, \|\cdot\|_{\operatorname{Sh}, s}}[\Theta] \ge O(1) \min\left[\frac{n^{\frac{1}{2} + \frac{1}{s}}}{\sqrt{K}}, 1\right].$$

When K is "large:" $K \ge n^{1+\frac{2}{s}}$, this lower bound matches, within an absolute constant factor, the upper bound (4.217) on the risk of the plug-in estimate, so that the latter estimate is near-optimal. When $K < n^{1+\frac{2}{s}}$, the lower risk bound becomes O(1), so that here a nearly optimal estimate is the trivial estimate $\hat{\theta}(\omega) \equiv I_n$.

Exercise 4.82. † [follow-up to Exercise 4.75]

1. Prove the following version of Proposition 4.77:

Proposition 4.83. In the situation of item 4.75.B and under assumptions 1) – 3) from this item, let

• $\|\cdot\|$ be a norm on \mathbf{R}^k such that

$$\|\theta\|_2 \le \kappa \|\theta\| \ \forall \theta \in \mathbf{R}^k$$

- $\Theta \subset \mathbf{R}^k$ be $\|\cdot\|$ -ball of radius r > 0,
- the family \mathcal{P} be such that $\mathcal{I}(\theta) \preceq \mathcal{J}$ for some $\mathcal{J} \succ 0$ and all $\theta \in \Theta$.

Then the minimax optimal risk

$$Risk_{\text{opt},\|\cdot\|} = \inf_{\widehat{\theta}(\cdot)} \left(\sup_{\theta \in \Theta} \mathbf{E}_{\omega \sim p(\cdot,\theta)} \left\{ \|\theta - \widehat{\theta}(\omega)\|^2 \right\} \right)^{1/2}$$

of recovering parameter $\theta \in \Theta$ from observation $\omega \sim p(\cdot, \theta)$ in the norm $\|\cdot\|$ satisfies the bound

$$Risk_{\text{opt},\|\cdot\|} \ge \frac{rk}{r\kappa\sqrt{\text{Tr}(\mathcal{J})} + k}.$$
(4.218)

In particular, when $\mathcal{J} = \alpha^{-1}I_k$, we get

$$Risk_{\text{opt},\|\cdot\|} \ge \frac{r\sqrt{\alpha k}}{r\kappa + \sqrt{\alpha k}}.$$
 (4.219)

- 2. Apply Proposition 4.83 to get lower bounds on the minimax $\|\cdot\|$ -risk in the following estimation problems:
- 2.1. Given indirect observation $\omega = A\theta + \sigma\xi$, $\xi \sim \mathcal{N}(0, I_m)$ of unknown vector θ known to belong to $\Theta = \{\theta \in \mathbf{R}^k : \|\theta\|_p \leq r\}$ with given A, Ker $A = \{0\}$, $p \in [2, \infty], r > 0$, we want to recover θ in $\|\cdot\|_p$.
- 2.2. Given indirect observation $\omega = L\theta R + \sigma \xi$, where θ is unknown $\mu \times \nu$ matrix known to belong to the Shatten norm ball $\Theta \in \mathbf{R}^{\mu \times \nu} : \|\theta\|_{\mathrm{Sh},p} \leq r$, we want to recover θ in $\|\cdot\|_{\mathrm{Sh},p}$. Here $L \in \mathbf{R}^{m \times \mu}$, Ker $L = \{0\}$ and $R = \mathbf{R}^{\nu \times n}$, Ker $R^T = \{0\}$ are given matrices, $p \in [2, \infty]$, and ξ is random Gaussian $m \times n$ matrix (i.e., the entries in ξ are independent of each other $\mathcal{N}(0, 1)$ random variables).
- 2.3. Given K-repeated observation $\omega^{K} = (\omega_{1}, ..., \omega_{K})$ with i.i.d. components $\omega_{t} \sim \mathcal{N}(0, \theta), 1 \leq t \leq K$, with unknown $\theta \in \mathbf{S}^{n}$ known to belong to the matrix box $\Theta = \{\theta : \beta_{-}I_{n} \leq \theta \leq \beta_{+}I_{n}\}$ with given $0 < \beta_{-} < \beta_{+} < \infty$, we want to recover θ in the spectral norm.

4.9.7 Around S-Lemma

Exercise 4.84. Proposition 4.6 provides us with upper bound on the quality of semidefinite relaxation as applied to the problem of upper-bounding the maximum of a homogeneous quadratic form over an ellitope. Extend the construction to the case when an inhomogeneous quadratic form is maximized over a shifted ellitope, so that quantity to upper-bound is

$$Opt = \max_{x \in X} \left[f(x) := x^T A x + 2b^T x + c \right], \ X = \{ x : \exists (y, t \in \mathcal{T}) : x = Py + p, y^T S_k y \le t_k, 1 \le k \le K \}$$

with our standard assumptions on S_k and \mathcal{T} .

<u>Note:</u> X is centered at p, and a natural upper bound on Opt is

$$Opt \le f(p) + \widehat{Opt},$$

where \widehat{Opt} is an upper bound on the quantity

$$\overline{\text{Opt}} = \max_{x \in X} \left[f(x) - f(p) \right];$$

what you are interested to upper-bound, is the ratio $\widehat{Opt}/\overline{Opt}$.

S-Lemma is a classical result of extreme importance in Semidefinite Optimization. Basically, Lemma states that when the ellitope \mathcal{X} in Proposition 4.6 is just an ellipsoid, (4.20) can be strengthen to Opt = Opt_{*}. In fact, *S*-Lemma is even stronger:

Lemma 4.85. [S-Lemma] Consider two quadratic forms $f(x) = x^T A x + 2a^T x + \alpha$, $g(x) = x^T B x + 2b^T x + \beta$ such that $g(\bar{x}) < 0$ for some \bar{x} . Then the implication

$$g(x) \le 0 \Rightarrow f(x) \le 0$$

takes place if and only if for some $\lambda \ge 0$ it holds $f(x) \le \lambda g(x)$ for all x, or, which is the same, if and only if Linear Matrix Inequality

$$\begin{bmatrix} \lambda B - A & \lambda b - a \\ \hline \lambda b^T - a^T & \lambda \beta - \alpha \end{bmatrix} \succeq 0$$

in scalar variable λ has a nonnegative solution.

Proof of \mathcal{S} -Lemma can be found, e.g., in [11, Section 3.5.2]

The goal of subsequent exercises is to get "tight" tractable outer approximations of sets obtained from ellitopes by quadratic lifting. We fix an ellitope

$$X = \{x \in \mathbf{R}^n : \exists t \in \mathcal{T} : x^T S_k x \le t_k, 1 \le k \le K\}$$

$$(4.220)$$

where, as always, S_k are positive semidefinite matrices with positive definite sum, and \mathcal{T} is a computationally tractable convex compact subset in \mathbf{R}^k_+ such that $t \in \mathcal{T}$ implies $t' \in \mathcal{T}$ whenever $0 \leq t' \leq t$ and \mathcal{T} contains a positive vector.

Exercise 4.86. Let us associate with ellitope X given by (4.220) the sets

 $\mathcal{X} = \operatorname{Conv}\{xx^T : x \in X\}, \ \widehat{\mathcal{X}} = \{Y \in \mathbf{S}^n : Y \succeq 0, \exists t \in \mathcal{T} : \operatorname{Tr}(S_k Y) \le t_k, 1 \le k \le K\},\$

411

so that $\mathcal{X}, \hat{\mathcal{X}}$ are convex compact sets containing the origin, and $\hat{\mathcal{X}}$ is computationally tractable along with \mathcal{T} . Prove that

- 1. When K = 1, we have $\mathcal{X} = \widehat{\mathcal{X}}$;
- 2. We always have $\mathcal{X} \subset \widehat{\mathcal{X}} \subset 4\ln(5K)\mathcal{X}$.

Exercise 4.87. For
$$x \in \mathbf{R}^n$$
 let $Z(x) = [x; 1][x; 1]^T$, $Z^o[x] = \left[\begin{array}{c|c} xx^T & x \\ \hline x^T & \end{array} \right]$. Let
$$C = \left[\begin{array}{c|c} \\ \hline \end{array} \right],$$

and let us associate with ellitope X given by (4.220) the sets

$$\mathcal{X}^+ = \operatorname{Conv}\{Z^o[x] : x \in X\},\$$

$$\widehat{\mathcal{X}}^+ = \{Y = \left[\frac{U \mid u}{u^T \mid}\right] \in \mathbf{S}^{n+1} : Y + C \succeq 0, \exists t \in \mathcal{T} : \operatorname{Tr}(S_k U) \le t_k, 1 \le k \le K\},\$$

so that \mathcal{X}^+ , $\hat{\mathcal{X}}^+$ are convex compact sets containing the origin, and $\hat{\mathcal{X}}^+$ is computationally tractable along with \mathcal{T} . Prove that

- 1. When K = 1, we have $\mathcal{X}^+ = \widehat{\mathcal{X}}^+$; 2. We always have $\mathcal{X}^+ \subset \widehat{\mathcal{X}}^+ \subset 4 \ln(5(K+1))\mathcal{X}^+$.

4.9.8Estimation by stochastic optimization

Exercise 4.88. Consider the following "multinomial" version of logistic regression problem from Section 4.7.1:

For k = 1, ..., K, we observe pairs

$$(\zeta_k, \ell_k) \in \mathbf{R}^n \times \{0, 1, \dots, m\}$$

$$(4.221)$$

drawn independently of each other from a probability distribution P_x parameterized by an unknown signal $x = [x^1; ...; x^m] \in \mathbf{R}^n \times ... \times \mathbf{R}^n$ in the following fashion:

- The probability distribution of regressor ζ induced by the distribution S_x of (ζ, ℓ) is a once for ever fixed independent of x distribution R on \mathbf{R}^n with finite second order moments and positive definite matrix $Z = \mathbf{E}_{\zeta \sim R} \{\zeta \zeta^T\}$ of second order moments;
- The conditional, ζ given, probability distribution of *label* ℓ induced by the distribution S_x of (ζ, ℓ) is the distribution of discrete random variable taking value $\iota \in \{0, 1, ..., m\}$ with probability

$$p_{\iota} = \begin{cases} \frac{\exp\{\zeta^{T}x^{\iota}\}}{1 + \sum_{i=1}^{m} \exp\{\zeta^{T}x^{i}\}}, & 1 \le \iota \le m \\ \frac{1}{1 + \sum_{i=1}^{m} \exp\{\zeta^{T}x^{i}\}}, & \iota = 0 \end{cases} \qquad [x = [x^{1}; ...; x^{m}]]$$

Given a nonempty convex compact set $\mathcal{X} \in \mathbf{R}^{mn}$ known to contain the (unknown) signal x underlying observations (4.221), we want to recover x. Note that the recovery problem associated with the standard logistic regression model is the case m = 1 of the just defined problem.

Your task is to process the above recovery problem via the approach developed in Section 4.7 and to compare the resulting SAA estimate with the Maximum 412

LECTURE 4

Likelihood estimate.

4.10 PROOFS

4.10.1 Preliminaries

4.10.1.1 Technical lemma

Lemma 4.89. Given basic spectratope

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : R_k^2[x] \leq t_k I_{d_k}, 1 \leq k \leq K \}$$

$$(4.222)$$

and a positive definite $n \times n$ matrix Q and setting $\Lambda_k = \mathcal{R}_k[Q]$, we get a collection of positive semidefinite matrices, and $\sum_k \mathcal{R}_k^*(\Lambda_k)$ is positive definite.

As a corollaries,

(i) whenever M_k , $k \leq K$, are positive definite matrices, the matrix $\sum_k \mathcal{R}_k^*[M_k]$ is positive definite;

(ii) the set $Q_T = \{Q \succeq 0 : \mathcal{R}_k[Q] \preceq TI_{d_k}, k \leq K\}$ is bounded for every T.

Proof. Let us prove the first claim, Assuming the opposite, we would be able to find a nonzero vector y such that $\sum_k y^T \mathcal{R}_k^*(\Lambda_k) y \leq 0$, whence

$$0 \ge \sum_{k} y^{T} \mathcal{R}_{k}^{*}(\Lambda_{k}) y = \sum_{k} \operatorname{Tr}(\mathcal{R}_{k}^{*}[\Lambda_{k}][yy^{T}]) = \sum_{k} \operatorname{Tr}(\Lambda_{k} \mathcal{R}_{k}[yy^{T}])$$

(we have used (4.28), (4.24)). Since $\Lambda_k = \mathcal{R}_k[Q] \succeq 0$ due to $Q \succeq 0$, see (4.25), it follows that $\operatorname{Tr}(\Lambda_k \mathcal{R}_k[yy^T]) = 0$ for all k. Now, the linear mapping $\mathcal{R}_k[\cdot]$ is \succeq -monotone, and Q is positive definite, implying that $Q \succeq r_k yy^T$ for some $r_k > 0$, whence $\Lambda_k \succeq r_k \mathcal{R}_k[yy^T]$, and therefore $\operatorname{Tr}(\Lambda_k \mathcal{R}_k[yy^T]) = 0$ implies that $\operatorname{Tr}(\mathcal{R}_k[yy^T]) = 0$, that is, $\mathcal{R}_k[yy^T] = \mathcal{R}_k^2[y] = 0$. Since $\mathcal{R}_k[\cdot]$ takes values in \mathbf{S}^{d_k} , we get $\mathcal{R}_k[y] = 0$ for al k, which is impossible due to $y \neq 0$ and property S.3, see Section 4.3.1.

To verify (i), note that when M_k are positive definite, we can find $\gamma > 0$ such that $\Lambda_k \preceq \gamma M_k$ for all $k \leq K$; invoking (4.29), we conclude that $\mathcal{R}_k^*[\Lambda_k] \preceq \gamma \mathcal{R}_k^*[M_k]$, whence $\sum_k \mathcal{R}_k^*[M_k]$ is positive definite along with $\sum_k \mathcal{R}_k^*[\Lambda_k]$.

To verify (ii), assume, on the contrary to what should be proved, that Q_T is unbounded. Since Q_T is closed and convex, it mush possess a nonzero recessive direction, that is, there should exist nonzero positive semidefinite matrix D such that $\mathcal{R}_k[D] \leq 0$ for all k. Selecting positive definite matrices M_k , the matrices $\mathcal{R}_k^*[M_k]$ are positive semidefinite (see Section 4.3.1), and their sum S is positive definite by (i). We have

$$0 \ge \sum_{k} \operatorname{Tr}(\mathcal{R}_{k}[D]M_{k}) = \sum_{k} \operatorname{Tr}(D\mathcal{R}_{k}^{*}[M_{k}]) = \operatorname{Tr}(DS),$$

where the first inequality is due to $M_k \succeq 0$, and the first equality is due to (4.28). The resulting inequality is impossible due to $0 \neq D \succeq 0$ and $S \succ 0$, which is a desired contradiction.

4.10.1.2 Noncommutative Khintchine Inequality

We will use deep result from Functional Analysis ("Noncommutative Khintchine Inequality") due to Lust-Piquard [105], Pisier [125] and Buchholz [25], see [140, Theorem 4.6.1]:

Theorem 4.90. Let $Q_i \in \mathbf{S}^n$, $1 \leq i \leq I$, and let ξ_i , i = 1, ..., I, be independent Rademacher (± 1 with probabilities 1/2) or $\mathcal{N}(0, 1)$ random variables. Then for all $t \geq 0$ one has

$$\operatorname{Prob}\left\{\left\|\sum_{i=1}^{I} \xi_{i} Q_{i}\right\| \geq t\right\} \leq 2n \exp\left\{-\frac{t^{2}}{2v_{Q}}\right\}$$

where $\|\cdot\|$ is the spectral norm, and $v_Q = \left\|\sum_{i=1}^{I} Q_i^2\right\|$.

We need the following immediate consequence of Theorem:

Lemma 4.91. Given spectratope (4.21), let $Q \in \mathbf{S}^n_+$ be such that

$$\mathcal{R}_k[Q] \preceq \rho t_k I_{d_k}, \ 1 \le k \le K, \tag{4.223}$$

for some $t \in \mathcal{T}$ and some $\rho \in (0, 1]$. Then

$$\operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q)} \{ \xi \notin \mathcal{X} \} \leq \min \left[2De^{-\frac{1}{2\rho}}, 1 \right], \ D := \sum_{k=1}^{K} d_k.$$

Proof. When setting $\xi = Q^{1/2}\eta$, $\eta \sim \mathcal{N}(0, I_n)$, we have

$$R_k[\xi] = R_k[Q^{1/2}\eta] =: \sum_{i=1}^n \eta_i \bar{R}^{ki} = \bar{R}_k[\eta]$$

with

$$\sum_{i} [\bar{R}^{ki}]^2 = \mathbf{E}_{\eta \sim \mathcal{N}(0, I_n)} \left\{ \bar{R}_k^2[\eta] \right\} = \mathbf{E}_{\xi \sim \mathcal{N}(0, Q)} \left\{ R_k^2[\xi] \right\} = \mathcal{R}_k[Q] \preceq \rho t_k I_{d_k}$$

due to (4.26). Hence, by Theorem 4.90

$$\operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q)}\{\|R_k[\xi]\|^2 \ge t_k\} = \operatorname{Prob}_{\eta \sim \mathcal{N}(0,I_n)}\{\|\bar{R}_k[\zeta]\|^2 \ge t_k\} \le 2d_k e^{-\frac{1}{2\rho}}.$$

We conclude that

$$\operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q)} \{ \xi \notin \mathcal{X} \} \le \operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q)} \{ \exists k : \|R_k[\xi]\|^2 > t_k \} \le 2De^{-\frac{1}{2\rho}}. \qquad \Box$$

The ellitopic version of Lemma 4.91 is as follows:

Lemma 4.92. Given ellitope (4.10), let $Q \in \mathbf{S}^n_+$ be such that

$$\operatorname{Tr}(R_k Q) \le \rho t_k, \ 1 \le k \le K, \tag{4.224}$$

for some $t \in \mathcal{T}$ and some $\rho \in (0, 1]$. Then

$$\operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q)} \{ \xi \notin \mathcal{X} \} \le 2K \exp\{-\frac{1}{3\rho}\}.$$

Proof. Observe that if $P \in \mathbf{S}^n_+$ satisfies $\operatorname{Tr}(R) \leq 1$, we have

$$\mathbf{E}_{\eta \sim \mathcal{N}(0,I_n)} \left\{ \exp\{\frac{1}{3}\eta^T P \eta\} \right\} \le \sqrt{3}.$$
(4.225)

Indeed, we lose nothing when assuming that $P = \text{Diag}\{\lambda_1, ..., \lambda_n\}$ with $\lambda_i \ge 0$, $\sum_i \lambda_i \le 1$. In this case

$$\mathbf{E}_{\eta \sim \mathcal{N}(0,I_n)} \left\{ \exp\{\frac{1}{3}\eta^T P \eta\} \right\} = f(\lambda) := \mathbf{E}_{\eta \sim \mathcal{N}(0,I_n)} \left\{ \exp\{\frac{1}{3}\sum_i \lambda_i \eta_i^2\} \right\}.$$

Function f is convex, so that its maximum on the simplex $\{\lambda \ge 0 : \sum_i \lambda_i \le 1\}$ is achieved at a vertex, that is,

$$f(\lambda) \leq \mathbf{E}_{\eta \sim \mathcal{N}(0,1)} \left\{ \exp\{\frac{1}{3}\eta^2\} \right\} = \sqrt{3};$$

(4.225) is proved. Note that (4.225) implies that

$$\operatorname{Prob}_{\eta \sim \mathcal{N}(0,I_n)}\left\{\eta : \eta^T P \eta > s\right\} < \sqrt{3} \exp\{-s/3\}, s \ge 0.$$

$$(4.226)$$

Now let Q and t satisfy lemma's premise. Setting $\xi = Q^{1/2}\eta$, $\eta \sim \mathcal{N}(0, I_n)$, for $k \leq K$ such that $t_k > 0$ we have

$$\xi^T R_k \xi = \rho t_k \eta^T P_k \eta, \ P_k := [\rho t_k]^{-1} Q^{1/2} R_k Q^{1/2} \succeq 0 \& \operatorname{Tr}(P_k) = [\rho t_k]^{-1} \operatorname{Tr}(QR_k) \le 1,$$

so that

$$\operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q)}\left\{\xi : \xi^T R_k \xi > s \rho t_k\right\} = \operatorname{Prob}_{\eta \sim \mathcal{N}(0,I_n)}\left\{\eta^T P_k \eta > s\right\} < \sqrt{3} \exp\{-s/3\},$$
(4.227)

where the inequality is due to (4.226). Relation (4.227) was established for k with $t_k > 0$; it is trivially true when $t_k = 0$, since in this case $Q^{1/2}R_kQ^{1/2} = 0$ due to $\operatorname{Tr}(QR_k) \leq 0$ and $R_k, Q \in \mathbf{S}^n_+$. Setting $s = 1/\rho$, we get from (4.227) that

$$\operatorname{Prob}_{x \sim \mathcal{N}(0,Q)} \left\{ \xi^T R_k \xi > t_k \right\} \le \sqrt{3} \exp\{-\frac{1}{3\rho}\}, \, k \le K,$$

and (4.226) follows due to the union bound.

