Course:

Optimization III Convex Analysis Nonlinear Programming Theory Nonlinear Programming Algorithms ISyE 6663 Spring 2008

Lecturer: Prof. Arkadi Nemirovski

Lecture Notes, Transparencies, Assignments: https://t-square.gatech.edu/portal/site/21746.200802

Grading Policy:

Assignments	5%
Midterm exam I	20%
Midterm exam II	25%
Final exam	50%

♣ To make decisions optimally is one of the most basic desires of a human being.

Whenever the candidate decisions, design restrictions and design goals can be properly quantified, optimal decision-making yields an *optimization problem*, most typically, a *Mathematical Programming* one:

minimize	
f(x)	[objective]
subject to	
$h_i(x) = 0, i = 1,, m$	equality constraints
$g_j(x) \leq 0, j = 1,, k$	inequality constraints
$x \in X$	[domain]
♣ In (MP),	(MP)

 \diamond a solution $x \in \mathbf{R}^n$ represents a candidate decision,

the constraints express restrictions on the meaningful decisions (balance and state equations, bounds on resources, etc.),

 the objective to be minimized represents the losses (minus profit) associated with a decision.



♣ To solve problem (MP) means to find its optimal solution x_* , that is, a feasible (i.e., satisfying the constraints) solution with the value of the objective \leq its value at any other feasible solution:

$$x_*: \begin{cases} h_i(x_*) = 0 \forall i \& g_j(x_*) \leq 0 \forall j \& x_* \in X \\ h_i(x) = 0 \forall i \& g_j(x) \leq 0 \forall j \& x \in X \\ \Rightarrow f(x_*) \leq f(x) \end{cases}$$

$$\min_{x} f(x) \\
s.t. \\
h_{i}(x) = 0, i = 1, ..., m \\
g_{j}(x) \leq 0, j = 1, ..., k \\
x \in X$$
(MP)

A In *Combinatorial* (or *Discrete*) Optimization, the domain X is a discrete set, like the set of all integral or 0/1 vectors.

In contrast to this, in *Continuous* Optimization we will focus on, X is a "continuum" set like the entire \mathbb{R}^n , a box $\{x : a \le x \le b\}$, or simplex $\{x \ge 0 : \sum_j x_j = 1\}$, etc., and the objective and the constraints are (at least) continuous on X.

A In Linear Programming, $X = \mathbf{R}^n$ and the objective and the constraints are linear functions of x.

In contrast to this, in *Nonlinear Continuous Optimization*, the objective and/or some of the constraints are nonlinear.

$$\min_{x} f(x) \\
s.t. \\
h_{i}(x) = 0, i = 1, ..., m \\
g_{j}(x) \leq 0, j = 1, ..., k \\
x \in X$$
(MP)

The goals of our course is to present

- basic theory of Continuous Optimization, with emphasis on existence and uniqueness of optimal solutions and their characterization (i.e., necessary and/or sufficient optimality conditions);
- traditional algorithms for building (approximate) optimal solutions to Continuous Optimization problems.

♣ Mathematical foundation of Optimization Theory is given by Convex Analysis – a specific combination of Real Analysis and Geometry unified by and focusing on investigating convexity-related notions.

Convex Sets

<u>Definition.</u> A set $X \subset \mathbb{R}^n$ is called *convex*, if X contains, along with every pair x, y of its points, the entire segment [x, y] with the endpoints x, y:

 $x, y \in X \Rightarrow \lambda x + (1 - \lambda)y \in X \ \forall \lambda \in [0, 1].$

<u>Note</u>: when λ runs through [0, 1], the point $x + \lambda(y - x) \equiv x + \lambda(y - x)$ runs through the segment [x, y].



Examples of convex sets, I: Affine sets

<u>Definition</u>: Affine set M in \mathbb{R}^n is a set which can be obtained as a shift of a *linear subspace* $L \subset \mathbb{R}^n$ by a vector $a \in \mathbb{R}^n$:

$$M = a + L = \{x = a + y : y \in L\}$$
(1)

<u>Note:</u> **I**. The linear subspace L is uniquely defined by affine subspace M and is the set of differences of vectors from M:

(1)
$$\Rightarrow L = M - M = \{y = x' - x'' : x', x'' \in M\}$$

II. The shift vector a is *not* uniquely defined by affine subspace M; in (1), one can take as a every vector from M (and only vector from M):

$$(1) \Rightarrow M = a' + L \ \forall a' \in M.$$

III. Generic example of affine subspace: the set of solutions of a *solvable* system of linear equations:

M is affine subspace in \mathbf{R}^n

 $\emptyset \neq M \equiv \{x \in \mathbf{R}^n : Ax = b\} \equiv \underbrace{a}_{Aa=b} + \underbrace{\{x : Ax = 0\}}_{\mathsf{Ker}A}$

By III, affine subspace is convex, due to <u>Proposition</u>. The solution set of an *arbitrary* (finite or infinite) system of linear inequalities is convex:

 $X = \{x \in \mathbf{R}^n : a_{\alpha}^T x \leq b_{\alpha}, \alpha \in \mathcal{A}\} \Rightarrow X \text{ is convex}$ In particular, every *polyhedral* set $\{x : Ax \leq b\}$ is convex.

Proof:

 $x, y \in X, \lambda \in [0, 1]$

 $\Leftrightarrow \ a_{\alpha}^{T}x \leq b_{\alpha}, \ a_{\alpha}^{T}y \leq b_{\alpha} \forall \alpha \in \mathcal{A}, \ \lambda \in [0, 1]$

$$\Rightarrow \underbrace{\lambda a_{\alpha}^{T} x + (1 - \lambda) a_{\alpha}^{T} y}_{a_{\alpha}^{T} [\lambda x + (1 - \lambda) y]} \leq \underbrace{\lambda b_{\alpha} + (1 - \lambda) b_{\alpha}}_{b_{\alpha}} \quad \forall \alpha \in \mathcal{A}$$

 $\Rightarrow [\lambda x + (1 - \lambda)y] \in X \ \forall \lambda \in [0, 1].$

<u>Remark:</u> Proposition remains valid when part of the nonstrict inequalities $a_{\alpha}^T x \leq b_{\alpha}$ are replaced with their strict versions $a_{\alpha}^T x < b_{\alpha}$.

Remark: The solution set

$$X = \{ x : a_{\alpha}^T x \le b_{\alpha}, \alpha \in \mathcal{A} \}$$

of a system of *nonstrict* inequalities is not only convex, it is closed (i.e., contains limits of all converging sequences $\{x_i \in X\}_{i=1}^{\infty}$ of points from X).

We shall see in the mean time that

Vice versa, every closed and convex set $X \subset \mathbf{R}^n$ is the solution set of an appropriate countable system of *nonstrict* linear inequalities:

 \boldsymbol{X} is closed and convex

 $\overset{\Downarrow}{X = \{x : a_i^T x \leq b_i, i = 1, 2, \ldots\}}$

Examples of convex sets, II: Unit balls of norms

<u>Definition</u>: A real-valued function ||x|| on \mathbb{R}^n is called *a norm*, if it possesses the following three properties:

♦ [positivity] $||x|| \ge 0$ for all x and ||x|| = 0 iff x = 0;

♦ [homogeneity] $\|\lambda x\| = |\lambda| \|x\|$ for all vectors x and reals λ ;

♦ [triangle inequality] $||x + y|| \le ||x|| + ||y||$ for all vectors x, y.

<u>Proposition</u>: Let $\|\cdot\|$ be a norm on \mathbb{R}^n . The unit ball of this norm – the set $\{x : \|x\| \le 1\}$, same as any other $\|\cdot\|$ -ball $\{x : \|x - a\| \le r\}$, is convex.

Proof:

$$||x - a|| \le r, ||y - a|| \le r, \lambda \in [0, 1]$$

$$\Rightarrow \underbrace{\lambda \|x-a\| + (1-\lambda)\|y-a\|}_{\geq \|\lambda(x-a)\| + \|(1-\lambda)(y-a)\|} \leq \underbrace{\lambda r + (1-\lambda)r}_{r}$$
$$\stackrel{\geq}{= \|\lambda(x-a) + (1-\lambda)(y-a)\|}_{\equiv \|[\lambda x + (1-\lambda)y] - a\|}$$

 $\Rightarrow \|[\lambda x + (1-\lambda)y] - a\| \le r \; \forall \lambda \in [0,1].$

Standard examples of norms on \mathbf{R}^n : ℓ_p -norms

$$\|x\|_p = \begin{cases} \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}, & 1 \le p < \infty \\ \max_i |x_i|, & p = \infty \end{cases}$$

<u>Note:</u> • $||x||_2 = \sqrt{\sum_i x_i^2}$ is the standard Euclidean norm:

• $||x||_1 = \sum_i |x_i|;$

• $||x||_{\infty} = \max_{i} |x_i|$ (uniform norm).

<u>Note</u>: except for the cases p = 1 and $p = \infty$, triangle inequality for $\|\cdot\|_p$ requires a nontrivial proof!

Proposition [characterization of $\|\cdot\|$ -balls] A set \mathcal{U} in \mathbb{R}^n is the unit ball of a norm iff \mathcal{U} is (a) convex and symmetric w.r.t. 0: V = -V, (b) bounded and closed, and

(c) contains a neighbourhood of the origin.

Examples of convex sets, III: Ellipsoid

<u>Definition</u>: An ellipsoid in \mathbb{R}^n is a set X given by

♦ positive definite and symmetric $n \times n$ matrix Q (that is, $Q = Q^T$ and $u^T Q u > 0$ whenever $u \neq 0$),

 \diamondsuit center $a \in \mathbf{R}^n$,

$$\diamond$$
 radius $r > 0$

via the relation

$$X = \{x : (x - a)^T Q(x - a) \le r^2\}.$$

Proposition: An ellipsoid is convex.

Proof: Since Q is symmetric positive definite, by Linear Algebra $Q = (Q^{1/2})^2$ for uniquely defined symmetric positive definite matrix $Q^{1/2}$. Setting $||x||_Q = ||Q^{1/2}x||_2$, we clearly get a norm on \mathbb{R}^n (since $|| \cdot ||_2$ is a norm and $Q^{1/2}$ is nonsingular). We have

$$(x-a)^T Q(x-a) = [(x-a)^T Q^{1/2}] [Q^{1/2}(x-a)] = \|Q^{1/2}(x-a)\|_2^2 = \|x-a\|_Q^2,$$

so that X is a $\|\cdot\|_Q$ -ball and is therefore a convex set.

Examples of convex sets, IV: ϵ -neighbourhood of convex set

<u>Proposition</u>: Let M be a nonempty convex set in \mathbb{R}^n , $\|\cdot\|$ be a norm, and $\epsilon \geq 0$. Then the set

$$X = \{x : \operatorname{dist}_{\|\cdot\|}(x, M) \equiv \inf_{y \in M} \|x - y\| \le \epsilon\}$$

is convex.

Proof: $x \in X$ if and only if for every $\epsilon' > \epsilon$ there exists $y \in M$ such that $||x - y|| \le \epsilon'$. We now have

 $x, y \in X, \lambda \in [0, 1]$

$$\Rightarrow \quad \forall \epsilon' > \epsilon \exists u, v \in M : ||x - u|| \le \epsilon', ||y - v|| \le \epsilon'$$

$$\Rightarrow \forall \epsilon' > \epsilon \exists u, v \in M :$$

$$\underbrace{\lambda \| x - u \| + (1 - \lambda) \| y - v \|}_{\geq \| [\lambda x + (1 - \lambda)y] - [\lambda u + (1 - \lambda)v] \|} \leq \epsilon' \forall \lambda \in [0, 1]$$

$$\Rightarrow \quad \forall \epsilon' > \epsilon \, \forall \lambda \in [0, 1] \exists w = \lambda u + (1 - \lambda) v \in M : \\ \| [\lambda x + (1 - \lambda) y] - w \| \le \epsilon'$$

 $\Rightarrow \lambda x + (1 - \lambda)y \in X \ \forall \lambda \in [0, 1]$

Convex Combinations and Convex Hulls

<u>Definition</u>: A convex combination of m vectors $x_1, ..., x_m \in \mathbf{R}^n$ is their linear combination

$$\sum_i \lambda_i x_i$$

with *nonnegative* coefficients and *unit sum of the coefficients*:

$$\lambda_i \ge 0 \ \forall i, \ \sum_i \lambda_i = 1.$$

<u>Proposition</u>: A set $X \subset \mathbb{R}^n$ is convex iff it is closed w.r.t. taking convex combinations of its points:

$$\begin{array}{c} X \text{ is convex} \\ \\ \\ x_i \in X, \lambda_i \geq 0, \\ \\ \sum_i \lambda_i = 1 \Rightarrow \\ \\ \\ i \\ \lambda_i x_i \in X. \end{array}$$

Proof, \Rightarrow : Assume that X is convex, and let us prove by induction in k that every kterm convex combination of vectors from X belongs to X. Base k = 1 is evident. Step $k \Rightarrow k + 1$: let $x_1, ..., x_{k+1} \in X$ and $\lambda_i \ge 0$, $\sum_{i=1}^k \lambda_i = 1$; we should prove that $\sum_{i=1}^{k+1} \lambda_i x_i \in X$. Assume w.l.o.g. that $0 \le \lambda_{k+1} < 1$. Then

$$\sum_{i=1}^{k+1} \lambda_i x_i = (1 - \lambda_{k+1}) \left(\sum_{\substack{i=1 \\ \in X}}^k \frac{\lambda_i}{1 - \lambda_{k+1}} x_i \right)$$
$$+ \lambda_{k+1} x_{k+1} \in X.$$

Proof, \Leftarrow : evident, since the definition of convexity of X is nothing but the requirement for every 2-term convex combination of points from X to belong to X.

<u>Proposition</u>: The intersection $X = \bigcap_{\alpha \in \mathcal{A}} X_{\alpha}$ of an *arbitrary* family $\{X_{\alpha}\}_{\alpha \in \mathcal{A}}$ of convex subsets of \mathbb{R}^{n} is convex.

Proof: evident.

<u>Corollary</u>: Let $X \subset \mathbb{R}^n$ be an arbitrary set. Then among convex sets containing X (which do exist, e.g. \mathbb{R}^n) there exists the smallest one, namely, the intersection of all convex sets containing X.

<u>Definition</u>: The smallest convex set containing X is called the *convex hull* Conv(X) of X.

<u>Proposition</u> [convex hull via convex combinations] For every subset X of \mathbb{R}^n , its convex hull Conv(X) is exactly the set \widehat{X} of all convex combinations of points from X.

Proof. 1) Every convex set which contains X contains every convex combination of points from X as well. Therefore $Conv(X) \supset \widehat{X}$.

2) It remains to prove that $Conv(X) \subset \widehat{X}$. To this end, by definition of Conv(X), it suffices to verify that the set \widehat{X} contains X (evident) and is convex. To see that \widehat{X} is convex, let $x = \sum_{i} \nu_{i} x_{i}, \ y = \sum_{i} \mu_{i} x_{i}$ be two points from \widehat{X} represented as convex combinations of points from X, and let $\lambda \in [0, 1]$. We have

$$\lambda x + (1 - \lambda)y = \sum_{i} [\lambda \nu_i + (1 - \lambda)\mu_i]x_i,$$

i.e., the left hand side vector is a convex combination of vectors from X.



Examples of convex sets, V: simplex

<u>Definition</u>: A collection of m+1 points x_i , i = 0, ..., m, in \mathbb{R}^n is called *affine independent*, if no nontrivial combination of the points *with zero sum of the coefficients* is zero:

 $x_0, ..., x_m$ are affine independent

<u>Motivation</u>: Let $X \subset \mathbf{R}^n$ be nonempty.

I. For every nonempty set $X \in \mathbb{R}^n$, the intersection of all affine subspaces containing X is an affine subspace. This clearly is the *smallest* affine subspace containing X; it is called the *affine hull* Aff(X) of X.

II. It is easily seen that Aff(X) is nothing but the set of all *affine combinations* of points from X, that is, linear combinations with unit sum of coefficients:

$$Aff(X) = \{x = \sum_{i} \lambda_{i} x_{i} : x_{i} \in X, \sum_{i} \lambda_{i} = 1\}.$$

III. m + 1 points $x_0, ..., x_m$ are affinely independent iff every point $x \in Aff(\{x_0, ..., x_m\})$ of their affine hull can be *uniquely represented* as an affine combination of $x_0, ..., x_m$:

$$\sum_{i} \lambda_{i} x_{i} = \sum_{i} \mu_{i} x_{i} \& \sum_{i} \lambda_{i} = \sum_{i} \mu_{i} = 1 \Rightarrow \lambda_{i} \equiv \mu_{i}$$

In this case, the coefficients λ_i in the representation

$$x = \sum_{i=0}^{m} \lambda_i x_i \qquad \qquad [\sum_i \lambda_i = 1]$$

of a point $x \in M = Aff(\{x_0, ..., x_m\})$ as an affine combination of $x_0, ..., x_m$ are called the *barycentric coordinates* of $x \in M$ *taken w.r.t. affine basis* $x_0, ..., x_m$ *of* M. <u>Definition:</u> *m*-dimensional simplex Δ with vertices $x_0, ..., x_m$ is the convex hull of m + 1 affine independent points $x_0, ..., x_m$:

$$\Delta = \Delta(x_0, ..., x_m) = \operatorname{Conv}(\{x_0, ..., x_m\}).$$

Examples: **A.** 2-dimensional simplex is given by 3 points not belonging to a line and is the triangle with vertices at these points.

B. Let $e_1, ..., e_n$ be the standard basic orths in \mathbb{R}^n . These n points are affinely independent, and the corresponding (n - 1)-dimensional simplex is the standard simplex $\Delta_n = \{x \in \mathbb{R}^n : x \ge 0, \sum_i x_i = 1\}.$

C. Adding to $e_1, ..., e_n$ the vector $e_0 = 0$, we get n+1 affine independent points. The corresponding *n*-dimensional simplex is

$$\Delta_n^+ = \{ x \in \mathbf{R}^n : x \ge 0, \sum_i x_i \le 1 \}.$$

Simplex with vertices $x_0, ..., x_m$ is convex (as a convex hull of a set), and every point from the simplex is a convex combination of the vertices with the coefficients uniquely defined by the point. Examples of convex sets, VI: cone

<u>Definition</u>: A nonempty subset K of \mathbb{R}^n is called *conic*, if it contains, along with every point x, the entire ray emanating from the origin and passing through x:

A convex conic set is called a cone.

Examples: A. Nonnegative orthant

$$\mathbf{R}^n_+ = \{ x \in \mathbf{R}^n : x \ge \mathbf{0} \}$$

B. Lorentz cone

$$\mathbf{L}^{n} = \{ x \in \mathbf{R}^{n} : x_{n} \ge \sqrt{x_{1}^{2} + \dots + x_{n-1}^{2}} \}$$

C. Semidefinite cone S^n_+ . This cone "lives" in the space S^n of $n \times n$ symmetric matrices and is comprised of all positive semidefinite symmetric $n \times n$ matrices **D.** The solution set $\{x : a_{\alpha}^T x \leq 0 \forall \alpha \in A\}$ of an arbitrary (finite or infinite) homogeneous system of nonstrict linear inequalities is a closed cone. In particular, so is a polyhedral cone $\{x : Ax \leq 0\}$.

<u>Note:</u> Every *closed* cone in \mathbb{R}^n is the solution set of a countable system of nonstrict linear inequalities.

<u>Proposition:</u> A nonempty subset K of \mathbf{R}^n is a cone iff

 $\diamondsuit K$ is conic: $x \in K, t \ge 0 \Rightarrow tx \in K$, and

 $\diamondsuit K$ is closed w.r.t. addition:

 $x, y \in K \Rightarrow x + y \in K.$

Proof, \Rightarrow : Let K be convex and $x, y \in K$, Then $\frac{1}{2}(x + y) \in K$ by convexity, and since K is conic, we also have $x + y \in K$. Thus, a convex conic set is closed w.r.t. addition.

Proof, \Leftarrow : Let *K* be conic and closed w.r.t. addition. In this case, a convex combination $\lambda x + (1 - \lambda)y$ of vectors x, y from *K* is the sum of the vectors λx and $(1 - \lambda)y$ and thus belongs to *K*, since *K* is closed w.r.t. addition. Thus, a conic set which is closed w.r.t. addition is convex.

♣ Cones form an extremely important class of convex sets with properties "parallel" to those of general convex sets. For example, ◇ Intersection of an arbitrary family of cones again is a cone. As a result, for every nonempty set X, among the cones containing X there exists the smallest cone Cone (X), called the conic hull of X.

♦ A nonempty set is a cone iff it is closed w.r.t. taking *conic* combinations of its elements (i.e., linear combinations with nonnegative coefficients).

 \diamond The conic hull of a nonempty set X is exactly the set of all conic combinations of elements of X.

"Calculus" of Convex Sets

<u>Proposition.</u> The following operations preserve convexity of sets:

1. Intersection: If $X_{\alpha} \subset \mathbb{R}^{n}$, $\alpha \in \mathcal{A}$, are convex sets, so is $\bigcap_{\alpha \in \mathcal{A}} X_{\alpha}$

2. Direct product: If $X_{\ell} \subset \mathbb{R}^{n_{\ell}}$ are convex sets, $\ell = 1, ..., L$, so is the set

$$X = X_1 \times ... \times X_L$$

$$\equiv \{x = (x^1, ..., x^L) : x^{\ell} \in X_{\ell}, 1 \le \ell \le L\}$$

$$\subset \mathbf{R}^{n_1 + ... + n_L}$$

3. Taking weighted sums: Let $X_1, ..., X_L$ be nonempty convex subsets in \mathbb{R}^n and $\lambda_1, ..., \lambda_L$ be reals. Then the set

$$\begin{split} \lambda_1 X_1 + \ldots + \lambda_L X_L \\ &\equiv \ \{x = \lambda_1 x_1 + \ldots + \lambda_L x_\ell : x_\ell \in X_\ell, 1 \leq \ell \leq L\} \\ \text{is convex.} \end{split}$$

4. Affine image: Let $X \subset \mathbf{R}^n$ be convex and $x \mapsto \mathcal{A}(x) = Ax + b$ be an affine mapping from \mathbf{R}^n to \mathbf{R}^k . Then the image of X under the mapping – the set

$$\mathcal{A}(X) = \{ y = Ax + b : x \in X \}$$

is convex.

5. Inverse affine image: Let $X \subset \mathbb{R}^n$ be convex and $y \mapsto \mathcal{A}(y) = Ay + b$ be an affine mapping from \mathbb{R}^k to \mathbb{R}^n . Then the inverse image of X under the mapping – the set

$$\mathcal{A}^{-1}(X) = \{ y : Ay + b \in X \}$$

is convex.

Application example: A point $x \in \mathbf{R}^n$ is

(a) "good", if it satisfies a given system of linear constraints $Ax \leq b$,

(b) "excellent", if it dominates a good point: $\exists y$: y is good and $x \ge y$,

(c) "semi-excellent", if it can be approximated, within accuracy 0.1 in the coordinatewise fashion, by excellent points:

 $\forall (i, \epsilon' > 0.1) \exists y : |y_i - x_i| \leq \epsilon' \& y \text{ is excellent}$

Question: Whether the set of semi-excellent points is convex?

Answer: Yes. Indeed,

• The set X_g of good points is convex (as a polyhedral set)

• \Rightarrow The set X_{exc} of excellent points is convex (as the sum of convex set X_{g} and the nonnegative orthant \mathbf{R}^{n}_{+} , which is convex) • \Rightarrow For every *i*, the set X^{i}_{exc} of *i*-th coordinates of excellent points is convex (as the projection of X_{exc} onto *i*-th axis; projection is an affine mapping) ⇒ For every *i*, the set Yⁱ on the axis which is the 0.1-neighbourhood of Xⁱ_{exc}, is convex (as 0.1-neighbourhood of a convex set)
⇒ The set of semi-excellent points, which is the direct product of the sets Y¹,..., Yⁿ, is convex (as direct product of convex sets).

Nice Topological Properties of Convex Sets

Recall that the set $X \subset \mathbf{R}^n$ is called \diamond *closed*, if X contains limits of all converging sequences of its points:

 $x_i \in X \And x_i \to x, i \to \infty \Rightarrow x \in X$

 \diamond open, if it contains, along with every of its points x, a ball of a positive radius centered at x:

 $x \in X \Rightarrow \exists r > 0 : \{y : \|y - x\|_2 \le r\} \subset X.$

E.g., the solution set of an arbitrary system of *nonstrict* linear inequalities $\{x : a_{\alpha}^T x \leq b_{\alpha}\}$ is closed; the solution set of *finite* system of *strict* linear inequalities $\{x : Ax < b\}$ is open.

<u>Facts</u>: **A.** X is closed iff $\mathbb{R}^n \setminus X$ is open

B. The intersection of an arbitrary family of closed sets and the union of a finite family of closed sets are closed

B'. The union of an arbitrary family of open sets and the intersection of a finite family of open sets are open \diamond From **B** it follows that the intersection of all closed sets containing a given set X is closed; this intersection, called the closure clX of X, is the smallest closed set containing X. clX is exactly the set of limits of all converging sequences of points of X:

$$\mathsf{cI}X = \{x : \exists x_i \in X : x = \lim_{i \to \infty} x_i\}.$$

 \diamond From **B**' it follows that the union of all open sets contained in a given set X is open; this union, called the interior intX of X, is the largest open set contained X. intX is exactly the set of all interior points of X – points x belonging to X along with balls of positive radii centered at the points:

int $X = \{x : \exists r > 0 : \{y : \|y - x\|_2 \le r\} \subset X\}.$ \diamondsuit Let $X \subset \mathbb{R}^n$. Then int $X \subset X \subset c | X$. The "difference" $\partial X = c | X \setminus int X$ is called the *boundary* of X; boundary always is closed (as the intersection of the closed sets c | X and the complement of int X).

$\mathsf{int}X \subset X \subset \mathsf{cl}X \tag{(*)}$

\clubsuit In general, the discrepancy between int*X* and cl*X* can be pretty large.

E.g., let $X \subset \mathbb{R}^1$ be the set of irrational numbers in [0,1]. Then $int X = \emptyset$, cl X = [0,1], so that int X and cl X differ dramatically.

Fortunately, a convex set is perfectly well approximated by its closure (and by interior, if the latter is nonempty).

Proposition: Let $X \subset \mathbb{R}^n$ be a nonempty *convex* set. Then (i) Both int*X* and cl*X* are convex (ii) If int*X* is nonempty, then int*X* is dense in cl*X*. Moreover,

 $x \in \text{int}X, \ y \in \text{cl}X \Rightarrow \\ \lambda x + (1 - \lambda)y \in \text{int}X \ \forall \lambda \in (0, 1]$ (!)