4.10.1.3 Anderson's Lemma

Below we use a simply-looking, but by far nontrivial, fact

Andesron's Lemma [1]. Let f be a nonnegative even $(f(x) \equiv f(-x))$ summable function on \mathbb{R}^N such that the level sets $\{x : f(x) \ge t\}$ are convex for all t. and let $X \subset \mathbb{R}^n$ be a symmetric w.r.t. the origin closed convex set. Then for every $y \in \mathbb{R}^n$

$$\int_{X+ty} f(z)dz$$

is a nonincreasing function of $t \geq 0$. In particular, if ζ is a zero mean n-dimensional

Gaussian random vector, then for every $y \in \mathbf{R}^n$

$$\operatorname{Prob}\{\zeta \notin y + X\} \ge \operatorname{Prob}\{\zeta \notin X\}$$

whence also for every norm $\|\cdot\|$ on \mathbf{R}^n it holds

$$\operatorname{Prob}\{\zeta : \|\zeta - y\| > \rho\} \ge \operatorname{Prob}\{\zeta : \|\zeta\| > \rho\} \ \forall (y \in \mathbf{R}^n, \rho \ge 0).$$

4.10.2 **Proof of Proposition 4.6**

We need the following

Lemma 4.93. Let S be positive semidefinite $\bar{n} \times \bar{n}$ matrix with trace ≤ 1 and ξ be \bar{n} -dimensional Rademacher random vector (i.e., the entries in ξ are independent and take values ± 1 with probabilities 1/2). Then

$$\mathbf{E}\left\{\exp\left\{\zeta^T S\zeta/3\right\}\right\} \le \sqrt{3},\tag{4.228}$$

implying that

$$\operatorname{Prob}\{\xi^T S\xi > s\} \le \sqrt{3} \exp\{-s/3\}, \ s \ge 0.$$

Proof. Let $S = \sum_i \sigma_i h^i [h^i]^T$ be the eigenvalue decomposition of S, so that $[h^i]^T h^i = 1, \sigma_i \ge 0$, and $\sum_i \sigma_i \le 1$. The function

$$F(\sigma_1, ..., \sigma_{\bar{n}}) = \mathbf{E} \left\{ e^{\frac{1}{3} \sum_i \sigma_i \xi^T h^i [h^i]^T \xi} \right\}$$

is convex on the simplex $\{\sigma \ge 0, \sum_i \sigma_i \le 1\}$ and thus attains it maximum over the simplex at a vertex, implying that for some $f = h^i$, $f^T f = 1$, it holds

$$\mathbf{E}\{\mathrm{e}^{\frac{1}{3}\xi^T S\xi}\} \le \mathbf{E}\{\mathrm{e}^{\frac{1}{3}(f^T\xi)^2}\}.$$

Let $\zeta \sim \mathcal{N}(0,1)$ be independent of ξ . We have

$$\begin{aligned} \mathbf{E}_{\xi} \left\{ \exp\{\frac{1}{3}(f^{T}\xi)^{2}\} \right\} &= \mathbf{E}_{\xi} \left\{ \mathbf{E}_{\zeta} \left\{ \exp\{\left[\sqrt{2/3}f^{T}\xi\right]\zeta\} \right\} \right\} \\ &= \mathbf{E}_{\zeta} \left\{ \mathbf{E}_{\xi} \left\{ \exp\{\left[\sqrt{2/3}f^{T}\xi\right]\zeta\} \right\} \right\} = \mathbf{E}_{\zeta} \left\{ \prod_{j=1}^{N} \mathbf{E}_{\xi} \left\{ \exp\{\sqrt{2/3}\zeta f_{j}\xi_{j}\} \right\} \right\} \\ &= \mathbf{E}_{\zeta} \left\{ \prod_{j=1}^{N} \cosh(\sqrt{2/3}\zeta f_{j}) \right\} \leq \mathbf{E}_{\zeta} \left\{ \prod_{j=1}^{N} \exp\{\zeta^{2}f_{j}^{2}/3\} \right\} \\ &= \mathbf{E}_{\zeta} \left\{ \exp\{\zeta^{2}/3\} \right\} = \sqrt{3} \end{aligned}$$

 2^{0} . The right inequality in (4.20) has been justified in Section 4.2.5. To prove the left inequality in (4.20), let **T** be the closed conic hull of \mathcal{T} (see Section 4.1.1), and consider the conic problem

$$Opt_* = \max_{Q,t} \left\{ Tr(P^T C P Q) : Q \succeq 0, Tr(QS_k) \le t_k \,\forall k \le K, [t;1] \in \mathbf{T} \right\}.$$
(4.229)

416

LECTURE 4

We claim that

$$Opt = Opt_*. (4.230)$$

Indeed, (4.229) clearly is a strictly feasible and bounded conic problem; so that its optimal value is equal to the one in its conic dual (Conic Duality Theorem). Taking into account that the cone \mathbf{T}_* dual to \mathbf{T} is $\{[g;s]: s \ge \phi_{\mathcal{T}}(-g)\}$, see Section 4.1.1, we therefore get

 Opt_*

$$= \min_{\lambda, [g;s],L} \left\{ s: \begin{array}{l} \operatorname{Tr}([\sum_{k} \lambda_{k} S_{k} - L]Q) - \sum_{k} [\lambda_{k} + g_{k}]t_{k} = \operatorname{Tr}(P^{T}CPQ) \ \forall (Q, t), \\ \lambda \geq 0, L \succeq 0, s \geq \phi_{\mathcal{T}}(-g) \end{array} \right\}$$
$$= \min_{\lambda, [g;s],L} \left\{ s: \begin{array}{l} \sum_{k} \lambda_{k} S_{k} - L = P^{T}CP, \ g = -\lambda, \\ \lambda \geq 0, L \succeq 0, s \geq \phi_{\mathcal{T}}(-g) \end{array} \right\}$$
$$= \min_{\lambda} \left\{ \phi_{\mathcal{T}}(\lambda) : \sum_{k} \lambda_{k} S_{k} \succeq P^{T}CP, \lambda \geq 0 \right\} = \operatorname{Opt}, \end{array}$$

as claimed.

3⁰. With Lemma 4.93 and (4.230) at our disposal, we can now complete the proof of Proposition 4.6 by adjusting the technique from [119]. Specifically, problem (4.229) clearly is solvable; let Q_*, t^* be an optimal solution to the problem. Next, let us set $R_* = Q_*^{1/2}, \bar{C} = R_* P^T CP R_*$, let $\bar{C} = U D U^T$ be the eigenvalue decomposition of \bar{C} , and let $\bar{S}_k = U^T R_* S_k R_* U$. Observe that

$$\operatorname{Tr}(D) = \operatorname{Tr}(R_*P^TCPR_*) = \operatorname{Tr}(Q_*P^TCP) = \operatorname{Opt}_* = \operatorname{Opt},$$

$$\operatorname{Tr}(\bar{S}_k) = \operatorname{Tr}(R_*S_kR_*) = \operatorname{Tr}(Q_*S_k) \le t_k^*.$$

Now let ξ be Rademacher random vector. For k with $t_k^* > 0$, applying Lemma 4.93 to matrices \bar{S}_k/t_k^* , we get for s > 0

$$\operatorname{Prob}\{\xi^T \bar{S}_k \xi > st_k^*\} \le \sqrt{3} \exp\{-s/3\}; \tag{4.231}$$

if k is such that $t_k^* = 0$, we have $\text{Tr}(\bar{S}_k) = 0$, that is, $\bar{S}_k = 0$, and (4.231) holds true as well. Now let

$$s_* = 3\ln(\sqrt{3K}),$$

so that $\sqrt{3} \exp\{-s/3\} < 1/K$ when $s > s_*$. The latter relation combines with (4.231) to imply that for every $s > s_*$ there exists a realization $\bar{\xi}$ of ξ such that

$$\bar{\xi}^T \bar{S}_k \bar{\xi} \le s t_k^* \,\forall k.$$

Let us set $\bar{y} = \frac{1}{\sqrt{s}} R_* U \bar{\xi}$. Then

$$\bar{y}^T S_k \bar{y} = s^{-1} \bar{\xi}^T U^T R_* S_k R_* U \bar{\xi} = s^{-1} \bar{\xi}^T \bar{S}_k \bar{\xi} \le t_k^* \quad \forall k$$

implying that $\bar{y} \in \bar{X}$, and

$$\bar{y}^T P^T C P \bar{y} = s^{-1} \bar{\xi}^T U^T R_* C R_* U \bar{\xi} = s^{-1} \bar{\xi}^T D \bar{\xi} = s^{-1} \operatorname{Tr}(D) = s^{-1} \operatorname{Opt}.$$

Thus, $\max_{y \in \bar{X}} y^T P^T C P y \ge s^{-1}$ Opt whenever $s > s_*$, which implies the left inequality in (4.20).
4.10.3 **Proof of Proposition 4.8**

Proof follows the lines of the proof of Proposition 4.6. First, passing from C to the matrix $\overline{C} = P^T C P$, the situation clearly reduces to the one where P = I, which we assume in the sequel. Second, from Lemma 4.89 and the fact that the level sets of $\phi_{\mathcal{T}}(\cdot)$ on the nonnegative orthant are bounded (since \mathcal{T} contains a positive vector) it immediately follows that problem (4.32) is feasible with bounded level sets of the objective, so that the problem is solvable. The left inequality in (4.33) was proved in Section 4.3.2. Thus, all we need is to prove the right inequality in (4.33).

1°. Let **T** be the closed conic hull of \mathcal{T} (see Section 4.1.1). Consider the conic problem

$$Opt_{\#} = \max_{Q,t} \left\{ Tr(\bar{C}Q) : Q \succeq 0, \mathcal{R}_k[Q] \preceq t_k I_{d_k} \,\forall k \le K, [t;1] \in \mathbf{T} \right\}.$$
(4.232)

This problem clearly is strictly feasible; by Lemma 4.89, the feasible set of the problem is bounded, so that the problem is solvable. We claim that

$$Opt_{\#} = Opt_*. \tag{4.233}$$

Indeed, (4.232) is a strictly feasible and bounded conic problem, so that its optimal value is equal to the one in its conic dual, x that is,

$$\begin{array}{lll}
\operatorname{Opt}_{\#} &= & \min_{\Lambda = \{\Lambda_k\}_{k \leq K}, [g;s], L} \left\{ s: & \operatorname{Tr}([\sum_k \mathcal{R}_k^*[\Lambda_k] - L]Q) - \sum_k [\operatorname{Tr}(\Lambda_k) + g_k] t_k \\ s: & = & \operatorname{Tr}(\bar{C}Q) \ \forall (Q, t), \\ \Lambda_k \succeq 0 \ \forall k, L \succeq 0, s \geq \phi_{\mathcal{T}}(-g) \\ &= & \min_{\Lambda, [g;s], L} \left\{ s: & \sum_k \mathcal{R}_k^*[\Lambda_k] - L = \bar{C}, g = -\lambda[\Lambda], \\ \Lambda_k \succeq 0 \ \forall k, L \succeq 0, s \geq \phi_{\mathcal{T}}(-g) \\ &= & \min_{\Lambda} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) : \sum_k \mathcal{R}_k^*[\Lambda_k] \succeq \bar{C}, \Lambda_k \succeq 0 \ \forall k \right\} = \operatorname{Opt}_*.
\end{array}$$

as claimed.

2°. Problem (4.232), as we already know, is solvable; let Q_*, t^* be an optimal solution to the problem. Next, let us set $R_* = Q_*^{1/2}$, $\hat{C} = R_*\bar{C}R_*$, and let $\hat{C} = UDU^T$ be the eigenvalue decomposition of \hat{C} , so that the matrix $D = U^T R_*\bar{C}R_*U$ is diagonal, and the trace of this matrix is $\operatorname{Tr}(R_*\bar{C}R_*) = \operatorname{Tr}(\bar{C}Q_*) = \operatorname{Opt}_{\#} = \operatorname{Opt}_*$, Now let $V = R_*U$, and let $\xi = V\eta$, where η is *n*-dimensional random Rademacher vector (independent entries taking values ± 1 with probabilities 1/2). We have

$$\xi^T \bar{C} \xi = \eta^T [V^T \bar{C} V] \eta = \eta^T [U^T R_* \bar{C} R_* U] \eta = \eta^T D \eta \equiv \text{Tr}(D) = \text{Opt}_*, \quad (4.234)$$

(recall that D is diagonal) and

$$\mathbf{E}_{\xi}\{\xi\xi^{T}\} = \mathbf{E}_{\eta}\{V\eta\eta^{T}V^{T}\} = VV^{T} = R_{*}UU^{T}R_{*} = R_{*}^{2} = Q_{*}.$$

From the latter relation,

$$\mathbf{E}_{\xi}\left\{R_{k}^{2}[\xi]\right\} = \mathbf{E}_{\xi}\left\{\mathcal{R}_{k}[\xi\xi^{T}]\right\} = \mathcal{R}_{k}[\mathbf{E}_{\xi}\{\xi\xi^{T}\}] = \mathcal{R}_{k}[Q_{*}] \preceq t_{k}^{*}I_{d_{k}}, 1 \le k \le K.$$
(4.235)

On the other hand, with properly selected symmetric matrices \bar{R}^{kj} we have

$$R_k[Vy] = \sum_i \bar{R}^{ki} y_i$$

identically in $y \in \mathbf{R}^n$, whence

$$\mathbf{E}_{\xi}\left\{R_{k}^{2}[\xi]\right\} = \mathbf{E}_{\eta}\left\{R_{k}^{2}[V\eta]\right\} = \mathbf{E}_{\eta}\left\{\left[\sum_{i}\eta_{i}\bar{R}^{ki}\right]^{2}\right\} = \sum_{i,j}\mathbf{E}_{\eta}\{\eta_{i}\eta_{j}\}\bar{R}^{ki}\bar{R}^{kj} = \sum_{i}[\bar{R}^{ki}]^{2}.$$

This combines with (4.235) to imply that

$$\sum_{i} [\bar{R}^{ki}]^2 \leq t_k^* I_{d_k}, \ 1 \leq k \leq K.$$
(4.236)

3°. Let us fix $k \leq K$. Assuming $t_k^* > 0$ and applying Theorem 4.90, we derive from (4.236) that

$$\operatorname{Prob}\{\eta: \|\bar{R}_k[\eta]\|^2 > t_k^*/\rho\} < 2d_k e^{-\frac{1}{2\rho}},$$

and recalling the relation between ξ and η , we arrive at

$$\operatorname{Prob}\{\xi : \|R_k[\xi]\|^2 > t_k^*/\rho\} < 2d_k e^{-\frac{1}{2\rho}} \quad \forall \rho \in (0, 1].$$
(4.237)

Note that when $t_k^* = 0$ (4.236) implies $\bar{R}^{ki} = 0$ for all *i*, so that $R_k[\xi] = \bar{R}_k[\eta] = 0$, and (4.237) holds for those *k* as well.

Now let us set $\rho = \frac{1}{2 \max[\ln(2D),1]}$. For this ρ , the sum over $k \leq K$ of the right hand sides in inequalities (4.237) is ≤ 1 , implying that there exists a realization $\bar{\xi}$ of ξ such that

$$||R_k[\bar{\xi}]||^2 \le t_k^*/\rho, \ \forall k$$

or, equivalently,

$$\bar{x} := \rho^{1/2} P \bar{\xi} \in \mathcal{X},$$

implying that

$$Opt \ge \bar{x}^T C \bar{x} = \rho \xi^T \bar{C} \xi = \rho Opt_*$$

(the concluding equality is due to (4.234)), and we arrive at the right inequality in (4.33).

4.10.4 **Proof of Lemma 4.17**

1°. Let us verify (4.65). When $Q \succ 0$, passing from variables (Θ, Υ) in problem (4.64) to the variables $(G = Q^{1/2} \Theta Q^{1/2}, \Upsilon)$, the problem becomes exactly the optimization problem in (4.65), implying that $\operatorname{Opt}[Q] = \overline{\operatorname{Opt}}[Q]$ when $Q \succ 0$. As it is easily seen, both sides in this equality are continuous in $Q \succeq 0$, and (4.65) follows.

2^o. Let us set $\zeta = Q^{1/2}\eta$ with $\eta \sim \mathcal{N}(0, I_N)$ and $Z = Q^{1/2}Y$. Let us show that when $\varkappa \geq 1$ one has

$$\operatorname{Prob}_{\eta}\{\|Z^{T}\eta\| \geq \bar{\delta}\} \geq \beta_{\varkappa} := 1 - \frac{e^{3/8}}{2} - 2Fe^{-\varkappa^{2}/2}, \quad \bar{\delta} = \frac{\operatorname{Opt}[Q]}{4\varkappa}, \quad (4.238)$$

where

$$\begin{bmatrix} \overline{\mathrm{Opt}}[Q] =] & \mathrm{Opt}[Q] := \min_{\Theta, \Upsilon = \{\Upsilon_{\ell}, \ell \leq L\}} & \left\{ \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \mathrm{Tr}(\Theta) : \\ & \Upsilon_{\ell} \succeq 0, \left[\frac{\Theta}{\frac{1}{2}M^{T}Z^{T}} \left| \frac{1}{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]} \right| \right] \succeq 0 \right\}$$

$$(4.239)$$

 $\mathbf{3}^{o}$. Let us represent Opt[Q] as the optimal value of a conic problem. Setting

 $\mathbf{K} = \mathrm{cl}\{[r;s] : s > 0, r/s \in \mathcal{R}\},\$

we ensure that

$$\mathcal{R} = \{r : [r;1] \in \mathbf{K}\}, \ \mathbf{K}_* = \{[g;s] : s \ge \phi_{\mathcal{R}}(-g)\},\$$

where \mathbf{K}_* is the cone dual to \mathbf{K} . Consequently, (4.239) reads

$$\operatorname{Opt}[Q] = \min_{\Theta, \Upsilon, \theta} \left\{ \theta + \operatorname{Tr}(\Theta) : \begin{array}{cc} \Upsilon_{\ell} \succeq 0, 1 \leq \ell \leq L & (a) \\ \frac{\Theta}{\frac{1}{2}ZM} & \frac{1}{2}ZM \\ \frac{1}{2}M^{T}Z^{T} & \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}] \\ [-\lambda[\Upsilon]; \theta] \in \mathbf{K}_{*} & (c) \end{array} \right\}.$$
(P)

 $\mathbf{4}^{o}$. Now let us prove that there exists matrix $W \in \mathbf{S}^{q}_{+}$ and $r \in \mathcal{R}$ such that

$$\mathcal{S}_{\ell}[W] \preceq r_{\ell} I_{f_{\ell}}, \, \ell \leq L, \tag{4.240}$$

and

$$\operatorname{Opt}[Q] \le \sum_{i} \sigma_{i}(ZMW^{1/2}), \qquad (4.241)$$

where $\sigma_1(\cdot) \ge \sigma_2(\cdot) \ge \dots$ are singular values.

To get the announced result, let us pass from problem (P) to its conic dual. Applying Lemma 4.89 we conclude that (P) is strictly feasible; in addition, (P) clearly is bounded, so that the dual to (P) problem (D) is solvable with optimal value Opt[Q]. Let us build (D). Denoting by $\Lambda_{\ell} \succeq 0, \ell \leq L, \left[\begin{array}{c|c} G & -R \\ \hline -R^T & W \end{array} \right] \succeq 0, [r; \tau] \in \mathbf{K}$ the Lagrange multipliers for the respective constraints in (P), and aggregating these constraints, the multipliers being the aggregation weights, we arrive at the following aggregated constraint:

$$\operatorname{Tr}(\Theta G) + \operatorname{Tr}(W \sum_{\ell} \mathcal{S}_{\ell}^*[\Upsilon_{\ell}]) + \sum_{\ell} \operatorname{Tr}(\Lambda_{\ell} \Upsilon_{\ell}) - \sum_{\ell} r_{\ell} \operatorname{Tr}(\Upsilon_{\ell}) + \theta \tau \ge \operatorname{Tr}(ZMR^T).$$

To get the dual problem, we impose on the Lagrange multipliers, in addition to the initial conic constraints like $\Lambda_{\ell} \succeq 0$, $1 \leq \ell \leq L$, the restriction that the left hand side in the aggregated constraint, identically in Θ , Υ_{ℓ} and θ , is equal to the objective of (P), that is,

$$G = I, \ \mathcal{S}_{\ell}[W] + \Lambda_{\ell} - r_{\ell}I_{f_{\ell}} = 0, \ 1 \ \leq \ell \leq L, \ \tau = 1,$$

and maximize, under the resulting restrictions, the right-hand side of the aggregated constraint. After immediate simplifications, we arrive at

$$Opt[Q] = \max_{W,R,r} \left\{ Tr(ZMR^T) : W \succeq R^T R, r \in \mathcal{R}, \mathcal{S}_{\ell}[W] \preceq r_{\ell} I_{f_{\ell}}, 1 \le \ell \le L \right\}$$

(note that $r \in \mathcal{R}$ is equivalent to $[r; 1] \in \mathbf{K}$, and $W \succeq R^T R$ is the same as $\begin{bmatrix} I & -R \\ \hline -R^T & W \end{bmatrix} \succeq 0$). Now, to say that $R^T R \preceq W$ is exactly the same as to say that $R = SW^{1/2}$ with the spectral norm $\|S\|_{\mathrm{Sh},\infty}$ of S not exceeding 1, so that

$$\operatorname{Opt}[Q] = \max_{W,S,r} \left\{ \underbrace{\operatorname{Tr}([ZM[SW^{1/2}]^T)]}_{=\operatorname{Tr}([ZMW^{1/2}]S^T)} : W \succeq 0, \|S\|_{\operatorname{Sh},\infty} \le 1, r \in \mathcal{R}, \mathcal{S}_{\ell}[W] \preceq r_{\ell}I_{f_{\ell}}, \, \ell \le L \right\}$$

and we can immediately eliminate the S-variable, using the well-known fact that for every $p \times q$ matrix J, it holds

$$\max_{S \in \mathbf{R}^{p \times q}, \|S\|_{\mathrm{Sh}, \infty} \le 1} \operatorname{Tr}(JS^T) = \|J\|_{\mathrm{Sh}, 1},$$

where $||J||_{\text{Sh},1}$ is the nuclear norm (the sum of singular values) of J. We arrive at

$$\operatorname{Opt}[Q] = \max_{W,r} \left\{ \|ZMW^{1/2}\|_{\operatorname{Sh},1} : r \in \mathcal{R}, W \succeq 0, \mathcal{S}_{\ell}[W] \preceq r_{\ell} I_{d_{\ell}}, \, \ell \leq L \right\}.$$

The resulting problem clearly is solvable, and its optimal solution W ensures the target relations (4.240) and (4.241).

5°. Given W satisfying (4.240) and (4.241), let $UJV = W^{1/2}M^TZ^T$ be the singular value decomposition of $W^{1/2}M^TZ^T$, so that U and V are, respectively, $q \times q$ and $N \times N$ orthogonal matrices, J is $q \times N$ matrix with diagonal $\sigma = [\sigma_1; ...; \sigma_p], p = \min[q, N]$, and zero off-diagonal entries; the diagonal entries $\sigma_i, 1 \leq i \leq p$ are the singular values of $W^{1/2}M^TZ^T$, or, which is the same, of $ZMW^{1/2}$. Therefore, we have

$$\sum_{i} \sigma_i \ge \operatorname{Opt}[Q]. \tag{4.242}$$

Now consider the following construction. Let $\eta \sim \mathcal{N}(0, I_N)$; we denote by v the vector comprised of the first p entries in $V\eta$; note that $v \sim \mathcal{N}(0, I_p)$, since V is orthogonal. We then augment, if necessary, v by q - p independent of each other and of $\eta \mathcal{N}(0, 1)$ random variables to obtain a q-dimensional normal vector $v' \sim \mathcal{N}(0, I_q)$, and set $\chi = Uv'$; because U is orthogonal we also have $\chi \sim \mathcal{N}(0, I_q)$. Observe that

$$\chi^T W^{1/2} M^T Z^T \eta = \chi^T U J V \eta = [\upsilon']^T J \upsilon = \sum_{i=1}^p \sigma_i \upsilon_i^2.$$
(4.243)

To continue we need two simple observations.

(i) One has

$$\alpha := \operatorname{Prob}\left\{\sum_{i=1}^{p} \sigma_{i} v_{i}^{2} < \frac{1}{4} \sum_{i=1}^{p} \sigma_{i}\right\} \le \frac{\mathrm{e}^{3/8}}{2} \ [= 0.7275...].$$
(4.244)

The claim is evident when $\sigma := \sum_i \sigma_i = 0$. Now let $\sigma > 0$, and let us apply the Cramer bounding scheme. Namely, given $\gamma > 0$, consider the random variable

$$\omega = \exp\left\{\frac{1}{4}\gamma \sum_{i} \sigma_{i} - \gamma \sum_{i} \sigma_{i} v_{i}^{2}\right\}.$$

Note that $\omega > 0$ a.s., and is > 1 when $\sum_{i=1}^{p} \sigma_i v_i^2 < \frac{1}{4} \sum_{i=1}^{p} \sigma_i$, so that $\alpha \leq \mathbf{E}\{\omega\}$, or, equivalently, thanks to $\upsilon \sim \mathcal{N}(0, I_p)$,

$$\ln(\alpha) \le \ln(\mathbf{E}\{\omega\}) = \frac{1}{4}\gamma \sum_{i} \sigma_{i} + \sum_{i} \ln\left(\mathbf{E}\{\exp\{-\gamma\sigma_{i}v_{i}^{2}\}\}\right) \le \frac{1}{4}\gamma\sigma - \frac{1}{2}\sum_{i} \ln(1+2\gamma\sigma_{i}).$$

Function $-\sum_{i} \ln(1+2\gamma\sigma_i)$ is convex in $[\sigma_1;...;\sigma_p] \ge 0$, therefore, its maximum over the simplex $\{\sigma_i \ge 0, i \le p, \sum_i \sigma_i = \sigma\}$ is attained at a vertex, and we get

$$\ln(\alpha) \le \frac{1}{4}\gamma\sigma - \frac{1}{2}\ln(1+2\gamma\sigma).$$

Minimizing the right hand side in $\gamma > 0$, we arrive at (4.244). (ii) Whenever $\varkappa \ge 1$, one has

$$\operatorname{Prob}\{\|MW^{1/2}\chi\|_* > \varkappa\} \le 2F \exp\{-\varkappa^2/2\}, \qquad (4.245)$$

with F given by (4.63).

Indeed, setting $\rho = 1/\varkappa^2 \leq 1$ and $\omega = \sqrt{\rho}W^{1/2}\chi$, we get $\omega \sim \mathcal{N}(0, \rho W)$. Let us apply Lemma 4.91 to $Q = \rho W$, \mathcal{R} in the role of \mathcal{T} , L in the role of K, and $\mathcal{S}_{\ell}[\cdot]$ in the role of $\mathcal{R}_k[\cdot]$. Denoting

$$\mathcal{Y} := \{ y : \exists r \in \mathcal{R} : S^2_{\ell}[y] \preceq r_{\ell} I_{f_{\ell}}, \ell \leq L \},\$$

we have $S_{\ell}[Q] = \rho S_{\ell}[W] \leq \rho r_{\ell} I_{f_{\ell}}, \ \ell \leq L$, with $r \in \mathcal{R}$ (see (4.240)), so we are under the premise of Lemma 4.91 (with \mathcal{Y} in the role of \mathcal{X} and therefore with Fin the role of D). Applying the lemma, we conclude that

$$\operatorname{Prob}\left\{\chi: \varkappa^{-1} W^{1/2} \chi \notin \mathcal{Y}\right\} \le 2F \exp\{-1/(2\rho)\} = 2F \exp\{-\varkappa^2/2\}.$$

Recalling that $\mathcal{B}_* = M\mathcal{Y}$, we see that $\operatorname{Prob}\{\chi : \varkappa^{-1}MW^{1/2}\chi \notin \mathcal{B}_*\}$ is indeed upper-bounded by the right hand side of (4.245), and (4.245) follows.

Now, for $\varkappa \geq 1$, let

$$E_{\varkappa} = \left\{ (\chi, \eta) : \| M W^{1/2} \chi \|_{\ast} \le \varkappa, \sum_{i} \sigma_{i} v_{i}^{2} \ge \frac{1}{4} \sum_{i} \sigma_{i} \right\}.$$

For $(\chi, \eta) \in E_{\varkappa}$ we have

$$\varkappa \|Z^{T}\eta\| \ge \|MW^{1/2}\chi\|_{*}\|Z^{T}\eta\| \ge \chi^{T}W^{1/2}M^{T}Z^{T}\eta = \sum_{i}\sigma_{i}v_{i}^{2} \ge \frac{1}{4}\sum_{i}\sigma_{i} \ge \frac{1}{4}\text{Opt}[Q],$$

(we have used (4.243) and (4.242)). On the other hand, due to (4.244) and (4.245),

$$\operatorname{Prob}\{E_{\varkappa}\} \ge \beta_{\varkappa},\tag{4.246}$$

and we arrive at (4.67). The latter relation clearly implies (4.68) which, in turn, implies the right inequality in (4.66).

4.10.5 **Proofs of Propositions 4.5, 4.16, 4.19**

Below, we focus on the proof of Proposition 4.16; Propositions 4.5, 4.19 will be derived from it in Sections 4.10.5.2, 4.10.6.2, respectively.

4.10.5.1 Proof of Proposition 4.16

In what follows, we use the assumptions and the notation of Proposition 4.16.

1⁰. Let

$$\Phi(H,\Lambda,\Upsilon,\Upsilon',\Theta;Q) = \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \operatorname{Tr}(Q\Theta) : \mathcal{M} \times \Pi \to \mathbf{R},$$

where

$$\mathcal{M} = \left\{ (H, \Lambda, \Upsilon, \Upsilon', \Theta) : \begin{array}{ll} \Lambda = \{\Lambda_k \succeq 0, k \le K\}, \\ \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \le L\}, \ \Upsilon' = \{\Upsilon'_\ell \succeq 0, \ell \le L\} \\ \left[\frac{\sum_k \mathcal{R}_k^* [\Lambda_k]}{\frac{1}{2} M^T [B - H^T A]} \left| \frac{\sum_\ell \mathcal{S}_\ell^* [\Upsilon_\ell]}{\sum_\ell \mathcal{S}_\ell^* [\Upsilon_\ell]} \right] \succeq 0 \\ \left[\frac{\Theta}{\frac{1}{2} M^T H^T} \left| \sum_\ell \mathcal{S}_\ell^* [\Upsilon'_\ell] \right| \ge 0 \end{array} \right\}$$

Looking at (4.50), we conclude immediately that the optimal value Opt in (4.50) is nothing but

$$Opt = \min_{(H,\Lambda,\Upsilon,\Upsilon',\Theta)\in\mathcal{M}} \left[\overline{\Phi}(H,\Lambda,\Upsilon,\Upsilon',\Theta) := \max_{Q\in\Pi} \Phi(H,\Lambda,\Upsilon,\Upsilon',\Theta;Q)\right]. \quad (4.247)$$

Note that the sets \mathcal{M} and Π are closed and convex, Π is compact, and Φ is a continuous convex-concave function on $\mathcal{M} \times \Pi$. In view of these observations, the fact that $\Pi \subset \operatorname{int} \mathbf{S}^m_+$ combines with the Sion-Kakutani Theorem to imply that Φ possesses saddle point $(H_*, \Lambda_*, \Upsilon_*, \Upsilon'_*, \Theta_*; Q_*)$ (min in $(H, \Lambda, \Upsilon, \Upsilon', \Theta)$, max in Q) on $\mathcal{M} \times \Pi$, whence Opt is the saddle point value of Φ by (4.247). We conclude that

for properly selected $Q_* \in \Pi$ it holds

$$\begin{aligned} \operatorname{Opt} &= \min_{(H,\Lambda,\Upsilon,\Upsilon',\Theta)\in\mathcal{M}} \Phi(H,\Lambda,\Upsilon,\Upsilon',\Theta;Q_*) \\ &= \min_{(H,\Lambda,\Upsilon,\Upsilon',\Theta)\in\mathcal{M}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \operatorname{Tr}(Q_*\Theta) : \\ \Lambda &= \{\Lambda_k \succeq 0, k \le K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \le L\}, \Upsilon' = \{\Upsilon'_\ell \succeq 0, \ell \le L\} \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T[B - H^T A]} \sum_{\ell} \mathcal{S}_\ell^*[\Upsilon_\ell]} \right] \succeq 0, \\ \left[\frac{\Theta}{\frac{1}{2}M^T H^T} \frac{\frac{1}{2}E^T}{\mathcal{L}_\ell} \mathcal{S}_\ell^*[\Upsilon'_\ell]} \right] \succeq 0, \\ \left[\frac{\Theta}{\frac{1}{2}M^T H^T} \frac{\frac{1}{2}E^T}{\mathcal{L}_\ell} \mathcal{S}_\ell^*[\Upsilon'_\ell]} \right] \succeq 0, \\ &= \min_{H,\Lambda,\Upsilon,\Upsilon',G} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \operatorname{Tr}(G) : \\ \Lambda &= \{\Lambda_k \succeq 0, k \le K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \le L\}, \Upsilon' = \{\Upsilon'_\ell \succeq 0, \ell \le L\} \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T H^T \mathcal{Q}_k^{1/2}} \frac{1}{2} \mathcal{O}_\ell} \mathcal{S}_\ell^*[\Upsilon_\ell]}{\mathcal{O}_\ell} \right] \succeq 0, \\ &= \min_{H,\Lambda,\Upsilon} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \overline{\Psi}(H) : \\ \Lambda &= \{\Lambda_k \succeq 0, k \le K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \le L\} \\ \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T (B - H^T A]} \frac{\sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell]}{\mathcal{O}_\ell} \right] \succeq 0 \\ \end{bmatrix}$$

where Opt is given by (4.50), and the equalities are due to (4.64) and (4.65).

From now on we assume that the observation noise ξ in observation (4.35) is $\xi \sim \mathcal{N}(0, Q_*)$. Besides this, we assume that $B \neq 0$, since otherwise the conclusion of Proposition 4.16 is evident.

2⁰. ϵ -risk. In Proposition 4.16, we are speaking about $\|\cdot\|$ -risk of an estimate – the maximal, over signals $x \in \mathcal{X}$, expected norm $\|\cdot\|$ of the error in recovering Bx; what we need to prove that the minimax optimal risk RiskOpt_{II,||}.|| \mathcal{X}] as given by (4.61) can be lower-bounded by a quantity "of order of" Opt. To this end, of course, it suffices to build such a lower bound for the quantity

$$\operatorname{RiskOpt}_{\|\cdot\|} := \inf_{\widehat{x}(\cdot)} \left[\sup_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim \mathcal{N}(0,Q_*)} \{ \|Bx - \widehat{x}(Ax + \xi)\| \} \right],$$

since this quantity is a lower bound on RiskOpt_{II,||·||}. Technically, it is more convenient to work with the ϵ -risk defined in terms of "||·||-confidence intervals" rather than in terms of the expected norm of the error. Specifically, in the sequel we will heavily use the minimax ϵ -risk defined as

$$\operatorname{RiskOpt}_{\epsilon} = \inf_{\widehat{x},\rho} \left\{ \rho : \operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q_*)} \{ \|Bx - \widehat{x}(Ax + \xi)\| \le \epsilon \ \forall x \in \mathcal{X} \right\}$$
(4.249)

When $\epsilon \in (0,1)$ is once for ever fixed (in the sequel, we use $\epsilon = \frac{1}{8}$), ϵ -risk lowerbounds RiskOpt_{||.||}, since by evident reasons

$$\operatorname{RiskOpt}_{\parallel \cdot \parallel} \ge \epsilon \operatorname{RiskOpt}_{\epsilon}. \tag{4.250}$$

424

LECTURE 4

Consequently, all we need in order to prove Proposition 4.16 is to lower-bound RiskOpt $_{\frac{1}{8}}$ by a "not too small" multiple of Opt, and this is what we are achieving below.