• Claim (i): Let X be convex. Then both int X and cl X are convex

Proof. (i) is nearly evident. Indeed, to prove that int*X* is convex, note that for every two points $x, y \in \text{int}X$ there exists a common r > 0 such that the balls B_x , B_y of radius r centered at x and y belong to X. Since X is convex, for every $\lambda \in [0, 1] X$ contains the set $\lambda B_x + (1 - \lambda)B_y$, which clearly is nothing but the ball of the radius r centered at $\lambda x + (1 - \lambda)y$. Thus, $\lambda x + (1 - \lambda)y \in \text{int}X$ for all $\lambda \in [0, 1]$.

Similarly, to prove that clX is convex, assume that $x, y \in clX$, so that $x = \lim_{i \to \infty} x_i$ and $y = \lim_{i \to \infty} y_i$ for appropriately chosen $x_i, y_i \in X$. Then for $\lambda \in [0, 1]$ we have

$$\lambda x + (1 - \lambda)y = \lim_{i \to \infty} \underbrace{[\lambda x_i + (1 - \lambda)y_i]}_{\in X},$$

so that $\lambda x + (1 - \lambda)y \in c X$ for all $\lambda \in [0, 1]$.

• <u>Claim (ii)</u>: Let X be convex and intX be nonempty. Then intX is dense in clX; moreover,

$$x \in \text{int}X, \ y \in \text{cl}X \Rightarrow \\ \lambda x + (1 - \lambda)y \in \text{int}X \ \forall \lambda \in (0, 1]$$
(!)

Proof. It suffices to prove (!). Indeed, let $\bar{x} \in \text{int}X$ (the latter set is nonempty). Every point $x \in \text{cl}X$ is the limit of the sequence $x_i = \frac{1}{i}\bar{x} + (1 - \frac{1}{i})x$. Given (!), all points x_i belong to intX, thus intX is dense in clX.

• Claim (ii): Let X be convex and intX be nonempty. Then

$$x \in \text{int}X, \ y \in \text{cl}X \Rightarrow \\ \lambda x + (1 - \lambda)y \in \text{int}X \ \forall \lambda \in (0, 1]$$
 (!)

Proof of (!): Let $x \in intX$, $y \in clX$, $\lambda \in (0,1]$. Let us prove that $\lambda x + (1-\lambda)y \in intX$. Since $x \in intX$, there exists r > 0 such that the ball B of radius r centered at x belongs to X. Since $y \in clX$, there exists a sequence $y_i \in X$ such that $y = \lim_{i \to \infty} y_i$. Now let

$$B^{i} = \lambda B + (1 - \lambda)y_{i}$$

= $\{\underbrace{z = [\lambda x + (1 - \lambda)y_{i}]}_{z_{i}} + \lambda d : ||d||_{2} \leq r\}$
= $\{z = z_{i} + \Delta : ||\delta|| \leq r' = \lambda d\}.$

Since $B \subset X$, $y_i \in X$ and X is convex, the sets B^i (which are balls of radius r' > 0 centered at z_i) contain in X. Since $z_i \to z = \lambda x + (1 - \lambda)y$ as $i \to \infty$, all these balls, starting with certain number, contain the ball B' of radius r'/2 centered at z. Thus, $B' \subset X$, i.e., $z \in int X$.

Let X be a convex set. It may happen that int $X = \emptyset$ (e.g., X is a segment in 3D); in this case, interior definitely does not approximate X and clX. What to do?

The natural way to overcome this difficulty is to pass to *relative* interior, which is nothing but the interior of X taken w.r.t. the affine hull Aff(X) of X rather than to \mathbb{R}^n . This affine hull, geometrically, is just certain \mathbb{R}^m with $m \leq n$; replacing, if necessary, \mathbb{R}^n with this \mathbb{R}^m , we arrive at the situation where intXis nonempty.

Implementation of the outlined idea goes through the following

<u>Definition</u>: [relative interior and relative boundary] Let X be a nonempty convex set and Mbe the affine hull of X. The *relative interior* rint X is the set of all point $x \in X$ such that a ball *in* M of a positive radius, centered at x, is contained in X:

rint $X = \{x : \exists r > 0 :$ $\{y \in Aff(X), \|y - x\|_2 \leq r\} \subset X\}.$ The *relative boundary* of X is, by definition, $c|X\setminus rint X.$ <u>Note</u>: An affine subspace M is given by a list of linear equations and thus is closed; as such, it contains the closure of every subset $Y \subset M$; this closure is nothing but the closure of Y which we would get when replacing the original "universe" \mathbb{R}^n with the affine subspace M (which, geometrically, is nothing but \mathbb{R}^m with certain $m \leq n$).

The essence of the matter is in the following fact:

<u>Proposition</u>: Let $X \subset \mathbf{R}^n$ be a nonempty *convex* set. Then rint $X \neq \emptyset$.
♣ Thus, replacing, if necessary, the original "universe" \mathbb{R}^n with a smaller geometrically similar universe, we can reduce investigating an arbitrary nonempty convex set X to the case where this set has a nonempty interior (which is nothing but the relative interior of X). In particular, our results for the "fulldimensional" case imply that For a nonempty convex set X, both rint X

and cIX are convex sets such that

 $\emptyset \neq \operatorname{rint} X \subset X \subset \operatorname{cl} X \subset \operatorname{Aff}(X)$

and rint X is dense in clX. Moreover, whenever $x \in \text{rint } X$, $y \in \text{cl} X$ and $\lambda \in (0, 1]$, one has

$$\lambda x + (1 - \lambda)y \in \operatorname{rint} X.$$

 $\emptyset \neq X$ is convex ?? \Rightarrow ?? rint $X \neq \emptyset$

Proof. 1⁰. Shifting X, we may assume w.l.o.g. that $0 \in X$, so that Aff(X) is merely a linear subspace L in \mathbb{R}^n .

2⁰. It may happen that $L = \{0\}$. Here $X = \{0\} = L$, whence the interior of X taken w.r.t. L is X and thus is nonempty.

3⁰. Now let $L \neq \{0\}$. Since $0 \in X$, L = Aff(X) (which always is the set of all *affine* combinations of vectors from X) is the same as the set of all *linear* combinations of vectors from X, i.e., is the linear span of X. Since dim L > 0, we can choose in L a linear basis $e_1, ..., e_m$ with $e_1, ..., e_m \in X$. Setting $e_0 = 0 \in X$, we get m + 1 affine independent vectors $e_0, ..., e_m$ in X. Since X is convex, the simplex Δ with vertices $e_0 = 0, e_1, ..., e_m$ is contained in X. Thus,

$$X \supset \Delta = \{x = \sum_{i=0}^{m} \lambda_i e_i : \lambda \ge 0, \sum_{i=0}^{m} \lambda_i = 1\}$$
$$\supset \{x = \sum_{i=1}^{m} \lambda_i e_i : \lambda_1, \dots, \lambda_m > 0, \sum_{i=1}^{m} \lambda_i < 1\}.$$

We see that X contains all vectors from L for which the coordinates w.r.t. the basis $\{e_1, ..., e_m\}$ are positive and with sum < 1. This set is open in L. Thus, the interior of X w.r.t. L is nonempty. **♣** Let X be convex and $\bar{x} \in \operatorname{rint} X$. As we know,

 $\lambda \in [0, 1], y \in C X \Rightarrow y_{\lambda} = \lambda \overline{x} + (1 - \lambda)y \in X.$

It follows that in order to pass from X to its closure clX, it suffices to pass to "radial closure":

For every direction $0 \neq d \in Aff(X) - \bar{x}$, let $T_d = \{t \ge 0 : \bar{x} + td \in X\}.$

<u>Note:</u> T_d is a convex subset of \mathbf{R}_+ which contains all small enough positive t's.

♦ If T_d is unbounded or is a bounded segment: $T_d = \{t : 0 \le t \le t(d) < \infty\}$, the intersection of clX with the ray $\{\bar{x} + td : t \ge 0\}$ is exactly the same as the intersection of X with the same ray.

 \diamond If T_d is a bounded half-segment: $T_d = \{t : 0 \le t < t(d) < \infty\}$, the intersection of clX with the ray $\{\bar{x} + td : t \ge 0\}$ is larger than the intersection of X with the same ray by exactly one point, namely, $\bar{x} + t(d)d$. Adding to X these "missing points" for all d, we arrive at clX.

Main Theorems on Convex Sets, I: Caratheodory Theorem

<u>Definition</u>: Let M be affine subspace in \mathbb{R}^n , so that M = a + L for a linear subspace L. The *linear* dimension of L is called the *affine* dimension dim M of M.

Examples: The affine dimension of a singleton is 0. The affine dimension of \mathbb{R}^n is n. The affine dimension of an affine subspace $M = \{x : Ax = b\}$ is $n - \operatorname{Rank}(A)$.

For a nonempty set $X \subset \mathbf{R}^n$, the *affine dimension* dim X of X is exactly the affine dimension of the affine hull Aff(X) of X.

Theorem [Caratheodory] Let $\emptyset \neq X \subset \mathbb{R}^n$. Then every point $x \in \text{Conv}(X)$ is a convex combination of *at most* dim(X) + 1 points of X. **Theorem** [Caratheodory] Let $\emptyset \neq X \subset \mathbb{R}^n$. Then every point $x \in \text{Conv}(X)$ is a convex combination of *at most* dim(X) + 1 points of X.

Proof. 1^0 . We should prove that if x is a convex combination of finitely many points $x_1, ..., x_k$ of X, then x is a convex combination of at most m + 1 of these points, where $m = \dim(X)$. Replacing, if necessary, \mathbb{R}^n with Aff(X), it suffices to consider the case of m = n.

2⁰. Consider a representation of x as a convex combination of $x_1, ..., x_k$ with minimum possible number of nonzero coefficients; it suffices to prove that this number is $\leq n+1$. Let, on the contrary, the "minimum representation" of x

$$x = \sum_{i=1}^{p} \lambda_i x_i \qquad [\lambda_i \ge 0, \sum_i \lambda_i = 1]$$

has p > nm + 1 terms.

3⁰. Consider the homogeneous system of linear equations in p variables δ_i

 $\begin{cases} (a) \sum_{i=1}^{p} \delta_{i}x_{i} = 0 \quad [n \text{ linear equations}] \\ (b) \sum_{i} \delta_{i} = 0 \quad [\text{single linear equation}] \end{cases}$ Since p > n + 1, this system has a nontrivial solution δ . Observe that for every $t \ge 0$ one has

$$x = \sum_{i=1}^{p} \underbrace{[\lambda_i + t\delta_i]}_{\lambda_i(t)} x_i \& \sum_i \lambda_i(t) = 1.$$

$$\delta: \qquad \delta \neq 0 \& \sum_{i} \delta_{i} = 0$$
$$\forall t \ge 0: \quad x = \sum_{i=1}^{p} \underbrace{[\lambda_{i} + t\delta_{i}]}_{\lambda_{i}(t)} x_{i} \& \sum_{i} \lambda_{i}(t) = 1.$$

 \diamond When t = 0, all coefficients $\lambda_i(t)$ are nonnegative

♦ When $t \to \infty$, some of the coefficients $\lambda_i(t)$ go to $-\infty$ (indeed, otherwise we would have $\delta_i \ge 0$ for all *i*, which is impossible since $\sum_i \delta_i = 0$ and not all δ_i are zeros). ↓ It follows that the quantity

 $t_* = \max \{t : t \ge 0 \& \lambda_i(t) \ge 0 \forall i\}$

is well defined; when $t = t_*$, all coefficients in the representation

$$x = \sum_{i=1}^{p} \lambda_i(t_*) x_i$$

are nonnegative, sum of them equals to 1, and at least one of the coefficients $\lambda_i(t_*)$ vanishes. This contradicts the assumption of minimality of the original representation of x as a convex combination of x_i . <u>Theorem</u> [Caratheodory, Conic Version.] Let $\emptyset \neq X \subset \mathbb{R}^n$. Then every vector $x \in \text{Cone}(X)$ is a conic combination of *at most* n vectors from X.

<u>Remark:</u> The bounds given by Caratheodory Theorems (usual and conic version) are sharp: \diamond for a simplex Δ with m+1 vertices $v_0, ..., v_m$ one has dim $\Delta = m$, and it takes all the vertices to represent the barycenter $\frac{1}{m+1} \sum_{i=0}^{m} v_i$ as a convex combination of the vertices; \diamond The conic hull of n standard basic orths in \mathbf{R}^n is exactly the nonnegative orthant \mathbf{R}^n_+ , and it takes all these vectors to get, as their conic combination, the n-dimensional vector of ones. <u>Problem:</u> Supermarkets sell 99 different herbal teas; every one of them is certain blend of 26 herbs A,...,Z. In spite of such a variety of marketed blends, John is not satisfied with any one of them; the only herbal tea he likes is their mixture, in the proportion

1:2:3:...:98:99

Once it occurred to John that in order to prepare his favorite tea, there is no necessity to buy all 99 marketed blends; a smaller number of them will do. With some arithmetics, John found a combination of 66 marketed blends which still allows to prepare his tea. Do you believe John's result can be improved? <u>Theorem</u> [Radon] Let $x_1, ..., x_m$ be $m \ge n+2$ vectors in \mathbb{R}^n . One can split these vectors into two nonempty and non-overlapping groups A, B such that

$\operatorname{Conv}(A) \cap \operatorname{Conv}(B) \neq \emptyset.$

Proof. Consider the homogeneous system of linear equations in m variables δ_i :

 $\begin{cases} \sum_{i=1}^{m} \delta_{i} x_{i} = 0 \quad [n \text{ linear equations}] \\ \sum_{i=1}^{m} \delta_{i} = 0 \quad [\text{single linear equation}] \end{cases}$

Since $m \ge n + 2$, the system has a nontrivial solution δ . Setting $I = \{i : \delta_i > 0\}, J = \{i : \delta_i \le 0\}$, we split the index set $\{1, ..., m\}$ into two *nonempty* (due to $\delta \ne 0, \sum_i u\delta_i = 0$)

groups such that

$$\sum_{i \in I} \delta_i x_i = \sum_{j \in J} [-\delta_j] x_j$$
$$\gamma = \sum_{i \in I} \delta_i = \sum_{j \in J} -\delta_j > 0$$

whence

$$\sum_{i \in I} \frac{\delta_i}{\gamma} x_i = \sum_{\substack{j \in J}} \frac{-\delta_j}{\gamma} x_j$$

$$\in \operatorname{Conv}(\{x_i : i \in A\}) = \operatorname{Conv}(\{x_j : j \in B\})$$

<u>Theorem</u> [Helley] Let $A_1, ..., A_M$ be convex sets in \mathbb{R}^n . Assume that every n+1 sets from the family have a point in common. Then all sets also have point in common.

Proof: induction in *M*. Base $M \le n+1$ is trivially true.

<u>Step</u>: Assume that for certain $M \ge n+1$ our statement hods true for every *M*-member family of convex sets, and let us prove that it holds true for M+1-member family of convex sets $A_1, ..., A_{M+1}$.

 \diamondsuit By inductive hypotheses, every one of the $M+1~{\rm sets}$

 $B_{\ell} = A_1 \cap A_2 \cap \dots \cap A_{\ell-1} \cap A_{\ell+1} \cap \dots \cap A_{M+1}$

is nonempty. Let us choose $x_{\ell} \in B_{\ell}$, $\ell = 1, ..., M + 1$. \diamond By Radon's Theorem, the collection $x_1, ..., x_{M+1}$ can be split in two sub-collections with intersecting convex hulls. W.I.o.g., let the split be $\{x_1, ..., x_{J-1}\} \cup$ $\{x_J, ..., x_{M+1}\}$, and let

 $z \in \text{Conv}(\{x_1, ..., x_{J-1}\}) \cap \text{Conv}(\{x_J, ..., x_{M+1}\}).$

<u>Situation:</u> x_j belongs to all sets A_ℓ except, perhaps, for A_j and

 $z \in \text{Conv}(\{x_1, ..., x_{J-1}\}) \bigcap \text{Conv}(\{x_J, ..., x_{M+1}\}).$

<u>Claim</u>: $z \in A_{\ell}$ for all $\ell \leq M + 1$ (Q.E.D.) Indeed, for $\ell \leq J-1$, the points $x_J, x_{J+1}, ..., x_{M+1}$ belong to the convex set A_{ℓ} , whence

 $z \in \mathsf{Conv}(\{x_J, ..., x_{M+1}\}) \subset A_{\ell}.$

For $\ell \geq J$, the points $x_1, ..., x_{J-1}$ belong to the convex set A_{ℓ} , whence

 $z \in \mathsf{Conv}(\{x_1, ..., x_{J-1}\}) \subset A_{\ell}.$

<u>Refinement:</u> Assume that $A_1, ..., A_M$ are convex sets in \mathbb{R}^n and that

 \diamondsuit the union $A_1 \cup A_2 \cup \ldots \cup A_M$ of the sets belongs to an affine subspace P of affine dimension m

 \diamondsuit every m+1 sets from the family have a point in common

Then all the sets have a point in common.

Proof. We can think of A_j as of sets in P, or, which is the same, as sets in \mathbb{R}^m and apply the Helley Theorem!

Exercise: We are given a function f(x) on a 7,000,000-point set $X \subset \mathbf{R}$. At every 7-point subset of X, this function can be approximated, within accuracy 0.001 at every point, by appropriate polynomial of degree 5. To approximate the function on the entire X, we want to use a spline of degree 5 (a piecewise polynomial function with pieces of degree 5). How many pieces do we need to get accuracy 0.001 at every point?

<u>Answer</u>: Just one. Indeed, let A_x , $x \in X$, be the set of coefficients of all polynomials of degree 5 which reproduce f(x) within accuracy 0.001:

$$A_x = \left\{ p = (p_0, ..., p_5) \in \mathbf{R}^6 : | f(x) - \sum_{i=0}^5 p_i x^i | \le 0.001 \right\}.$$

The set A_x is polyhedral and therefore convex, and we know that every 6 + 1 = 7 sets from the family $\{A_x\}_{x \in X}$ have a point in common. By Helley Theorem, all sets A_x , $x \in X$, have a point in common, that is, there exists a *single* polynomial of degree 5 which approximates f within accuracy 0.001 at *every* point of X.

Exercise: We should design a factory which, mathematically, is described by the following Linear Programming model:

```
\begin{array}{rcl} Ax & \geq & d & [d_1, ..., d_{1000}: \text{ demands}] \\ Bx & \leq & f & [f_1, ..., f_{10}: \text{ facility capacities}] \\ Cx & \leq & c & [\text{other constraints}] \end{array}
```

(F) The data A, B, C, c are given in advance. We should create in advance facility capacities f_i , i = 1, ..., 10, in such a way that the factory will be capable to satisfy all demand scenarios d from a given finite set D, that is, (F) should be feasible for every $d \in D$. Creating capacity f_i of *i*-th facility costs us $a_i f_i$.

It is known that in order to be able to satisfy every single demand from D, it suffices to invest \$1 in creating the facilities.

How large should be investment in facilities in the cases when ${\cal D}$ contains

- \diamond just one scenario?
- ♦ 3 scenarios?
- ♦ 10 scenarios?
- ♦ 2004 scenarios?

<u>Answer</u>: $D = \{d_1\} \Rightarrow \1 is enough $D = \{d_1, d_2, d_3\} \Rightarrow \3 is enough $D = \{d_1, ..., d_{10}\} \Rightarrow \10 is enough $D = \{d_1, ..., d_{2004}\} \Rightarrow \11 is enough! Indeed, for $d \in D$ let F_d be the set of all $f \in \mathbf{R}^{10}$, $f \ge 0$ which cost at most \$11 and result in solvable system

$$\begin{array}{rcl} Ax & \geq & d \\ Bx & \leq & f \\ Cx & \leq & c \end{array} \qquad (F[d])$$

in variables x. The set F_d is convex (why?), and every 11 sets of this type have a common point. Indeed, given 11 scenarios $d^1, ..., d^{11}$ from D, we can "materialize" d^i with appropriate $f^i \ge 0$ at the cost of \$1; therefore we can "materialize" every one of the 11 scenarios $d^1, ..., d^{11}$ by a single vector of capacities $f^1 + ... + f^{11}$ at the cost of \$11, and therefore this vector belongs to $F_{d^1}, ..., F_{d^{11}}$. Since every 11 of 2004 convex sets $F_d \subset \mathbb{R}^{10}$, $d \in D$, have a point in common, all these sets

have a point f in common; for this f, every one of the systems $(F[d]), d \in D$, is solvable.

Exercise: Consider an optimization program

$$c_* = \left\{ c^T x : g_i(x) \le 0, \ i = 1, ..., 2004 \right\}$$

with 11 variables $x_1, ..., x_{11}$. Assume that the constraints are convex, that is, every one of the sets

$$X_i = \{x : g_i(x) \le 0\}, i = 1, ..., 2004$$

is convex. Assume also that the problem is solvable with optimal value 0.

Clearly, when dropping one or more constraints, the optimal value can only decrease or remain the same.

♦ Is it possible to find a constraint such that dropping it, we preserve the optimal value? Two constraints which can be dropped simultaneously with no effect on the optimal value? Three of them? <u>Answer:</u> You can drop as many as 2004 - 11 = 1993 appropriately chosen constraints without varying the optimal value!

Assume, on the contrary, that every 11constraint relaxation of the original problem has negative optimal value. Since there are finitely many such relaxations, there exists $\epsilon < 0$ such that every problem of the form

$$\min_{x} \{ c^{T} x : g_{i_{1}}(x) \le 0, ..., g_{i_{11}}(x) \le 0 \}$$

has a feasible solution with the value of the objective $< -\epsilon$. Since this problem has a feasible solution with the value of the objective equal to 0 (namely, the optimal solution of the original problem) and its feasible set is convex, the problem has a feasible solution x with $c^T x = -\epsilon$. In other words, every 11 of the 2004 sets

 $Y_i = \{x : c^T x = -\epsilon, g_i(x) \le 0\}, i = 1, ..., 2004$ have a point in common.

Every 11 of the 2004 sets

 $Y_i = \{x : c^T x = -\epsilon, g_i(x) \le 0\}, i = 1, ..., 2004$

have a point in common!

The sets Y_i are convex (as intersections of convex sets X_i and an affine subspace). If $c \neq 0$, then these sets belong to affine subspace of affine dimension 10, and since every 11 of them intersect, all 2004 intersect; a point x from their intersection is a feasible solution of the original problem with $c^T x < 0$, which is impossible.

When c = 0, the claim is evident: we can drop all 2004 constraints without varying the optimal value! <u>Helley Theorem II:</u> Let A_{α} , $\alpha \in \mathcal{A}$, be a family of convex sets in \mathbb{R}^n such that every n+1sets from the family have a point in common. Assume, in addition, that

 \diamondsuit the sets A_{α} are closed

 \diamond one can find finitely many sets $A_{\alpha_1}, ..., A_{\alpha_M}$ with a bounded intersection.

Then all sets A_{α} , $\alpha \in \mathcal{A}$, have a point in common.

Proof. By the Helley Theorem, every finite collection of the sets A_{α} has a point in common, and it remains to apply the following standard fact from Analysis:

Let B_{α} be a family of closed sets in \mathbf{R}^n such that

 every finite collection of the sets has a nonempty intersection;

♦ in the family, there exists finite collection with bounded intersection.

Then all sets from the family have a point in common.

Proof of the Standard Fact is based upon the following fundamental property of \mathbf{R}^n :

Every closed and bounded subset of \mathbf{R}^n is a compact set.

Recall two equivalent definitions of a compact set:

A subset X in a metric space M is called compact, if from every sequence of points of X one can extract a sub-sequence converging to a point from X

A subset X in a metric space M is called compact, if from every open covering of X(i.e., from every family of open sets such that every point of X belongs to at least one of them) one can extract a finite sub-covering. Now let B_{α} be a family of closed sets in \mathbb{R}^n such that every finite sub-family of the sets has a nonempty intersection and at least one of these intersection, let it be B, is bounded. Let us prove that all sets B_{α} have a point in common.

• Assume that it is not the case. Then for every point $x \in B$ there exists a set B_{α} which does not contain x. Since B_{α} is closed, it does not intersect an appropriate open ball V_x centered at x. Note that the system $\{V_x :$ $x \in B\}$ forms an open covering of B.

• By its origin, B is closed (as intersection of closed sets) and bounded and *thus is a compact set.* Therefore one can find a *finite* collection $V_{x_1}, ..., V_{x_M}$ which covers B. For every $i \leq M$, there exists a set B_{α_i} in the family which does not intersect V_{x_i} ; therefore $\bigcap_{i=1}^{M} B_{\alpha_i}$ does not intersect B. Since B itself is i=1

the intersection of finitely many sets B_{α} , we see that the intersection of finitely many sets B_{α} (those participating in the description of B and the sets $B_{\alpha_1}, \ldots, B_{\alpha_M}$) is empty, which is a contradiction.

Theory of Systems of Linear Inequalities, I Homogeneous Farkas Lemma

♣ Question: When a vector $a \in \mathbb{R}^n$ is a conic combination of given vectors $a_1, ..., a_m$? <u>Answer:</u> [Homogeneous Farkas Lemma] *A* vector $a \in \mathbb{R}^n$ can be represented as a conic combination $\sum_{i=1}^m \lambda_i a_i, \lambda_i \ge 0$, of given vectors $a_1, ..., a_m \in \mathbb{R}^n$ iff the homogeneous linear inequality

$$a^T x \ge 0 \tag{1}$$

is a consequence of the system of homogeneous linear inequalities

$$a_i^T x \ge 0, \ i = 1, ..., m$$
 (S)

i.e., iff the following implication is true:

$$a_i^T x \ge 0, i = 1, ..., m \Rightarrow a^T x \ge 0.$$
 (*)

Proof, \Rightarrow : If $a = \sum_{i=1}^{m} \lambda_i a_i$ with $\lambda_i \ge 0$, then, of course, $a^T x = \sum_i \lambda_i a_i^T x$ for all x, and (*) clearly is true.

$$a_i^T x \ge 0, i = 1, ..., m \Rightarrow a^T x \ge 0.$$
 (*)

Proof, \Leftarrow : Assume that (*) is true, and let us prove that *a* is a conic combination of $a_1, ..., a_m$. The case of a = 0 is trivial, thus assume that $a \neq 0$.

1⁰. Let

$$A_i = \left\{ x : a^T x = -1, a_i^T x \ge 0 \right\}, i = 1, ..., m.$$

Note that A_i are convex sets, and their intersection is empty by (*). Consider a minimal, in number of sets, sub-family of the family $A_1, ..., A_m$ with empty intersection; w.l.o.g. we may assume that this sub-family is $A_1, ..., A_k$. Thus, the k sets $A_1, ..., A_k$ do not intersect, while every k-1 sets from this family do intersect.

2⁰. <u>Claims:</u> **A:** Vector a is a linear combination of vectors $a_1, ..., a_k$

B: Vectors $a_1, ..., a_k$ are linearly independent

Situation: The sets $A_i = \{x : a^T x = -1, a_i^T x \ge 0\}$, i = 1, ..., k, are such that the k sets $A_1, ..., A_k$ do not intersect, while every k - 1 sets from this family do intersect.