3°. Let W be a positive semidefinite $n \times n$ matrix, let $\eta \sim \mathcal{N}(0, W)$ be random signal, and let $\xi \sim \mathcal{N}(0, Q_*)$ be independent of η ; vectors (η, ξ) induce random vector

 $\omega = A\eta + \xi \sim \mathcal{N}(0, AWA^T + Q_*).$

Consider the Bayesian version of the estimation problem where given ω we are interested to recover $B\eta$. Recall that, because $[\omega; B\eta]$ is zero mean Gaussian, the conditional expectation $\mathbf{E}_{|\omega}\{B\eta\}$ of $B\eta$ given ω is linear in ω : $\mathbf{E}_{|\omega}\{B\eta\} = \bar{H}^T \omega$ for some \bar{H} depending on W only⁸⁵. Therefore, denoting by $P_{|\omega}$ conditional, ω given, probability distribution, for any $\rho > 0$ and estimate $\hat{x}(\cdot)$ one has

$$\begin{aligned} \operatorname{Prob}_{\eta,\xi}\{\|B\eta - \widehat{x}(A\eta + \xi)\| \ge \rho\} &= \mathbf{E}_{\omega}\{\operatorname{Prob}_{|\omega}\{\|B\eta - \widehat{x}(\omega)\| \ge \rho\}\}\\ &\ge \mathbf{E}_{\omega}\{\operatorname{Prob}_{|\omega}\{\|B\eta - \mathbf{E}_{|\omega}\{B\eta\}\| \ge \rho\}\}\\ &= \operatorname{Prob}_{\eta,\xi}\{\|B\eta - \overline{H}^{T}(A\eta + \xi)\| \ge \rho\}, \end{aligned}$$

$$(4.251)$$

with the inequality given by the Anderson Lemma as applied to the shift of the Gaussian distribution $P_{|\omega}$ by its mean. Applying the Anderson Lemma again we get

$$\operatorname{Prob}_{\eta,\xi}\{\|B\eta - \bar{H}^T(A\eta + \xi)\| \ge \rho\} = \mathbf{E}_{\xi}\left\{\operatorname{Prob}_{\eta}\{\|(B - \bar{H}^T A)\eta - \bar{H}^T \xi\| \ge \rho\}\right\}$$
$$\ge \operatorname{Prob}_{\eta}\{\|(B - \bar{H}^T A)\eta\| \ge \rho\},$$

and, by "symmetric" reasoning,

$$\operatorname{Prob}_{\eta,\xi}\{\|B\eta - \bar{H}^T(A\eta + \xi)\| \ge \rho\} \ge \operatorname{Prob}_{\xi}\{\|\bar{H}^T\xi\| \ge \rho\}.$$

We conclude that for any $\widehat{x}(\cdot)$

$$\begin{aligned} \operatorname{Prob}_{\eta,\xi}\{\|B\eta - \widehat{x}(\omega)\| \geq \rho\} \\ \geq \max\left\{\operatorname{Prob}_{\eta}\{\|(B - \bar{H}^T A)\eta\| \geq \rho\}, \operatorname{Prob}_{\xi}\{\|\bar{H}^T\xi\| \geq \rho\}\right\}. \end{aligned} \tag{4.252}$$

4°. Let H be $m \times \nu$ matrix. Applying Lemma 4.17 to N = m, $Y = \overline{H}$, $Q = Q_*$, we get from (4.67)

$$\operatorname{Prob}_{\xi \sim \mathcal{N}(0,Q_*)}\{ \|H^T \xi\| \ge [4\varkappa]^{-1} \overline{\Psi}(\overline{H}) \} \ge \beta_{\varkappa}$$

$$(4.253)$$

where $\overline{\Psi}(H)$ is defined by (4.248). Similarly, applying Lemma 4.17 to N = n, $Y = (B - \overline{H}^T A)^T$, Q = W, we obtain

$$\operatorname{Prob}_{\eta \sim \mathcal{N}(0,W)} \{ \| (B - \bar{H}^T A) \eta \| \ge [4\varkappa]^{-1} \Phi(W,\bar{H}) \} \ge \beta_{\varkappa}$$

$$(4.254)$$

⁸⁵We have used the following standard fact: let $\zeta = [\omega, \eta] \sim \mathcal{N}(0, S)$, the covariance matrix of the marginal distribution of ω being nonsingular. Then the conditional, ω given, distribution of η is Gaussian with mean linearly depending on ω and covariance matrix independent of ω .

where

$$\overline{\Phi}(W,H) = \min_{\Upsilon = \{\Upsilon_{\ell},\ell \leq L\},\Theta} \left\{ \operatorname{Tr}(W\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) : \Upsilon_{\ell} \succeq 0 \,\forall \ell, \\ \left[\frac{\Theta}{\frac{1}{2}M^{T}[B - H^{T}A]} \mid \frac{1}{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]} \right] \succeq 0 \right\}.$$

$$(4.255)$$

Let us put $\rho(W, \bar{H}) = [8\varkappa]^{-1}[\overline{\Psi}(\bar{H}) + \overline{\Phi}(W, \bar{H})]$; when combining (4.254) with (4.253) we conclude that

$$\max\left\{\operatorname{Prob}_{\eta}\{\|(B-\bar{H}^{T}A)\eta\| \ge \rho(W,\bar{H})\}, \operatorname{Prob}_{\xi}\{\|\bar{H}^{T}\xi\| \ge \rho(W,\bar{H})\}\right\} \ge \beta_{\varkappa},$$

and the same inequality holds if $\rho(W, \bar{H})$ is replaced with the smaller quantity

$$\bar{\rho}(W) = [8\varkappa]^{-1} \inf_{H} [\overline{\Psi}(H) + \overline{\Phi}(W,H)].$$

Now, the latter bound combines with (4.252) to imply the following result:

Lemma 4.94. Let W be a positive semidefinite $n \times n$ matrix, and $\varkappa \geq 1$. Then for any estimate $\hat{x}(\cdot)$ of $B\eta$ given observation $\omega = A\eta + \xi$, one has

$$\operatorname{Prob}_{\eta,\xi}\{\|B\eta - \widehat{x}(\omega)\| \ge [8\varkappa]^{-1} \inf_{H}[\overline{\Psi}(H) + \overline{\Phi}(W,H)]\} \ge \beta_{\varkappa} = 1 - \frac{e^{3/8}}{2} - 2Fe^{-\varkappa^2/2}$$

where $\overline{\Psi}(H)$ and $\overline{\Phi}(W, H)$ are defined, respectively, by (4.248) and (4.255). In particular, for

$$\varkappa = \bar{\varkappa} := \sqrt{2\ln F} + 10\ln 2 \tag{4.256}$$

the latter probability is > 3/16.

5°. For $0 < \kappa \leq 1$, let us set

(a)
$$\mathcal{W}_{\kappa} = \{ W \in \mathbf{S}_{+}^{n} : \exists t \in \mathcal{T} : \mathcal{R}_{k}[W] \leq \kappa t_{k} I_{d_{k}}, 1 \leq k \leq K \},$$

(b) $\mathcal{Z} = \left\{ (\Upsilon = \{\Upsilon_{\ell}, \ell \leq L\}, \Theta, H) : \begin{bmatrix} \Theta & \left| \frac{1}{2} [B^{T} - A^{T} H] M \right| \\ \frac{1}{2} M^{T} [B - H^{T} A] & \sum_{\ell} \mathcal{S}_{\ell}^{*} [\Upsilon_{\ell}] \end{bmatrix} \succeq 0 \right\}.$

$$(4.257)$$

Note that \mathcal{W}_{κ} is a nonempty convex compact (by Lemma 4.89) set such that $\mathcal{W}_{\kappa} = \kappa \mathcal{W}_1$, and \mathcal{Z} is a nonempty closed convex set. Consider the parametric saddle point problem

$$Opt(\kappa) = \max_{W \in \mathcal{W}_{\kappa}} \min_{(\Upsilon,\Theta,H) \in \mathcal{Z}} \left[E(W;\Upsilon,\Theta,H) := Tr(W\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \overline{\Psi}(H) \right].$$
(4.258)

This problem is convex-concave; utilizing the fact that \mathcal{W}_{κ} is compact and contains positive definite matrices, it is immediately seen that the Sion-Kakutani theorem ensures the existence of a saddle point whenever $\kappa \in (0, 1]$. We claim that

$$0 < \kappa \le 1 \Rightarrow \operatorname{Opt}(\kappa) \ge \sqrt{\kappa} \operatorname{Opt}(1). \tag{4.259}$$

426

Indeed, \mathcal{Z} is invariant w.r.t. scalings

$$(\Upsilon = \{\Upsilon_{\ell}, \ell \leq L\}, \Theta, H) \mapsto (\theta\Upsilon := \{\theta\Upsilon_{\ell}, \ell \leq L\}, \theta^{-1}\Theta, H), \qquad [\theta > 0]$$

When taking into account that $\phi_{\mathcal{R}}(\lambda[\theta \Upsilon]) = \theta \phi_{\mathcal{R}}(\lambda[\Upsilon])$, we get

$$\begin{split} \underline{E}(W) &:= \min_{\substack{(\Upsilon,\Theta,H)\in\mathcal{Z}}} E(W;\Upsilon,\Theta,H) = \min_{\substack{(\Upsilon,\Theta,H)\in\mathcal{Z}\\ (\Upsilon,\Theta,H)\in\mathcal{Z}}} \inf_{\theta>0} E(W;\theta\Upsilon,\theta^{-1}\Theta,H) \\ &= \min_{\substack{(\Upsilon,\Theta,H)\in\mathcal{Z}}} \left[2\sqrt{\operatorname{Tr}(W\Theta)\phi_{\mathcal{R}}(\lambda[\Upsilon])} + \overline{\Psi}(H) \right]. \end{split}$$

Because $\overline{\Psi}$ is nonnegative we conclude that whenever $W \succeq 0$ and $\kappa \in (0, 1]$, one has

$$\underline{E}(\kappa W) \ge \sqrt{\kappa} \underline{E}(W),$$

which combines with $\mathcal{W}_{\kappa} = \kappa \mathcal{W}_1$ to imply that

$$\operatorname{Opt}(\kappa) = \max_{W \in \mathcal{W}_{\kappa}} \underline{E}(W) = \max_{W \in \mathcal{W}_{1}} \underline{E}(\kappa W) \ge \sqrt{\kappa} \max_{W \in \mathcal{W}_{1}} \underline{E}(W) = \sqrt{\kappa} \operatorname{Opt}(1),$$

and (4.259) follows.

6^o**.** We claim that

$$Opt(1) = Opt, (4.260)$$

where Opt is given by (4.50) (and, as we have seen, by (4.248) as well). Note that (4.260) combines with (4.259) to imply that

$$0 < \kappa \le 1 \Rightarrow \operatorname{Opt}(\kappa) \ge \sqrt{\kappa} \operatorname{Opt}.$$
 (4.261)

Verification of (4.260) is given by the following computation. By the Sion-Kakutani Theorem,

$$\begin{aligned} \operatorname{Opt}(1) &= \max_{W \in \mathcal{W}_1} \min_{(\Upsilon, \Theta, H) \in \mathcal{Z}} \left\{ \operatorname{Tr}(W\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \overline{\Psi}(H) \right\} \\ &= \min_{(\Upsilon, \Theta, H) \in \mathcal{Z}} \max_{W \in \mathcal{W}_1} \left\{ \operatorname{Tr}(W\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \overline{\Psi}(H) \right\} \\ &= \min_{(\Upsilon, \Theta, H) \in \mathcal{Z}} \left\{ \overline{\Psi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \max_{W} \left\{ \operatorname{Tr}(\Theta W) : \\ W \succeq 0, \exists t \in \mathcal{T} : \mathcal{R}_k[W] \preceq t_k I_{d_k}, k \leq K \right\} \right\} \\ &= \min_{(\Upsilon, \Theta, H) \in \mathcal{Z}} \left\{ \overline{\Psi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \max_{W, t} \left\{ \operatorname{Tr}(\Theta W) : \\ W \succeq 0, [t; 1] \in \mathbf{T}, \mathcal{R}_k[W] \preceq t_k I_{d_k}, k \leq K \right\} \right\}, \end{aligned}$$

where \mathbf{T} is the closed conic hull of \mathcal{T} . Now, using Conic Duality combined with the

. .

fact that $\mathbf{T}_* = \{[g; s] : s \ge \phi_{\mathcal{T}}(-g)\}$ we obtain

$$\max_{W,t} \{ \operatorname{Tr}(\Theta W) : W \succeq 0, [t;1] \in \mathbf{K}[\mathcal{T}], \mathcal{R}_k[W] \preceq t_k I_{d_k}, k \leq K \}$$

$$= \min_{Z, [g;s], \Lambda = \{\Lambda_k\}} \begin{cases} s : \begin{cases} Z \succeq 0, [g;s] \in (\mathbf{K}[\mathcal{T}])_*, \Lambda_k \succeq 0, k \leq K \\ -\operatorname{Tr}(ZW) - g^T t + \sum_k \operatorname{Tr}(\mathcal{R}_k^*[\Lambda_k]W) \\ -\sum_k t_k \operatorname{Tr}(\Lambda_k) = \Theta \end{cases}$$

$$= \min_{Z, [g;s], \Lambda = \{\Lambda_k\}} \{ s : \begin{cases} Z \succeq 0, s \geq \phi_{\mathcal{T}}(-g), \Lambda_k \succeq 0, k \leq K \\ \Theta = \sum_k \mathcal{R}_k^*[\Lambda_k] - Z, g = -\lambda[\Lambda] \end{cases} \}$$

$$= \min_{\Lambda} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) : \Lambda = \{\Lambda_k \succeq 0, k \leq K\}, \Theta \preceq \sum_k \mathcal{R}_k^*[\Lambda_k] \right\},$$

and we arrive at

$$\begin{split} \operatorname{Opt}(1) &= \min_{\Upsilon,\Theta,H,\Lambda} & \left\{ \overline{\Psi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{T}}(\lambda[\Lambda]) : \\ \Upsilon &= \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \Lambda = \{\Lambda_k \succeq 0, k \leq K\}, \\ \Theta \preceq \sum_k \mathcal{R}_k^*[\Lambda_k], \\ & \left[\frac{\Theta}{\frac{1}{2}[B^T - A^T H]M} \right] \frac{1}{2}[B^T - A^T H]M}{\frac{1}{2}M^T[B - H^T A]} \right] \succeq 0 \\ &= \min_{\Upsilon,H,\Lambda} & \left\{ \overline{\Psi}(H) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{T}}(\lambda[\Lambda]) : \\ \Upsilon &= \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \Lambda = \{\Lambda_k \succeq 0, k \leq K\} \\ & \left[\frac{\sum_k \mathcal{R}_k^*[\Lambda_k]}{\frac{1}{2}M^T[B - H^T A]} \right] \frac{1}{2}[B^T - A^T H]M}{\sum_\ell \mathcal{S}_\ell^*[\Upsilon_\ell]} \right] \succeq 0 \\ &= \operatorname{Opt} \quad [\operatorname{see} (4.248)]. \end{split}$$

7°. Now we can complete the proof. For $\kappa \in (0,1]$, let W_{κ} be the W-component of a saddle point solution to the saddle point problem (4.258). Then, by (4.261),

$$\sqrt{\kappa} \operatorname{Opt} \leq \operatorname{Opt}(\kappa) = \min_{\substack{(\Upsilon,\Theta,H)\in\mathcal{Z}\\H}} \left\{ \operatorname{Tr}(W_{\kappa}\Theta) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \overline{\Psi}(H) \right\}
= \min_{H} \left\{ \overline{\Phi}(W_{\kappa},H) + \overline{\Psi}(H) \right\}$$
(4.262)

On the other hand, when applying Lemma 4.91 to $Q = W_{\kappa}$ and $\rho = \kappa$, we obtain, in view of relations $0 < \kappa \leq 1, W_{\kappa} \in \mathcal{W}_{\kappa}$,

$$\delta(\kappa) := \operatorname{Prob}_{\zeta \sim \mathcal{N}(0, I_n)} \{ W_{\kappa}^{1/2} \zeta \notin \mathcal{X} \} \le 2De^{-\frac{1}{2\kappa}}, \qquad (4.263)$$

with D given by (4.63). In particular, when setting

$$\bar{\kappa} = \frac{1}{2\ln D + 10\ln 2} \tag{4.264}$$

we obtain $\delta_{\kappa} \leq 1/16$. Therefore,

$$\operatorname{Prob}_{\eta \sim \mathcal{N}(0, W_{\bar{\kappa}})} \{ \eta \notin \mathcal{X} \} \le \frac{1}{16}.$$

$$(4.265)$$

Now let

$$\varrho_* := \frac{\text{Opt}}{8\sqrt{(2\ln F + 10\ln 2)(2\ln D + 10\ln 2)}}.$$
(4.266)

All we need in order to achieve our goal, that is, to justify (4.62), is to show that

$$\operatorname{RiskOpt}_{\frac{1}{2}} \ge \varrho_*, \tag{4.267}$$

since given the latter relation, (4.62) will be immediately given by (4.250) as applied with $\epsilon = \frac{1}{8}$.

To prove (4.267), assume, on the contrary to what should be proved, that the $\frac{1}{8}$ -risk is $\langle \varrho_*$, and let $\bar{x}(\cdot)$ be an estimate with $\frac{1}{8}$ -risk $\leq \varrho_*$. We can utilize \bar{x} to estimate $B\eta$, in the Bayesian problem of recovering $B\eta$ from observation $\omega = A\eta + \xi$, $(\eta, \xi) \sim \mathcal{N}(0, \Sigma)$ with $\Sigma = \text{Diag}\{W_{\bar{\kappa}}, Q_*\}$. From (4.265) we conclude that

$$\begin{aligned} \operatorname{Prob}_{(\eta,\xi)\sim\mathcal{N}(0,\Sigma)}\{\|B\eta - \bar{x}(A\eta + \xi)\| > \varrho_*\} \\ &\leq \operatorname{Prob}_{(\eta,\xi)\sim\mathcal{N}(0,\Sigma)}\{\|B\eta - \bar{x}(A\eta + \xi)\| > \varrho_*, \ \eta \in \mathcal{X}\} + \operatorname{Prob}_{\eta\sim\mathcal{N}(0,W_{\bar{\kappa}})}\{\eta \notin \mathcal{X}\} \\ &\leq \frac{1}{8} + \frac{1}{16} = \frac{3}{16}. \end{aligned}$$

$$(4.268)$$

On the other hand, by (4.262) we have

$$\min_{H} \left[\overline{\Phi}(W_{\bar{\kappa}}, H) + \overline{\Psi}(H) \right] = \operatorname{Opt}(\bar{\kappa}) \ge \sqrt{\bar{\kappa}} \operatorname{Opt} = [8\bar{\varkappa}] \varrho_*$$

with $\bar{\varkappa}$ given by (4.256), so by Lemma 4.94, for any estimate $\hat{x}(\cdot)$ of $B\eta$ via observation $\omega = Ax + \xi$ it holds

$$\operatorname{Prob}_{\eta,\xi}\{\|B\eta - \widehat{x}(A\eta + \xi)\| \ge \varrho_*\} \ge \beta_{\bar{\varkappa}} > 3/16;$$

in particular, this relation should hold true for $\hat{x}(\cdot) \equiv \bar{x}(\cdot)$, but the latter is impossible: the $\frac{1}{8}$ -risk of \bar{x} is $\leq \varrho_*$, see (4.268).

4.10.5.2 Proof of Proposition 4.5

We shall extract Proposition 4.5 from the following result, meaningful by its own right (it can be considered as "ellitopic refinement" of Proposition 4.16):

Proposition 4.95. Consider recovery of the linear image $Bx \in \mathbf{R}^{\nu}$ of unknown signal x known to belong to a given signal set $\mathcal{X} \subset \mathbf{R}^n$ from noisy observation

$$\omega = Ax + \xi \in \mathbf{R}^m \qquad [\xi \sim \mathcal{N}(0, \Gamma), \, \Gamma \succ 0]$$

the recovery error being measured in a norm $\|\cdot\|$ on \mathbf{R}^{ν} . Assume that \mathcal{X} and the unit ball \mathcal{B}_* of the norm $\|\cdot\|_*$ conjugate to $\|\cdot\|$ are ellitopes:

$$\mathcal{X} = \{ x \in \mathbf{R}^n : \exists t \in \mathcal{T} : x^T R_k x \le t_k, k \le K \}$$

$$\mathcal{B}_* = \{ y \in \mathbf{R}^\nu : \exists (r \in \mathcal{R}, y) : u = My, y^T S_\ell y \le r_\ell, \ell \le L \}$$
(4.269)

with our standard restrictions on \mathcal{T} , \mathcal{R} , R_k and S_ℓ (as always, we lose nothing when assuming that the ellitope \mathcal{X} is basic).

Consider the optimization problem

$$Opt_{\#} = \min_{\substack{\Theta, H, \lambda, \mu, \mu' \\ l \ge 0, \mu \ge 0, \mu' \ge 0}} \left\{ \phi_{\mathcal{T}}(\lambda) + \phi_{\mathcal{R}}(\mu) + \phi_{\mathcal{R}}(\mu') + Tr(\Gamma\Theta) : \\ \lambda \ge 0, \mu \ge 0, \mu' \ge 0 \\ \left[\frac{\sum_{k} \lambda_{k} R_{k}}{\frac{1}{2} [B - H^{T}A]^{T} M} \right] \sum_{\ell} \mu_{\ell} S_{\ell} \right] \succeq 0$$

$$\left[\frac{\Theta}{\frac{1}{2} M^{T} H^{T}} \left| \sum_{\ell} \mu_{\ell}' S_{\ell} \right] \ge 0$$

$$(4.270)$$

The problem is solvable, and the linear estimate $\hat{x}_{H_*}(\omega) = H_*^T \omega$ yielded by the *H*-component of an optimal solution to the problem satisfies the risk bound

$$\operatorname{Risk}_{\Gamma, \|\cdot\|} [\widehat{x}_{H_*} | \mathcal{X}] := \max_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim \mathcal{N}(0, \Gamma)} \left\{ \|Bx - \widehat{x}_{H_*} (Ax + \xi)\| \right\} \le \operatorname{Opt}_{\#}.$$

Besides this, the estimate is near-optimal:

$$Opt_{\#} \le 64\sqrt{(3\ln K + 15\ln 2)(3\ln L + 15\ln 2)}RiskOpt, \qquad (4.271)$$

where RiskOpt is the minimax optimal risk:

RiskOpt =
$$\inf_{\widehat{x}} \sup_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim \mathcal{N}(0,\Gamma)} \{ \|Bx - \widehat{x}(Ax + \xi)\| \},\$$

the infimum being taken w.r.t. all estimates.

Proposition 4.95 \Rightarrow **Proposition 4.5:** The situation considered in Proposition 4.5 is the special case of the situation considered in Proposition 4.95, namely, the case when \mathcal{B}_* is the standard Euclidean ball:

$$\mathcal{B}_* = \{ u \in \mathbf{R}^\nu : u^T u \le 1 \}.$$

In this case, problem (4.270) reads

Comparing the resulting representation of $\operatorname{Opt}_{\#}$ with (4.13), we see that the upper bound $\sqrt{\operatorname{Opt}}$ on the risk of the linear estimate \hat{x}_{H_*} appearing in (4.16) is $\leq \operatorname{Opt}_{\#}$. Combining this observation with (4.271) and the evident relation

$$\operatorname{RiskOpt} = \inf_{\widehat{x}} \sup_{x \in \mathcal{X}} \mathbf{E}_{x \sim \mathcal{N}(0,\Gamma)} \{ \|Bx - \widehat{x}(Ax + \xi)\|_2 \}$$

$$\leq \inf_{\widehat{x}} \sqrt{\sup_{x \in \mathcal{X}} \mathbf{E}_{x \sim \mathcal{N}(0,\Gamma)}} \{ \|Bx - \widehat{x}(Ax + \xi)\|_2 \} = \operatorname{Risk}_{\operatorname{opt}},$$

(recall that we are in the case of $\|\cdot\| = \|\cdot\|_2$), we arrive at (4.16) and thus justify Proposition 4.5.

Proof of Proposition 4.95. It is immediately seen that problem (4.270) is nothing but problem (4.50) in the case when the spectratopes $\mathcal{X}, \mathcal{B}_*$ and the set Π participating in Proposition 4.14 are, respectively, the ellitopes given by (4.269), and the singleton $\{\Gamma\}$. Thus, Proposition 4.95 is, essentially, the special case of Proposition 4.16. The only refinement in Proposition 4.95 as compared to Proposition 4.16 is the form of the logarithmic "non-optimality" factor in (4.271); similar factor in Proposition 4.16 is expressed in terms of spectratopic sizes D, F of \mathcal{X} and \mathcal{B} (the total ranks of matrices R_k , $k \leq K$, and S_{ℓ} , $\ell \leq L$, in the case of (4.269)), while in (4.271) the nonoptimality factor is expressed in terms of ellitopic sizes K, L of \mathcal{X} and \mathcal{B}_* . Strictly speaking, to arrive at this (slight – the sizes in question are under logs) refinement, we were supposed to reproduce, with minimal modifications, the reasoning of items $2^0 - 7^0$ of Section 4.10.5.1, with Γ in the role of Q_* , and slightly refine Lemma 4.17 underlying this reasoning. It would be counter-productive to carry out this course of actions literally; instead, we intend to indicate "local modifications" to be made in the proof of Proposition 4.16 in order to prove Proposition 4.95. Here are the modifications:

A. The collections of matrices $\Lambda = \{\Lambda_k \succeq 0, k \leq K\}, \Upsilon = \{\Upsilon_\ell \succeq 0, \ell \leq L\}$ should be substituted by collections of nonnegative reals $\lambda \in \mathbf{R}_+^K$, resp., $\mu \in \mathbf{R}_+^L$, and

vectors $\lambda[\Lambda]$, $\lambda[\Upsilon]$ – with vectors λ , resp. μ . Expressions like $\mathcal{R}_k[W]$, $\mathcal{R}_k^*[\Lambda_k]$, $\mathcal{S}_{\ell}^*[\Upsilon_{\ell}]$ should be replaced, respectively, with $\operatorname{Tr}(R_k W)$, $\lambda_k R_k$, $\mu_{\ell} S_{\ell}$. Finally, Q_* should be replaced with Γ , and scalar matrices, like $t_k I_{d_k}$, should be replaced with the corresponding reals, like t_k .

B. The role of Lemma 4.17 is now played by

Lemma 4.96. Let Y be an $N \times \nu$ matrix, let $\|\cdot\|$ be a norm on \mathbf{R}^{ν} such that the unit ball \mathcal{B}_* of the conjugate norm is the ellitope (4.269):

$$\mathcal{B}_* = \{ y \in \mathbf{R}^{\nu} : \exists (r \in \mathcal{R}, y) : u = My, y^T S_\ell y \le r_\ell, \, \ell \le L \},$$
(4.269)

and let $\zeta \sim \mathcal{N}(0, Q)$ for some positive semidefinite $N \times N$ matrix Q. Then the best upper bound on $\psi_Q(Y) := \mathbf{E}\{||Y^T\zeta||\}$ yielded by Lemma 4.11, that is, the optimal value $\operatorname{Opt}[Q]$ in the convex optimization problem (cf. (4.48))

$$\operatorname{Opt}[Q] = \min_{\Theta,\mu} \left\{ \phi_{\mathcal{R}}(\mu) + \operatorname{Tr}(Q\Theta) : \mu \ge 0, \left[\frac{\Theta}{\frac{1}{2}M^{T}Y^{T}} \middle| \sum_{\ell} \mu_{\ell} R_{\ell} \right] \ge 0 \right\}$$
(4.272)

satisfies the identity

$$\forall (Q \succeq 0) :$$

$$\operatorname{Opt}[Q] = \overline{\operatorname{Opt}}[Q] := \min_{G,\mu} \left\{ \phi_{\mathcal{R}}(\mu) + \operatorname{Tr}(G) : \mu \ge 0, \left[\frac{G}{\frac{1}{2}M^T Y^T Q^{1/2}} \middle| \frac{1}{\sum_{\ell} \mu_{\ell} R_{\ell}} \right] \succeq 0 \right\}$$

$$(4.273)$$

and is a tight bound on $\psi_Q(Y)$, namely,

$$\psi_Q(Y) \le \text{Opt}[Q] \le 22\sqrt{3\ln L + 15\ln 2}\psi_Q(Y),$$
 (4.274)

where L is the size of the ellitope \mathcal{B}_* , see (4.269). Besides this, for all $\varkappa \geq 1$ one has

$$\operatorname{Prob}_{\zeta}\left\{\|Y^{T}\zeta\| \geq \frac{\operatorname{Opt}[Q]}{4\varkappa}\right\} \geq \beta_{\varkappa} := 1 - \frac{e^{3/8}}{2} - 2Le^{-\varkappa^{2}/3}.$$
(4.275)

In particular, when selecting $\varkappa = \sqrt{3 \ln L + 15 \ln 2}$, we obtain

$$\operatorname{Prob}_{\zeta} \left\{ \|Y^{T}\zeta\| \ge \frac{\operatorname{Opt}[Q]}{4\sqrt{3\ln L + 15\ln 2}} \right\} \ge \beta_{\varkappa} = 0.2100 > \frac{3}{16}.$$
(4.276)

The proof of Lemma 4.96 follows the one of Lemma 4.17, with Lemma 4.92 substituting Lemma 4.91.

Proof of Lemma 4.96

 1^o .

Relation (4.273) can be verified exactly in the same fashion as in the case of Lemma 4.17.

 2^{o} .

432

Let us set $\zeta = Q^{1/2}\eta$ with $\eta \sim \mathcal{N}(0, I_N)$ and $Z = Q^{1/2}Y$. Let us show that when $\varkappa \geq 1$ one has

$$\operatorname{Prob}_{\eta}\{\|Z^{T}\eta\| \geq \bar{\delta}\} \geq \beta_{\varkappa} := 1 - \frac{e^{3/8}}{2} - 2Le^{-\varkappa^{2}/3}, \quad \bar{\delta} = \frac{\operatorname{Opt}[Q]}{4\varkappa}, \quad (4.277)$$

where

$$[\overline{\mathrm{Opt}}[Q] =] \quad \mathrm{Opt}[Q] := \min_{\Theta,\mu} \left\{ \phi_{\mathcal{R}}(\mu) + \mathrm{Tr}(\Theta) : \mu \ge 0, \left[\begin{array}{c|c} \Theta & \frac{1}{2}ZM \\ \hline \frac{1}{2}M^TZ^T & \sum_{\ell} \mu_{\ell}R_{\ell} \\ \end{array} \right] \ge 0 \right\}$$
(4.278)

 3^{o} .

Let us represent Opt[Q] as the optimal value of a conic problem. Setting

$$\mathbf{K} = \mathrm{cl}\{[r; s] : s > 0, r/s \in \mathcal{R}\},\$$

we ensure that

$$\mathcal{R} = \{r : [r;1] \in \mathbf{K}\}, \ \mathbf{K}_* = \{[g;s] : s \ge \phi_{\mathcal{R}}(-g)\},\$$

where \mathbf{K}_* is the cone dual to \mathbf{K} . Consequently, (4.278) reads

$$\operatorname{Opt}[Q] = \min_{\Theta, \Upsilon, \theta} \left\{ \theta + \operatorname{Tr}(\Theta) : \begin{array}{cc} \mu \ge 0 & (a) \\ \left[\frac{\Theta}{\frac{1}{2}M^T Z^T} \middle| \sum_{\ell} \mu_{\ell} S_{\ell} \right] \ge 0 & (b) \\ \left[-\mu; \theta \right] \in \mathbf{K}_* & (c) \end{array} \right\}.$$
(P)

 4^{o} .

Now let us prove that there exists matrix $W \in \mathbf{S}^q_+$ and $r \in \mathcal{R}$ such that

$$\operatorname{Tr}(WS_{\ell}) \le r_{\ell}, \ \ell \le L, \tag{4.279}$$

and

$$\operatorname{Opt}[Q] \le \sum_{i} \sigma_{i}(ZMW^{1/2}), \qquad (4.280)$$

where $\sigma_1(\cdot) \ge \sigma_2(\cdot) \ge \dots$ are singular values.

To get the announced result, let us pass from problem (P) to its conic dual. Applying Lemma 4.89 we conclude that (P) is strictly feasible; in addition, (P) clearly is bounded, so that the dual to (P) problem (D) is solvable with optimal value Opt[Q]. Let us build (D). Denoting by $\lambda_{\ell} \geq 0, \ell \leq L, \left[\begin{array}{c|c} G & -R \\ \hline -R^T & W \end{array} \right] \succeq 0, [r; \tau] \in \mathbf{K}$ the Lagrange multipliers for the respective constraints in (P), and aggregating these constraints, the multipliers being the aggregation weights, we arrive at the following aggregated constraint:

$$\operatorname{Tr}(\Theta G) + \operatorname{Tr}(W \sum_{\ell} \mu_{\ell} S_{\ell}) + \sum_{\ell} \lambda_{\ell} \mu_{\ell} - \sum_{\ell} r_{\ell} \mu_{\ell} + \theta \tau \ge \operatorname{Tr}(ZMR^{T}).$$

To get the dual problem, we impose on the Lagrange multipliers, in addition to the initial conic constraints like $\lambda_{\ell} \geq 0$, $1 \leq \ell \leq L$, the restriction that the left hand side in the aggregated constraint, identically in Θ , μ_{ℓ} and θ , is equal to the objective of (P), that is,

$$G = I, \operatorname{Tr}(WS_{\ell}) + \lambda_{\ell} - r_{\ell} = 0, \ 1 \leq \ell \leq L, \ \tau = 1,$$

and maximize, under the resulting restrictions, the right-hand side of the aggregated constraint. After immediate simplifications, we arrive at

$$Opt[Q] = \max_{W,R,r} \left\{ Tr(ZMR^T) : W \succeq R^T R, r \in \mathcal{R}, Tr(WS_\ell) \le r_\ell, 1 \le \ell \le L \right\}$$

(note that $r \in \mathcal{R}$ is equivalent to $[r; 1] \in \mathbf{K}$, and $W \succeq R^T R$ is the same as $\begin{bmatrix} I & | & -R \\ \hline & -R^T & | & W \end{bmatrix} \succeq 0$). Now, to say that $R^T R \preceq W$ is exactly the same as to say that $R = SW^{1/2}$ with the spectral norm $\|S\|_{\mathrm{Sh},\infty}$ of S not exceeding 1, so that

$$\operatorname{Opt}[Q] = \max_{W,S,r} \left\{ \underbrace{\operatorname{Tr}([ZM[SW^{1/2}]^T)]}_{=\operatorname{Tr}([ZMW^{1/2}]S^T)} : W \succeq 0, \|S\|_{\operatorname{Sh},\infty} \le 1, r \in \mathcal{R}, \operatorname{Tr}(WS_{\ell}) \le r_{\ell}, \, \ell \le L \right\}$$

and we can immediately eliminate the S-variable, using the well-known fact that for every $p \times q$ matrix J, it holds

$$\max_{S \in \mathbf{R}^{p \times q}, \|S\|_{\mathrm{Sh}, \infty} \le 1} \operatorname{Tr}(JS^T) = \|J\|_{\mathrm{Sh}, 1},$$

where $||J||_{\text{Sh},1}$ is the nuclear norm (the sum of singular values) of J. We arrive at

$$\operatorname{Opt}[Q] = \max_{W,r} \left\{ \|ZMW^{1/2}\|_{\operatorname{Sh},1} : r \in \mathcal{R}, W \succeq 0, \operatorname{Tr}(WS_{\ell}) \le r_{\ell}, \, \ell \le L \right\}.$$

The resulting problem clearly is solvable, and its optimal solution W ensures the target relations (4.279) and (4.280).

 5^{o} .

Given W satisfying (4.279) and (4.280), let $UJV = W^{1/2}M^TZ^T$ be the singular value decomposition of $W^{1/2}M^TZ^T$, so that U and V are, respectively, $q \times q$ and $N \times N$ orthogonal matrices, J is $q \times N$ matrix with diagonal $\sigma = [\sigma_1; ...; \sigma_p], p = \min[q, N]$, and zero off-diagonal entries; the diagonal entries $\sigma_i, 1 \leq i \leq p$ are the singular values of $W^{1/2}M^TZ^T$, or, which is the same, of $ZMW^{1/2}$. Therefore, we have

$$\sum_{i} \sigma_i \ge \operatorname{Opt}[Q]. \tag{4.281}$$

Now consider the following construction. Let $\eta \sim \mathcal{N}(0, I_N)$; we denote by v the vector comprised of the first p entries in $V\eta$; note that $v \sim \mathcal{N}(0, I_p)$, since V is orthogonal. We then augment, if necessary, v by q - p independent of each other and of $\eta \mathcal{N}(0, 1)$ random variables to obtain a q-dimensional normal vector $v' \sim \mathcal{N}(0, I_q)$, and set $\chi = Uv'$; because U is orthogonal we also have $\chi \sim \mathcal{N}(0, I_q)$.

Observe that

$$\chi^{T} W^{1/2} M^{T} Z^{T} \eta = \chi^{T} U J V \eta = [\upsilon']^{T} J \upsilon = \sum_{i=1}^{p} \sigma_{i} \upsilon_{i}^{2}.$$
(4.282)

To continue we need two simple observations.