<u>Claim A:</u> Vector a is a linear combination of vectors $a_1, ..., a_k$

Proof of A: Assuming that $a \notin \text{Lin}(\{a_1, ..., a_k\})$, there exists a vector h with $a^T h \neq 0$ and $a_i^T h = 0$, i = 1, ..., k (you can take as h the projection of a onto the orthogonal complement of $a_1, ..., a_k$). Setting $x = -(a^T h)^{-1}h$, we get $a^T x = -1$ and $a_i^T x = 0$, i = 1, ..., k, so that $x \in A_i$, i = 1, ..., k, which is a contradiction.

Situation: The sets $A_i = \{x : a^T x = -1, a_i^T x \ge 0\}$, i = 1, ..., k, are such that the k sets $A_1, ..., A_k$ do not intersect, while every k - 1 sets from this family do intersect.

<u>Claim B:</u> Vectors $a_1, ..., a_k$ are linearly independent

Proof of B: The case of k = 1 is evident. Indeed, in this case we should prove that $a_1 \neq 0$; if $a_1 = 0$ and k = 1, then by Claim A also a = 0, which is not the case. Thus, let us prove Claim B in the case of k > 1. Assume, on the contrary, that

 $E = Lin(\{a_1, ..., a_k\})$

is of dimension r < k, and let

 $B_i = \{x \in E : a^T x = -1, a_i^T x \ge 0\}.$

♦ By Claim A, $0 \neq a \in E$, so that B_i belong to (r-1)-dimensional affine subspace $M = \{x \in E : a^T x = -1\}$. We claim that every r of the sets $B_1, ..., B_k$ have a point in common. Indeed, to prove, e.g., that $B_1, ..., B_r$ have a point in common, note that r < k and therefore the sets $A_1, ..., A_r$ have a point in common; projecting this point onto E, we clearly get a point from $B_1 \cap ... \cap B_r$. ♦ $B_1, ..., B_k$ are convex subsets in (m - 1)dimensional affine subspace, and every mof them have a common point. By Helley Theorem, all k sets $B_1, ..., B_k$ have a common point, whence $A_1, ..., A_k$ have a common point, which is a contradiction. **3**⁰. By Claim A,

$$a = \sum_{i=1}^{k} \lambda_i a_i$$

with certain coefficients λ_i . All we need is to prove that $\lambda_i \ge 0$. To this end assume that certain λ_i , say, λ_1 , is < 0. By claim B, the vectors $a_1, ..., a_k$ are linearly independent, and therefore there exists a vector x such that

$$a_1^T x = 1, a_i^T x = 0, i = 2, ..., k.$$

Since $\lambda_1 < 0$, it follows that

$$a^T x = \lambda_1 a_1^T x < 0,$$

and the vector

$$\bar{x} = \frac{x}{|a^T x|}$$

clearly satisfies

 $a^T \bar{x} = -1, \ a_1^T \bar{x} > 0, \ a_i^T \bar{x} = 0, \ i = 2, ..., k,$

that is, $\overline{x} \in A_1 \cap ... \cap A_k$, which is a contradiction.

Theory of Systems of Linear Inequalities, II Theorem on Alternative

A general (finite!) system of linear inequalities with unknowns $x \in \mathbb{R}^n$ can be written down as

$$\begin{array}{ll}
a_i^T x &> b_i, \ i = 1, ..., m_{\mathsf{S}} \\
a_i^T x &\geq b_i, \ i = m_{\mathsf{S}} + 1, ..., m
\end{array} \tag{S}$$

Question: How to certify that (S) is solvable? Answer: A solution is a certificate of solvability!

Question: How to certify that S is <u>not</u> solvable?

Answer: ???

$$\begin{array}{ll}
a_i^T x &> b_i, \, i = 1, ..., m_{\mathsf{S}} \\
a_i^T x &\geq b_i, \, i = m_{\mathsf{S}} + 1, ..., m
\end{array} \tag{S}$$

<u>Question:</u> How to certify that S is <u>not</u> solvable?

<u>Conceptual sufficient insolvability condition:</u> If we can lead the assumption that x solves (S) to a contradiction, then (S) has no solutions.

"Contradiction by linear aggregation": Let us associate with inequalities of (S) nonnegative weights λ_i and sum up the inequalities with these weights. The resulting inequality

$$\begin{bmatrix} m \\ \sum_{i=1}^{m} \lambda_{i} a_{i} \end{bmatrix}^{T} x \begin{cases} > \sum_{i} \lambda_{i} b_{i}, & \sum_{i=1}^{m_{s}} \lambda_{s} > 0 \\ \ge \sum_{i} \lambda_{i} b_{i}, & \sum_{i=1}^{m_{s}} \lambda_{s} = 0 \end{cases}$$
(C)

by its origin is a consequence of (S), that is, it is satisfied at every solution to (S). Consequently, if there exist $\lambda \ge 0$ such that (C) has no solutions at all, then (S) has no solutions! Question: When a linear inequality

$$d^T x \left\{ \begin{array}{l} > \\ \ge \end{array} e \right.$$

has no solutions at all?

<u>Answer</u>: This is the case if and only if d = 0 and

— either the sign is ">", and $e \ge 0$,

— or the sign is " \geq ", and e > 0.

<u>Conclusion:</u> Consider a system of linear inequalities

$$\begin{array}{lll}
a_{i}^{T}x &> b_{i}, \, i = 1, ..., m_{\mathsf{S}} \\
a_{i}^{T}x &\geq b_{i}, \, i = m_{\mathsf{S}} + 1, ..., m
\end{array} \tag{S}$$

in variables x, and let us associate with it two systems of linear inequalities in variables λ :

$$\mathcal{T}_{\mathrm{I}}: \begin{cases} \lambda \geq 0 \\ \sum \limits_{i=1}^{m} \lambda_{i}a_{i} = 0 \\ \sum \limits_{i=1}^{m_{s}} \lambda_{i} > 0 \\ \sum \limits_{i=1}^{m} \lambda_{i}b_{i} \geq 0 \end{cases} \qquad \mathcal{T}_{\mathrm{II}}: \begin{cases} \lambda \geq 0 \\ \sum \limits_{i=1}^{m} \lambda_{i}a_{i} = 0 \\ \sum \limits_{i=1}^{m_{s}} \lambda_{i}b_{i} \geq 0 \\ \sum \limits_{i=1}^{m} \lambda_{i}b_{i} > 0 \end{cases}$$

<u>If</u> one of the systems T_{I} , T_{II} is solvable, <u>then</u> (S) is unsolvable.

Note: If \mathcal{T}_{II} is solvable, then already the system

$$a_i^T x \ge b_i, \ i = m_{\rm S} + 1, ..., m$$

is unsolvable!

<u>General Theorem on Alternative:</u> *A system of linear inequalities*

$$\begin{array}{ll}
a_{i}^{T}x &> b_{i}, \, i = 1, ..., m_{\mathsf{S}} \\
a_{i}^{T}x &\geq b_{i}, \, i = m_{\mathsf{S}} + 1, ..., m
\end{array} \tag{S}$$

is <u>un</u>solvable <u>iff</u> one of the systems

$$\mathcal{T}_{\mathrm{I}}: \begin{cases} \lambda \geq 0 \\ \sum \limits_{i=1}^{m} \lambda_{i} a_{i} = 0 \\ \sum \limits_{i=1}^{m_{s}} \lambda_{i} > 0 \\ \sum \limits_{i=1}^{m_{s}} \lambda_{i} b_{i} \geq 0 \end{cases} \qquad \mathcal{T}_{\mathrm{II}}: \begin{cases} \lambda \geq 0 \\ \sum \limits_{i=1}^{m} \lambda_{i} a_{i} = 0 \\ \sum \limits_{i=1}^{m_{s}} \lambda_{i} b_{i} \geq 0 \\ \sum \limits_{i=1}^{m} \lambda_{i} b_{i} > 0 \end{cases}$$

is solvable. <u>Note:</u> *The subsystem*

$$a_i^T x \ge b_i, \ i = m_{\mathsf{S}} + 1, ..., m$$

of (S) is <u>un</u>solvable <u>iff</u> \mathcal{T}_{II} is solvable!

Proof. We already know that solvability of one of the systems \mathcal{T}_{I} , \mathcal{T}_{II} is a sufficient condition for <u>un</u>solvability of (S). All we need to prove is that if (S) is unsolvable, then one of the systems \mathcal{T}_{I} , \mathcal{T}_{II} is solvable.

Assume that the system

$$\begin{array}{rcl}
a_{i}^{T}x &> b_{i}, \, i = 1, ..., m_{\mathsf{S}} \\
a_{i}^{T}x &\geq b_{i}, \, i = m_{\mathsf{S}} + 1, ..., m
\end{array} \tag{S}$$

in variables x has no solutions. Then every solution x, τ, ϵ to the homogeneous system of inequalities

has $\epsilon \leq 0$.

Indeed, in a solution with $\epsilon > 0$ one would also have $\tau > 0$, and the vector $\tau^{-1}x$ would solve (S). <u>Situation:</u> Every solution to the system of homogeneous inequalities

$$\begin{aligned}
\tau & -\epsilon &\geq 0 \\
a_i^T x & -b_i \tau & -\epsilon &\geq 0, \ i = 1, ..., m_{\mathsf{S}} \\
a_i^T x & -b_i \tau &\geq 0, \ i = m_{\mathsf{S}} + 1, ..., m
\end{aligned} \tag{U}$$

has $\epsilon \leq 0$, i.e., the homogeneous inequality

$$-\epsilon \ge 0$$
 (I)

is a consequence of system (U) of homogeneous inequalities. By Homogeneous Farkas Lemma, the vector of coefficients in the left hand side of (I) is a conic combination of the vectors of coefficients in the left hand sides of (U):

$$\exists \lambda \ge 0, \nu \ge 0 :$$

$$\sum_{\substack{i=1\\i=1}^{m} \lambda_i a_i = 0$$

$$-\sum_{\substack{i=1\\m_{s}\\i=1}}^{m} \lambda_i b_i + \nu = 0$$

$$-\sum_{\substack{i=1\\i=1}}^{m_{s}} \lambda_i - \nu = -1$$

Assuming that $\lambda_1 = ... = \lambda_{m_s} = 0$, we get $\nu = 1$, and therefore λ solves \mathcal{T}_{II} . In the case of $\sum_{i=1}^{m_s} \lambda_i > 0$, λ clearly solves \mathcal{T}_{I} .
Corollaries of GTA

Principle A: A finite system of linear inequalities has no solutions iff one can lead it to a contradiction by linear aggregation, i.e., an appropriate weighted sum of the inequalities with "legitimate" weights is either a contradictory inequality

$$\mathbf{0}^T x > a \qquad \qquad [a \ge \mathbf{0}]$$

or a contradictory inequality

$$\mathbf{0}^T x \ge a \qquad \qquad [a > \mathbf{0}]$$

Principle B: [Inhomogeneous Farkas Lemma] A linear inequality

 $a^T x \leq b$

is a consequence of <u>solvable</u> system of linear inequalities

$$a_i^T x \le b_i, \ i = 1, ..., m$$

iff the target inequality can be obtained from the inequalities of the system <u>*and*</u> *the identically true inequality*

 $\mathbf{0}^T x \leq \mathbf{1}$

by linear aggregation, that is, iff there exist <u>nonnegative</u> $\lambda_0, \lambda_1, ..., \lambda_m$ such that

$$a = \sum_{i=1}^{m} \lambda_{i} a_{i}$$

$$b = \lambda_{0} + \sum_{i=1}^{m} \lambda_{i} b_{i}$$

$$a = \sum_{\substack{i=1 \\ m}}^{m} \lambda_{i} a_{i}$$

$$b \ge \sum_{\substack{i=1 \\ m}}^{m} \lambda_{i} b_{i}$$

Linear Programming Duality Theorem

The origin of the LP dual of a Linear Programming program

Opt(P) = $\min_{x} \{c^T x : Ax \ge b\}$ (P) is the desire to get a systematic way to bound from below the optimal value in (P). The conceptually simplest bounding scheme is *linear aggregation of the constraints*: Observation: For every vector λ of nonnegative weights, the constraint

 $[A^T \lambda]^T x \equiv \lambda^T A x \ge \lambda^T b$

is a consequence of the constraints of (P)and as such is satisfied at every feasible solution of (P).

Corollary: For every vector $\lambda \ge 0$ such that $\overline{A^T \lambda} = c$, the quantity $\lambda^T b$ is a lower bound on Opt(P).

\clubsuit The problem dual to (P) is nothing but the problem

 $Opt(D) = \max_{\lambda} \left\{ b^T \lambda : \lambda \ge 0, A^T \lambda = c \right\} \quad (D)$

of maximizing the lower bound on Opt(P) given by Corollary.

The origin of (D) implies the following Weak Duality Theorem: The value of the primal objective at every feasible solution of the primal problem

$$Opt(P) = \min_{x} \left\{ c^T x : Ax \ge b \right\}$$
 (P)

is \geq the value of the dual objective at every feasible solution to the dual problem

 $\label{eq:opt} {\rm Opt}(D) = \max_{\lambda} \left\{ b^T \lambda : \lambda \geq 0, A^T \lambda = c \right\} \quad (D)$ that is,

 $\left. \begin{array}{c} x \text{ is feasible for } (P) \\ \lambda \text{ is feasible for } (D) \end{array} \right\} \Rightarrow c^T x \geq b^T \lambda$

In particular,

$$Opt(P) \ge Opt(D).$$

LP Duality Theorem: *Consider an LP pro*gram along with its dual:

 $Opt(P) = \min_{x} \left\{ c^{T}x : Ax \ge b \right\}$ (P)

 $Opt(D) = \max_{\lambda} \left\{ b^T \lambda : A^T \lambda = c, \lambda \ge 0 \right\} (D)$

Then

♦ Duality is symmetric: the problem dual to dual is (equivalent to) the primal

 \diamond The value of the dual objective at every dual feasible solution is \leq the value of the primal objective at every primal feasible solution

♦ The following 5 properties are equivalent to each other:

- (*i*) (*P*) is feasible and bounded (below)
- (*ii*) (*D*) is feasible and bounded (above)
- (iii) (P) is solvable
- (iv) (D) is solvable

(v) both (P) and (D) are feasible and whenever they take place, one has Opt(P) = Opt(D).

$$Opt(P) = \min_{x} \left\{ c^{T}x : Ax \ge b \right\}$$
(P)
$$Opt(D) = \max_{\lambda} \left\{ b^{T}\lambda : A^{T}\lambda = c, \lambda \ge 0 \right\}$$
(D)

 \diamond Duality is symmetric

Proof: Rewriting (D) in the form of (P), we arrive at the problem

$$\min_{\lambda} \left\{ -b^T \lambda : \begin{bmatrix} A^T \\ -A^T \\ I \end{bmatrix} \lambda \ge \begin{bmatrix} c \\ -c \\ 0 \end{bmatrix} \right\},\$$

with the dual being

$$\max_{u,v,w} \left\{ c^T u - c^T v + 0^T w : \begin{array}{c} u \ge 0, v \ge 0, w \ge 0, \\ Au - Av + w = -b \end{array} \right\}$$

$$\lim_{x=v-u,w} \left\{ -c^T x : w \ge 0, Ax = b + w \right\}$$

$$\lim_{x} \left\{ c^T x : Ax \ge b \right\}$$

 \diamondsuit The value of the dual objective at every dual feasible solution is \leq the value of the primal objective at every primal feasible solution

This is Weak Duality

♦ The following 5 properties are equivalent to each other:

(P) is feasible and bounded below (i) \Downarrow

(D) is solvable (iv)

Indeed, by origin of Opt(P), the inequality

$$c^T x \ge \mathsf{Opt}(P)$$

is a consequence of the (solvable!) system of inequalities

$$Ax \ge b.$$

By Principle B, the inequality is a linear consequence of the system:

$$\exists \lambda \geq 0 : A^T \lambda = c \& b^T \lambda \geq \mathsf{Opt}(P).$$

Thus, the dual problem has a feasible solution with the value of the dual objective $\geq Opt(P)$. By Weak Duality, this solution is optimal, and Opt(D) = Opt(P). ♦ The following 5 properties are equivalent to each other: (D) is solvable ↓ (D) is feasible and bounded above (ii) Evident

The following 5 properties are equivalent to each other: (D) is feasible and bounded above (ii) ψ (P) is solvable (iii) Implied by already proved relation (P) is feasible and bounded below (i) ψ (D) is solvable (iv) in view of primal-dual symmetry

We proved that

$$(\mathsf{i}) \Leftrightarrow (\mathsf{ii}) \Leftrightarrow (\mathsf{iii}) \Leftrightarrow (\mathsf{iv})$$

and that when these 4 equivalent properties take place, one has

$$Opt(P) = Opt(D)$$

It remains to prove that properties (i) – (iv) are equivalent to

both
$$(P)$$
 and (D) are feasible (v)

♦ In the case of (v), (P) is feasible and below bounded (Weak Duality), so that (v)⇒(i) ♦ in the case of (i)≡(ii), both (P) and (D) are feasible, so that (i)⇒(v) <u>Theorem:</u> Consider a primal-dual pair of feasible LP programs

 $Opt(P) = \min_{x} \left\{ c^{T}x : Ax \ge b \right\}$ (P) $Opt(D) = \max_{\lambda} \left\{ b^{T}\lambda : A^{T}\lambda = c, \lambda \ge 0 \right\}$ (D)

and let x, λ be <u>feasible</u> solutions to the respective programs. These solutions are optimal for the respective problems \diamondsuit iff $c^T x - b^T \lambda = 0$ ["zero duality gap"] as well as

 $iff [Ax - b]_i \cdot \lambda_i = 0$ for all i ["complementary slackness"]

Proof: Under Theorem's premise, Opt(P) = Opt(D), so that

$$c^T x - b^T \lambda = \underbrace{c^T x - \operatorname{Opt}(P)}_{\geq 0} + \underbrace{\operatorname{Opt}(D) - b^T \lambda}_{\geq 0}$$

Thus, duality gap $c^T x - b^T \lambda$ is always nonnegative and is zero iff x, λ are optimal for the respective problems. The complementary slackness condition is given by the identity

$$c^T x - b^T \lambda = (A^T \lambda)^T x - b^T \lambda = [Ax - b]^T \lambda$$

Since both [Ax-b] and λ are nonnegative, duality gap is zero iff the complementary slackness holds true.

♣ Every linear form f(x) on \mathbb{R}^n is representable via inner product:

$$f(x) = f^T x$$

for appropriate vector $f \in \mathbb{R}^n$ uniquely defined by the form. Nontrivial (not identically zero) forms correspond to nonzero vectors f.

A level set

$$M = \left\{ x : f^T x = a \right\} \tag{(*)}$$

of a *nontrivial* linear form on \mathbb{R}^n is affine subspace of affine dimension n - 1; vice versa, every affine subspace M of affine dimension n - 1 in \mathbb{R}^n can be represented by (*) with appropriately chosen $f \neq 0$ and a; f and aare defined by M up to multiplication by a common nonzero factor.

(n-1)-dimensional affine subspaces in \mathbb{R}^n are called *hyperplanes*.

$$M = \left\{ x : f^T x = a \right\} \tag{(*)}$$

Level set (*) of nontrivial linear form splits \mathbb{R}^n into two parts:

$$M_{+} = \{x : f^{T}x \ge a\} \\ M_{-} = \{x : f^{T}x \le a\}$$

called *closed half-spaces* given by (f, a); the hyperplane M is the common boundary of these half-spaces. The interiors M_{++} of M_{+} and M_{--} of M_{-} are given by

$$M_{++} = \{x : f^T x > a\} \\ M_{--} = \{x : f^T x < a\}$$

and are called *open half-spaces* given by (f, a). We have

 $\mathbb{R}^n = M_- \bigcup M_+ \quad [M_- \bigcap M_+ = M]$

and

$$\mathbb{R}^n = M_{--} \bigcup M \bigcup M_{++}$$

♣ <u>Definition.</u> Let T, S be two nonempty sets in \mathbb{R}^n .

(i) We say that a hyperplane

$$M = \{x : f^T x = a\} \tag{(*)}$$

separates S and T, if

 $\diamondsuit S \subset M_{-}, T \subset M_{+}$ ("S does not go above M, and T does not go below M") and

 $\diamondsuit \ S \cup T \not\subset M.$

(ii) We say that a nontrivial linear form $f^T x$ separates S and T if, for properly chosen a, the hyperplane (*) separates S and T.

Examples: The linear form x_1 on \mathbb{R}^2 1) separates the sets







<u>Observation:</u> A linear form f^Tx separates nonempty sets S, T iff

$$\begin{split} \sup_{\substack{x \in S \\ x \in S}} f^T x &\leq \inf_{\substack{y \in T \\ y \in T}} f^T y \\ \sup_{x \in S} f^T x &< \sup_{y \in T} f^T y \\ \end{split} \tag{*}$$

In the case of (*), the associated with f hyperplanes separating S and T are exactly the hyperplanes

 $\{x : f^T x = a\}$ with $\sup_{x \in S} f^T x \le a \le \inf_{y \in T} f^T y.$

Separation Theorem: Two nonempty <u>convex</u> sets S, T can be separated iff their relative interiors do not intersect.

<u>Note:</u> In this statement, convexity of both S and T is crucial!

S. Т

Proof, \Rightarrow : (!) <u>If</u> nonempty convex sets S, Tcan be separated, <u>then</u> rint $S \cap$ rint $T = \emptyset$ <u>Lemma.</u> Let X be a convex set, $f(x) = f^T x$ be a linear form and $a \in$ rint X. Then

$$f^T a = \max_{x \in X} f^T x \Leftrightarrow f(\cdot) \Big|_X = \text{const.}$$

♣ Lemma ⇒ (!): Let $a \in \operatorname{rint} S \cap \operatorname{rint} T$. Assume, on contrary to what should be proved, that $f^T x$ separates S, T, so that

$$\sup_{x \in S} f^T x \le \inf_{y \in T} f^T y.$$

♦ Since $a \in T$, we get $f^T a \ge \sup_{x \in S} f^T x$, that is, $f^T a = \max_{x \in S} f^T x$. By Lemma, $f^T x = f^T a$ for all $x \in S$. ♦ Since $a \in S$, we get $f^T a \le \inf_{y \in T} f^T y$, that is, $f^T a = \min_{y \in T} f^T y$. By Lemma, $f^T y = f^T a$ for all $y \in T$. Thus,

$$z \in S \cup T \Rightarrow f^T z \equiv f^T a,$$

so that f does <u>not</u> separate S and T, which is a contradiction.

<u>Lemma.</u> Let X be a convex set, $f(x) = f^T x$ be a linear form and $a \in \text{rint } X$. Then

1

$$f^T a = \max_{x \in X} f^T x \Leftrightarrow f(\cdot) \Big|_X = \text{const.}$$

Proof. Shifting X, we may assume a = 0. Let, on the contrary to what should be proved, $f^T x$ be non-constant on X, so that there exists $y \in X$ with $f^T y \neq f^T a = 0$. The case of $f^T y > 0$ is impossible, since $f^T a = 0$ is the maximum of $f^T x$ on X. Thus, $f^T y < 0$. The line $\{ty : t \in \mathbb{R}\}$ passing through 0 and through y belongs to Aff(X); since $0 \in \text{rint } X$, all points $z = -\epsilon y$ on this line belong to X, provided that $\epsilon > 0$ is small enough. At every point of this type, $f^T z > 0$, which contradicts the fact that $\max_{x \in X} f^T x = f^T a = 0$. **Proof,** \Leftarrow : Assume that *S*, *T* are nonempty convex sets such that rint $S \cap$ rint $T = \emptyset$, and let us prove that *S*, *T* can be separated. **Step 1: Separating a point and a convex hull of a finite set.** Let $S = \text{Conv}(\{b_1, ..., b_m\})$ and $T = \{b\}$ with $b \notin S$, and let us prove that *S* and *T* can be separated. $\mathbf{1}^0$. Let

$$\beta_i = \begin{bmatrix} b_i \\ 1 \end{bmatrix}, \ \beta = \begin{bmatrix} b \\ 1 \end{bmatrix}.$$

Observe that β is *not* a conic combination of $\beta_1, ..., \beta_m$:

$$\beta_i = \begin{bmatrix} b_i \\ 1 \end{bmatrix}, \ \beta = \begin{bmatrix} b \\ 1 \end{bmatrix}.$$

2⁰. Since β is not conic combination of β_i , by Homogeneous Farkas Lemma there exists $h = \begin{bmatrix} f \\ -a \end{bmatrix}$ such that

 $f^Tb-a\equiv h^T\beta>0\geq h^T\beta_i\equiv f^Tb_i-a,\,i=1,...,m$ that is,

$$f^T b > \max_{i=1,...,m} f^T b_i = \max_{x \in S = \text{Conv}(\{b_1,...,b_m\})} f^T x.$$

Note: We have used the evident fact that

$$\max_{x \in \text{Conv}(\{b_1, \dots, b_m\})} f^T x \equiv \max_{\lambda \ge 0, \sum_i \lambda_i = 1} \frac{f^T[\sum_i \lambda_i b_i]}{i}$$
$$= \max_{\lambda \ge 0, \sum_i \lambda_i = 1} \sum_i \lambda_i [f^T b_i]$$
$$= \max_i f^T b_i.$$

Step 2: Separating a point and a convex set which does not contain the point. Let S be a nonempty convex set and $T = \{b\}$ with $b \notin S$, and let us prove that S and T can be separated.

1⁰. Shifting *S* and *T* by -b (which clearly does not affect the possibility of separating the sets), we can assume that $T = \{0\} \not\subset S$. 2⁰. Replacing, if necessary, \mathbb{R}^n with Lin(*S*), we may further assume that $\mathbb{R}^n = \text{Lin}(S)$.

Lemma: Every nonempty subset S in \mathbb{R}^n is separable: one can find a sequence $\{x_i\}$ of points from S which is dense in S, i.e., is such that every point $x \in S$ is the limit of an appropriate subsequence of the sequence. **Lemma** \Rightarrow **Separation:** Let $\{x_i \in S\}$ be a sequence which is dense in *S*. Since *S* is convex and does not contain 0, we have

$$0 \not\in \mathsf{Conv}(\{x_1, ..., x_i\}) \ \forall i$$

whence

$$\exists f_i : 0 = f_i^T 0 > \max_{1 \le j \le i} f_i^T x_j.$$
 (*)

By scaling, we may assume that $||f_i||_2 = 1$. The sequence $\{f_i\}$ of unit vectors possesses a converging subsequence $\{f_{i_s}\}_{s=1}^{\infty}$; the limit f of this subsequence is, of course, a unit vector. By (*), for every fixed j and all large enough s we have $f_{i_s}^T x_j < 0$, whence

$$f^T x_j \le 0 \ \forall j. \tag{**}$$

Since $\{x_j\}$ is dense in S, (**) implies that $f^T x \leq 0$ for all $x \in S$, whence

$$\sup_{x \in S} f^T x \le 0 = f^T 0.$$

Situation: (a) $Lin(S) = \mathbb{R}^n$ (b) $T = \{0\}$ (c) We have built a unit vector f such that $\sup f^T x < 0 = f^T 0.$ (!)

$$\sup_{x \in S} f^T x \le 0 = f^T 0. \tag{!}$$

By (!), all we need to prove that f separates $T = \{0\}$ and S is to verify that

$$\inf_{x \in S} f^T x < f^T 0 = 0.$$

Assuming the opposite, (!) would say that $f^T x = 0$ for all $x \in S$, which is impossible, since $\text{Lin}(S) = \mathbb{R}^n$ and f is nonzero.