(i) One has

$$\alpha := \operatorname{Prob}\left\{\sum_{i=1}^{p} \sigma_{i} v_{i}^{2} < \frac{1}{4} \sum_{i=1}^{p} \sigma_{i}\right\} \le \frac{e^{3/8}}{2} \ [= 0.7275...].$$
(4.283)

The claim is evident when $\sigma := \sum_i \sigma_i = 0$. Now let $\sigma > 0$, and let us apply the Cramer bounding scheme. Namely, given $\gamma > 0$, consider the random variable

$$\omega = \exp\left\{\frac{1}{4}\gamma \sum_{i} \sigma_{i} - \gamma \sum_{i} \sigma_{i} v_{i}^{2}\right\}.$$

Note that $\omega > 0$ a.s., and is > 1 when $\sum_{i=1}^{p} \sigma_i v_i^2 < \frac{1}{4} \sum_{i=1}^{p} \sigma_i$, so that $\alpha \leq \mathbf{E}\{\omega\}$, or, equivalently, thanks to $v \sim \mathcal{N}(0, I_p)$,

$$\ln(\alpha) \le \ln(\mathbf{E}\{\omega\}) = \frac{1}{4}\gamma \sum_{i} \sigma_i + \sum_{i} \ln\left(\mathbf{E}\{\exp\{-\gamma\sigma_i v_i^2\}\}\right) \le \frac{1}{4}\gamma\sigma - \frac{1}{2}\sum_{i} \ln(1+2\gamma\sigma_i).$$

Function $-\sum_{i} \ln(1 + 2\gamma \sigma_i)$ is convex in $[\sigma_1; ...; \sigma_p] \ge 0$, therefore, its maximum over the simplex $\{\sigma_i \ge 0, i \le p, \sum_i \sigma_i = \sigma\}$ is attained at a vertex, and we get

$$\ln(\alpha) \le \frac{1}{4}\gamma\sigma - \frac{1}{2}\ln(1+2\gamma\sigma).$$

Minimizing the right hand side in $\gamma > 0$, we arrive at (4.283). (ii) Whenever $\varkappa \ge 1$, one has

$$\operatorname{Prob}\{\|MW^{1/2}\chi\|_* > \varkappa\} \le 2L \exp\{-\varkappa^2/3\}, \qquad (4.284)$$

with L coming from (4.269).

Indeed, setting $\rho = 1/\varkappa^2 \leq 1$ and $\omega = \sqrt{\rho}W^{1/2}\chi$, we get $\omega \sim \mathcal{N}(0, \rho W)$. Let us apply Lemma 4.92 to $Q = \rho W$, \mathcal{R} in the role of \mathcal{T} , L in the role of K, and S_{ℓ} 's in the role of R_k 's. Denoting

$$\mathcal{Y} := \{ y : \exists r \in \mathcal{R} : y^T S_\ell y \preceq r_\ell, \ell \leq L \},\$$

we have $\operatorname{Tr}(QS_{\ell}) = \rho \operatorname{Tr}(WS_{\ell}) = \rho \operatorname{Tr}(WS_{\ell}) \leq \rho r_{\ell}, \ \ell \leq L$, with $r \in \mathcal{R}$ (see (4.279)), so we are under the premise of Lemma 4.92 (with \mathcal{Y} in the role of \mathcal{X} and therefore with L in the role of K). Applying the lemma, we conclude that

$$\operatorname{Prob}\left\{\chi: \varkappa^{-1} W^{1/2} \chi \notin \mathcal{Y}\right\} \le 2L \exp\{-1/(3\rho)\} = 2L \exp\{-\varkappa^2/3\}.$$

Recalling that $\mathcal{B}_* = M\mathcal{Y}$, we see that $\operatorname{Prob}\{\chi : \varkappa^{-1}MW^{1/2}\chi \notin \mathcal{B}_*\}$ is indeed upper-bounded by the right hand side of (4.284), and (4.284) follows.

Now, for $\varkappa \geq 1$, let

$$E_{\varkappa} = \left\{ (\chi, \eta) : \| M W^{1/2} \chi \|_{\ast} \le \varkappa, \sum_{i} \sigma_{i} v_{i}^{2} \ge \frac{1}{4} \sum_{i} \sigma_{i} \right\}.$$

For $(\chi, \eta) \in E_{\varkappa}$ we have

$$\varkappa \|Z^{T}\eta\| \ge \|MW^{1/2}\chi\|_{*}\|Z^{T}\eta\| \ge \chi^{T}W^{1/2}M^{T}Z^{T}\eta = \sum_{i}\sigma_{i}\upsilon_{i}^{2} \ge \frac{1}{4}\sum_{i}\sigma_{i} \ge \frac{1}{4}\text{Opt}[Q],$$

(we have used (4.282) and (4.281)). On the other hand, due to (4.283) and (4.284),

$$\operatorname{Prob}\{E_{\varkappa}\} \ge \beta_{\varkappa},\tag{4.285}$$

and we arrive at (4.275). The latter relation clearly implies (4.276) which, in turn, implies the right inequality in (4.274).

C. As a result of substituting Lemma 4.17 with Lemma 4.96, the analogy of Lemma 4.94 used in item 4^0 of the proof of Proposition 4.16 now reads as follows:

Lemma 4.97. Let W be a positive semidefinite $n \times n$ matrix, and $\varkappa \geq 1$. Then for any estimate $\hat{x}(\cdot)$ of $B\eta$ given observation $\omega = A\eta + \xi$, one has

$$\operatorname{Prob}_{\eta,\xi}\{\|B\eta - \widehat{x}(\omega)\| \ge [8\varkappa]^{-1} \inf_{H}[\overline{\Psi}(H) + \overline{\Phi}(W,H)]\} \ge \beta_{\varkappa} = 1 - \frac{e^{3/8}}{2} - 2Le^{-\varkappa^2/3}$$

where $\overline{\Psi}(H)$ and $\overline{\Phi}(W,H)$ are defined, respectively, by (4.248) and (4.255). In particular, for

$$\varkappa = \bar{\varkappa} := \sqrt{3\ln K} + 15\ln 2 \tag{4.286}$$

the latter probability is > 3/16.

- D. Reference to Lemma 4.91 in item 7^0 of the proof should be substituted with reference to Lemma 4.92, resulting in replacing
 - relation (4.263 with the relation

$$\delta(\kappa) := \operatorname{Prob}_{\zeta \sim \mathcal{N}(0, I_n)} \{ W_{\kappa}^{1/2} \zeta \notin \mathcal{X} \} \leq 3K e^{-\frac{1}{3\kappa}},$$

• relation (4.264) with the relation

$$\bar{\kappa} = \frac{1}{3\ln K + 15\ln 2};$$

• relation (4.266) – with the relation

$$\varrho_* := \frac{\text{Opt}}{8\sqrt{(3\ln L + 15\ln 2)(3\ln K + 15\ln 2)}}.$$

4.10.6 Proofs of Propositions 4.18, 4.19, and justification of Remark 4.20

4.10.6.1 Proof of Proposition 4.18

The only claim in Proposition which is not an immediate consequence of Proposition 4.8 is that problem (4.73) is solvable; let us justify this claim. Let F = ImA. Clearly, feasibility of a candidate solution (H, Λ, Υ) to the problem depends solely on the restriction of the linear mapping $z \mapsto H^T z$ onto F, so that adding to the constraints of the problem the requirement that the restriction of this linear mapping on the orthogonal complement of F in \mathbb{R}^m is identically zero, we get an equivalent problem. It is immediately seen that in the resulting problem, the feasible solutions with the value of the objective $\leq a$ for every $a \in \mathbb{R}$ form a compact set, so that the latter problem (and thus – the original one) indeed is solvable. \Box

4.10.6.2 Proof of Proposition 4.19

We are about to derive Proposition 4.19 from Proposition 4.16. Observe that in the situation of the latter Proposition, setting formally $\Pi = \{0\}$, problem (4.50) becomes problem (4.73), so that Proposition 4.19 looks as the special case $\Pi = \{0\}$ of Proposition 4.16. However, the premise of the latter Proposition forbids specializing Π as $\{0\}$ – this would violate the regularity assumption \mathbf{R} which is part of the premise. The difficulty, however, can be easily resolved. Assume w.l.o.g. that the image space of A is the entire \mathbf{R}^m (otherwise we could from the very beginning replace \mathbf{R}^m with the image space of A, and let us pass our current noiseless recovery problem of interest (!) to its "noisy modification," the differences with (!) being

- noisy observation $\omega = Ax + \sigma\xi$, $\sigma > 0$, $\xi \sim \mathcal{N}(0, I_1m)$;
- risk quantification of a candidate estimate $\hat{x}(\cdot)$ according to

$$\operatorname{Risk}_{\|\cdot\|}^{\sigma}[\widehat{x}(Ax+\sigma\xi)|\mathcal{X}] = \sup_{x\in\mathcal{X}} \mathbf{E}_{\xi\sim\mathcal{N}(0,I_m)}\left\{\|Bx-\widehat{x}(Ax+\sigma\xi)\}\right\},$$

the corresponding minimax optimal risk being

$$\operatorname{RiskOpt}_{\|\cdot\|}^{\sigma}[\mathcal{X}] = \inf_{\widehat{x}(\cdot)} \operatorname{Risk}_{\|\cdot\|}^{\sigma}[\widehat{x}(Ax + \sigma\xi)|\mathcal{X}]$$

Proposition 4.16 does apply to the modified problem – it suffices to specify Π as $\{\sigma^2 I_m\}$; according to this Proposition, the quantity

$$\begin{aligned} \operatorname{Opt}[\sigma] &= \min_{\substack{H,\Lambda,\Upsilon,\Upsilon',\Theta\\ H,\Lambda,\Upsilon,\Upsilon',\Theta}} \left\{ \phi_{\mathcal{T}}(\lambda[\Lambda]) + \phi_{\mathcal{R}}(\lambda[\Upsilon]) + \phi_{\mathcal{R}}(\lambda[\Upsilon']) + \sigma^{2}\operatorname{Tr}(\Theta) : \\ \Lambda &= \{\Lambda_{k} \succeq 0, k \leq K\}, \ \Upsilon = \{\Upsilon_{\ell} \succeq 0, \ell \leq L\}, \ \Upsilon' &= \{\Upsilon'_{\ell} \succeq 0, \ell \leq L\}, \\ \left[\frac{\sum_{k} \mathcal{R}_{k}^{*}[\Lambda_{k}]}{\frac{1}{2}M^{T}[B - H^{T}A]} \middle| \frac{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon_{\ell}]}{\sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon'_{\ell}]} \right] \succeq 0, \\ \left[\frac{\Theta}{\frac{1}{2}M^{T}H^{T}} \middle| \sum_{\ell} \mathcal{S}_{\ell}^{*}[\Upsilon'_{\ell}]} \right] \succeq 0 \end{aligned} \end{aligned}$$

satisfies the relation

$$\operatorname{Opt}[\sigma] \le O(1) \ln(D) \operatorname{RiskOpt}_{\parallel \cdot \parallel}^{\sigma} [\mathcal{X}]$$
(4.287)

with D defined in (4.74). Looking at problem (4.73) we immediately conclude that $Opt_{\#} \leq Opt[\sigma]$. Thus, all we need in order to extract the target relation (4.74) from

(4.287) is to prove that the minimax optimal risk $\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}]$ defined in Proposition 4.19 satisfies the relation

$$\liminf_{\sigma \to +0} \operatorname{RiskOpt}_{\|\cdot\|}^{\sigma}[\mathcal{X}] \le \operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}].$$
(4.288)

To prove this relation let us fix $r > \text{Risk}_{\text{opt}}[\mathcal{X}]$, so that for some Borel estimate $\hat{x}(\cdot)$ it holds

$$\sup_{x \in \mathcal{X}} \|Bx - \widehat{x}(Ax)\| < r.$$
(4.289)

Were we able to ensure that $\hat{x}(\cdot)$ is bounded and continuous, we would be done, since in this case, due to compactness of \mathcal{X} , it clearly holds

$$\begin{split} &\lim \inf_{\sigma \to +0} \operatorname{RiskOpt}_{\|\cdot\|}^{\sigma} [\mathcal{X}] \\ &\leq \lim \inf_{\sigma \to +0} \sup_{x \in \mathcal{X}} \mathbf{E}_{\xi \sim \mathcal{N}(0, I_m)} \left\{ \|Bx - \widehat{x}(Ax + \sigma\xi)\| \right\} \\ &\leq \sup_{x \in \mathcal{X}} \|Bx - \widehat{x}(Ax)\| < r, \end{split}$$

and since $r > \text{Risk}_{\text{opt}}[\mathcal{X}]$ is arbitrary, (4.288) would follow. Thus, all we need to do is to verify that given Borel estimate $\hat{x}(\cdot)$ satisfying (4.289), we can update it into a bounded and continuous estimate satisfying the same relation. Verification is as follows:

1. Setting $\beta = \max_{x \in \mathcal{X}} \|Bx\|$ and replacing estimate \hat{x} with its truncation

$$\widetilde{x}(\omega) = \begin{cases} \widehat{x}(\omega), & \|\widehat{x}(\omega)\| \le 2\beta \\ 0, & \text{otherwise} \end{cases}$$

we only reduce the norm of the recovery error, whatever be a signal from \mathcal{X} ; at the same time, \tilde{x} is Borel and bounded. Thus, we lose nothing when assuming in the rest of the proof that $\hat{x}(\cdot)$ is Borel and bounded.

2. For $\epsilon > 0$, let $\hat{x}_{\epsilon}(\omega) = (1 + \epsilon)\hat{x}(\omega/(1 + \epsilon))$ and let $\mathcal{X}_{\epsilon} = (1 + \epsilon)\mathcal{X}$. Observe that

$$\begin{aligned} \sup_{x \in \mathcal{X}_{\epsilon}} \|Bx - \widehat{x}_{\epsilon}(Ax)\| &= \sup_{y \in \mathcal{X}} \|B[1+\epsilon]y - \widehat{x}_{\epsilon}(A[1+\epsilon]y)\| \\ &= \sup_{y \in \mathcal{X}} \|B[1+\epsilon]y - [1+\epsilon]\widehat{x}(Ay)\| = [1+\epsilon] \sup_{y \in \mathcal{X}} \|By - \widehat{x}(Ay)\|, \end{aligned}$$

implying, in view of (4.289), that for small enough positive ϵ we have

$$\bar{r} := \sup_{x \in \mathcal{X}_{\epsilon}} \|Bx - \hat{x}_{\epsilon}(Ax)\| < r.$$
(4.290)

3. Finally, let A^{\dagger} be the pseudoinverse of A, so that $AA^{\dagger}z = z$ for every $z \in \mathbf{R}^m$ (recall that the image space of A is the entire \mathbf{R}^m). Given $\rho > 0$, let $\theta_{\rho}(\cdot)$ be a nonnegative smooth function on \mathbf{R}^m with integral 1 such that θ_{ρ} vanishes outside of the ball of radius ρ centered at the origin, and let

$$\widehat{x}_{\epsilon,\rho}(\omega) = \int_{\mathbf{R}^m} \widehat{x}_{\epsilon}(\omega - z)\theta_{\rho}(z)dz$$

be the convolution of \hat{x}_{ϵ} and θ_{ρ} ; since $\hat{x}_{\epsilon}(\cdot)$ is Borel and bounded, this convolution is well-defined smooth function on \mathbf{R}^{m} . Since \mathcal{X} contains a neighbourhood of the origin, for all small enough $\rho > 0$, all z from the support of θ_{ρ} and all $x \in \mathcal{X}$

the point $x - A^{\dagger}z$ belongs to \mathcal{X}_{ϵ} . For such a ρ and any $x \in \mathcal{X}$ we have

$$\begin{aligned} \|Bx - \widehat{x}_{\epsilon}(Ax - z)\| &= \|Bx - \widehat{x}_{\epsilon}(A[x - A^{\dagger}z])\| \\ &\leq \|BA^{\dagger}z\| + \|B[x - A^{\dagger}z] - \widehat{x}_{\epsilon}(A[x - A^{\dagger}z])\| \\ &\leq C\rho + \bar{r} \end{aligned}$$

with properly selected constant C independent of ρ (we have used (4.290); note that for our ρ and x, $x - A^{\dagger}z \in \mathcal{X}_{\epsilon}$). We conclude that for properly selected $r' < r, \rho > 0$ and all $x \in \mathcal{X}$ we have

$$\|Bx - \widehat{x}_{\epsilon}(Ax - z)\| \le r' \,\forall (z \in \operatorname{supp} \theta_{\rho}),$$

implying, by construction of $\hat{x}_{\epsilon,\rho}$, that

$$\forall (x \in \mathcal{X}) : \|Bx - \widehat{x}_{\epsilon,\rho}(Ax)\| \le r' < r.$$

The resulting estimate $\hat{x}_{\epsilon,\rho}$ is the continuous and bounded estimate satisfying (4.289) we were looking for.

4.10.6.3 Justification of Remark 4.20

Justification of Remark is given by repeating word by word the proof of Proposition 4.19, with Proposition 4.95 in the role of Proposition 4.16.

4.10.7 **Proof of (4.96)**

Let $h \in \mathbf{R}^m$, and let ω be random vector with independent across *i* entries $\omega_i \sim \text{Poisson}(\mu_i)$. Taking into account that ω_i are independent across *i*, we have

$$\mathbf{E} \left\{ \exp\{\gamma h^T \omega\} \right\} = \prod_i \mathbf{E} \left\{ \gamma h_i \omega_i \right\} = \prod_i \exp\{ [\exp\{\gamma h_i\} - 1] \mu_i \}$$

=
$$\exp\{ \sum_i [\exp\{\gamma h_i\} - 1] \mu_i \},$$

whence by Tschebyshev inequality for $\gamma \geq 0$ it holds

$$\begin{aligned} &\operatorname{Prob}\{h^{T}\omega > h^{T}\mu + t\} = \operatorname{Prob}\{\gamma h^{T}\omega > \gamma h^{T}\mu + \gamma t\} \\ &\leq \mathbf{E}\left\{\exp\{\gamma h^{T}\omega\}\right\}\exp\{-\gamma h^{T}\mu - \gamma t\} \leq \exp\{\sum_{i}[\exp\{\gamma h_{i}\} - 1]\mu_{i} - \gamma h^{T}\mu - \gamma t\}, \end{aligned}$$

$$\begin{aligned} & (4.291) \end{aligned}$$

Now, it is easily seen that when $|s| \le 2/3$, one has $e^s \le 1 + s + \frac{3}{4}s^2$, which combines with (4.291) to imply that

$$0 \le \gamma \le \frac{2}{3\|h\|_{\infty}} \Rightarrow \ln\left(\operatorname{Prob}\{h^T\omega > h^T\mu + t\}\right) \le \frac{3}{4}\gamma^2 \sum_i h_i^2\mu_i - \gamma t.$$
(4.292)

Minimizing the right hand side in this inequality in $\gamma \in [0, \frac{2}{3} ||h||_{\infty}]$, we get

$$\operatorname{Prob}\left\{h^{T}\omega > h^{T}\mu + t\right\} \leq \exp\{-\frac{t^{2}}{3[\sum_{i}h_{i}^{2}\mu_{i} + \|h\|_{\infty}t]}\}.$$

This inequality combines with the same inequality applied to -h in the role of h to imply (4.96).

439

4.10.8 **Proof of Lemma 4.28**

(i): When $p(\operatorname{Col}_{\ell}[H]) \leq 1$ for all ℓ and $\lambda \geq 0$, denoting by $[h]^2$ the vector comprised of squares of the entries in h, we have

$$\begin{aligned} \phi(\operatorname{dg}(H\operatorname{Diag}\{\lambda\}H^T)) &= \phi(\sum_{\ell} \lambda_{\ell}[\operatorname{Col}_{\ell}[H]]^2) \leq \sum_{\ell} \lambda_{\ell} \phi([\operatorname{Col}_{\ell}[H]]^2) \\ &= \sum_{\ell} \lambda_{\ell} p^2(\operatorname{Col}_{\ell}[H]) \leq \sum_{\ell} \lambda_{\ell}, \end{aligned}$$

implying that $(H^T \text{Diag}\{\lambda\} H^T, \varkappa \sum_{\ell} \lambda_{\ell})$ belongs to **H**.

(ii): Let Θ, μ, Q, V be as stated in (ii); there is nothing to prove when $\mu = 0$, thus assume that $\mu > 0$. Let $d = dg(\Theta)$, so that

$$d_i = \sum_j Q_{ij}^2 \& \varkappa \phi(d) \le \mu \tag{4.293}$$

(the second relation is due to $(\Theta, \mu) \in \mathbf{H}$). (4.127) is evident. We have

$$[H_{\chi}]_{ij} = \sqrt{m/\mu} [G_{\chi}]_{ij}, G_{\chi} = Q \text{Diag}\{\chi\} V = \left[\sum_{k=1}^{m} Q_{ik} \chi_k V_{kj}\right]_{i,j}.$$

We claim that

$$\forall \gamma > 0 : \operatorname{Prob}\left\{ [G_{\chi}]_{ij}^2 > 3\gamma d_i/m \right\} \le \sqrt{3} \exp\{-\gamma/2\}.$$
 (4.294)

Indeed, there is nothing to prove when $d_i = 0$, since in this case $Q_{ij} = 0$ for all j and therefore $[G_{\chi}]_{ij} \equiv 0$. When $d_i > 0$, by homogeneity in Q it suffices to verify (4.294) when $d_i/m = 1/3$. Assuming that this is the case, let $\eta \sim \mathcal{N}(0, 1)$ be independent of χ . We have

$$\mathbf{E}_{\eta} \{ \mathbf{E}_{\chi} \{ \exp\{\eta[G_{\chi}]_{ij} \} \} = \mathbf{E}_{\eta} \{ \prod_{k} \cosh(\eta Q_{ik} V_{kj}) \} \leq \mathbf{E}_{\eta} \{ \prod_{k} \exp\{\frac{1}{2} \eta^{2} Q_{ik}^{2} V_{kj}^{2} \} \}$$

$$= \mathbf{E}_{\eta} \{ \exp\{\frac{1}{2} \eta^{2} \underbrace{\sum_{k} Q_{ik}^{2} V_{kj}^{2} }_{\leq 2d_{i}/m} \} \} \leq \mathbf{E}_{\eta} \{ \eta^{2} d_{i}/m \} = \mathbf{E}_{\eta} \{ \exp\{\eta^{2}/3\} \} = \sqrt{3},$$

and

$$\mathbf{E}_{\chi}\left\{\mathbf{E}_{\eta}\left\{\exp\{\eta[G_{\chi}]_{ij}\}\right\}\right\} = \mathbf{E}_{\chi}\left\{\exp\{\frac{1}{2}[G_{\chi}]_{ij}^{2}\}\right\}$$

implying that

$$\mathbf{E}_{\chi}\left\{\exp\{\frac{1}{2}[G_{\chi}]_{ij}^{2}\}\right\} \leq \sqrt{3}.$$

Therefore in the case of $d_i/m = 1/3$ for all s > 0 it holds

$$\operatorname{Prob}\{\chi: [G_{\chi}]_{ij}^2 > s\} \le \sqrt{3} \exp\{-s/2\},\$$

and (4.294) follows. Recalling the relation between H and G, we get from (4.294) that

 $\forall \gamma > 0: \operatorname{Prob}\{\chi: [H_{\chi}]_{ij}^2 > 3\gamma d_i/\mu\} \leq \sqrt{3} \exp\{-\gamma/2\},$

440

LECTURE 4

By this inequality, with \varkappa given by (4.126) the probability of the event

$$\forall i, j : [H_{\chi}]_{ij}^2 \le \varkappa \frac{d_i}{\mu}$$

is at least 1/2. Let this event take place; in this case we have $[\operatorname{Col}_{\ell}[H]]^2 \leq \varkappa d/\mu$, whence, due to what the norm $p(\cdot)$ is, $p^2(\operatorname{Col}_{\ell}[H]) \leq \kappa \phi(d)/\mu \leq 1$ (see the second relation in (4.293)). Thus, the probability of the event (4.128) is at least 1/2. \Box

4.10.9 Justification of (4.141)

Given $s \in [2,\infty]$ and setting $\bar{s} = s/2$, $s_* = \frac{s}{s-1}$, $\bar{s}_* = \frac{\bar{s}}{\bar{s}-1}$, we want to prove that

$$\{ (V,\tau) \in \mathbf{S}_{+}^{N} \times \mathbf{R}_{+} : \exists (W \in \mathbf{S}^{N}, w \in \mathbf{R}_{+}^{N}) : V \preceq W + \text{Diag}\{w\} \& \|W\|_{s_{*}} + \|w\|_{\bar{s}_{*}} \leq \tau \}$$

= $\{ (V,\tau) \in \mathbf{S}_{+}^{N} \times \mathbf{R}_{+} : \exists w \in \mathbf{R}_{+}^{N} : V \preceq \text{Diag}\{w\}, \|w\|_{\bar{s}_{*}} \leq \tau \}.$

To this end it clearly suffices to check that whenever $W \in \mathbf{S}^N$, there exists $w \in \mathbf{R}^N$ satisfying

$$W \preceq \text{Diag}\{w\}, \|w\|_{\bar{s}_*} \leq \|W\|_{s_*}$$

The latter claim is nothing but the claim that whenever $W \in \mathbf{S}^N$ and $\|W\|_{s_*} \le 1$, the conic optimization problem

$$Opt = \min_{t \ w} \{ t : t \ge \|w\|_{\bar{s}_*}, Diag\{w\} \succeq W \}$$
(4.295)

is solvable (which is evident) with optimal value ≤ 1 . To see that the latter indeed is the case, note that the problem clearly is strictly feasible, whence its optimal value is the same as the optimal value in the conic problem

$$Opt = \max_{P} \left\{ Tr(PW) : P \succeq 0, \| dg\{P\} \|_{\bar{s}_{*}/(\bar{s}_{*}-1)} \le 1 \right\}$$
$$[dg\{P\} = [P_{11}; P_{22}; ...; P_{NN}]]$$

dual to (4.295). Since $\operatorname{Tr}(PW) \leq \|P\|_{s_*/(s_*-1)} \|W\|_{s_*} \leq \|P\|_{s_*/(s_*-1)}$, recalling what s_* and \bar{s}_* are, our task boils down to verifying that when a matrix $P \succeq 0$ satisfies $\|\mathrm{dg}\{P)\|_{s/2} \leq 1$, one has also $\|P\|_s \leq 1$, which is immediate: since P is positive semidefinite, we have $|P_{ij}| \leq P_{ii}^{1/2} P_{jj}^{1/2}$, whence, assuming $s < \infty$,

$$||P||_{s}^{s} = \sum_{i,j} |P_{ij}|^{s} \le \sum_{i,j} P_{ii}^{s/2} P_{jj}^{s/2} = \left(\sum_{i} P_{ii}^{s/2}\right)^{2} \le 1.$$

When $s = \infty$, the same argument says that $||P||_{\infty} = \max_{i,j} |P_{ij}| = \max_i |P_{ii}| = ||dg\{P\}||_{\infty}$.

Bibliography

- T. W. Anderson. The integral of a symmetric unimodal function over a symmetric convex set and some probability inequalities. *Proceedings of the American Mathematical Society*, 6(2):170–176, 1955.
- [2] A. Antoniadis and I. Gijbels. Detecting abrupt changes by wavelet methods. *Journal of Nonparametric Statistics*, 14(1-2):7–29, 2002.
- [3] B. F. Arnold and P. Stahlecker. Another view of the kuks–olman estimator. Journal of statistical planning and inference, 89(1):169–174, 2000.
- [4] T. Augustin and R. Hable. On the impact of robust statistics on imprecise probability models: a review. *Structural Safety*, 32(6):358–365, 2010.
- [5] R. Bakeman and J. Gottman. Observing Interaction: An Introduction to Sequential Analysis. Cambridge University Press, 1997.
- [6] M. Basseville. Detecting changes in signals and systems a survey. Automatica, 24(3):309–326, 1988.
- [7] M. Basseville and I. Nikiforov. Detection of Abrupt Changes: Theory and Application. Prentice-Hall, Englewood Cliffs, N.J., 1993.
- T. Bednarski. Binary experiments, minimax tests and 2-alternating capacities. The Annals of Statistics, 10(1):226-232, 1982.
- [9] D. Belomestny and A. Goldenschluger. Nonparametric density estimation from observations with multiplicative measurement errors. *arXiv preprint arXiv:1709.00629*, 2017.
- [10] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust optimization*, volume 28. Princeton University Press, 2009.
- [11] A. Ben-Tal and A. Nemirovski. *Lectures on modern convex optimization: analysis, algorithms, and engineering applications,* volume 2. Siam, 2001.
- [12] A. Ben-Tal and A. Nemirovski. Lectures on modern convex optimization: analysis, algorithms, and engineering applications, volume 2. Siam, 2001.
- [13] M. Bertero and P. Boccacci. Application of the OS-EM method to the restoration of LBT images. Astronomy and Astrophysics Supplement Series, 144(1):181–186, 2000.
- [14] M. Bertero and P. Boccacci. Image restoration methods for the large binocular telescope (LBT). Astronomy and Astrophysics Supplement Series, 147(2):323–333, 2000.

- [15] E. Betzig, G. H. Patterson, R. Sougrat, O. W. Lindwasser, S. Olenych, J. S. Bonifacino, M. W. Davidson, J. Lippincott-Schwartz, and H. F. Hess. Imaging intracellular fluorescent proteins at nanometer resolution. *Science*, 313(5793):1642–1645, 2006.
- [16] L. Birgé. Approximation dans les spaces métriques et théorie de l'estimation: inégalités de Cràmer-Chernoff et théorie asymptotique des tests. PhD thesis, Université Paris VII, 1980.
- [17] L. Birgé. Vitesses maximales de décroissance des erreurs et tests optimaux associés. Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete, 55(3):261–273, 1981.
- [18] L. Birgé. Sur un théorème de minimax et son application aux tests. Probab. Math. Stat., 3:259–282, 1982.
- [19] L. Birgé. Approximation dans les espaces métriques et théorie de l'estimation. Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete, 65(2):181– 237, 1983.
- [20] L. Birgé. Robust testing for independent non identically distributed variables and Markov chains. In *Specifying Statistical Models*, pages 134–162. Springer, 1983.
- [21] L. Birgé. Model selection via testing: an alternative to (penalized) maximum likelihood estimators. In Annales de l'Institut Henri Poincare (B) Probability and Statistics, volume 42:3, pages 273–325. Elsevier, 2006.
- [22] S. P. Boyd, L. El Ghaoui, E. Feron, and V. Balakrishnan. *Linear matrix inequalities in system and control theory*, volume 15. SIAM, 1994.
- [23] E. Brodsky and B. S. Darkhovsky. Nonparametric methods in change point problems, volume 243. Springer Science & Business Media, 2013.
- [24] E. Brunel, F. Comte, and V. Genon-Catalot. Nonparametric density and survival function estimation in the multiplicative censoring model. *Test*, 25(3):570–590, 2016.
- [25] A. Buchholz. Operator khintchine inequality in non-commutative probability. Mathematische Annalen, 319(1):1–16, 2001.
- [26] A. Buja. On the huber-strassen theorem. Probability Theory and Related Fields, 73(1):149–152, 1986.
- [27] M. Burnashev. On the minimax detection of an imperfectly known signal in a white noise background. *Theory Probab. Appl.*, 24:107–119, 1979.
- [28] M. Burnashev. Discrimination of hypotheses for gaussian measures and a geometric characterization of the gaussian distribution. *Math. Notes*, 32:757– 761, 1982.
- [29] T. T. Cai and M. G. Low. A note on nonparametric estimation of linear functionals. *The Annals of Statistics*, pages 1140–1153, 2003.
- [30] T. T. Cai and M. G. Low. Minimax estimation of linear functionals over

nonconvex parameter spaces. The Annals of Statistics, 32(2):552–576, 2004.

- [31] T. T. Cai and M. G. Low. On adaptive estimation of linear functionals. The Annals of Statistics, 33(5):2311–2343, 2005.
- [32] E. Candes and T. Tao. The dantzig selector: statistical estimation when p is much larger than n. *The Annals of Statistics*, pages 2313–2351, 2007.
- [33] E. J. Candes. The restricted isometry property and its implications for compressed sensing. *Comptes rendus mathematique*, 346(9-10):589–592, 2008.
- [34] E. J. Candès et al. Compressive sampling. In Proceedings of the international congress of mathematicians, volume 3, pages 1433–1452. Madrid, Spain, 2006.
- [35] E. J. Candes, J. K. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences, 59(8):1207–1223, 2006.
- [36] E. J. Candes and T. Tao. Decoding by linear programming. *IEEE transactions on information theory*, 51(12):4203–4215, 2005.
- [37] E. J. Candes and T. Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE transactions on information* theory, 52(12):5406-5425, 2006.
- [38] J. Chen and A. Gupta. *Parametric statistical change point analysis: with applications to genetics, medicine, and finance.* Boston: Birkhäuser, 2012.
- [39] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pages 493–507, 1952.
- [40] H. Chernoff. Sequential Analysis and Optimal Design. SIAM, 1972.
- [41] N. Christopeit and K. Helmes. Linear minimax estimation with ellipsoidal constraints. Acta Applicandae Mathematica, 43(1):3–15, 1996.
- [42] I. Dattner, A. Goldenshluger, and A. Juditsky. On deconvolution of distribution functions. *The Annals of Statistics*, 39(5):2477–2501, 2011.
- [43] D. Donoho and R. Liu. Geometrizing rate of convergence I. Technical report, Tech. Report 137a, Dept. of Statist., University of California, Berkeley, 1987.
- [44] D. L. Donoho. Statistical estimation and optimal recovery. The Annals of Statistics, 22(1):238–270, 1994.
- [45] D. L. Donoho. De-noising by soft-thresholding. IEEE Transactions on Information Theory, 41(3):613–627, 1995.
- [46] D. L. Donoho. Neighborly polytopes and sparse solutions of underdetermined linear equations. 2005.
- [47] D. L. Donoho, M. Elad, and V. N. Temlyakov. Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Transactions on*

information theory, 52(1):6–18, 2006.

- [48] D. L. Donoho and X. Huo. Uncertainty principles and ideal atomic decomposition. Information Theory, IEEE Transactions on, 47(7):2845–2862, 2001.
- [49] D. L. Donoho and I. M. Johnstone. Minimax estimation via wavelet shrinkage. The Annals of Statistics, 26(3):879–921, 1998.
- [50] D. L. Donoho and R. C. Liu. Geometrizing rates of convergence, II. The Annals of Statistics, pages 633–667, 1991.
- [51] D. L. Donoho and R. C. Liu. Geometrizing rates of convergence, iii. The Annals of Statistics, pages 668–701, 1991.
- [52] D. L. Donoho, R. C. Liu, and B. MacGibbon. Minimax risk over hyperrectangles, and implications. *The Annals of Statistics*, pages 1416–1437, 1990.
- [53] D. L. Donoho and M. G. Low. Renormalization exponents and optimal pointwise rates of convergence. *The Annals of Statistics*, pages 944–970, 1992.
- [54] H. Drygas. Spectral methods in linear minimax estimation. Acta Applicandae Mathematica, 43(1):17–42, 1996.
- [55] S. Efromovich. Nonparametric curve estimation: methods, theory, and applications. Springer Science & Business Media, 2008.
- [56] S. Efromovich and M. Pinsker. Sharp-optimal and adaptive estimation for heteroscedastic nonparametric regression. *Statistica Sinica*, pages 925–942, 1996.
- [57] J. Fan. On the optimal rates of convergence for nonparametric deconvolution problems. *The Annals of Statistics*, pages 1257–1272, 1991.
- [58] R. Fano. Transmission of information: a statistical theory of communications. MIT Press, Cambridge, MA, 1968.
- [59] G. Fellouris and G. Sokolov. Second-order asymptotic optimality in multisensor sequential change detection. *IEEE Transactions on Information Theory*, 62(6):3662–3675, 2016.
- [60] J.-J. Fuchs. On sparse representations in arbitrary redundant bases. *IEEE transactions on Information theory*, 50(6):1341–1344, 2004.
- [61] J.-J. Fuchs. Recovery of exact sparse representations in the presence of bounded noise. *IEEE Transactions on Information Theory*, 51(10):3601–3608, 2005.
- [62] W. R. Gaffey. A consistent estimator of a component of a convolution. The Annals of Mathematical Statistics, 30(1):198–205, 1959.
- [63] N. H. Gholson and R. L. Moose. Maneuvering target tracking using adaptive state estimation. *IEEE Transactions on Aerospace and Electronic Systems*, 13(3):310–317, 1977.
- [64] A. Goldenshluger, A. Juditsky, and A. Nemirovski. Hypothesis testing by

convex optimization. *Electronic Journal of Statistics*, 9(2):1645–1712, 2015.

- [65] A. Goldenshluger, A. Juditsky, A. Tsybakov, and A. Zeevi. Change-point estimation from indirect observations. 1. minimax complexity. Ann. Inst. Henri Poincare Probab. Stat., 44:787–818, 2008.
- [66] A. Goldenshluger, A. Juditsky, A. Tsybakov, and A. Zeevi. Change-point estimation from indirect observations. 2. adaptation. Ann. Inst. H. Poincare Probab. Statist, 44(5):819–836, 2008.
- [67] Y. K. Golubev, B. Y. Levit, and A. B. Tsybakov. Asymptotically efficient estimation of analytic functions in gaussian noise. *Bernoulli*, pages 167–181, 1996.
- [68] L. Gordon and M. Pollak. An efficient sequential nonparametric scheme for detecting a change of distribution. *The Annals of Statistics*, pages 763–804, 1994.
- [69] M. Grant and S. Boyd. The CVX Users Guide. Release 2.1, 2014. http: //web.cvxr.com/cvx/doc/CVX.pdf.
- [70] V. Guigues, A. Juditsky, and A. Nemirovski. Hypothesis testing via euclidean separation. arXiv preprint arXiv:1705.07196, 2017.
- [71] V. Guigues, A. Juditsky, A. Nemirovski, Y. Cao, and Y. Xie. Change detection via affine and quadratic detectors. arXiv preprint arXiv:1608.00524, 2016.
- [72] F. Gustafsson. Adaptive filtering and change detection, volume 1. Wiley New York, 2000.
- [73] S. W. Hell. Toward fluorescence nanoscopy. *Nature biotechnology*, 21(11):1347, 2003.
- [74] S. W. Hell. Microscopy and its focal switch. *Nature methods*, 6(1):24, 2009.
- [75] S. W. Hell and J. Wichmann. Breaking the diffraction resolution limit by stimulated emission: stimulated-emission-depletion fluorescence microscopy. *Optics letters*, 19(11):780–782, 1994.
- [76] S. T. Hess, T. P. Girirajan, and M. D. Mason. Ultra-high resolution imaging by fluorescence photoactivation localization microscopy. *Biophysical journal*, 91(11):4258–4272, 2006.
- [77] J.-B. Hiriart-Urruty and C. Lemarechal. Convex analysis and minimization algorithms i: Fundamentals (grundlehren der mathematischen wissenschaften). 1993.
- [78] P. J. Huber. A robust version of the probability ratio test. The Annals of Mathematical Statistics, 36(6):1753–1758, 1965.
- [79] P. J. Huber and V. Strassen. Minimax tests and the neyman-pearson lemma for capacities. *The Annals of Statistics*, pages 251–263, 1973.
- [80] P. J. Huber and V. Strassen. Note: Correction to minimax tests and the

neyman-pearson lemma for capacities. The Annals of Statistics, 2(1):223–224, 1974.