<u>Lemma:</u> Every nonempty subset S in \mathbb{R}^n is separable: one can find a sequence $\{x_i\}$ of points from S which is dense in S, i.e., is such that every point $x \in S$ is the limit of an appropriate subsequence of the sequence.

Proof. Let $r_1, r_2, ...$ be the countable set of all rational vectors in \mathbb{R}^n . For every positive integer t, let $X_t \subset S$ be the countable set given by the following construction:

We look, one after another, at the points $r_1, r_2, ...$ and for every point r_s check whether there is a point z in S which is at most at the distance 1/t away from r_s . If points z with this property exist, we take one of them and add it to X_t and then pass to r_{s+1} , otherwise directly pass to r_{s+1} .

Is is clear that

(*) Every point $x \in S$ is at the distance at most 2/t from certain point of X_t .

Indeed, since the rational vectors are dense in \mathbb{R}^n , there exists s such that r_s is at the distance $\leq \frac{1}{t}$ from x. Therefore, when processing r_s , we definitely add to X_t a point zwhich is at the distance $\leq 1/t$ from r_s and thus is at the distance $\leq 2/t$ from x.

By construction, the countable union $\bigcup_{t=1}^{\infty} X_t$ of countable sets $X_t \subset S$ is a countable set in S, and by (*) this set is dense in S. Step 3: Separating two non-intersecting nonempty convex sets. Let S, T be nonempty convex sets which do not intersect; let us prove that S,T can be separated. Let $\hat{S} = S - T$ and $\hat{T} = \{0\}$. The set \hat{S} clearly is convex and does not contain 0 (since $S \cap T = \emptyset$). By Step 2, \hat{S} and $\{0\} = \hat{T}$ can be separated: there exists f such that

$$\begin{cases} \sup_{x \in S} f^T s - \inf_{y \in T} f^T y \\ \underbrace{\sup_{x \in S, y \in T} [f^T x - f^T y]}_{x \in S, y \in T} &\leq 0 = \inf_{z \in \{0\}} f^T z \\ \underbrace{\inf_{x \in S, y \in T} [f^T x - f^T y]}_{\inf_{x \in S} f^T x - \sup_{y \in T} f^T y} &< 0 = \sup_{z \in \{0\}} f^T z \\ \underbrace{\inf_{x \in S} f^T x - \sup_{y \in T} f^T y}_{y \in T} &\leq 0 = \sup_{z \in \{0\}} f^T z \end{cases}$$

whence

$$\sup_{x \in S} f^T x \leq \inf_{y \in T} f^T y \\ \inf_{x \in S} f^T x < \sup_{y \in T} f^T y \\ y \in T$$

Step 4: Completing the proof of Separation Theorem. Finally, let S, T be nonempty convex sets with non-intersecting relative interiors, and let us prove that S, T can be separated.

As we know, the sets $S' = \operatorname{rint} S$ and $T' = \operatorname{rint} T$ are convex and nonempty; we are in the situation when these sets do not intersect. By Step 3, S' and T' can be separated: for properly chosen f, one has

$$\sup_{x \in S'} f^T x \leq \inf_{y \in T'} f^T y \\ \inf_{x \in S'} f^T x < \sup_{y \in T'} f^T y \\ (*)$$

Since S' is dense in S and T' is dense in T, inf's and sup's in (*) remain the same when replacing S' with S and T' with T. Thus, fseparates S and T. Alternative proof of Separation Theorem starts with separating a point $T = \{a\}$ and a *closed* convex set S, $a \notin S$, and is based on the following fact:

Let S be a nonempty closed convex set and let $a \notin S$. There exists a unique closest to a point in S:

 $\operatorname{Proj}_{S}(a) = \operatorname{argmin}_{x \in S} \|a - x\|_{2}$ and the vector $e = a - \operatorname{Proj}_{S}(a)$ sepa-

rates a and S:

 $\max_{x \in S} e^T x = e^T \operatorname{Proj}_S(a) = e^T a - ||e||_2^2 < e^T a.$



Proof: 1^0 . The closest to *a* point in *S* does exist. Indeed, let $x_i \in S$ be a sequence such that

$$\|a - x_i\|_2 \rightarrow \inf_{x \in S} \|a - x\|_2, \ , \ i \rightarrow \infty$$

The sequence $\{x_i\}$ clearly is bounded; passing to a subsequence, we may assume that $x_i \rightarrow \overline{}$ as $i \to \infty$. Since S is closed, we have $\overline{x} \in S$, and

$$||a - \bar{x}||_2 = \lim_{i \to \infty} ||a - x_i||_2 = \inf_{x \in S} ||a - x||_2.$$

 2^{0} . The closest to a point in S is unique. Indeed, let x, y be two closest to a points in S, so that $||a - x||_2 = ||a - y||_2 = d$. Since S is convex, the point $z = \frac{1}{2}(x+y)$ belongs to S; therefore $||a - z||_2 \ge d$. We now have

$$\underbrace{ = \|2(a-z)\|_{2}^{2} \ge 4d^{2}}_{\|[a-x] + [a-y]\|_{2}^{2}} + \underbrace{ = \|x-y\|^{2}}_{\|[a-x] - [a-y]\|_{2}^{2}} = \underbrace{ = \|x-y\|_{2}^{2} + 2\|a-y\|_{2}^{2}}_{4d^{2}}$$
There $\|x-y\|_{2}^{2} = 0$

whence $||x - y||_2 = 0$.
3⁰. Thus, the closest to a point in S exists and is unique. With $e = a - \operatorname{Proj}_{S}(a)$, we have

$$x \in S, f = x - \operatorname{Proj}_{S}(a)$$

$$\downarrow$$

$$\phi(t) \equiv \|e - tf\|_{2}^{2}$$

$$= \|a - [\operatorname{Proj}_{S}(a) + t(x - \operatorname{Proj}_{S}(a))]\|_{2}^{2}$$

$$\geq \|a - \operatorname{Proj}_{S}(a)\|_{2}^{2}$$

$$= \phi(0), 0 \leq t \leq 1$$

$$\downarrow$$

$$0 \leq \phi'(0) = -2e^{T}(x - \operatorname{Proj}_{S}(a))$$

$$\downarrow$$

$$\forall x \in S : e^{T}x \leq e^{T}\operatorname{Proj}_{S}(a) = e^{T}a - \|e\|_{2}^{2}.$$

♣ Separation of sets S, T by linear form $f^T x$ is called *strict*, if

$$\sup_{x \in S} f^T x < \inf_{y \in T} f^T y$$

<u>Theorem:</u> Let S, T be nonempty convex sets. These sets can be strictly separated iff they are at positive distance:

dist(S,T) =
$$\inf_{x \in S, y \in T} ||x - y||_2 > 0.$$

Proof, \Rightarrow : Let f strictly separate S, T; let us prove that S, T are at positive distance. Otherwise we could find sequences $x_i \in S$, $y_i \in T$ with $||x_i - y_i||_2 \to 0$ as $i \to \infty$, whence $f^T(y_i - x_i) \to 0$ as $i \to \infty$. It follows that the sets on the axis

$$\widehat{S} = \{a = f^T x : x \in S\}, \widehat{T} = \{b = f^T y : y \in T\}$$

are at zero distance, which is a contradiction with

$$\sup_{a \in \widehat{S}} a < \inf_{b \in \widehat{T}} b.$$

Proof, \Leftarrow : Let *T*, *S* be nonempty convex sets which are at positive distance 2δ :

$$2\delta = \inf_{x \in S, y \in T} \|x - y\|_2 > 0.$$

Let

$$S^+ = S + \{z : \|z\|_2 \le \delta\}$$

The sets S^+ and T are convex and do not intersect, and thus can be separated:

$$\sup_{x_+ \in S^+} f^T x_+ \le \inf_{y \in T} f^T y \qquad [f \neq 0]$$

Since

$$\sup_{x_+ \in S^+} f^T x_+ = \sup_{\substack{x \in S, \|z\|_2 \le \delta}} [f^T x + f^T z]$$
$$= [\sup_{x \in S} f^T x] + \delta \|f\|_2,$$

we arrive at

$$\sup_{x \in S} f^T x < \inf_{y \in T} f^T y$$

<u>Exercise</u> Below S is a nonempty convex set and $T = \{a\}$.

Statement	True?
If T and S can be separated	
then $a \not\in S$	
If $a \not\in S$, then T and S can be	
separated	
If T and S can be strictly	
separated, then $a \not\in S$	
If $a \not\in S$, then T and S can be	
strictly separated	
If S is closed and $a \not\in S$, then T	
and S can be strictly separated	

Supporting Planes and Extreme Points

♣ <u>Definition</u>. Let Q be a *closed* convex set in \mathbb{R}^n and \overline{x} be a point from the relative boundary of Q. A hyperplane

$$\Pi = \{ x : f^T x = a \} \qquad [a \neq 0]$$

is called supporting to Q at the point \overline{x} , if the hyperplane separates Q and $\{\overline{x}\}$:

$$\sup_{\substack{x \in Q \\ x \in Q}} f^T x \leq f^T \bar{x}$$
$$\inf_{x \in Q} f^T x < f^t \bar{x}$$

Equivalently: Hyperplane $\Pi = \{x : f^T x = a\}$ supports Q at \bar{x} iff the linear form $f^T x$ attains its maximum on Q, equal to a, at the point \bar{x} and the form is non-constant on Q. <u>Proposition.</u> Let Q be a convex closed set in \mathbb{R}^n and \bar{x} be a point from the relative boundary of Q. Then

 \diamond There exist at least one hyperplane Π which supports Q at \bar{x} ;

 \diamond For every such hyperplane Π , the set $Q \cap \Pi$ has dimension less than the one of Q.

Proof: Existence of supporting plane is given by Separation Theorem. This theorem is applicable since

 $\bar{x} \notin \operatorname{rint} Q \Rightarrow \{\bar{x}\} \equiv \operatorname{rint} \{\bar{x}\} \cap \operatorname{rint} Q = \emptyset.$

Further,

 $Q \nsubseteq \Pi \Rightarrow \mathsf{Aff}(Q) \nsubseteq \Pi \Rightarrow \mathsf{Aff}(\Pi \cap Q) \subsetneqq \mathsf{Aff}(Q),$

and if two *distinct* affine subspaces are embedded one into another, then the dimension of the embedded subspace is strictly less than the dimension of the embedding one.

Extreme Points

♣ <u>Definition</u>. Let Q be a convex set in \mathbb{R}^n and \overline{x} be a point of X. The point is called *extreme*, if it is not a convex combination, with positive weights, of two points of X distinct from \overline{x} :

$$\begin{aligned} \bar{x} \in \mathsf{Ext}(Q) \\ & \updownarrow \\ \{ \bar{x} \in Q \} & \& \left\{ \begin{array}{l} u, v \in Q, \lambda \in (0, 1) \\ \bar{x} = \lambda u + (1 - \lambda)v \end{array} \right\} \Rightarrow u = v = \bar{x} \end{aligned}$$

Equivalently: A point $\overline{x} \in Q$ is extreme iff it is <u>not</u> the midpoint of a nontrivial segment in Q:

 $x \pm h \in Q \Rightarrow h = 0.$

Equivalently: A point $\bar{x} \in Q$ is extreme iff the set $Q \setminus \{\bar{x}\}$ is convex.

Examples:

- 1. Extreme points of [x, y] are ...
- 2. Extreme points of $\triangle ABC$ are ...
- 3. Extreme points of the ball $\{x : ||x||_2 \le 1\}$ are ...

<u>Theorem</u> [Krein-Milman] Let Q be a closed convex and nonempty set in \mathbb{R}^n . Then

 $\diamondsuit Q$ possess extreme points iff Q does not contain lines;

 \diamondsuit If Q is bounded, then Q is the convex hull of its extreme points:

 $Q = \operatorname{Conv}(\operatorname{Ext}(Q))$

so that every point of Q is convex combination of extreme points of Q.

<u>Note:</u> If Q = Conv(A), then $\text{Ext}(Q) \subset A$. Thus, extreme points of a *closed convex bounded* set Q give the *minimal* representation of Q as Conv(...). **Proof.** 1⁰: If closed convex set Q does not contain lines, then $Ext(Q) \neq \emptyset$ Important lemma: Let S be a closed convex set and $\Pi = \{x : f^T x = a\}$ be a hyperplane which supports S at certain point. Then

$\mathsf{Ext}(\Pi \cap S) \subset \mathsf{Ext}(S).$

Proof of Lemma. Let $\bar{x} \in \text{Ext}(\Pi \cap S)$; we should prove that $\bar{x} \in \text{Ext}(S)$. Assume, on the contrary, that \bar{x} is a midpoint of a nontrivial segment $[u, v] \subset S$. Then $f^T \bar{x} =$ $a = \max_{x \in S} f^T x$, whence $f^T \bar{x} = \max_{x \in [u,v]} f^T x$. A linear form can attain its maximum on a segment at the midpoint of the segment iff the form is constant on the segment; thus, $a = f^T \bar{x} = f^T u = f^T v$, that is, $[u, v] \subset \Pi \cap S$. But \bar{x} is an extreme point of $\Pi \cap S$ – contradiction! Let Q be a nonempty closed convex set which does not contain lines. In order to build an extreme point of Q, apply the *Purification algorithm*:

Initialization: Set $S_0 = Q$ and choose $x_0 \in Q$. Step t: Given a nonempty closed convex set S_t which does not contain lines and is such that $\text{Ext}(S_t) \subset \text{Ext}(Q)$ and $x_t \in S_t$,

1) check whether S_t is a singleton $\{x_t\}$. If it is the case, terminate: $x_t \in \text{Ext}\{S_t\} \subset \text{Ext}(Q)$. 2) if S_t is not a singleton, find a point x_{t+1} on the relative boundary of S_t and build a hyperplane Π_t which supports S_t at x_{t+1} .

To find x_{t+1} , take a direction $h \neq 0$ parallel to Aff (S_t) . Since S_t does not contain lines, when moving from x_t either in the direction h, or in the direction -h, we eventually leave S_t , and thus cross the relative boundary of S_t . The intersection point is the desired x_{t+1} . 3) Set $S_{t+1} = S_t \cap \Pi_t$, replace t with t + 1 and loop to 1).



Justification: By Important Lemma,

 $\mathsf{Ext}(S_{t+1}) \subset \mathsf{Ext}(S_t),$

so that

 $\mathsf{Ext}(S_t) \subset \mathsf{Ext}(Q) \ \forall t.$

Besides this, dim $(S_{t+1} < \dim (S_t))$, so that Purification algorithm does terminate. the algorithm

<u>Note</u>: Assume you are given a linear form $g^T x$ which is bounded from above on Q. Then in the Purification algorithm one can easily ensure that $g^T x_{t+1} \ge g^T x_t$. Thus,

If Q is a nonempty closed set in \mathbb{R}^n which does not contain lines and $f^T x$ is a linear form which is bounded above on Q, then for every point $x_0 \in Q$ there exists (and can be found by Purification) a point $\overline{x} \in \text{Ext}(Q)$ such that $g^T \overline{x} \ge g^T x_0$. In particular, if $g^T x$ attains its maximum on Q, then the maximizer can be found among extreme points of Q. **Proof, 2**⁰ If a closed convex set Q contains lines, it has no extreme points.

Another Important Lemma: Let S be a closed convex set such that $\{\bar{x} + th : t \ge 0\} \subset S$ for certain \bar{x} . Then

 $\{x + th : t \ge 0\} \subset S \ \forall x \in S.$

Proof: For every $s \ge 0$ and $x \in S$ we have

$$x + sh = \lim_{i \to \infty} \underbrace{\left[(1 - s/i)x + (s/i)[\bar{x} + th] \right]}_{\in S}.$$

<u>Note</u>: The set of all directions $h \in \mathbb{R}^n$ such that $\{x + th : t \ge 0\} \subset S$ for some (and then, for all) $x \in S$, is called the *recessive cone* Rec(S) of closed convex set S. Rec(S) indeed is a cone, and

$$S + \operatorname{Rec}(S) = S.$$

<u>Corollary</u>: If a closed convex set Q contains a line ℓ , then the parallel lines, passing through points of Q, also belong to Q. In particular, Q possesses no extreme points. **Proof, 3**⁰: If a nonempty closed convex set Q is bounded, then Q = Conv(Ext(Q)).

The inclusion $Conv(Ext(Q)) \subset Q$ is evident. Let us prove the opposite inclusion, i.e., prove that every point of Q is a convex combination of extreme points of Q.

Induction in $k = \dim Q$. Base k = 0 (Q is a singleton) is evident.

Step $k \mapsto k + 1$: Given (k + 1)-dimensional closed and bounded convex set Q and a point $x \in Q$, we, as in the Purification algorithm, can represent x as a convex combination of two points x_+ and x_- from the relative boundary of Q. Let Π_+ be a hyperplane which supports Q at x_+ , and let $Q_+ = \Pi_+ \cap Q$. As we know, Q_+ is a closed convex set such that

dim $Q_+ < \dim Q$, $Ext(Q_+) \subset Ext(Q)$, $x_+ \in Q_+$. Invoking inductive hypothesis,

 $x_+ \in \text{Conv}(\text{Ext}(Q_+)) \subset \text{Conv}(\text{Ext}(Q)).$ Similarly, $x_- \in \text{Conv}(\text{Ext}(Q))$. Since $x \in [x_-, x_+]$, we get $x \in \text{Conv}(\text{Ext}(Q))$. Structure of Polyhedral Sets

♣ <u>Definition</u>: A *polyhedral* set Q in \mathbb{R}^n is a *nonempty* subset in \mathbb{R}^n which is a solution set of a finite system of nonstrict inequalities:

Q is polyhedral $\Leftrightarrow Q = \{x : Ax \ge b\} \neq \emptyset$.

• Every polyhedral set is convex and closed.

Question: When a polyhedral set $Q = \{x : Ax \ge b\}$ contains lines? What are these lines, if any?

<u>Answer:</u> Q contains lines iff A has a nontrivial nullspace:

 $\mathsf{Null}(A) \equiv \{h : Ah = 0\} \neq \{0\}.$

Indeed, a line $\ell = \{x = \bar{x} + th : t \in \mathbb{R}\}, h \neq 0$, belongs to Q iff

$$\forall t : A(\bar{x} + th) \ge b$$

$$\Leftrightarrow \quad \forall t : tAh \ge b - A\bar{x}$$

$$\Leftrightarrow \quad Ah = 0 \& \bar{x} \in Q.$$

<u>Fact</u>: A polyhedral set $Q = \{x : Ax \ge b\}$ always can be represented as

 $Q = Q_* + L,$

where Q_* is a polyhedral set which does not contain lines and L is a linear subspace. In this representation,

 $\diamond L$ is uniquely defined by Q and coincides with Null(A),

 $\Diamond Q_*$ can be chosen, e.g., as

 $Q_* = Q \cap L^{\perp}$

Structure of polyhedral set which does <u>not</u> contain lines

♣ <u>Theorem.</u> Let

$$Q = \{x : Ax \ge b\} \neq \emptyset$$

be a polyhedral set which does not contain lines (or, which is the same, $Null(A) = \{0\}$). Then the set Ext(Q) of extreme points of Qis *nonempty* and *finite*, and

 $Q = \text{Conv}(\text{Ext}(Q)) + \text{Cone}\{r_1, ..., r_S\}$ (*)

for properly chosen vectors $r_1, ..., r_S$. <u>Note:</u> Cone $\{r_1, ..., r_s\}$ is exactly the recessive cone of Q:

Cone
$$\{r_1, ..., r_S\}$$

= $\{r : x + tr \in Q \ \forall (x \in Q, t \ge 0)\}$
= $\{r : Ar \ge 0\}.$

This cone is the trivial cone $\{0\}$ iff Q is a *bounded* polyhedral set (called *polytope*).

Combining the above theorems, we come to the following results:

A polyhedral set Q always can be represented in the form

$$Q = \left\{ x = \sum_{i=1}^{I} \lambda_i v_i + \sum_{j=1}^{J} \mu_j w_j : \begin{array}{c} \lambda \ge 0, \mu \ge 0\\ \sum \lambda_i = 1 \end{array} \right\}$$
(!)

where I, J are positive integers and $v_1, ..., v_I$, $w_1, ..., w_J$ are appropriately chosen points and directions.

Vice versa, every set Q of the form (!) is a polyhedral set.

<u>Note:</u> Polytopes (bounded polyhedral sets) are *exactly* the sets of form (!) with "trivial w-part": $w_1 = ... = w_J = 0$.

Exercise 1: Is it true that the intersection of two polyhedral sets, if nonempty, is a polyhedral set?

Exercise 2: Is it true that the affine image $\{y = Px + p : x \in Q\}$ of a polyhedral set Q is a polyhedral set?

Applications to Linear Programming

Consider a *feasible* Linear Programming program

 $\min_{x} c^{T} x \text{ s.t. } x \in Q = \{x : Ax \ge b\} \qquad (\mathsf{LP})$

<u>Observation</u>: We lose nothing when assuming that $Null(A) = \{0\}$. Indeed, we have

$$Q = Q_* + \operatorname{Null}(A),$$

where Q_* is a polyhedral set not containing lines. If c is not orthogonal to Null(A), then (LP) clearly is unbounded. If c is orthogonal to Null(A), then (LP) is equivalent to the LP program

$$\min_{x} c^{T}x \text{ s.t. } x \in Q_{*},$$

and now the matrix in a representation $Q_* = \{x : \widetilde{A}x \ge \widetilde{b}\}$ has trivial nullspace.

Assuming $Null(A) = \{0\}$, let (LP) be bounded (and thus solvable). Since Q is convex, closed and does not contain lines, among the (nonempty!) set of minimizers the objective on Q there is an extreme point of Q. $\min_{x} c^{T}x \text{ s.t. } x \in Q = \{x : Ax \ge b\}$ (LP) We have proved

<u>Proposition.</u> Assume that (LP) is feasible and bounded (and thus is solvable) and that Null(A) = {0}. Then among optimal solutions to (LP) there exists at least one which is an extreme point of Q.

Question: How to characterize extreme points of the set

$$Q = \{x : Ax \ge b\} \neq \emptyset$$

provided that A is $m \times n$ matrix with Null(A) = $\{0\}$?

<u>Answer:</u> Extreme points \overline{x} of Q are fully characterized by the following two properties:

 $\diamondsuit \ A\bar{x} \ge b$

 \diamond Among constraints $Ax \ge b$ which are *active* at \overline{x} (i.e., are satisfied as equalities), there are n linearly independent. Justification of the answer, \Rightarrow : If \bar{x} is an extreme point of Q, then among the constraints $Ax \ge b$ active at \bar{x} there are n linearly independent.

W.l.o.g., assume that the constraints active at \bar{x} are the first k constraints

$$a_i^T x \ge b_i, \ i = 1, \dots, k.$$

We should prove that among *n*-dimensional vectors $a_1, ..., a_k$, there are *n* linearly independent. Assuming otherwise, there exists a nonzero vector *h* such that $a_i^T h = 0$, i = 1, ..., k, that is,

$$a_i^T[\bar{x} \pm \epsilon h] = a_i^T \bar{x} = b_i, \ i = 1, ..., k$$

for all $\epsilon > 0$. Since the remaining constraints $a_i^T x \ge b_i$, i > k, are strictly satisfied at \bar{x} , we conclude that

$$a_i^T[\bar{x} \pm \epsilon h] \ge b_i, \ i = k + 1, ..., m$$

for all small enough values of $\epsilon > 0$. We conclude that $\overline{x} \pm \epsilon h \in Q = \{x : Ax \ge b\}$ for all small enough $\epsilon > 0$. Since $h \neq 0$ and \overline{x} is an extreme point of Q, we get a contradiction. Justification of the answer, \Leftarrow : If $\bar{x} \in Q$ makes equalities n of the constraints $a_i^T x \ge b_i$ with linearly independent vectors of coefficients, then $\bar{x} \in \text{Ext}(Q)$.

W.I.o.g., assume that n active at \bar{x} constraints with linearly independent vectors of coefficients are the first n constraints

$$a_i^T x \ge b_i, i = 1, ..., n$$

We should prove that if h is such that $\bar{x} \pm h \in Q$, then h = 0. Indeed, we have

 $\bar{x} \pm h \in Q \Rightarrow a_i^T [\bar{x} \pm h] \ge b_i, i = 1, ..., n;$ since $a_i^T \bar{x} = b_i$ for $i \le n$, we get

 $a_i^T \bar{x} \pm a_i^T h = a_i^T [\bar{x} \pm h] \ge a_i^T \bar{x}, \ i = 1, ..., n,$

whence

$$a_i^T h = 0, \ i = 1, ..., n.$$
 (*)

Since *n*-dimensional vectors $a_1, ..., a_n$ are linearly independent, (*) implies that h = 0, Q.E.D.

Convex Functions

Definition: Let f be a real-valued function defined on a nonempty subset Dom f in \mathbb{R}^n . f is called *convex*, if $\Diamond Dom f$ is a convex set \Diamond for all $x, y \in Dom f$ and $\lambda \in [0, 1]$ one has

$$f(\lambda x + (1 - \lambda)y) \le \lambda f(x) + (1 - \lambda)f(y)$$

Equivalent definition: Let f be a real-valued function defined on a nonempty subset Domfin \mathbb{R}^n . The function is called convex, if its *epigraph* – the set

$$\mathsf{Epi}{f} = {(x,t) \in \mathbb{R}^{n+1} : f(x) \le t}$$

is a convex set in \mathbb{R}^{n+1} .

What does the definition of convexity actually mean?

The inequality

 $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$ (*) where $x, y \in \text{Dom} f$ and $\lambda \in [0, 1]$ is automatically satisfied when x = y or when $\lambda = 0/1$. Thus, it says something only when the points x, y are distinct from each other and the point $z = \lambda x + (1 - \lambda)y$ is a (relative) interior point of the segment [x, y]. What does (*) say in this case?

Observe that $z = \lambda x + (1 - \lambda)y = x + (1 - \lambda)(y - x)$, whence

 $||y - x|| : ||y - z|| : ||z - x|| = 1 : \lambda : (1 - \lambda)$ Therefore

Similarly,

<u>Conclusion:</u> f is convex iff for every three distinct points x, y, z such that $x, y \in \text{Dom} f$ and $z \in [x, y]$, we have $z \in \text{Dom} f$ and

$$\frac{f(z) - f(x)}{\|z - x\|} \le \frac{f(y) - f(x)}{\|y - x\|} \le \frac{f(y) - f(z)}{\|y - z\|} \quad (*)$$

Note: From 3 inequalities in (*):

$$\frac{f(z) - f(x)}{\|z - x\|} \le \frac{f(y) - f(x)}{\|y - x\|}$$
$$\frac{f(y) - f(x)}{\|y - x\|} \le \frac{f(y) - f(z)}{\|y - z\|}$$
$$\frac{f(z) - f(x)}{\|z - x\|} \le \frac{f(y) - f(z)}{\|y - z\|}$$

every single one implies the other two.



<u>Jensen's Inequality:</u> Let f(x) be a convex function. Then

$$x_i \in \mathsf{Dom} f, \lambda_i \ge 0, \sum_i \lambda_i = 1 \Rightarrow f(\sum_i \lambda_i x_i) \le \sum_i \lambda_i f(x_i)$$

<u>Proof</u>: The points $(x_i, f(x_i))$ belong to Epi $\{f\}$. Since this set is convex, the point

$$(\sum_{i} \lambda_{i} x_{i}, \sum_{i} \lambda_{i} f(x_{i})) \in \mathsf{Epi}\{f\}.$$

By definition of the epigraph, it follows that

$$f(\sum_{i} \lambda_i x_i) \leq \sum_{i} \lambda_i f(x_i).$$

Extension: Let f be convex, Domf be closed and f be continuous on Domf. Consider a probability distribution $\pi(dx)$ supported on Domf. Then

 $f(\mathbf{E}_{\pi}\{x\}) \leq \mathbf{E}_{\pi}\{f(x)\}.$

Examples:

 \diamond Functions convex on \mathbb{R} : • x^2 , x^4 , x^6 ,...

• $\exp\{x\}$

Nonconvex functions on \mathbb{R} : • x^3 • sin(x)

 \diamond Functions convex on \mathbb{R}_+ : • x^p , $p \ge 1$

• $-x^p$, $0 \le p \le 1$ • $x \ln x$

 \diamond Functions convex on \mathbb{R}^n : • affine function $f(x) = f^T c$

A norm $\|\cdot\|$ on \mathbb{R}^n is a convex function:

 $\begin{aligned} \|\lambda x + (1-\lambda)y\| &\leq \|\lambda x\| + \|(1-\lambda)y\| \\ & \text{[Triangle inequality]} \\ &= \lambda \|x\| + (1-\lambda)\|y\| \\ & \text{[homogeneity]} \end{aligned}$

Application of Jensen's Inequality: Let $p = \{p_i > 0\}_{i=1}^n$, $q = \{q_i > 0\}_{i=1}^n$ be two discrete probability distributions.

<u>Claim:</u> The Kullback-Liebler distance

$$\sum_i p_i \ln \frac{p_i}{q_i}$$

between the distributions is ≥ 0 .

Indeed, the function $f(x) = -\ln x$, $\text{Dom}f = \{x > 0\}$, is convex. Setting $x_i = q_i/p_i$, $\lambda_i = p_i$ we have

$$0 = -\ln\left(\sum_{i} q_{i}\right) = f(\sum_{i} p_{i} x_{i})$$

$$\leq \sum_{i} p_{i} f(x_{i}) = \sum_{i} p_{i} (-\ln q_{i}/p_{i})$$

$$= \sum_{i} p_{i} \ln(p_{i}/q_{i})$$

What is the value of a convex function outside its domain?

<u>Convention.</u> To save words, it is convenient to think that a convex function f is defined *everywhere* on \mathbb{R}^n and takes real values *and value* $+\infty$. With this interpretation, f "remembers" its domain:

> $Dom f = \{x : f(x) \in \mathbb{R}\}\$ $x \notin Dom f \Rightarrow f(x) = +\infty$

and the definition of convexity becomes

$$\begin{split} f(\lambda x + (1-\lambda)y) &\leq \lambda f(x) + (1-\lambda)f(y) \; \forall \; \substack{x,y \in \mathbb{R}^n \\ \lambda \in [0,1]} \\ \text{where the arithmetics of } +\infty \text{ and reals is} \\ \text{given by the rules} \end{split}$$

$$a \in \mathbb{R} \Rightarrow a + (+\infty) = (+\infty) + (+\infty) = +\infty$$
$$0 \cdot (+\infty) = +\infty$$
$$\lambda > 0 \Rightarrow \lambda \cdot (+\infty) = +\infty$$

<u>Note</u>: Operations like $(+\infty) - (+\infty)$ or $(-5) \cdot (+\infty)$ are undefined!

Convexity-preserving operations:

Taking conic combinations: If $f_i(x)$ are convex function on \mathbb{R}^n and $\lambda_i \ge 0$, then the function $\sum_i \lambda_i f_i(x)$ is convex

 \diamond Affine substitution of argument: If f(x) is convex function on \mathbb{R}^n and x = Ay + b is an affine mapping from \mathbb{R}^k to \mathbb{R}^n , then the function g(y) = f(Ax + b) is convex on \mathbb{R}^m \diamond Taking supremum: If $f_{\alpha}(x)$, $\alpha \in \mathcal{A}$, is a family of convex function on \mathbb{R}^n , then the function $\sup_{\alpha \in \mathcal{A}} f_{\alpha}(x)$ is convex.

<u>Proof:</u> Epi{sup $f_{\alpha}(\cdot)$ } = \bigcap_{α} Epi{ $f_{\alpha}(\cdot)$ }, and intersections of convex sets are convex. \diamond Superposition Theorem: Let $f_i(x)$ be convex functions on \mathbb{R}^n , i = 1, ..., m, and $F(y_1, ..., y_m)$

be a convex and *monotone* function on \mathbb{R}^m . Then the function

 $g(x) = \begin{cases} F(f_1(x), ..., f_m(x)) & , x \in \text{Dom}f_i, \forall i \\ +\infty & , \text{otherwise} \end{cases}$

is convex.

 \diamond Partial minimization: Let f(x, y) be a convex function of $z = (x, y) \in \mathbb{R}^n$, and let

$$g(x) = \inf_{y} f(x, y)$$

be $> -\infty$ for all x. Then the function g(x) is convex.

Proof: g clearly takes real values and value $+\infty$. Let us check the Convexity Inequality

$$egin{aligned} g(\lambda x' + (1-\lambda)x'') &\leq \lambda g(x') + (1-\lambda)g(x'') \ & [\lambda \in [0,1]] \end{aligned}$$

There is nothing to check when $\lambda = 0$ or $\lambda = 1$, so let $0 < \lambda < 1$. In this case, there is nothing to check when g(x') or g(x'') is $+\infty$, so let $g(x') < +\infty$, $g(x'') < +\infty$. Since $g(x') < +\infty$, for every $\epsilon > 0$ there exists y' such that $f(x', y') \leq g(x') + \epsilon$. Similarly, there exists y'' such that $f(x'', y'') \leq g(x'') + \epsilon$. Now,

$$egin{aligned} &g(\lambda x'+(1-\lambda)x'')\ &\leq f(\lambda x'+(1-\lambda)x'',\lambda y'+(1-\lambda)y'')\ &\leq \lambda f(x',y')+(1-\lambda)f(x'',y'')\ &\leq \lambda (g(x')+\epsilon)+(1-\lambda)(g(x'')+\epsilon)\ &= \lambda g(x')+(1-\lambda)g(x'')+\epsilon \end{aligned}$$

Since $\epsilon > 0$ is arbitrary, we get

 $g(\lambda x' + (1-\lambda)x'') \leq \lambda g(x') + (1-\lambda)g(x'').$

How to detect convexity?

 $\frac{\text{Convexity is one-dimensional property: A set}}{X \subset \mathbb{R}^n \text{ is convex iff the set}}$

$$\{t : a + th \in X\}$$

is, for every (a, h), a convex set on the axis A function f on \mathbb{R}^n is convex iff the function

$$\phi(t) = f(a + th)$$

is, for every (a, h), a convex function on the axis.

A When a function ϕ on the axis is convex? Let ϕ be convex and finite on (a, b). This is exactly the same as

$$\frac{\phi(z) - \phi(x)}{z - x} \le \frac{\phi(y) - \phi(x)}{y - x} \le \frac{\phi(y) - \phi(z)}{y - z}$$

when a < x < z < y < b. Assuming that $\phi'(x)$ and $\phi'(y)$ exist and passing to limits as $z \to x + 0$ and $z \to y - 0$, we get

$$\phi'(x) \le \frac{\phi(y) - \phi(x)}{y - x} \le \phi'(y)$$

that is, $\phi'(x)$ is nondecreasing on the set of points from (a, b) where it exists.

The following conditions are necessary and sufficient for convexity of a univariate function:

 $\Delta = (a, b)$, possibly with added endpoint(s) (provided that the corresponding endpoint(s) is/are finite)

 $\Diamond \phi$ should be continuous on (a,b) and differentiable everywhere, except, perhaps, a countable set, and the derivative should be monotonically non-decreasing

 \diamond at endpoint of (a, b) which belongs to Dom ϕ , ϕ is allowed to "jump up", but not to jump down.
♣ Sufficient condition for convexity of a univariate function ϕ : Dom ϕ is convex, ϕ is continuous on Dom ϕ and is twice differentiable, with nonnegative ϕ'' , on intDom ϕ .

Indeed, we should prove that under the condition, if x < z < y are in Dom ϕ , then

$$\frac{\phi(z) - \phi(x)}{z - x} \le \frac{\phi(y) - \phi(z)}{y - z}$$

By Lagrange Theorem, the left ratio is $\phi'(\xi)$ for certain $\xi \in (x, z)$, and the right ratio is $\phi'(\eta)$ for certain $\eta \in (z, y)$. Since $\phi''(\cdot) \ge 0$ and $\eta > \xi$, we have $\phi'(\eta) \ge \phi'(\xi)$, Q.E.D. Sufficient condition for convexity of a multivariate function f: Domf is convex, f is continuous on Domf and is twice differentiable, with positive semidefinite Hessian matrix f'', on intDomf.

Instructive example: The function $f(x) = \prod_{i=1}^{n} \exp\{x_i\}$ is convex on \mathbb{R}^n . Indeed,

$$h^{T}f'(x) = \frac{\sum_{i} \exp\{x_{i}\}h_{i}}{\sum_{i} \exp\{x_{i}\}}$$
$$h^{T}f''(x)h = -\frac{\left(\sum_{i} \exp\{x_{i}\}h_{i}\right)^{2}}{\left(\sum_{i} \exp\{x_{i}\}\right)^{2}} + \frac{\sum_{i} \exp\{x_{i}\}h_{i}^{2}}{\sum_{i} \exp\{x_{i}\}}$$

$$h^T f''(x)h = -\left(\frac{\sum_{i} \exp\{x_i\}h_i}{\sum_{i} \exp\{x_i\}}\right)^2 + \frac{\sum_{i} \exp\{x_i\}h_i^2}{\sum_{i} \exp\{x_i\}}$$

Setting
$$p_i = \frac{\exp\{x_i\}}{\sum\limits_j \exp\{x_j\}}$$
, we have

$$h^{T}f''(x)h = \sum_{i} p_{i}h_{i}^{2} - \left(\sum_{i} p_{i}h_{i}\right)^{2}$$

$$= \sum_{i} p_{i}h_{i}^{2} - \left(\sum_{i} \sqrt{p_{i}}(\sqrt{p_{i}}h_{i})\right)^{2}$$

$$\geq \sum_{i} p_{i}h_{i}^{2} - \left(\sum_{i} (\sqrt{p_{i}})^{2}\right) \left(\sum_{i} (\sqrt{p_{i}}h_{i})^{2}\right)$$

$$= \sum_{i} p_{i}h_{i}^{2} - \left(\sum_{i} p_{i}h_{i}^{2}\right) = 0$$

(note that $\sum_i p_i = 1$)

<u>Corollary</u>: When $c_i > 0$, the function

$$g(y) = \ln\left(\sum_{i} c_{i} \exp\{a_{i}^{T}y\}\right)$$

is convex.

Indeed,

$$g(y) = \ln\left(\sum_{i} \exp\{\ln c_i + a_i^T y\}\right)$$

is obtained from the convex function

$$\ln\left(\sum_{i} \exp\{x_i\}\right)$$

by affine substitution of argument.

Gradient Inequality

Proposition: Let f be a function, x be an interior point of the domain of f and Q, $x \in Q$, be a convex set such that f is convex on Q. Assume that f is differentiable at x. Then

 $\forall y \in Q : f(y) \ge f(x) + (y - x)^T f'(x). \quad (*)$

Proof. Let $y \in Q$. There is nothing to prove when y = x or $f(y) = +\infty$, thus, assume that $f(y) < \infty$ and $y \neq x$. Let is set $z_{\epsilon} =$ $x + \epsilon(y - x)$, $0 < \epsilon < 1$. Then z_{ϵ} is an interior point of the segment [x, y]. Since f is convex, we have

$$\frac{f(y) - f(x)}{\|y - x\|} \ge \frac{f(z_{\epsilon}) - f(x)}{\|z_{\epsilon} - x\|} = \frac{f(x + \epsilon(y - x)) - f(x)}{\epsilon \|y - x\|}$$

Passing to limit as $\epsilon \to +0$, we arrive at

$$\frac{f(y) - f(x)}{\|y - x\|} \ge \frac{(y - x)^T f'(x)}{\|y - x\|},$$

as required by (*).

Lipschitz continuity of a convex function

<u>Proposition</u>: Let f be a convex function, and let K be a *closed* and *bounded* set belonging to relative interior of the domain of f. Then f is Lipschitz continuous on K, that is, there exists a constant $L < \infty$ such that

 $|f(x) - f(y)| \le L ||x - y||_2 \quad \forall x, y \in K.$

<u>Note</u>: All three assumptions on K are essential, as is shown by the following examples: $\langle f(x) = -\sqrt{x}, \text{ Dom} f = \{x \ge 0\}, K = [0, 1].$ Here $K \subset \text{Dom} f$ is closed and bounded, but is not contained in the relative interior of Dom f, and f is *not* Lipschitz continuous on K

 $\oint f(x) = x^2$, Dom $f = K = \mathbb{R}$. Here K is closed and belongs to rint Domf, but is unbounded, and f is *not* Lipschitz continuous on K

 $\oint f(x) = \frac{1}{x}$, Dom $f = \{x > 0\}$, K = (0, 1]. Here K is bounded and belongs to rint *Domf*, but is not closed, and f is *not* Lipschitz continuous on K Maxima and Minima of Convex Functions

(!) <u>Proposition</u> ["unimodality"] Let f be a convex function and x_* be a local minimizer of f:

$$x_* \in \mathsf{Dom} f$$

&
 $\exists r > 0 : f(x) \ge f(x_*) \ \forall (x : ||x - x_*|| \le r).$

Then x_* is a global minimizer of f:

 $f(x) \ge f(x_*) \ \forall x.$

Proof: All we need to prove is that if $x \neq x_*$ and $x \in \text{Dom} f$, then $f(x) \ge f(x_*)$. To this end let $z \in (x_*, x)$. By convexity we have

$$\frac{f(z) - f(x_*)}{\|z - x_*\|} \le \frac{f(x) - f(x_*)}{\|x - x_*\|}$$

When $z \in (x_*, x)$ is close enough to x_* , we have $\frac{f(z)-f(x_*)}{\|z-x_*\|} \ge 0$, whence $\frac{f(x)-f(x_*)}{\|x-x_*\|} \ge 0$, that is, $f(x) \ge f(x_*)$.

<u>Proposition</u> Let f be a convex function. The set of X_* of global minimizers is convex.

Proof: This is an immediate corollary of important

<u>Lemma:</u> Let f be a convex function. Then the level sets of f, that is, the sets

$$X_a = \{x : f(x) \le a\}$$

where a is a real, are convex.

Proof of Lemma: If $x, y \in X_a$ and $\lambda \in [0, 1]$, then

$$\begin{array}{rcl} f(\lambda x + (1-\lambda)y) &\leq & \lambda f(x) + (1-\lambda)f(y) \\ &\leq & \lambda a + (1-\lambda)a = a. \end{array}$$

Thus, $[x, y] \subset X_a$.

When the minimizer of a convex function is unique?

<u>Definition:</u> A convex function is called *strictly convex*, if

 $f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$

whenever $x \neq y$ and $\lambda \in (0, 1)$.

<u>Note:</u> If a convex function f has open domain and is twice continuously differentiable on this domain with

 $h^T f''(x)h > 0 \quad \forall (x \in \mathsf{Dom} f, h \neq 0),$

then f is strictly convex.

<u>Proposition:</u> For a strictly convex function f a minimizer, if it exists, is unique.

Proof. Assume that $X_* = \operatorname{Argmin} f$ contains two distinct points x', x''. By strong convexity,

$$f(\frac{1}{2}x' + \frac{1}{2}x'') < \frac{1}{2}\left[f(x') + f(x'')\right] = \inf_{x} f_{x},$$

which is impossible.

<u>Theorem</u> [Optimality conditions in convex minimization] Let f be a function which is differentiable at a point x_* and is convex on a convex set $Q \subset \text{Dom} f$ which contains x_* . A necessary and sufficient condition for f to attain its minimum on Q at x_* is

$$(x-x_*)^T f'(x_*) \ge 0 \quad \forall x \in Q.$$
 (*)

Proof, \Leftarrow : Assume that (*) is valid, and let us verify that $f(x) \ge f(x_*)$ for every $x \in Q$. There is nothing to prove when $x = x_*$, thus, let $f(x) < \infty$ and $x \ne x_*$. For $z_{\lambda} = x_* + \lambda(x - x_*)$ we have

$$\frac{f(z_{\lambda})-f(x_*)}{\|z_{\lambda}-x_*\|} \leq \frac{f(x)-f(x_*)}{\|x-x_*\|} \quad \forall \lambda \in (0,1)$$

or, which is the same,

$$\frac{f(x_* + \lambda[x - x_*]) - f(x_*)}{\lambda \|x - x_*\|} \le \frac{f(x) - f(x_*)}{\|x - x_*\|} \forall \lambda \in (0, 1)$$

As $\lambda \to +0$, the left ratio converges to $(x-x_*)^T f'(x_*)/||x-x_*|| \ge 0$; thus, $\frac{f(x)-f(x_*)}{||x-x_*||} \ge 0$, whence $f(x) \ge f(x_*)$.

"Let f be a function which is differentiable at a point x_* and is convex on a convex set $Q \subset \text{Dom} f$ which contains x_* . A necessary and sufficient condition for f to attain its minimum on Q at x_* is

$$(x-x_*)^T f'(x_*) \ge 0 \quad \forall x \in Q.''$$

Proof, \Rightarrow : Given that $x_* \in \operatorname{Argmin}_{y \in Q} f(y)$, let $x \in Q$. Then

$$0 \leq rac{f(x_*+\lambda[x-x_*])-f(x_*)}{\lambda} \quad orall \lambda \in (0,1),$$
 whence $(x-x_*)^T f'(x_*) \geq 0.$

♣ Equivalent reformulation: Let f be a function which is differentiable at a point x_* and is convex on a convex set $Q \subset \text{Dom}f$, $x_* \in Q$. Consider the radial cone of Q at x_* :

 $T_Q(x_*) = \{h : \exists t > 0 : x_* + th \in Q\}$

<u>Note</u>: $T_Q(x_*)$ is indeed a cone which is comprised of all vectors of the form $s(x - x_*)$, where $x \in Q$ and $s \ge 0$.

f attains its minimum on Q at x_* iff

 $h^T f'(x_*) \ge 0 \ \forall h \in T_Q(x_*),$

or, which is the same, iff

 $f'(x_*) \in \underbrace{N_Q(x_*) = \{g : g^T h \ge 0 \forall h \in T_Q(x_*)\}}_{normal \ cone \ of \ Q \ at \ x_*}$ (*)

Example I: $x_* \in intQ$. Here $T_Q(x_*) = \mathbb{R}^n$, whence $N_Q(x_*) = \{0\}$, and (*) becomes the Fermat equation

 $f'(x_*) = 0$

Example II: $x_* \in \operatorname{rint} Q$. Let $\operatorname{Aff}(Q) = x_* + L$, where L is a linear subspace in \mathbb{R}^n . Here $T_Q(x_*) = L$, whence $N_Q(x_*) = L^{\perp}$. (*) becomes the condition

 $f'(x_*)$ is orthogonal to L.

Equivalently: Let Aff(Q) = $\{x : Ax = b\}$. Then $L = \{x : Ax = 0\}, L^{\perp} = \{y = A^T\lambda\},\$ and the optimality condition becomes

 $\frac{\text{Example III: } Q = \{x : Ax - b \le 0\} \text{ is polyhedral.}}{\text{In this case,}}$

 $= \left\{ \begin{aligned} T_Q(x_*) \\ h : a_i^T h &\leq 0 \ \forall i \in I(x_*) = \{i : a_i^T x_* - b_i = 0\} \right\}. \end{aligned}$

By Homogeneous Farkas Lemma,

$$N_Q(x_*) \equiv \{y : a_i^T h \le 0, i \in I(x_*) \Rightarrow y^T h \ge 0\}$$
$$= \{y = -\sum_{i \in I(x_*)} \lambda_i a_i : \lambda_i \ge 0\}$$

and the optimality condition becomes

$$\exists (\lambda_i^* \ge 0, i \in I(x_*)) : f'(x_*) + \sum_{i \in I(x_*)} \lambda_i^* a_i = 0$$

or, which is the same:

$$\exists \lambda^* \ge 0 : \begin{cases} f'(x_*) + \sum_{i=1}^m \lambda_i^* a_i = 0\\ \lambda_i^* (a_i^T x_* - b_i) = 0, i = 1, ..., m \end{cases}$$

The point is that in the *convex* case these conditions are necessary *and sufficient* for x_* to be a minimizer of f on Q. Example: Let us solve the problem

$$\min_{x} \left\{ c^{T} x + \sum_{i=1}^{m} x_{i} \ln x_{i} : x \ge 0, \sum_{i} x_{i} = 1 \right\}.$$

The objective is convex, the domain $Q = \{x \ge 0, \sum_{i} x_i = 1\}$ is convex (and even polyhedral). Assuming that the minimum is achieved at a point $x_* \in \operatorname{rint} Q$, the optimality condition becomes

$$x_i = \frac{\exp\{-c_i\}}{\sum\limits_j \exp\{-c_j\}}.$$

At this point, the optimality condition is satisfied, so that the point indeed is a minimizer.

Maxima of convex functions

Proposition. Let f be a convex function. Then

f f attains its maximum over Domf at a point $x^* \in \text{rint Dom} f$, then f is constant on Domf

fightharpoinds fightharpoinds for the state of the stat

fightharpoondown fightharpoondown for a polyhedral and f is bounded from above on Dom f, then f attains its maximum on Dom f.

Subgradients of convex functions

Let f be a convex function and $\bar{x} \in$ intDom f. If f differentiable at \bar{x} , then, by Gradient Inequality, there exists an affine function, specifically,

$$h(x) = f(\bar{x}) + (x - \bar{x})^T f'(\bar{x}),$$

such that

$$f(x) \ge h(x) \forall x \& f(\bar{x}) = h(\bar{x})$$
 (*)

Affine function with property (*) may exist also in the case when f is *not* differentiable at $\bar{x} \in \text{Dom} f$. (*) implies that

$$h(x) = f(\bar{x}) + (x - \bar{x})^T g$$
 (**)

for certain g. Function (**) indeed satisfies (*) if and only if g is such that

$$f(x) \ge f(\bar{x}) + (x - \bar{x})^T g \quad \forall x$$
 (!)

<u>Definition.</u> Let f be a convex function and $\overline{x} \in \text{Dom} f$. Every vector g satisfying

$$f(x) \ge f(\bar{x}) + (x - \bar{x})^T g \quad \forall x \qquad (!)$$

is called a *subgradient* of f at \bar{x} . The set of all subgradients, if any, of f at \bar{x} is called *subdifferential* $\partial f(\bar{x})$ of f at \bar{x} .

<u>Example I</u>: By Gradient Inequality, if convex function f is differentiable at \bar{x} , then $\nabla f(\bar{x}) \in \partial f(\bar{x})$. If, in addition, $\bar{x} \in \text{intDom} f$, then $\nabla f(\bar{x})$ is the *unique* element of $\partial f(\bar{x})$.

Example II: Let f(x) = |x| ($x \in \mathbb{R}$). When $\overline{x} \neq 0$, f is differentiable at \overline{x} , whence $\partial f(\overline{x}) = f'(\overline{x})$. When $\overline{x} = 0$, subgradients g are given by

$$|x| \ge \mathbf{0} + gx = gx \ \forall x,$$

that is, $\partial f(0) = [-1, 1]$.

<u>Note:</u> In the case in question, f has directional derivative

$$Df(x)[h] = \lim_{t \to +0} \frac{f(x+th) - f(x)}{t}$$

at every point $x \in \mathbb{R}$ along every direction $h \in \mathbb{R}$, and this derivative is nothing but

$$Df(x)[h] = \max_{g \in \partial f(x)} g^T h$$

Proposition: Let f be convex. Then \diamond For every $x \in \text{Dom}f$, the subdifferential $\partial f(x)$ is closed convex set \diamond If $x \in \text{rint Dom}f$, then $\partial f(x)$ is nonempty. \diamond If $x \in \text{rint Dom}f$, then, for every $h \in \mathbb{R}^n$,

$$\exists Df(x)[h] \equiv \lim_{t \to +0} \frac{f(x+th) - f(x)}{t} = \max_{g \in \partial f(x)} g^T h.$$

 $Assume that \ \bar{x} \in Dom f$ is represented as $\lim_{i \to \infty} x_i$ with $x_i \in Dom f$ and that

$$f(\bar{x}) \leq \lim \inf_{i \to \infty} f(x_i)$$

If a sequence $g_i \in \partial f(x_i)$ converges to certain vector g, then $g \in \partial f(\bar{x})$.

 \Diamond The multi-valued mapping $x \mapsto \partial f(x)$ is locally bounded at every point $\overline{x} \in \text{intDom} f$, that is, whenever $\overline{x} \in \text{intDom} f$, there exist r > 0 and $R < \infty$ such that

$$||x - \bar{x}||_2 \le r, g \in \partial f(x) \Rightarrow ||g||_2 \le R.$$

Selected proof: "If $\bar{x} \in \text{rint Dom} f$, then $\partial f(\bar{x})$ is nonempty."

W.l.o.g. let Domf be full-dimensional, so that $\bar{x} \in intDomf$. Consider the convex set

 $T = \mathsf{Epi}\{f\} = \{(x, t) : t \ge f(x)\}.$

Since f is convex, it is continuous on intDomf, whence T has a nonempty interior. The point $(\bar{x}, f(\bar{x}))$ clearly does not belong to this interior, whence $S = \{(\bar{x}, f(\bar{x}))\}$ can be separated from T: there exists $(\alpha, \beta) \neq 0$ such that

 $\alpha^T \bar{x} + \beta f(\bar{x}) \le \alpha^T x + \beta t \quad \forall (x, t \ge f(x)) \quad (*)$

Clearly $\beta \ge 0$ (otherwise (*) will be impossible when $x = \overline{x}$ and $t > f(\overline{x})$ is large).