- [81] I. A. Ibragimov and R. Z. Has' Minskii. Statistical estimation: asymptotic theory, volume 16. Springer Science & Business Media, 2013.
- [82] I. A. Ibragimov and R. Z. Khas' minskii. On nonparametric estimation of the value of a linear functional in gaussian white noise. *Theory of Probability & Its Applications*, 29(1):18–32, 1985.
- [83] I. A. Ibragimov and R. Z. Khas minskii. Estimation of linear functionals in gaussian noise. Theory of Probability & Its Applications, 32(1):30–39, 1988.
- [84] Y. Ingster and I. A. Suslina. Nonparametric goodness-of-fit testing under Gaussian models, volume 169 of Lecture Notes in Statistics. Springer, 2002.
- [85] A. Juditsky, F. K. Karzan, A. Nemirovski, B. Polyak, et al. Accuracy guaranties for l₁-recovery of block-sparse signals. *The Annals of Statistics*, 40(6):3077–3107, 2012.
- [86] A. Juditsky and A. Nemirovski. Nonparametric estimation by convex programming. The Annals of Statistics, 37(5a):2278–2300, 2009.
- [87] A. Juditsky and A. Nemirovski. On sequential hypotheses testing via convex optimization. Avtomatika i Telemekhanika (Engl. translation: Automation and Remote Control), (5 (transl: 76:5)):100–120 (transl: 809–825), 2015.
- [88] A. Juditsky and A. Nemirovski. Estimating linear and quadratic forms via indirect observations. arXiv preprint arXiv:1612.01508, 2016.
- [89] A. Juditsky and A. Nemirovski. Hypothesis testing via affine detectors. *Electronic Journal of Statistics*, 10:2204–2242, 2016.
- [90] A. Juditsky and A. Nemirovski. Near-optimality of linear recovery from indirect observations. *Mathematical Statistics and Learning*, 1(2):101–110, 2018.
- [91] A. Juditsky and A. Nemirovski. Near-optimality of linear recovery from indirect observations. *Mathematical Statistics and Learning*, 1(2):171–225, 2018.
- [92] A. Juditsky, A. Nemirovski, et al. Near-optimality of linear recovery in gaussian observation scheme under ||·||²₂-loss. The Annals of Statistics, 46(4):1603– 1629, 2018.
- [93] J. Kiefer and J. Wolfowitz. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*, 23(3), 1952.
- [94] C. Kraft. Some conditions for consistency and uniform consistency of statistical procedures. Univ. of California Publ. Statist., 2:493–507, 1955.
- [95] J. A. Kuks and W. Olman. Minimax linear estimation of regression coefficients
 (i). Iswestija Akademija Nauk Estonskoj SSR, 20:480–482, 1971.
- [96] J. A. Kuks and W. Olman. Minimax linear estimation of regression coefficients (ii). Iswestija Akademija Nauk Estonskoj SSR, 21:66–72, 1972.

- [98] T. L. Lai. Sequential changepoint detection in quality control and dynamical systems. Journal of the Royal Statistical Society. Series B (Methodological), pages 613–658, 1995.
- [99] A. Lakhina, M. Crovella, and C. Diot. Diagnosing network-wide traffic anomalies. In ACM SIGCOMM Computer Communication Review, volume 34:4, pages 219–230. ACM, 2004.
- [100] L. Le Cam. On the assumptions used to prove asymptotic normality of maximum likelihood estimates. The Annals of Mathematical Statistics, pages 802–828, 1970.
- [101] L. Le Cam. Convergence of estimates under dimensionality restrictions. The Annals of Statistics, pages 38–53, 1973.
- [102] L. Le Cam. On local and global properties in the theory of asymptotic normality of experiments. Stochastic processes and related topics, 1:13–54, 1975.
- [103] L. Le Cam. Asymptotic Methods in Statistical Decision Theory. Springer Series in Statistics. Springer, 1986.
- [104] G. Lorden. Procedures for reacting to a change in distribution. The Annals of Mathematical Statistics, pages 1897–1908, 1971.
- [105] F. Lust-Piquard. Inégalités de khintchine dans c^p (1 . CR Acad. Sci. Paris, 303:289–292, 1986.
- [106] L. Mackey, M. I. Jordan, R. Y. Chen, B. Farrell, and J. A. Tropp. Matrix concentration inequalities via the method of exchangeable pairs. *The Annals* of *Probability*, 42(3):906–945, 2014.
- [107] L. Mackey, M. I. Jordan, R. Y. Chen, B. Farrell, J. A. Tropp, et al. Matrix concentration inequalities via the method of exchangeable pairs. *The Annals* of *Probability*, 42(3):906–945, 2014.
- [108] E. Mazor, A. Averbuch, Y. Bar-Shalom, and J. Dayan. Interacting multiple model methods in target tracking: a survey. *IEEE Transactions on Aerospace* and Electronic Systems, 34(1):103–123, 1998.
- [109] Y. Mei. Asymptotic optimality theory for decentralized sequential hypothesis testing in sensor networks. *IEEE Transactions on Information Theory*, 54(5):2072–2089, 2008.
- [110] A. Meister. Deconvolution problems in nonparametric statistics, volume 193. Springer, 2009.
- [111] A. Mosek. The MOSEK optimization toolbox for MATLAB manual. Version 8.0, 2015. http://docs.mosek.com/8.0/toolbox/.
- [112] G. V. Moustakides. Optimal stopping times for detecting changes in distributions. *The Annals of Statistics*, pages 1379–1387, 1986.

- [113] H.-G. Müller and U. Stadtmüller. Discontinuous versus smooth regression. The Annals of Statistics, 27(1):299–337, 1999.
- [114] A. Nemirovski. Topics in non-parametric statistics. In P. Bernard, editor, Lectures on Probability Theory and Statistics, Ecole dEté de Probabilités de Saint-Flour, volume 28, pages 87–285. Springer, 2000.
- [115] A. Nemirovski. Interior point polynomial time methods in convex programming. lecture notes, 2005.
- [116] A. Nemirovski. Introduction to linear optimization. lecture notes, 2015.
- [117] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on opti*mization, 19(4):1574–1609, 2009.
- [118] A. Nemirovski, S. Onn, and R. U. Accuracy certificates for computational problems with convex structure. *Mathematics of Operations Research*, 35(1):52–78, 2010.
- [119] A. Nemirovski, C. Roos, and T. Terlaky. On maximization of quadratic form over intersection of ellipsoids with common center. *Mathematical Programming*, 86(3):463–473, 1999.
- [120] Y. Nesterov and A. Nemirovskii. Interior-point polynomial algorithms in convex programming, volume 13. Siam, 1994.
- [121] M. H. Neumann. Optimal change-point estimation in inverse problems. Scandinavian Journal of Statistics, 24(4):503–521, 1997.
- [122] F. Osterreicher. On the construction of least favourable pairs of distributions. Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete, 43(1):49– 55, 1978.
- [123] J. Pilz. Minimax linear regression estimation with symmetric parameter restrictions. Journal of Statistical Planning and Inference, 13:297–318, 1986.
- [124] M. Pinsker. Optimal filtration of square-integrable signals in gaussian noise. Prob. Info. Transmission, 16(2):120–133, 1980.
- [125] G. Pisier. Non-commutative vector valued l_p -spaces and completely *p*-summing maps. *Astérisque*, 247, 1998.
- [126] M. Pollak. Optimal detection of a change in distribution. The Annals of Statistics, pages 206–227, 1985.
- [127] M. Pollak. Average run lengths of an optimal method of detecting a change in distribution. *The Annals of Statistics*, pages 749–779, 1987.
- [128] H. V. Poor and O. Hadjiliadis. *Quickest detection*, volume 40. Cambridge University Press Cambridge, 2009.
- [129] C. R. Rao. Linear statistical inference and its applications, volume 22. John Wiley & Sons, 1973.

- [130] C. R. Rao. Estimation of parameters in a linear model. The Annals of Statistics, pages 1023–1037, 1976.
- [131] H. Rieder. Least favorable pairs for special capacities. The Annals of Statistics, pages 909–921, 1977.
- [132] M. J. Rust, M. Bates, and X. Zhuang. Sub-diffraction-limit imaging by stochastic optical reconstruction microscopy (storm). *Nature methods*, 3(10):793, 2006.
- [133] A. Shapiro, D. Dentcheva, and Ruszczyński. Lectures on Stochastic Programming: Modeling and Theory, Second Edition. SIAM, 2014.
- [134] A. N. Shiryaev. On optimum methods in quickest detection problems. Theory of Probability & Its Applications, 8(1):22–46, 1963.
- [135] D. Siegmund. Sequential Analysis: Tests and Confidence Intervals. Springer Science & Business Media, 1985.
- [136] D. Siegmund and B. Yakir. The statistics of gene mapping. Springer Science & Business Media, 2007.
- [137] A. Tartakovsky, I. Nikiforov, and M. Basseville. Sequential analysis: Hypothesis testing and changepoint detection. CRC Press, 2014.
- [138] A. G. Tartakovsky and V. V. Veeravalli. Change-point detection in multichannel and distributed systems. Applied Sequential Methodologies: Real-World Examples with Data Analysis, 173:339–370, 2004.
- [139] A. G. Tartakovsky and V. V. Veeravalli. Asymptotically optimal quickest change detection in distributed sensor systems. *Sequential Analysis*, 27(4):441–475, 2008.
- [140] J. A. Tropp. An introduction to matrix concentration inequalities. Foundations and Trends in Machine Learning, 8(1-2):1–230, 2015.
- [141] A. B. Tsybakov. Introduction to nonparametric estimation. Revised and extended from the 2004 French original. Translated by Vladimir Zaiats. Springer Series in Statistics. Springer, New York, 2009.
- [142] Y. Vardi, L. Shepp, and L. Kaufman. A statistical model for positron emission tomography. *Journal of the American statistical Association*, 80(389):8–20, 1985.
- [143] A. Wald. Sequential tests of statistical hypotheses. The Annals of Mathematical Statistics, 16(2):117–186, 1945.
- [144] A. Wald. Sequential Analysis. John Wiley and Sons, NY, 1947.
- [145] Y. Wang. Jump and sharp cusp detection by wavelets. *Biometrika*, 82(2):385– 397, 1995.
- [146] L. Wasserman. All of nonparametric statistics. Springer Science & Business Media, 2006.

- [147] A. S. Willsky. Detection of abrupt changes in dynamic systems. Springer, 1985.
- [148] H. Wolkowicz, R. Saigal, and L. Vandenberghe. Handbook of semidefinite programming: theory, algorithms, and applications, volume 27. Springer Science & Business Media, 2012.
- [149] Y. Xie and D. Siegmund. Sequential multi-sensor change-point detection. Annals of Statistics, 41(2):670–692, 2013.
- [150] Y. Yin. Detection of the number, locations and magnitudes of jumps. Communications in Statistics. Stochastic Models, 4(3):445–455, 1988.
- [151] C.-H. Zhang. Fourier methods for estimating mixing densities and distributions. The Annals of Statistics, pages 806–831, 1990.

Index

 $A^*, 2$ O(1), 3Diag, 1 Risk, 46 $\mathbf{E}_{\xi}\{ \}, \mathbf{E}_{\xi \sim P}\{ \}, \mathbf{E}\{ \}, 2$ $\mathbf{Q}_q(s,\kappa)$ -condition, 17 links with RIP, 29tractability when $q = \infty, 27$ verifiable sufficient conditions for, 25 $\mathbf{R}^n, \mathbf{R}^{m \times n}, \mathbf{1}$ $R_+, R_+^n, 2$ $\mathbf{S}^n, \mathbf{1}$ $S_{+}^{n}, 2$ $\mathcal{N}(\mu,\Theta), \mathbf{2}$ $\mathcal{R}_k[\cdot], \mathcal{R}_k^*[\cdot], \mathcal{S}_\ell[\cdot], \mathcal{S}_\ell^*[\cdot], ..., 280$ dg. 1 ℓ_1 minimization, see Compressed Sensing $\int_{\Omega} f(\xi) \Pi(d\xi), \ \mathbf{3}$ $\tilde{\lambda}[\cdot], 281$ $\succeq,\succ,\preceq,\prec,\,\mathbf{2}$ $\xi \sim P, \frac{1}{2}$ s-goodness, 13 Bisection estimate, 204 near-optimality of, 207 closeness relation, 65

Compressed Sensing, 7–10 via ℓ_1 minimization, 10–22 imperfect, 16 validity of, 12verifiable sufficient validity conditions, 22-31 verifiable sufficient validity conditions, limits of performance, 29 via penalized ℓ_1 recovery, 18 via regular ℓ_1 recovery, 18 conditional quantile, 203 cone dual, 267 Lorentz, 267 regular, 267 semidefinite, 267

conic problem, 268 dual of, 268 programming, 267, 270 Conic Duality Theorem, 269 conic hull, 268 contrast matrix, see nullspace property quantification Cramer-Rao risk bound, 398-401, 405 detector, 70 affine, 131 in simple observation schemes, 90 quadratic, 147risks of, 70 structural properties, 71 ellitope, 270-271 calculus of, 345-349 estimation of N-convex functions, 199-216of linear form, 190, 216-227 from repeated observations, 220-222 of sub-Gaussianity parameters, 222-227 of sub-Gaussianity parameters, direct product case, 224 of quadratic form, 227-237 Gaussian case, 227-232 Gaussian case, consistency, 231 Gaussian case, construction, 228 sub-Gaussian case, 232, 237 sub-Gaussian case, construction, 233

family of distributions regular/simple, 132–139 calculus of, 134 examples of, 132 spherical, 58 cap of, 59 function N-convex, 202 examples of, 203

StatOpt^{*}LN^{*}NS January 21, 2019 7x10

norm

452

Gaussian mixtures, 59 Hellinger affinity, 91 Hypothesis Testing change detection via quadratic lifting, 159–166 of multiple hypotheses, 64–70 in simple observation schemes, 95-113 up to closeness, 65, 99 via Euclidean separation, 67–70 via repeated observations, 102 of unions, 95 problem's setting, 46 sequential, 113-121 test. 47detector-based, 70 deterministic, 47 partial risks of, 50randomized, 47 simple, 47total risk of, 50two-point lower risk bound, 51 via affine detectors, 139-147 via Euclidean separation, 54-64 and repeated observations, 60majority test, 61 multiple hypotheses case, 67–70 pairwise, 55 via quadratic lifting, 147 Gaussian case, 147–153 sub-Gaussian case, 153-159 via repeated observations, 47

inequality Cramer-Rao, 400

lemma on Schur Complement, see Schur Complement Lemma LMI, 2 logistic regression, 332–333 matrices notation, 1 sensing, 5 MD, see Measurement Design Measurement Design, 121–131 simple case discrete o.s., 125 Gaussian o.s., 129

Poisson o.s., 128

Mutual Incoherence, 28

INDEX

conjugate, 284 Shatten, 352 Wasserstein, 386 Nullspace property, 13, 15 quantification of, 16 o.s., see observation scheme observation scheme discrete, 82, 93 Gaussian, 79, 92 Poisson, 80, 92 simple, 77-95 K-th power of, 93definition of, 78 direct product of, 82PET, see Positron Emission Tomography Poisson Imaging, 80 polyhedral estimate, 304-331 Positron Emission Tomography, 80 Rademacher random vector, 415 regular data, 132 repeated observations quasi-stationary, 49 semi-stationary, 48 stationary, 47 Restricted Isometry Property, 25 RIP, see Restricted Isometry Property risk $Risk(\mathcal{T}|H_1, ..., H_L), 50$ $\operatorname{RiskOpt}_{\Pi, \|\cdot\|}[\mathcal{X}], 295$ $\operatorname{Risk}[\widehat{x}(\cdot)|\mathcal{X}], 265$ $\operatorname{Risk}_{\epsilon}^{*}, 239$ $\operatorname{Risk}^{\mathcal{C}}(\mathcal{T}|H_1,...,H_L),$ 65 $\operatorname{Risk}_{\ell}^{\mathcal{C}}(\mathcal{T}|H_1,...,H_L), 65$ $\operatorname{Risk}_{\epsilon}^{\operatorname{opt}}(K), 225$ $\operatorname{Risk}_{\ell}(\mathcal{T}|H_1,...,H_L), 50$ $\operatorname{Risk}_{\epsilon}(\widehat{g}(\cdot)|G,\mathcal{X},\upsilon,\mathcal{A},\mathcal{H},\mathcal{M},\Phi),$ 217 $\operatorname{Risk}_{\pm}[\phi|\mathcal{P}], \operatorname{Risk}[\phi|\mathcal{P}_1, \mathcal{P}_2], 70$ $\operatorname{Risk}_{\operatorname{tot}}(\mathcal{T}|H_1,...,H_L), 50$ $\operatorname{Risk}_{\mathcal{H}}[\widehat{x}_*|\mathcal{X}], 303$ $\operatorname{Risk}_{\operatorname{opt}}[\mathcal{X}], 265$ $\operatorname{Risk}_{\mathcal{H},\|\cdot\|}[\widehat{x}|\mathcal{X}], \ \mathbf{300}$ C-, 65 \mathcal{H} -, 300 ϵ -, 217 in Hypothesis Testing partial, 50 total, 50up to closeness, 65
INDEX

minimax, 271 ϵ -, 225 of detector, 70 of simple test, 50

saddle point convex-concave saddle point problem, 84 Sample Average Approximation, 337-339 Schur Complement Lemma, 270 semidefinite relaxation on ellitope tightness of, 278 on spectratope tightness of, 282signal estimation, see signal recovery signal recovery linear, 272on ellitope, 272–276 on ellitope, near-optimality of,

276-278

on spectratope, 282–296 on spectratope under uncertainbut-bounded noise, 296-302on spectratope under uncertainbut-bounded noise, near-optimality of, 301 on spectratope, near-optimality of, 282, 295 problem setting, 5, 265 sparsity, s-sparsity, 7 spectratope, 279 calculus of, 345-349 examples of, 281Stochastic Approximation, 339–341 test, see Hypothesis Testing test theorem Sion-Kakutani, 86

vectors notation, 1 453

[111, 1, 2, 3, 4, 5, 6, 7, 8, 9, 12, 11, 10, 14, 13, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 35, 34, 36, 37, 32, 33, 71, 38, 39, 40, 41, 42, 43, 52, 50, 44, 51, 53, 45, 49, 48, 46, 47, 54, 55, 56, 57, 59, 60, 61, 62, 63, 65, 66, 64, 67, 68, 69, 70, 72, 75, 73, 74, 76, 77, 78, 79, 80, 81, 82, 83, 84, 86, 85, 87, 89, 92, 91, 88, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 107, 108, 109, 110, 112, 113, 119, 114, 117, 116, 115, 120, 121, 122, 123, 124, 125, 126, 127, 128, 130, 131, 132, 148, 134, 135, 136, 138, 139, 137, 106, 140, 142, 143, 144, 145, 146, 147, 149, 150, 151]