<u>Claim:</u> $\beta > 0$. Indeed, with $\beta = 0$, (*) implies

$$\alpha^T \bar{x} \le \alpha^T x \; \forall x \in \mathsf{Dom} f \tag{**}$$

Since $(\alpha, \beta) \neq 0$ and $\beta = 0$, we have $\alpha \neq 0$; but then (**) contradicts $\overline{x} \in \text{intDom} f$. $\diamond \text{Since } \beta > 0$, (*) implies that if $g = \beta^{-1} \alpha$, then

 $g^T\bar{x} + f(\bar{x}) \leq g^Tx + f(x) \; \forall x \in \mathrm{Dom} f,$ that is,

$$f(x) \ge f(\bar{x}) + (x - \bar{x})^T g \ \forall x.$$

Elementary Calculus of Subgradients

$$\begin{aligned} & \diamondsuit \text{If } g_i \in \partial f_i(x) \text{ and } \lambda_i \geq 0, \text{ then} \\ & \sum_i \lambda_i g_i \in \partial (\sum_i \lambda_i f_i)(x) \\ & \diamondsuit \text{If } g_\alpha \in \partial f_\alpha(x), \ \alpha \in \mathcal{A}, \\ & f(\cdot) = \sup_{\alpha \in \mathcal{A}} f_\alpha(\cdot) \end{aligned}$$

and

$$f(x) = f_{\alpha}(x), \ \alpha \in \mathcal{A}_*(x) \neq \emptyset,$$

then every convex combination of vectors g_{α} , $\alpha \in \mathcal{A}_*(x)$, is a subgradient of f at x $\langle If g_i \in \text{dom} f_i(x), i = 1, ..., m$, and $F(y_1, ..., y_m)$ is convex and monotone and $0 \leq d \in \partial F(f_1(x), ..., f_m(x))$, then the vector

$$\sum_i d_i g_i$$

is a subgradient of $F(f_1(\cdot), ..., f_m(\cdot))$ at x.

Convex Programming Lagrange Duality Saddle Points

Mathematical Programming program is

 $f_* = \min_{x} \left\{ \begin{array}{ll} g(x) \equiv (g_1(x), ..., g_m(x))^T &\leq 0\\ f(x): \ h(x) = (h_1(x), ..., h_k(x))^T &= 0\\ & x \in X \\ (P) \end{array} \right\}$

 $\diamondsuit x$ is the *design vector*. Values of x are called *solutions* to (P)

 $\Diamond f(x)$ is the objective

 $\displaystyle \diamondsuit g(x) \equiv (g_1(x), ..., g_m(x))^T \leq 0 - inequality constraints$

 $\diamondsuit h(x) = (h_1(x), ..., h_k(x))^T = 0 - equality con$ straints

 $\Diamond X \subset \mathbb{R}^n$ – domain. We always assume that the objective and the constraints are well-defined on X.

$$f_{*} = \min_{x} \left\{ \begin{array}{cc} g(x) \equiv (g_{1}(x), ..., g_{m}(x))^{T} \leq 0\\ f(x) : h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0\\ x \in X \end{array} \right\}$$

Solution x is called *feasible*, if it satisfies all the constraints. Problem which has feasible solutions is called *feasible*.

♣ If the objective is (below) bounded on the set of feasible solutions, (P) is called *bounded*.

\clubsuit The optimal value f_* is

 $f_* = \begin{cases} \inf_x \{f(x) : x \text{ is feasible}\}, & (P) \text{ is feasible} \\ +\infty, & \text{otherwise} \end{cases}$

 f_* is a real for feasible and bounded problem, is $-\infty$ for feasible unbounded problem, and is $+\infty$ for infeasible problem.

♣ Optimal solution of (P) is a feasible solution x_* such that $f(x_*) = f_*$. Problem which has optimal solutions is called *solvable*.

$$f_{*} = \min_{x} \left\{ \begin{array}{cc} g(x) \equiv (g_{1}(x), ..., g_{m}(x))^{T} \leq 0\\ f(x) : h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0\\ x \in X\\ (P) \end{array} \right\}$$

♣ Problem (P) is called *convex*, if $\Diamond X$ is a convex subset of \mathbb{R}^n $\Diamond f(\cdot), g_1(\cdot), \dots, g_m(\cdot)$ are *convex real-valued* functions on X

♦ There are no equality constraints [we could allow *linear* equality constraints, but this does not add generality]

Preparing tools for Lagrange Duality: Convex Theorem on Alternative

Question: How to certify insolvability of the system

$$\begin{array}{rcl}
f(x) &< c \\
g_j(x) &\leq 0, \, j = 1, ..., m \\
x &\in X
\end{array} (I)$$

Answer: Assume that there exist <u>nonnegative</u> weights λ_j , j = 1, ..., m, such that the inequality

$$f(x) + \sum_{j=1}^{m} \lambda_j g_j(x) < c$$

has no solutions in X:

$$\exists \lambda_j \ge 0: \quad \inf_{x \in X} [f(x) + \sum_{j=1}^m \lambda_j g_j(x)] \ge c.$$

Then (I) is insolvable.

A Convex Theorem on Alternative: Consider a system of constraints on x

$$\begin{array}{rcl}
f(x) &< c \\
g_j(x) &\leq 0, \, j = 1, ..., m \\
x &\in X
\end{array} (I)$$

along with system of constraints on λ :

$$\inf_{x \in X} [f(x) + \sum_{j=1}^{m} \lambda_j g_j(x)] \geq c$$
$$\lambda_j \geq 0, j = 1, ..., m$$
(II)

 $\langle [Trivial part] If (II) is solvable, then (I) is insolvable$

O[Nontrivial part] If (I) is insolvable and system (I) is convex:

- X is convex set

— $f,\ g_1,...,g_m$ are real-valued convex functions on X

and the subsystem

$$g_j(x) < 0, j = 1, ..., m, \ x \in X$$

is solvable [Slater condition], then (II) is solvable.

$$\begin{array}{rcl}
f(x) &< c \\
g_j(x) &\leq 0, \, j = 1, ..., m \\
x &\in X
\end{array} (I)$$

Proof of Nontrivial part: Assume that (*I*) has no solutions. Consider two sets in \mathbb{R}^{m+1} :

$$\underbrace{\begin{cases}
 T \\
 f(x) \leq u_0 \\
 g_1(x) \leq u_1 \\
 \\
 g_m(x) \leq u_m
\end{cases}}_{g_m(x) \leq u_m}$$

$$\underbrace{\left\{u \in \mathbb{R}^{m+1} : u_0 < c, u_1 \leq 0, ..., u_m \leq 0\right\}}_{S}$$

<u>Observations:</u> $\diamondsuit S$, T are convex and nonempty $\diamondsuit S$, T do not intersect (otherwise (I) would have a solution)

<u>Conclusion:</u> S and T can be separated:

$$\exists (a_0, ..., a_m) \neq 0 : \inf_{u \in T} a^T u \ge \sup_{u \in S} a^T u$$

$$\underbrace{\begin{cases}
 T \\
 f(x) \leq u_0 \\
 g_1(x) \leq u_1 \\
 \\
 g_m(x) \leq u_m
\end{cases}}_{g_m(x) \leq u_m}$$

$$\underbrace{\left\{u \in \mathbb{R}^{m+1} : u_0 < c, u_1 \leq 0, ..., u_m \leq 0\right\}}_{S}$$

$$\exists (a_0, ..., a_m) \neq 0 : \\ \inf_{\substack{x \in X \\ u_0 \geq f(x) \\ u_1 \geq g_1(x) \\ \vdots \\ u_m \geq g_m(x) \\ ext{index} \\ sup_{index} [a_0u_0 + a_1u_1 + ... + a_mu_m] \\ isometry_{index} \\ u_0 < c \\ u_1 \leq 0 \\ \vdots \\ u_m \leq 0 \end{cases}$$

<u>Conclusion:</u> $a \ge 0$, whence

 $\inf_{x \in X} [a_0 f(x) + a_1 g_1(x) + \dots + a_m g_m(x)] \ge a_0 c.$

Summary:

 $\exists a \geq 0, a \neq 0$: $\inf_{x \in X} [a_0 f(x) + a_1 g_1(x) + \dots + a_m g_m(x)] \ge a_0 c$ <u>Observation</u>: $a_0 > 0$. Indeed, otherwise $0 \neq (a_1, ..., a_m) \geq 0$ and $\inf_{x \in X} [a_1 g_1(x) + \dots + a_m g_m(x)] \ge 0,$ while $\exists \bar{x} \in X : g_j(\bar{x}) < 0$ for all j. <u>Conclusion</u>: $a_0 > 0$, whence 2.

$$\inf_{x \in X} \left[f(x) + \sum_{j=1}^{m} \underbrace{\left[\frac{a_j}{a_0} \right]}_{\lambda_j \ge 0} g_j(x) \right] \ge c$$

Lagrange Function

Consider optimization program
Opt(P) = min { $f(x) : g_j(x) \le 0, j \le m, x \in X$ (P)
and associate with it Lagrange function $L(x, \lambda) = f(x) + \sum_{j=1}^m \lambda_j g_j(x)$ along with the Lagrange Dual problem
Opt(D) = max $\underline{L}(\lambda), \ \underline{L}(\lambda) = \inf_{x \in X} L(x, \lambda)$ (D)

♣ Convex Programming Duality Theorem: ◊[Weak Duality] For every λ ≥ 0, L(λ) ≤ Opt(P). In particular,

 $Opt(D) \leq Opt(P)$

 \Diamond [Strong Duality] If (P) is convex and below bounded and satisfies Slater condition, then (D) is solvable, and

Opt(D) = Opt(P).

Weak Duality: " $Opt(D) \leq Opt(P)$ ": There is nothing to prove when (P) is infeasible, that is, when $Opt(P) = \infty$. If x is feasible for (P) and $\lambda \geq 0$, then $L(x, \lambda) \leq f(x)$, whence

$$\begin{array}{rcl} \lambda \geq 0 \Rightarrow \underline{L}(\lambda) &\equiv & \inf_{x \in X} L(x,\lambda) \\ &\leq & \inf_{x \in X \text{ is feasible}} L(x,\lambda) \\ &\leq & \inf_{x \in X \text{ is feasible}} f(x) \\ &= & \operatorname{Opt}(P) \\ &\Rightarrow \operatorname{Opt}(D) &= & \sup_{\lambda \geq 0} \underline{L}(\lambda) \leq \operatorname{Opt}(D). \end{array}$$

Strong Duality: "If (P) is convex and below bounded and satisfies Slater condition, then (D) is solvable and Opt(D) = Opt(P)": The system

 $f(x) < Opt(P), g_j(x) \le 0, j = 1, ..., m, x \in X$ has no solutions, while the system

$$g_j(x) < 0, j = 1, ..., m, x \in X$$

has a solution. By CTA,

$$\exists \lambda^* \ge 0 : f(x) + \sum_j \lambda_j^* g_j(x) \ge \operatorname{Opt}(P) \ \forall x \in X,$$

whence

$$\underline{L}(\lambda^*) \ge \operatorname{Opt}(P). \tag{*}$$

Combined with Weak Duality, (*) says that

$$Opt(D) = \underline{L}(\lambda^*) = Opt(P).$$

<u>Note:</u> The Lagrange function "remembers", up to equivalence, both (P) and (D). Indeed,

$$Opt(D) = \sup_{\lambda \ge 0} \inf_{x \in X} L(x, \lambda)$$

is given by the Lagrange function. Now consider the function

$$\overline{L}(x) = \sup_{\lambda \ge 0} L(x, \lambda) = \begin{cases} f(x), & g_j(x) \le 0, \ j \le m \\ +\infty, & \text{otherwise} \end{cases}$$

(P) clearly is equivalent to the problem of minimizing $\overline{L}(x)$ over $x \in X$:

$$Opt(P) = \inf_{x \in X} \sup_{\lambda \ge 0} L(x, \lambda)$$

Saddle Points

Let $X \subset \mathbb{R}^n$, $\Lambda \subset \mathbb{R}^m$ be nonempty sets, and let $F(x, \lambda)$ be a real-valued function on $X \times \Lambda$. This function gives rise to two optimization problems

$$Opt(P) = \inf_{\substack{x \in X \\ \lambda \in \Lambda}} \underbrace{\sup_{\lambda \in \Lambda} \overline{F(x,\lambda)}}_{K \in X} (P)$$
$$Opt(D) = \sup_{\substack{x \in X \\ \lambda \in \Lambda}} \inf_{\substack{x \in X \\ \underline{F}(\lambda)}} F(x,\lambda) (D)$$

$$Opt(P) = \inf_{x \in X} \underbrace{\sup_{\lambda \in \Lambda} \overline{F(x,\lambda)}}_{X \in X} (P)$$
$$Opt(D) = \sup_{\lambda \in \Lambda} \inf_{x \in X} F(x,\lambda) (D)$$
$$\underbrace{\underline{F(\lambda)}}_{\underline{F(\lambda)}} (D)$$

Game interpretation: Player I chooses $x \in X$, player II chooses $\lambda \in \Lambda$. With choices of the players x, λ , player I pays to player II the sum of $F(x, \lambda)$. What should the players do to optimize their wealth?

 \Diamond If Player I chooses x first, and Player II knows this choice when choosing λ , II will maximize his profit, and the loss of I will be $\overline{F}(x)$. To minimize his loss, I should solve (P), thus ensuring himself loss Opt(P) or less.

 \Diamond If Player II chooses λ first, and Player I knows this choice when choosing x, I will minimize his loss, and the profit of II will be $\underline{F}(\lambda)$. To maximize his profit, II should solve (D), thus ensuring himself profit Opt(D) or more.


<u>Observation</u>: For Player I, second situation seems better, so that it is natural to guess that his anticipated loss in this situation is \leq his anticipated loss in the first situation:

 $Opt(D) \equiv \sup_{\lambda \in \Lambda} \inf_{x \in X} F(x, \lambda) \leq \inf_{x \in X} \sup_{\lambda \in \Lambda} F(x, \lambda) \equiv Opt(P).$

This indeed is true: assuming $Opt(P) < \infty$ (otherwise the inequality is evident),

$$\begin{aligned} \forall (\epsilon > 0) : \quad \exists x_{\epsilon} \in X : \sup_{\lambda \in \Lambda} F(x_{\epsilon}, \lambda) \leq \operatorname{Opt}(P) + \epsilon \\ \Rightarrow \forall \lambda \in \Lambda : \underline{F}(\lambda) = \inf_{x \in X} F(x, \lambda) \leq F(x_{\epsilon}, \lambda) \leq \operatorname{Opt}(P) + \epsilon \\ \Rightarrow \operatorname{Opt}(D) \equiv \sup_{\lambda \in \Lambda} \underline{F}(\lambda) \leq \operatorname{Opt}(P) + \epsilon \\ \Rightarrow \operatorname{Opt}(D) \leq \operatorname{Opt}(P). \end{aligned}$$

$$Opt(P) = \inf_{x \in X} \underbrace{\sup_{\lambda \in \Lambda} \overline{F(x,\lambda)}}_{X \in \Lambda} (P)$$
$$Opt(D) = \sup_{\lambda \in \Lambda} \underbrace{\inf_{x \in X} F(x,\lambda)}_{\underline{F(\lambda)}} (D)$$

What should the players do when making their choices simultaneously?

<u>A "good case"</u> when we can answer this question -F has a saddle point.

<u>Definition</u>: We call a point $(x_*, \lambda_*) \in X \times \Lambda$ a saddle point of F, if

 $F(x, \lambda_*) \ge F(x_*, \lambda_*) \ge F(x_*, \lambda) \ \forall (x \in X, \lambda \in \Lambda).$ In game terms, a saddle point is an *equilibrium* – no one of the players can improve his wealth, provided the adversary keeps his choice unchanged.

<u>Proposition</u>: F has a saddle point if and only if both (P) and (D) are solvable with equal optimal values. In this case, the saddle points of F are exactly the pairs (x_*, λ_*) , where x_* is an optimal solution to (P), and λ_* is an optimal solution to (D).

$$Opt(P) = \inf_{x \in X} \underbrace{\sup_{\lambda \in \Lambda} \overline{F(x,\lambda)}}_{X \in \Lambda} (P)$$
$$Opt(D) = \sup_{\lambda \in \Lambda} \underbrace{\inf_{x \in X} F(x,\lambda)}_{F(\lambda)} (D)$$

Proof, \Rightarrow : Assume that (x_*, λ_*) is a saddle point of F, and let us prove that x_* solves (P), λ_* solves (D), and Opt(P) = Opt(D). Indeed, we have

$$F(x, \lambda_*) \ge F(x_*, \lambda_*) \ge F(x_*, \lambda) \ \forall (x \in X, \lambda \in \Lambda)$$

whence

$$Opt(P) \leq \overline{F}(x_*) = \sup_{\lambda \in \Lambda} F(x_*, \lambda) = F(x_*, \lambda_*)$$
$$Opt(D) \geq \underline{F}(\lambda_*) = \inf_{x \in X} F(x, \lambda_*) = F(x_*, \lambda_*)$$

Since $Opt(P) \ge Opt(D)$, we see that all inequalities in the chain

 $Opt(P) \leq \overline{F}(x_*) = F(x_*, \lambda_*) = \underline{F}(\lambda_*) \leq Opt(D)$ are equalities. Thus, x_* solves (P), λ_* solves (D) and Opt(P) = Opt(D).

$$Opt(P) = \inf_{x \in X} \underbrace{\sup_{\lambda \in \Lambda} \overline{F(x, \lambda)}}_{\lambda \in \Lambda} (P)$$
$$Opt(D) = \sup_{\lambda \in \Lambda} \underbrace{\inf_{x \in X} F(x, \lambda)}_{\underline{F(\lambda)}} (D)$$

Proof, \Leftarrow . Assume that (P), (D) have optimal solutions x_*, λ_* and Opt(P) = Opt(D), and let us prove that (x_*, λ_*) is a saddle point. We have

$$Opt(P) = \overline{F}(x_*) = \sup_{\lambda \in \Lambda} F(x_*, \lambda) \ge F(x_*, \lambda_*)$$
$$Opt(D) = \underline{F}(\lambda_*) = \inf_{x \in X} F(x, \lambda_*) \le F(x_*, \lambda_*)$$
(*)

Since Opt(P) = Opt(D), all inequalities in (*) are equalities, so that

$$\sup_{\lambda \in \Lambda} F(x_*, \lambda) = F(x_*\lambda_*) = \inf_{x \in X} F(x, \lambda_*).$$

<u>Theorem</u> [Saddle Point form of Optimality Conditions in Convex Programming] Let $x_* \in X$.

 $[Sufficient optimality condition] If x_* can be extended, by a <math>\lambda^* \ge 0$, to a saddle point of the Lagrange function on $X \times \{\lambda \ge 0\}$:

 $L(x, \lambda^*) \ge L(x_*, \lambda^*) \ge L(x_*, \lambda) \ \forall (x \in X, \lambda \ge 0),$ <u>then</u> x_* is optimal for (P).

O(Necessary optimality condition) If x_* is optimal for (P) and (P) is convex and satisfies the Slater condition, then x_* can be extended, by a $\lambda^* \ge 0$, to a saddle point of the Lagrange function on $X \times \{\lambda \ge 0\}$.

Proof, \Rightarrow : "Assume $x_* \in X$ and $\exists \lambda^* \geq 0$: $L(x, \lambda^*) \geq L(x_*, \lambda^*) \geq L(x_*, \lambda) \ \forall (x \in X, \lambda \geq 0).$ <u>Then</u> x_* is optimal for (P)." Clearly, $\sup_{\lambda \geq 0} L(x_*, \lambda) = \begin{cases} +\infty, & x_* \text{ is infeasible} \\ f(x_*), & \text{otherwise} \end{cases}$ Thus, $\lambda^* \geq 0 \& L(x_*, \lambda^*) \geq L(x_*, \lambda) \ \forall \lambda \geq 0$ is equivalent to

 $g_j(x_*) \leq 0 \forall j \& \lambda_j^* g_j(x_*) = 0 \forall j.$ Consequently, $L(x_*, \lambda^*) = f(x_*)$, whence

$$L(x,\lambda^*) \ge L(x_*,\lambda^*) \ \forall x \in X$$

reads as

$$L(x,\lambda^*) \ge f(x_*) \ \forall x. \tag{(*)}$$

Since for $\lambda \ge 0$ one has $f(x) \ge L(x, \lambda)$ for all feasible x, (*) implies that

x is feasible
$$\Rightarrow f(x) \ge f(x_*)$$
.

Proof, \Leftarrow : <u>Assume</u> x_* is optimal for convex problem (*P*) satisfying the Slater condition. Then $\exists \lambda^* \ge 0$:

 $L(x,\lambda^*) \ge L(x_*,\lambda^*) \ge L(x_*,\lambda) \ \forall (x \in X,\lambda \ge 0).$

By Lagrange Duality Theorem, $\exists \lambda^* \geq 0$:

$$f(x_*) = \underline{L}(\lambda^*) \equiv \inf_{x \in X} \left[f(x) + \sum_j \lambda_j^* g_j(x) \right]. \quad (*)$$

Since x_* is feasible, we have

$$\inf_{x \in X} \left[f(x) + \sum_{j} \lambda_j^* g_j(x) \right] \leq f(x_*) + \sum_{j} \lambda_j^* g_j(x_*) \leq f(x_*).$$

By (*), the last " \geq " here is " = ", which with $\lambda^* \geq 0$ is possible iff $\lambda_j^* g_j(x_*) = 0 \forall j$

$$\Rightarrow f(x_*) = L(x_*, \lambda^*) \ge L(x_*, \lambda) \ \forall \lambda \ge 0.$$

Now (*) reads $L(x, \lambda^*) \ge f(x_*) = L(x_*, \lambda^*)$.

<u>Theorem</u> [Karush-Kuhn-Tucker Optimality Conditions in Convex Programming] Let (*P*) be a convex program, let x^* be its feasible solution, and let the functions f, g_1, \ldots, g_m be differentiable at x^* . Then \diamond The Karush-Kuhn-Tucker condition: Exist Lagrange multipliers $\lambda^* \ge 0$ such that

$$abla f(x_*) + \sum_{j=1}^m \lambda_j^* \nabla g_j(x_*) \in N_X^*(x_*)$$

 $\lambda_j^* g_j(x_*) = 0, \ j \le m$ [complementary slackness]

is <u>sufficient</u> for x_* to be optimal. \Diamond If (P) satisfies restricted Slater condition: $\exists \overline{x} \in \text{rint } X : g_j(\overline{x}) \leq 0$ for all constraints and $g_j(\overline{x}) < 0$ for all <u>nonlinear</u> constraints, then the KKT is <u>necessary and sufficient</u> for x_* to be optimal.

Proof, \Rightarrow : Let (*P*) be convex, x_* be feasible, and *f*, g_j be differentiable at x_* . Assume also that the KKT holds:

Exist Lagrange multipliers $\lambda^* \geq 0$ such that

(a)
$$\nabla f(x_*) + \sum_{j=1}^m \lambda_j^* \nabla g_j(x_*) \in N_X^*(x_*)$$

(b)
$$\lambda_j^* g_j(x_*) = 0, j \le m$$
 [complementary slackness]

Then x_* is optimal.

Indeed, complementary slackness plus $\lambda^* \geq 0$ ensure that

$$L(x_*,\lambda^*)\geq L(x_*,\lambda) \quad orall \lambda\geq 0.$$

Further, $L(x, \lambda^*)$ is convex in $x \in X$ and differentiable at $x_* \in X$, so that (a) implies that

$$L(x, \lambda^*) \ge L(x_*, \lambda^*) \quad \forall x \in X.$$

Thus, x_* can be extended to a saddle point of the Lagrange function and therefore is optimal for (P).

Proof, \Leftarrow **[under Slater condition]** Let (*P*) be convex and satisfy the Slater condition, let x_* be optimal and *f*, g_j be differentiable at x_* . Then

Exist Lagrange multipliers $\lambda^* \geq 0$ such that

(a) $\nabla f(x_*) + \sum_{j=1}^m \lambda_j^* \nabla g_j(x_*) \in N_X^*(x_*)$

(b) $\lambda_j^* g_j(x_*) = 0, \ j \le m$ [complementary slackness]

By Saddle Point Optimality condition, from optimality of x_* it follows that $\exists \lambda^* \geq 0$ such that (x_*, λ^*) is a saddle point of $L(x, \lambda)$ on $X \times \{\lambda \geq 0\}$. This is equivalent to

$$\lambda_j^* g_j(x_*) = 0 \ \forall j \ \& \ \underbrace{\min_{x \in X} L(x, \lambda^*) = L(x_*, \lambda^*)}_{(*)}$$

Since the function $L(x, \lambda^*)$ is convex in $x \in X$ and differentiable at $x_* \in X$, relation (*) implies (*a*).

Application example: Assuming $a_i > 0$, $p \ge 1$, let us solve the problem

$$\min_{x} \left\{ \sum_{i} \frac{a_i}{x_i} : x > 0, \sum_{i} x_i^p \le 1 \right\}$$

Assuming $x_* > 0$ is a solution such that $\sum_i (x_i^*)^p = 1$, the KKT conditions read

$$\nabla_x \left\{ \sum_i \frac{a_i}{x_i} + \lambda \left(\sum_i x_i^p - 1 \right) \right\} = 0 \Leftrightarrow \frac{a_i}{x_i^2} = p \lambda x_i^{p-1}$$
$$\sum_i x_i^p = 1$$

whence $x_i = c(\lambda)a_i^{\frac{1}{p+1}}$. Since $\sum_i x_i^p$ should be 1, we get

$$x_{i}^{*} = \frac{a_{i}^{\frac{1}{p+1}}}{\left(\sum_{j} a_{j}^{\frac{p}{p+1}}\right)^{\frac{1}{p}}}.$$

This point is feasible, problem is convex, KKT at the point is satisfied $\Rightarrow x^*$ is optimal!

Existence of Saddle Points

♣ <u>Theorem</u> [Sion-Kakutani] Let $X \subset \mathbb{R}^n$, $\Lambda \subset \mathbb{R}^m$ be nonempty convex closed sets and $F(x,\lambda) : X \times \Lambda \to \mathbb{R}$ be a continuous function which is convex in $x \in X$ and concave in $\lambda \in \Lambda$.

Assume that X is compact, and that there exists $\bar{x} \in X$ such that all the sets

$$\Lambda_a : \{\lambda \in \Lambda : F(\bar{x}, \lambda) \ge a\}$$

are bounded (e.g., Λ is bounded).

Then *F* possesses a saddle point on $X \rightarrow \Lambda$. <u>Proof:</u>

<u>MiniMax Lemma:</u> Let $f_i(x)$, i = 1, ..., m, be convex continuous functions on a convex compact set $X \subset \mathbb{R}^n$. Then there exists $\mu^* \geq 0$ with $\sum_i \mu_i^* = 1$ such that

$$\min_{x \in X} \max_{1 \le i \le m} f_i(x) = \min_{x \in X} \sum_i \mu_i^* f_i(x)$$

<u>Note:</u> When $\mu \ge 0, \sum_i \mu_i = 1$, one has

$$\max_{\substack{1 \le i \le m}} f_i(x) \ge \sum_i \mu_i f_i(x)$$

$$\Rightarrow \min_{x \in X} \max_i f_i(x) \ge \min_{x \in X} \sum_i \mu_i f_i(x)$$

Proof of MinMax Lemma: Consider the optimization program

$$\min_{t,x} \left\{ t : f_i(x) - t \le 0, \ i \le m, (t,x) \in X_+ \right\},$$
$$X_+ = \{(t,x) : x \in X\}$$
(P)

This program clearly is convex, solvable and satisfies the Slater condition, whence there exists $\lambda^* \ge 0$ and an optimal solution (x_*, t_*) to (P) such that (x_*, λ^*) is the saddle point of the Lagrange function on $X^+ \times \{\lambda \ge 0\}$:

$$\min_{x \in X, t} \left\{ t + \sum_{i} \lambda_{i}^{*}(f_{i}(x) - t) \right\} = t_{*} + \sum_{i} \lambda_{i}^{*}(f_{i}(x_{*}) - t_{*}) \quad (a)$$

$$\max_{\lambda \ge 0} \left\{ t + \sum_{i} \lambda_i (f_i(x) - t) \right\} = t_* + \sum_{i} \lambda_i^* (f_i(x_*) - t_*) \quad (b)$$

(b) implies that $t_* + \sum_i \lambda_i^* (f_i(x_*) - t_*) = t_*$. (a) implies that $\sum_i \lambda_i^* = 1$ and therefore implies that

$$\min_{x \in X} \sum_{i} \lambda_i^* f_i(x) = t_* = \min_{x \in X} \max_{i} f_i(x).$$

Proof of Sion-Kakutani Theorem: We should prove that problems

$$Opt(P) = \inf_{x \in X} \underbrace{\sup_{\lambda \in \Lambda} \overline{F(x,\lambda)}}_{X \in X} (P)$$
$$Opt(D) = \sup_{\lambda \in \Lambda} \inf_{x \in X} F(x,\lambda) (D)$$
$$\underbrace{\underline{F(\lambda)}}_{\underline{F(\lambda)}} (D)$$

are solvable with equal optimal values. $\mathbf{1}^{0}$. Since X is compact and $F(x, \lambda)$ is continuous on $X \times \lambda$, the function $\underline{F}(\lambda)$ is continuous on Λ . Besides this, the sets

$$\Lambda^a = \{\lambda \in \Lambda : \underline{F}(\lambda) \ge a\}$$

are contained in the sets

$$\Lambda_a = \{\lambda \in \Lambda : F(\bar{x}, \lambda) \ge a\}$$

and therefore are bounded. Finally, Λ is closed, so that the *continuous* function $\underline{F}(\cdot)$ with *bounded* level sets Λ^a attains it maximum on a *closed* set Λ . Thus, (D) is solvable; let λ^* be an optimal solution to (D).

 2^0 . Consider the sets

$$X(\lambda) = \{ x \in X : F(x, \lambda) \le \operatorname{Opt}(D) \}.$$

These are closed convex subsets of a compact set X. Let us prove that every finite collection of these sets has a nonempty intersection. Indeed, assume that

$$X(\lambda^1) \cap ... \cap X(\lambda^N) = \emptyset.$$

so that

$$\max_{j=1,\ldots,N} F(x,\lambda^j) > \operatorname{Opt}(D).$$

By MinMax Lemma, there exist weights $\mu_j \ge 0$, $\sum_{j} \mu_j = 1$, such that

$$\min_{x \in X} \underbrace{\sum_{j} \mu_{j} F(x, \lambda^{j})}_{\geq F(x, \sum_{j} \mu_{j} \lambda^{j})} > Opt(D)$$

which is impossible.

3⁰. Since every finite collection of closed convex subsets $X(\lambda)$ of a compact set has a nonempty intersection, all those sets have a nonempty intersection:

 $\exists x_* \in X : F(x_*, \lambda) \leq \mathsf{Opt}(D) \ \forall \lambda.$

Due to $Opt(P) \ge Opt(D)$, this is possible iff x_* is optimal for (P) and Opt(P) = Opt(D).

Optimality Conditions in Mathematical Programming

Situation: We are given a Mathematical Programming problem

$$\min_{x} \left\{ \begin{aligned} g_{1}(x), g_{2}(x), ..., g_{m}(x)) &\leq 0\\ f(x): & (h_{1}(x), ..., h_{k}(x)) = 0\\ & x \in X \end{aligned} \right\}. \tag{P}$$

Question of interest: Assume that we are given a feasible solution x_* to (P). What are the conditions (necessary, sufficient, necessary and sufficient) for x_* to be optimal? Fact: Except for convex programs, there are no verifiable local sufficient conditions for global optimality. There exist, however, \diamond verifiable local necessary conditions for *local* (and thus – for global) optimality \diamond verifiable local sufficient conditions for *local* optimality Fact: Existing conditions for local optimality assume that $x_* \in intX$, which, from the viewpoint of local optimality of x_* , is exactly the same as to say that $X = \mathbb{R}^n$. Situation: We are given a Mathematical Programming problem

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\}. \tag{P}$$

<u>and</u> a feasible solution x_* to the problem, and are interested in necessary/sufficient conditions for *local* optimality of x_* :

There exists r > 0 such that for every feasible x with $||x - x_*|| \le r$ one has

$$f(x) \geq f(x_*).$$

Default assumption: The objective and all the constraints are continuously differentiable in a neighbourhood of x_* .

$$\min_{x} \left\{ f(x): \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\}. \tag{P}$$

First Order Optimality Conditions are expressed via values and gradients of the objective and the constraints at x_* . Except for convex case, only <u>necessary</u> First Order conditions are known.

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), ..., g_{m}(x)) \leq 0\\ (h_{1}(x), ..., h_{k}(x)) = 0 \end{array} \right\}. \tag{P}$$

The idea:

 \diamond Assume that x_* is locally optimal for (*P*). Let us approximate (*P*) around x_* by a Linear Programming program

$$\min_{x} f(x_{*}) + (x - x_{*})^{T} f'(x_{*})$$
s.t.

$$\underbrace{0}{g_{j}(x_{*})} + (x - x_{*})^{T} g'_{j}(x_{*}) \leq 0, \ j \in J(x_{*})$$

$$\underbrace{h_{i}(x_{*})}{0} + (x - x_{*})^{T} h'_{i}(x_{*}) = 0, \ 1 \leq i \leq k$$

$$\begin{bmatrix}J(x_{*}) = \{j : g_{j}(x_{*}) = 0\}\end{bmatrix}$$
(LP)

<u>Note</u>: Since all $g_j(\cdot)$ are continuous at x_* , the *non-active at* x_* inequality constraints (those with $g_j(x_*) < 0$) do not affect (*LP*).

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} \quad (P) \\
\Rightarrow \min_{x} \left\{ \begin{array}{c} (x - x_{*})^{T} f'(x_{*}) : \begin{array}{c} j \in J(x_{*})\\ (x - x_{*})^{T} h'_{i}(x_{*}) = 0, \\ i = 1, \dots, k \end{array} \right\} \quad (LP) \\
J(x_{*}) = \{j : g_{j}(x_{*}) = 0\}$$

 \diamond It is natural *to guess* that if x_* is locally optimal for (*P*), then x_* is locally optimal for (*LP*) as well.

LP is a *convex* program with *affine* constraints, whence the KKT conditions are necessary and sufficient for optimality:

$$x_* \text{ is optimal for } (LP)$$

$$\downarrow \\ \exists (\lambda_j^* \ge 0, j \in J(x_*), \mu_i) :$$

$$f'(x_*) + \sum_{j \in J(x_*)} \lambda_j^* g'_j(x_*) + \sum_{i=1}^k \mu_i h'_i(x_*) = 0$$

$$\downarrow \\ \exists (\lambda_j^* \ge 0, \mu_i^*) :$$

$$f'(x_*) + \sum_j \lambda_j^* g'_j(x_*) + \sum_i \mu_i^* h'_i(x_*) = 0$$

$$\lambda_j^* g_j(x_*) = 0, j = 1, ..., m$$

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), ..., g_{m}(x)) \leq 0\\ (h_{1}(x), ..., h_{k}(x)) = 0 \end{array} \right\}$$
(P)

<u>Proposition</u>. Let x_* be a locally optimal solution of (P).

Assume that x_* remains locally optimal when passing from (P) to the linearized problem

$$\min_{x} \left\{ (x - x_{*})^{T} f'(x_{*}) : \begin{array}{l} (x - x_{*})^{T} g'_{j}(x_{*}) \leq 0, \\ j \in J(x_{*}) \\ (x - x_{*})^{T} h'_{i}(x_{*}) = 0, \\ i = 1, \dots, k \end{array} \right\} \tag{LP}$$

Then at x_* the KKT condition holds:

$$\exists (\lambda_j^* \ge 0, \mu_i^*) : \\ f'(x_*) + \sum_j \lambda_j^* g'_j(x_*) + \sum_i \mu_i^* h'_i(x_*) = 0 \\ \lambda_j^* g_j(x_*) = 0, \ j = 1, ..., m$$

$$\min_{x} \left\{ f(x): \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} \quad (P)$$

$$\min_{x} \left\{ (x - x_{*})^{T} f'(x_{*}): \begin{array}{c} (x - x_{*})^{T} g'_{j}(x_{*}) \leq 0,\\ j \in J(x_{*})\\ (x - x_{*})^{T} h'_{i}(x_{*}) = 0,\\ i = 1, \dots, k \end{array} \right\} \quad (LP)$$

To make Proposition useful, we need a verifiable sufficient condition for " x_* remains locally optimal when passing from (P) to (LP)".

The most natural form of such a condition is *regularity*:

Gradients, taken at x_* , of all constraints active at x_* are linearly independent.

Of course, all equality constraints by definition are active at every feasible solution.

$$\min_{x} \left\{ f(x): \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} \quad (P)$$

$$\min_{x} \left\{ (x - x_{*})^{T} f'(x_{*}): \begin{array}{c} (x - x_{*})^{T} g'_{j}(x_{*}) \leq 0,\\ j \in J(x_{*})\\ (x - x_{*})^{T} h'_{i}(x_{*}) = 0,\\ i = 1, \dots, k \end{array} \right\} \quad (LP)$$

<u>Proposition</u>: Let x_* be a locally optimal *reg*ular solution of (P). Then x_* is optimal for (LP) and, consequently, the KKT conditions take place at x_* .

<u>Proof</u> is based on an important fact of Analysis – a version of Implicit Function Theorem. <u>Theorem</u>: Let $x_* \in \mathbb{R}^n$ and let $p_{\ell}(x)$, $\ell = 1, ..., L$, be real-valued functions such that $\Diamond p_{\ell}$ are $\kappa \geq 1$ times continuously differentiable in a neighbourhood of x_*

$$\Diamond p_\ell(x_*) = 0$$

 \diamond vectors $\nabla p_{\ell}(x_*)$, $\ell = 1, ..., L$, are linearly independent.

Then there exists substitution of variables

$$y \mapsto x = \Phi(y)$$

defined in a neighbourhood V of the origin and mapping V, in a one-to-one manner, onto a neighbourhood B of x_* , such that $\Diamond x_* = \Phi(0)$

(both Φ : $V \to B$ and its inverse mapping Φ^{-1} : $B \to V$ are κ times continuously differentiable

 \Diamond in coordinates y, the functions p_{ℓ} become just the coordinates:

$$y \in V \Rightarrow p_{\ell}(\Phi(y)) \equiv y_{\ell}, \ \ell = 1, ..., L.$$

$$\min_{x} \left\{ (x - x_{*})^{T} f'(x_{*}) : \begin{array}{l} (x - x_{*})^{T} g'_{j}(x_{*}) \leq 0, \\ j \in J(x_{*}) \\ (x - x_{*})^{T} h'_{i}(x_{*}) = 0, \\ i = 1, ..., k \end{array} \right\} \quad (LP)$$

Let x_* be a regular locally optimal solution to (P); assume, on the contrary to what should be proven, that x_* is not an optimal solution to (LP), and let us lead this to contradiction. $\mathbf{1}^0$. Since $x = x_*$ is not an optimal solution to (LP), there exists a feasible solution $x' = x_* + d$ to the problem with $(x' - x_*)^T f'(x_*) = d^T f'(x_*) < 0$, so that

$$d^T f'(x_*) < 0, \underbrace{d^T h'_i(x_*) = 0}_{\forall i}, \underbrace{d^T g'_j(x_*) \le 0}_{\forall j \in J(x_*)}$$

$$d^T f'(x_*) < 0, \underbrace{d^T h'_i(x_*) = 0}_{\forall i}, \underbrace{d^T g'_j(x_*) < 0}_{\forall j \in J(x_*)}$$

2⁰. W.I.o.g., assume that $J(x_*) = \{1, ..., \ell\}$. By Theorem, there exist continuously differentiable local substitution of argument

$$x = \Phi(y) \qquad \qquad [\Phi(0) = x_*]$$

with a continuously differentiable in a neighbourhood of x_* inverse $y = \Psi(x)$ such that in a neighbourhood of origin one has

 $h_i(\Phi(y)) \equiv y_i, \ g_j(\Phi(y)) = y_{k+j}, \ j = 1, ..., \ell.$ Since $\Psi(\Phi(y)) \equiv y_i$, we have $\Psi'(x_*)\Phi'(0) = I$, whence

$$\exists e : \Phi'(0)e = d.$$

<u>Situation</u>: We have found a smooth local substitution of argument $x = \Phi(y)$ (y = 0 corresponds to $x = x_*$) and a direction e such that in a neighbourhood of y = 0 one has

Consider the differentiable curve

$$x(t) = \Phi(te).$$

We have

$$te_{i} \equiv h_{i}(\Phi(te)) \Rightarrow e_{i} = [\Phi'(0)e]^{T}h'_{i}(x_{*}) = 0$$

$$te_{k+j} \equiv g_{j}(\Phi(te)) \Rightarrow e_{k+j} = [\Phi'(0)e]^{T}g'_{j}(x_{*}) < 0$$

$$\Rightarrow \underbrace{h_{i}(x(t)) = te_{i} = 0}_{\forall i}, \underbrace{g_{j}(x(t)) = te_{k+j} \leq 0}_{\forall j \in J(x_{*})}$$

Thus, x(t) is feasible for all small $t \ge 0$. But:

$$\frac{d}{dt}\Big|_{t=0} f(x(t)) = [\Phi'(0)e]^T f'(x_*) < 0,$$

whence $f(x(t)) < f(x(0)) = f(x_*)$ for all small enough t > 0, which is a contradiction with local optimality of x_* . Second Order Optimality Conditions

In the case of unconstrained minimization problem

$$\min_{x} f(x) \tag{P}$$

with continuously differentiable objective, the KKT conditions reduce to Fermat Rule: If x_* is locally optimal for (P), then $\nabla f(x_*) = 0$. Fermat Rule is the "first order" part of Second Order Necessary Optimality Condition in unconstrained minimization:

<u>If</u> x_* is locally optimal for (P) and f is twice differentiable in a neighbourhood of x_* , <u>then</u>

 $\nabla f(x_*) = 0 \& \nabla^2 f(x_*) \succeq 0 \Leftrightarrow d^T \nabla^2 f(x_*) d \ge 0 \forall d$

Indeed, let x_* be locally optimal for (P); then for appropriate $r_d > 0$

$$0 \leq t \leq r_d$$

$$\Rightarrow 0 \leq f(x_* + td) - f(x_*)$$

$$= t \underbrace{d^T \nabla f(x_*)}_{=0} + \frac{1}{2} t^2 d^T \nabla^2 f(x_*) d + t^2 \underbrace{R_d(t)}_{\stackrel{\to 0,}{t \to 0}}$$

$$\Rightarrow \frac{1}{2} d^T \nabla^2 f(x_*) d + R_d(t) \geq 0 \Rightarrow d^T \nabla^2 f(x_*) d \geq 0$$

$$\min_{x} f(x) \tag{P}$$

The *necessary* Second Order Optimality condition in unconstrained minimization can be strengthened to

Second Order Sufficient Optimality Condition in unconstrained minimization: Let f be twice differentiable in a neighbourhood of x_* . If

 $\nabla f(x_*) = 0, \nabla^2 f(x_*) \succ 0 \Leftrightarrow d^T \nabla^2 f(x_*) d > 0 \forall d \neq 0$

<u>then</u> x_* is locally optimal for (*P*). **Proof:** Since $d^T \nabla^2 f(x_*) d > 0$ for all d > 0, then there exists $\alpha > 0$ such that $d^T \nabla^2 f(x_*) d \ge \alpha d^T d$ for all d.

By differentiability, for every $\epsilon>0$ there exists $r_\epsilon>0$ such that

mizer of f.

We are given a Mathematical Programming problem

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), ..., g_{m}(x)) \leq 0\\ (h_{1}(x), ..., h_{k}(x)) = 0 \end{array} \right\} (P) \\
\downarrow \\ L(x; \lambda, \mu) = f(x) + \sum_{j} \lambda_{j} g_{j}(x) + \sum_{i} \mu_{i} h_{i}(x)$$

♣ In Optimality Conditions for a constrained problem (P), the role of $\nabla^2 f(x_*)$ is played by the Hessian of the Lagrange function: Second Order Necessary Optimality Condition: Let x_* be a regular feasible solution of (P) such that the functions f, g_j, h_i are twice continuously differentiable in a neighbourhood of x_* . If x_* is locally optimal, then \Diamond There exist uniquely defined Lagrange multipliers $\lambda_j^* \ge 0$, μ_i^* such that the KKT conditions hold:

$$\nabla_x L(x_*; \lambda^*, \mu^*) = 0$$

 $\lambda_j^* g_j(x_*) = 0, \ j = 1, ..., m$

 \diamond For every d orthogonal to the gradients, taken at x_* , of all active at x_* equality and inequality constraints, one has

$$d^T \nabla_x^2 L(x_*; \lambda^*, \mu^*) d \ge 0.$$

Proof. 1^0 . Constraints which are non-active at x_* clearly do not affect neither local optimality of x_* , nor the conclusion to be proven. Removing these constraints, we reduce the situation to one where *all* constraints in the problem

$$\min_{x} \left\{ f(x): \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} (P)$$

are active at x_* .

 2^{0} . Applying Implicit Function Theorem, we can find a local change of variables

$$x = \Phi(y) \Leftrightarrow y = \Psi(x)$$
$$[\Phi(0) = x_*, \Psi(x_*) = 0]$$
with locally twice continuously differentiable

 Φ, Ψ such that

$$g_j(\Phi(y)) \equiv y_j, \ j \leq m, h_i(\Phi(y)) \equiv y_{m+i}, \ i \leq k.$$

In variables y, problem (P) becomes

$$\min_{y} \left\{ \underbrace{f(\Phi(y))}_{\phi(y)} : y_j \le 0, \ j \le m, y_{k+i} = 0, \ i \le k \right\}.$$
(P')

$$\min_{x} \left\{ f(x): \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\}$$
(P)

$$\lim_{y} \left\{ \underbrace{f(\Phi(y))}_{\phi(y)}: y_{j} \leq 0, \ j \leq m, y_{k+i} = 0, \ i \leq k \right\}$$
(P)

$$\left\{ \underbrace{f(\Phi(y))}_{\phi(y)}: y_{j} \leq 0, \ j \leq m, y_{k+i} = 0, \ i \leq k \right\}$$
(P)

$$\left\{ \underbrace{f(\Phi(y))}_{\phi(y)}: y_{j} \leq 0, \ j \leq m, y_{k+i} = 0, \ i \leq k \right\}$$
(P)

Our plan is as follows:

 \diamond Since Φ is a smooth one-to-one mapping of a neighbourhood of x_* onto a neighbourhood of $y_* = 0$, x_* is locally optimal for (P) iff $y_* = 0$ is locally optimal for (P').

 \diamond We intend to build necessary/sufficient conditions for $y_* = 0$ to be locally optimal for (P'); "translated" to x-variables, these conditions will imply necessary/sufficient conditions for local optimality of x_* for (P).

3⁰. Since $x_* = \Phi(0)$ is locally optimal for (P), $y_* = 0$ is locally optimal for (P'). In particular, if e_i is *i*-th basic orth, then for appropriate $\epsilon > 0$:

$$\begin{aligned} j &\leq m \Rightarrow y(t) = -te_j \text{ is feasible for } (P') \text{ when} \\ \epsilon &\geq t \geq 0 \Rightarrow -\frac{\partial \phi(0)}{\partial y_t} = \frac{d}{dt} \Big|_{t=0} \phi(y(t)) \geq 0 \\ \Rightarrow \lambda_j^* \equiv -\frac{\partial \phi(0)}{\partial y_i} \geq 0 \end{aligned}$$

and

$$\begin{split} s > m + k &\Rightarrow y(t) = te_s \text{ is feasible for } (P') \text{ when } \\ \epsilon \ge t \ge \epsilon &\Rightarrow \left. \frac{\partial \phi(0)}{\partial y_s} = \frac{d}{dt} \right|_{t=0} \phi(y(t)) = 0 \\ \text{Setting } \mu_i^* = -\frac{\partial \phi(0)}{\partial y_{m+i}}, \ i = 1, \dots, k, \text{ we get} \\ \lambda^* \ge 0 \& \nabla_y M(0; \lambda^*, \mu^*) = 0. \quad (\text{KKT}) \end{split}$$

Note that the condition $\nabla_y M(0; \lambda^*, \mu^*) = 0$ defines λ^* , μ^* are in a unique fashion. $\mathbf{4}^0$. We have seen that for (P'), the first order part of the Necessary Second Order Optimality condition holds true. Let us prove the second order part of the condition, which reads

$$egin{aligned} &orall (d : d^T
abla_y y_\ell = 0, \, \ell \leq m+k): \ & d^T
abla_y^2 M(0; \lambda^*, \mu^*) d \geq 0. \end{aligned}$$
<u>Situation</u>: $y_* = 0$ is locally optimal solution to the problem

$$\min_{y} \left\{ \phi(y) \equiv f(\Phi(y)) : \begin{array}{l} y_j \leq 0, \ j \leq m \\ y_m + i = 0, \ i \leq k \end{array} \right\} \\ (P')$$

Claim:

$$orall (d: d^T
abla_y y_\ell = 0, \ \ell \leq m+k): \ d^T
abla_y^2 M(0; \lambda^*, \mu^*) d \geq 0.$$

This is evident: since $M(y; \lambda^*, \mu^*) = \phi(y) + \sum_{j=1}^{m} \lambda_j^* y_j + \sum_{i=1}^{k} \mu_i^* y_{m+i}$, we have

$$\nabla_y^2 M(0; \lambda^*, \mu^*) = \nabla^2 \phi(0).$$

Claim therefore states that $d^T \nabla^2 \phi(0) d \ge 0$ for every vector d from the linear subspace $L = \{d : d_1 = ... = d_{m+k} = 0\}$. But this subspace is feasible for (P'), so that ϕ , restricted onto L, should attain unconstrained local minimum at the origin. By Necessary Second Order Optimality condition for unconstrained minimization,

$$d^T \nabla^2 \phi(\mathbf{0}) d \ge \mathbf{0} \ \forall d \in L.$$

5⁰. We have seen that if x_* is locally optimal for (*P*), then there exist uniquely defined $\lambda^* \geq 0$, μ^* such that

$$\nabla_y M(0; \lambda^*, \mu^*) = 0,$$

and one has

 $d^T \nabla_y y_\ell = 0, \ \ell \leq m + k \Rightarrow d^T \nabla_y^2 M(0; \lambda^*, \mu^*) d \geq 0.$

Let us prove that then

$$\nabla_x L(x_*; \lambda^*, \mu^*) = 0 \qquad (*)$$

and

$$e^{T}g'_{j}(x_{*}) = 0, \ j \le m \\ e^{T}h'_{i}(x_{*}) = 0, \ i \le k \ \ \} \Rightarrow e^{T}\nabla_{x}^{2}L(x_{*}); \lambda^{*}, \mu^{*})e \ge 0.$$

$$(**)$$

Given:

$$abla_y M(0;\lambda^*,\mu^*) = 0 \ d^T
abla_y y_\ell = 0, \ \ell \leq m+k \Rightarrow d^T
abla_y^2 M(0;\lambda^*,\mu^*) d \geq 0.$$

Should prove:

$$\begin{array}{c} \nabla_{x}L(x_{*};\lambda^{*},\mu^{*}) = 0 & (*) \\ e^{T}g'_{j}(x_{*}) = 0, \ j \leq m \\ e^{T}h'_{i}(x_{*}) = 0, \ i \leq k \end{array} \right\} \Rightarrow e^{T}\nabla_{x}^{2}L(x_{*});\lambda^{*},\mu^{*})e \geq 0 \quad (**)$$

 \diamondsuit Setting $\mathcal{L}(x) = L(x; \lambda^*, \mu^*), \ \mathcal{M}(y) = M(y; \lambda^*, \mu^*),$ we have

$$\mathcal{L}(x) = \mathcal{M}(\Psi(x))$$

$$\Rightarrow \nabla_x \mathcal{L}(x_*) = [\Psi'(x_*)]^T \nabla_y \mathcal{M}(y_*) = 0,$$

as required in (*).

 \diamond Let *e* satisfy the premise in (**), and let $d = [\Phi'(0)]^{-1}e$. Then

$$\begin{aligned} \frac{\frac{d}{dt}\Big|_{t=0}^{td_j}}{\frac{d}{dt}\Big|_{t=0}^{t=0}} &= [g'_j(x_*)]^T \overbrace{[\Phi'(0)]d}^e \\ \Rightarrow d_j &= e^T g'_j(x_*) = 0, \ j \le m \\ \frac{d}{dt}\Big|_{t=0}^{td_m(\Phi(td))} &= [h'_i(x_*)]^T \underbrace{[\Phi'(0)]d}_e \\ \xrightarrow{\frac{d}{dt}\Big|_{t=0}^{td_m+i}} \\ \Rightarrow d_{m+i} &= e^T h'_i(x_*) = 0, \ i \le k \end{aligned}$$

We have

$$e^{T}\nabla^{2}\mathcal{L}(x_{*})e = \frac{d^{2}}{dt^{2}}\bigg|_{t=0}\mathcal{L}(x_{*}+te) = \frac{d^{2}}{dt^{2}}\bigg|_{t=0}\mathcal{M}(\Psi(x_{*}+te))$$

$$= \frac{d}{dt}\bigg|_{t=0}\left[e^{T}[\Psi'(x_{*}+te)]^{T}\nabla\mathcal{M}(\Psi(x_{*}+te))\right]$$

$$= e^{T}[\Psi'(x_{*})]^{T}\nabla^{2}\mathcal{M}(0)[\Psi'(x_{*})e]$$

$$+e^{T}[\frac{d}{dt}_{t=0}\Psi'(x_{*}+te)]^{T}\underbrace{\nabla\mathcal{M}(0)}_{=0}$$

$$= d^{T}\nabla^{2}\mathcal{M}d \ge 0,$$

Thus, whenever e is orthogonal to the gradients of all constraints active at x_* , we have $e^T \nabla^2 \mathcal{L} e \ge 0$.

Second Order Sufficient Condition for Local Optimality

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} (P) \\ \downarrow \\ L(x; \lambda, \mu) = f(x) + \sum_{j} \lambda_{j} g_{j}(x) + \sum_{i} \mu_{i} h_{i}(x) \end{array}$$

Second Order Sufficient Optimality Condition: Let x_* be a regular feasible solution of (P) such that the functions f, g_j, h_i are twice continuously differentiable in a neighbourhood of x_* . If there exist Lagrange multipliers $\lambda_j^* \ge 0$, μ_i^* such that

♦ the KKT conditions hold:

$$\nabla_x L(x_*; \lambda^*, \mu^*) = 0$$

 $\lambda_j^* g_j(x_*) = 0, \ j = 1, ..., m$

 \diamond For every $d \neq 0$ orthogonal to the gradients, taken at x_* , of all active at x_* equality constraints and those active at x_* inequality constraints for which $\lambda_i^* > 0$, one has

$$d^T \nabla^2_x L(x_*; \lambda^*, \mu^*) d > 0$$

<u>then</u> x_* is locally optimal for (P).

<u>Note:</u> Difference between Sufficient and Necessary optimality conditions is in their "second order" parts and is twofold:

 \Diamond [minor difference] Necessary condition states positive <u>semi</u>definiteness of $\nabla_x^2 L(x_*; \lambda^*, \mu^*)$ along linear subspace:

$$\forall d \in T = \{d : \overbrace{d^T h'_i(x_*) = 0}^{\forall i \leq k}, \overbrace{d^T g'_j(x_*) = 0}^{\forall j \in J(x_*)} \} : d^T \nabla_x^2 L(x_*; \lambda^*, \mu^*) d \geq 0$$

while Sufficient condition requires positive definiteness of $\nabla_x^2 L(x_*; \lambda^*, \mu^*)$ along linear subspace:

$$\forall 0 \neq d \in T^+ = \{d : \overbrace{d^T h'_i(x_*) = 0}^{\forall i \leq k}, \overbrace{d^T g'_j(x_*) = 0}^{\forall j : \lambda_j^* > 0} \} : d^T \nabla_x^2 L(x_*; \lambda^*, \mu^*) d > 0$$

 $\langle [major difference]$ The linear subspaces in question are different, and $T \subset T^+$; the subspaces are equal to each other iff *all* active at x_* inequality constraints have positive Lagrange multipliers λ_i^* . <u>Note:</u> This "gap" is essential, as is shown by example

$$\begin{split} \min_{x_1,x_2} \left\{ f(x) = x_2^2 - x_1^2 : g_1(x) = x_1 \leq 0 \right\} \\ & [x_* = (0,0)^T] \\ \text{Here the Necessary Second Order Optimal-} \\ & \text{ity condition is satisfied "strictly": } L(x;\lambda) = \\ & x_2^2 - x_1^2 + \lambda x_1, \text{ whence} \end{split}$$

$$\lambda^* = 0 \Rightarrow \nabla_x L(x_*; \lambda^*) = 0,$$

$$T = \{d : d^T g'_1(0) = 0\} = \{d : d_1 = 0\},$$

$$0 \neq d \in T \Rightarrow d^T \nabla_x^2 L(x_*; \lambda^*) d = d_2^2 > 0\}$$

while x_* is <u>not</u> a local solution.

Proof of Sufficient Second Order Optimality Condition. 1⁰. As in the case of Second Order Necessary Optimality Condition, we can reduce the situation to one where \Diamond All inequality constraints are active at x_* \Diamond The problem is of the special form

$$\min_{y} \left\{ \phi(y) : \begin{array}{l} y_j \leq 0, \ j \leq m \\ y_{m+i} = 0, \ i \leq k \end{array} \right\} \qquad (P')$$

2⁰. In the case of (P'), Sufficient condition reads: $\exists \lambda^* \geq 0, \mu^*$:

$$\begin{aligned} \nabla_{y} \Big|_{y=0} \left\{ \phi(y) + \sum_{j=1}^{m} \lambda_{j}^{*} y_{j} + \sum_{i=1}^{k} \mu_{i}^{*} y_{m+i} \right\} \\ d_{j} &= 0, j \in J, d \neq 0 \Rightarrow d^{T} \nabla^{2} \phi(0) d > 0 \\ \left[J &= \{ j \leq m : \lambda_{j}^{*} > 0 \} \cup \{ m+1, ..., m+k \} \right] \end{aligned}$$
(*)

Assuming w.l.o.g. $\{j : \lambda_j^* > 0\} = \{1, ..., q\}$, (*) reads:

$rac{\partial \phi(0)}{\partial y_\ell}$	<	$0,\ell=1,,q$
$\frac{\partial \phi(0)}{\partial u_{\ell}}$	=	$0, \ell = q+1,, m$
$\frac{\partial \phi(0)}{\partial u_{\ell}}$	=	$0, \ell = m + k + 1,, n$
$0 \neq d$	\in	$T^+ = \{d : d_\ell = 0, \ell \in \{1,, q, m+1,, m+k\}\}$
	\Rightarrow	$d^T abla^2 \phi(0) d > 0$

Our goal is to derive from this assumption local optimality of $y_* = 0$ for (P').

2⁰. The feasible set of (P') is the closed cone $K = \{d : d_{\ell} \le 0, \ell = 1, ..., m, d_{\ell} = 0, \ell = m+1, ..., m+k\}$ Lemma: For $0 \ne d \in K$ one has $d^T \nabla \phi(0) \ge 0$

Lemma: For $0 \neq d \in K$ one has $d^T \nabla \phi(0) \ge 0$ and

 $d^T \nabla \phi(\mathbf{0}) = \mathbf{0} \Rightarrow d^T \nabla^2 \phi(\mathbf{0}) d > \mathbf{0}.$

Situation:

$$\begin{aligned} \frac{\partial \phi(0)}{\partial y_{\ell}} &< 0, \ \ell = 1, ..., q \\ \frac{\partial \phi(0)}{\partial y_{\ell}} &= 0, \ \ell = q + 1, ..., m \\ \frac{\partial \phi(0)}{\partial y_{\ell}} &= 0, \ \ell = m + k + 1, ..., n \\ 0 \neq d &\in T^{+} = \{d : d_{\ell} = 0, \ \ell \in \{1, ..., q, m + 1, ..., m + k\}\}: \\ &\Rightarrow \ d^{T} \nabla^{2} \phi(0) d > 0 \end{aligned}$$

$$(*)$$

 $K = \{ d : d_{\ell} \le 0, \, \ell = 1, ..., m, d_{\ell} = 0, \, \ell = m+1, ..., m+k \}$ (**)

<u>Claim</u>: For $0 \neq d \in K$ one has $d^T \nabla \phi(0) \ge 0$ and

$$d^T \nabla \phi(\mathbf{0}) = \mathbf{0} \Rightarrow d^T \nabla^2 \phi(\mathbf{0}) d > \mathbf{0}.$$

Proof: For $d \in K$, we have

$$d^T \nabla \phi(0) = \sum_{\ell=1}^n \frac{\partial \phi(0)}{\partial y_\ell} d_\ell$$

By (*) - (**), the first q terms in this sum are nonnegative, and the remaining are 0. Thus, the sum always is ≥ 0 . For $d \neq 0$, the only possibility for the sum to vanish is to have $d \in T^+$, and in this case $d^T \phi''(0) d > 0$. <u>Situation:</u> (P') is the problem

$$\min_{y \in K} \phi(y), \tag{!}$$

K is a closed cone, ϕ is twice continuously differentiable in a neighbourhood of the origin and is such that

$$d \in K \Rightarrow d^T \nabla \phi(\mathbf{0}) \ge \mathbf{0}$$

$$d \in K \setminus \{\mathbf{0}\}, d^T \nabla \phi(\mathbf{0}) = \mathbf{0} \Rightarrow d^T \nabla^2 \phi(\mathbf{0}) d > \mathbf{0}$$

<u>Claim:</u> In the situation in question, 0 is a locally optimal solution to (!).

Proof: Let $M = \{d \in K : ||d||_2 = 1\}$, and let $M_0 = \{d \in M : d^T \nabla \phi(0) = 0\}$. Since K is closed, both M and M_0 are compact sets. We know that $d^T \nabla^2 \phi(0) d > 0$ for $d \in M_0$. Since M_0 is a compact set, there exists a neighbourhood V of M_0 and $\alpha > 0$ such that

$$d \in V \Rightarrow d^T \nabla^2 \phi(\mathbf{0}) d \ge \alpha.$$

The set $V_1 = M \setminus V$ is compact and $d^T \nabla \phi(0) > 0$ when $d \in V_1$; thus, there exists $\beta > 0$ such that

$$d \in V_1 \Rightarrow d^T \nabla \phi(\mathbf{0}) \ge \beta.$$

<u>Situation</u>: *K* is a cone, and the set $M = \{d \in K : ||d||_2 = 1\}$ is partitioned into two subsets $V_0 = V \cap M$ and V_1 in such a way that

$$\begin{aligned} d \in V_0 &\Rightarrow d^T \nabla \phi(0) \geq 0, d^T \nabla^2 \phi(0) d \geq \alpha > 0 \\ d \in V_1 &\to d^T \nabla \phi(0) \geq \beta > 0 \end{aligned}$$

<u>Goal</u>: To prove that 0 is local minimizer of ϕ on K, or, which is the same, that

$$\exists r > 0: \\ \phi(0) \leq \phi(td) \ \forall (d \in M, 0 \leq t \leq r).$$
Proof: Let $d \in M, t \geq 0$. When $d \in V_0$, we have
$$\phi(td) - \phi(0) \geq td^T \nabla \phi(0) + \frac{1}{2}t^2 d^T \nabla^2 \phi(0) d - t^2 \underbrace{R(t)}_{t \to 0} \\ \geq \frac{1}{2}t^2(\alpha - 2R(t)) \\ \Rightarrow \exists r_0 > 0: \qquad \phi(td) - \phi(0) \geq \frac{1}{4}t^2\alpha \geq 0 \ \forall t \leq r_0$$
When $d \in V_1$, we have
$$\phi(td) - \phi(0) \geq td^T \nabla \phi(0) + \frac{1}{2}t^2 d^T \nabla^2 \phi(0) d - t^2 \underbrace{R(t)}_{t \to 0} \\ \geq \beta t - Ct^2 - t^2 R(t) \\ \Rightarrow \exists r_1 > 0: \qquad \phi(td) - \phi(0) \geq \frac{\beta}{2}t \geq 0 \ \forall t \leq r_1$$
Thus, $\phi(td) - \phi(0) \geq 0$ for all $t \leq \min[r_0, r_1], d \in M.$

Sensitivity Analysis

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} (P) \\ \downarrow \\ L(x; \lambda, \mu) = f(x) + \sum_{j} \lambda_{j} g_{j}(x) + \sum_{i} \mu_{i} h_{i}(x)$$

<u>Definition</u>: Let x_* be a feasible solution to (P) such that the functions f, g_j, h_i are $\ell \ge 2$ times continuously differentiable in a neighbourhood of x_* .

 x_* is called a *nondegenerate locally optimal* solution to (P), if

 $\diamond x_*$ is a regular solution (i.e., gradients of active at x_* constraints are linearly independent)

♦ at x_* , Sufficient Second Order Optimality condition holds $\exists (\lambda^* \ge 0, \mu^*)$:

$$\begin{cases} \nabla_{x}L(x_{*};\lambda^{*},\mu^{*}) = 0 \\ \lambda_{j}^{*}g_{j}(x_{*}) = 0, \ j = 1,...,m \\ d^{T}\nabla g_{j}(x_{*}) = 0 \ \forall (j:\lambda_{j}^{*} > 0) \\ d^{T}\nabla h_{i}(x_{*}) = 0 \ \forall i \\ d \neq 0 \end{cases} \right\} \Rightarrow d^{T}\nabla_{x}^{2}L(x_{*};\lambda^{*},\mu^{*}) > 0$$

 \diamond for all active at x_* inequality constraints, Lagrange multipliers are positive:

$$g_j(x_*) = 0 \Rightarrow \lambda_j^* > 0.$$

$$\min_{x} \left\{ f(x) : \begin{array}{c} (g_{1}(x), g_{2}(x), \dots, g_{m}(x)) \leq 0\\ (h_{1}(x), \dots, h_{k}(x)) = 0 \end{array} \right\} (P)$$

<u>Theorem:</u> Let x_* be a nondegenerate locally optimal solution to (P). Let us embed (P) into the parametric family of problems

$$\min_{x} \left\{ f(x): \begin{array}{l} g_{1}(x) \leq a_{1}, \dots, g_{m}(x) \leq a_{m} \\ h_{1}(x) = b_{1}, \dots, h_{k}(x) = b_{k} \end{array} \right\} \quad (P[a, b])$$

so that (P) is (P[0,0]).

There exists a neighbourhood V_x of x_* and a neighbourhood $V_{a,b}$ of the point a = 0, b = 0in the space of parameters a, b such that $\Diamond \forall (a,b) \in V_{a,b}$, in V_v there exists a unique KKT point $x_*(a,b)$ of (P[a,b]), and this point is a nondegenerate locally optimal solution to (P[a,b]); moreover, $x_*(a,b)$ is optimal solution for the optimization problem

$$Opt_{loc}(a,b) = \min_{x} \left\{ f(x): \begin{array}{l} g_{1}(x) \leq a_{1}, ..., g_{m}(x) \leq a_{m} \\ h_{1}(x) = b_{1}, ..., h_{k}(x) = b_{k} \\ x \in V_{x} \end{array} \right\}$$
$$(P_{loc}[a,b])$$

 \diamond both $x_*(a,b)$ and the corresponding Lagrange multipliers $\lambda^*(a,b)$, $\mu^*(a,b)$ are $\ell - 1$ times continuously differentiable functions of $(a,b) \in V_{a,b}$, and

$$\frac{\partial \mathsf{Opt}_{\mathsf{loc}}(a,b)}{\partial a_j} = \frac{\partial f(x_*(a,b))}{\partial a_j} = -\lambda_j^*(a,b)$$
$$\frac{\partial \mathsf{Opt}_{\mathsf{loc}}(a,b)}{\partial b_i} = \frac{\partial f(x_*(a,b))}{\partial b_i} = -\mu_i^*(a,b)$$

Simple example: Existence of Eigenvalue

Consider optimization problem

$$Opt = \min_{x \in \mathbb{R}^n} \left\{ f(x) = x^T A x : h(x) = 1 - x^T x = 0 \right\}$$
(P)

where $A = A^T$ is an $n \times n$ matrix. The problem clearly is solvable. Let x_* be its optimal solution. What can we say about x_* ? <u>Claim:</u> x_* is a regular solution to (*P*). Indeed, we should prove that the gradients of active at x_* constraints are linearly independent. There is only one constraint, and its gradient at the feasible set is nonzero. Since x_* is a regular globally (and therefore locally) optimal solution, at x_* the Necessary Second Order Optimality condition should hold: $\exists \mu^*$:

$$\nabla_{x} \underbrace{\left[x^{T}Ax + \mu^{*}(1 - x^{T}x)\right]}_{\Leftrightarrow d^{T}\nabla_{x}h(x_{*}) = 0} \Rightarrow \underbrace{d^{T}\nabla_{x}^{2}L(x_{*}; \mu^{*})d \ge 0}_{\Leftrightarrow d^{T}x_{*} = 0} \Rightarrow \underbrace{d^{T}\nabla_{x}^{2}L(x_{*}; \mu^{*})d \ge 0}_{\Leftrightarrow d^{T}(A - \mu^{*}I)d \ge 0}$$

$$Opt = \min_{x \in \mathbb{R}^n} \left\{ f(x) = x^T A x : g(x) = 1 - x^T x = 0 \right\}$$
(P)

<u>Situation</u>: If x_* is optimal, then $\exists \mu^*$:

$$Ax_* = \mu^* x_* \qquad (A)$$

$$d^T x_* = 0 \Rightarrow d^T (A - \mu^* I) d \ge 0 \quad (B)$$

♣ (A) says that $x_* \neq 0$ is an eigenvector of A with eigenvalue μ^* ; in particular, we see that a symmetric matrix always has a real eigenvector

♣ (B) along with (A) says that $y^T (A - \mu^* I) y \ge$ 0 for all y. Indeed, every $y \in \mathbb{R}^n$ can be represented as $y = tx_* + d$ with $d^T x_* = 0$. We now have

$$y^{T}[A - \mu^{*}I]y = (tx_{*} + d)^{T}[A - \mu^{*}I](tx_{*} + d)$$

= $t^{2}x_{*}^{T} \underbrace{[A - \mu^{*}I]x_{*}}_{=0} + 2td^{T} \underbrace{d^{T}[A - \mu^{*}I]x_{*}}_{=0}$
+ $\underbrace{d^{T}[A - \mu^{*}I]d}_{\geq 0} \ge 0$

$$Opt = \min_{x \in \mathbb{R}^n} \left\{ f(x) = x^T A x : g(x) = 1 - x^T x = 0 \right\}$$
(P)

<u>Note:</u> In the case in question, Necessary Second Order Optimality conditions can be rewritten equivalently as $\exists \mu^*$:

$$[A - \mu^* I] x_* = 0 y^T [A - \mu^* I] y \ge 0 \,\forall y$$
 (*)

and are not only necessary, but also sufficient for feasible solution x_* to be globally optimal. To prove sufficiency, let x_* be feasible, and μ^* be such that (*) holds true. For every feasible solution x, one has

 $0 \leq x^{T}[A-\mu^{*}I]x = x^{T}Ax-\mu^{*}x^{T}x = x^{T}Ax-\mu^{*},$ whence $x^{T}Ax \geq \mu^{*}$. For $x = x_{*}$, we have $0 = x_{*}^{T}[A-\mu^{*}I]x_{*} = x_{*}^{T}Ax_{*}-\mu^{*}x_{*}^{T}x_{*} = x_{*}^{T}Ax_{*}-\mu^{*},$ whence $x_{*}^{T}Ax_{*} = \mu^{*}$. Thus, x_{*} is globally optimal for (P), and μ^{*} is the optimal value in (P). Extension: S-Lemma. Let A, B be symmetric matrices, and let B be such that

$$\exists \bar{x} : \bar{x}^T B \bar{x} > 0. \tag{*}$$

Then the inequality

$$x^T A x \ge 0 \tag{A}$$

is a consequence of the inequality

$$x^T B x \ge 0 \tag{B}$$

iff (A) is a "linear consequence" of (B): there exists $\lambda \ge 0$ such that

$$x^T [A - \lambda B] x \ge 0 \forall x \tag{C}$$

that is, (A) is a weighted sum of (B) (weight $\lambda \ge 0$) and identically true inequality (C). Sketch of the proof: The only nontrivial statement is that "If (A) is a consequence of (B), then t e exists $\lambda \ge 0$ such that ...". To prove this statement, assume that (A) is a consequence of (B). Situation:

$$\exists \bar{x} : \bar{x}^T B \bar{x} > 0; \ \underbrace{x^T B x \ge 0}_{(B)} \Rightarrow \underbrace{x^T A x \ge 0}_{(A)}$$

Consider optimization problem

Opt = min_x
$$\left\{ x^T A x : h(x) \equiv 1 - x^T B x = 0 \right\}$$
.

Problem is feasible by (*), and Opt ≥ 0 . Assume that an optimal solution x_* exists. Then, same as above, x_* is regular, and at x_* the Second Order Necessary condition holds true: $\exists \mu^*$:

$$\nabla_x \Big|_{x=x_*} \underbrace{ \begin{bmatrix} x^T A x + \mu^* [1 - x^T B x] \end{bmatrix} = 0 \Leftrightarrow [A - \mu^* B] x_* = 0}_{\substack{d^T \nabla_x \Big|_{x=x_*} h(x) = 0\\ \Leftrightarrow d^T B x_* = 0}} \Rightarrow d^T [A - \mu^* B] d \ge 0$$

We have $0 = x_*^T [A - \mu^* B] x_*$, that is, $\mu_* = Opt \ge 0$. Representing $y \in \mathbb{R}^n$ as $tx_* + d$ with $d^T B x_* = 0$ (that is, $t = x_*^T B y$), we get

$$y^{T}[A - \mu^{*}B]y = t^{2}x_{*}^{T} \underbrace{[A - \mu^{*}B]x_{*}}_{=0} + 2td^{T} \underbrace{[A - \mu^{*}B]x_{*}}_{=0} + \underbrace{d^{T}[A - \mu^{*}B]d}_{\geq 0} \ge 0,$$

Thus, $\mu^* \ge 0$ and $y^T [A - \mu^* B] y \ge 0$ for all y, Q.E.D.

Introduction to Optimization Algorithms

Goal: Approximate numerically solutions to Mathematical Programming problems

$$\min_{x} \left\{ f(x) : \begin{array}{l} g_{j}(x) \leq 0, \ j = 1, ..., m \\ h_{i}(x) = 0, \ i = 1, ..., k \end{array} \right\} \quad (P)$$

♣ Traditional MP algorithms to be considered in the Course do *not* assume the analytic structure of (*P*) to be known in advance (and do not know how to use the structure when it is known). These algorithms are *black-box-oriented*: when solving (*P*), method generates a sequence of *iterates* x_1 , x_2 ,... in such a way that x_{t+1} depends solely on local information of (*P*) gathered along the preceding *iterates* x_1 , ..., x_t .

Information on (P) obtained at x_t usually is comprised of the values and the first and the second derivatives of the objective and the constraints at x_t .

♠ <u>Note:</u> In optimization, there exist algorithms which do exploit problem's structure. Traditional methods of this type – Simplex method and its variations – do not go beyond Linear Programming and Linearly Constrained Quadratic Programming.

Recently, new efficient ways to exploit problem's structure were discovered (Interior Point methods). The resulting algorithms, however, do not go beyond Convex Programming. ♣ Except for very specific and relatively simple problem classes, like Linear Programming or Linearly Constrained Quadratic Programming, optimization algorithms cannot guarantee finding exact solution – local or global – in finite time. The best we can expect from these algorithms is *convergence* of approximate solutions generated by algorithms to the exact solutions.

♠ Even in the case when "finite" solution methods do exist (Simplex method in Linear Programming), no reasonable complexity bounds for these methods are known, therefore in reality the ability of a method to generate the exact solution in finitely many steps is neither necessary, nor sufficient to justify the method. Aside of Convex Programming, traditional optimization methods are unable to guarantee convergence to a globally optimal solution. Indeed, in the non-convex case there is no way to conclude from *local* information whether a given point is/is not globally optimal:



"looking" at problem around x', we get absolutely no hint that the trues global optimal solution is x''.

♠ In order to guarantee approximating global solution, it seems unavoidable to "scan" a dense set of the values of x in order to be sure that the globally optimal solution is not missed. Theoretically, such a possibility exists; however, the complexity of "exhaustive search" methods blows up exponentially with the dimension of the decision vector, which makes these methods completely impractical.

♣ Traditional optimization methods do *not* incorporate exhaustive search and, as a result, can*not* guarantee convergence to a global solution.

♠ A typical theoretical result on a traditional the optimization method as applied to a general (not necessary convex) problem sounds like:

Assume that problem (P) possesses the following properties:

• • •

Then the sequence of approximate solutions generated by method X is bounded, and all its limiting points are KKT points of the problem.

or

Assume that x_* is a nondegenerate local solution to (P). Then method X, started close enough to x_* , converges to x_* . Classification of MP Algorithms

♣ There are two major traditional classifications of MP algorithms:

♦ Classification by application fields, primarily into

- algorithms for unconstrained optimization
- algorithms for constrained optimization

♦ Classification by information used by the algorithms, primarily into

zero order methods which use only the values of the objective and the constraints

- first order methods (use both values and first order derivatives)
- second order methods (use values, firstand second order derivatives).

Rate of Convergence of MP Algorithms

There is a necessity to quantify the convergence properties of MP algorithms. Traditionally, this is done via asymptotical rate of convergence defined as follows:

<u>Step 1.</u> We introduce an appropriate *error measure* – a nonnegative function $\text{Error}_P(x)$ of approximate solution and of the problem we are solving which is zero exactly at the set X_* of solutions to (P) we intend to approximate.

Examples: (i) Distance to the set X_* :

$$\operatorname{Error}_{P}(x) = \inf_{x_{*} \in X_{*}} \|x - x_{*}\|_{2}$$

(ii) Residual in terms of the objective and the constraints

$$\mathsf{Error}_{P}(x) = \max \left[f(x) - \mathsf{Opt}(P), \\ [g_{1}(x)]_{+}, ..., [g_{m}(x)]_{+}, \\ |h_{1}(x)|, ..., |h_{k}(x)| \right]$$

<u>Step 2.</u> Assume that we have established *convergence* of our method, that is, we know that if x_t^* are approximate solutions generated by method as applied to a problem (*P*) from a given family, then

 $\operatorname{Error}_P(t) \equiv \operatorname{Error}_P(x_t^*) \to 0, t \to \infty$

We then roughly quantify the *rate* at which the sequence $\text{Error}_P(t)$ of nonnegative reals converges to 0. Specifically, we say that \diamond the method converges *sublinearly*, if the error goes to zero less rapidly than a geometric progression, e.g., as 1/t or $1/t^2$; \diamond the method converges *linearly*, if there exist

 $C < \infty$ and $q \in (0, 1)$ such that

 $\operatorname{Error}_{(P)}(t) \leq Cq^t$

q is called the *convergence ratio*. E.g.,

 $\operatorname{Error}_P(t) \asymp e^{-at}$

exhibits linear convergence with ratio e^{-a} . Sufficient condition for linear convergence with ratio $q \in (0, 1)$ is that

$$\frac{1}{t \to \infty} \frac{\mathsf{Error}_P(t+1)}{\mathsf{Error}_P(t)} < q$$

Othe method converges superlinearly, if the sequence of errors converges to 0 faster than every geometric progression:

$$\forall q \in (0, 1) \exists C : \mathsf{Error}_P(t) \leq Cq^t$$

For example,

$$\mathsf{Error}_P(t) \asymp e^{-at^2}$$

corresponds to superlinear convergence. <u>Sufficient condition</u> for superlinear convergence is

$$\lim_{t \to \infty} \frac{\operatorname{Error}_P(t+1)}{\operatorname{Error}_P(t)} = 0$$

 \diamondsuit the method exhibits convergence of order p>1, if

$$\exists C : \operatorname{Error}_P(t+1) \leq C (\operatorname{Error}_P(t))^p$$

Convergence of order 2 is called *quadratic*. For example,

$$\mathsf{Error}_P(t) = e^{-ap^t}$$

converges to 0 with order p.

Informal explanation: When the method converges, $\operatorname{Error}_P(t)$ goes to 0 as $t \to \infty$, that is, eventually the decimal representation of $\operatorname{Error}_P(t)$ has zero before the decimal dot and more and more zeros after the dot; the number of zeros following the decimal dot is called the number of accuracy digits in the corresponding approximate solution. Traditional classification of rates of convergence is based on how many steps, asymptotically, is required to add a new accuracy digit to the

is required to add a new accuracy digit to the existing ones.

 \diamond With *sublinear* convergence, the "price" of accuracy digit grows with the position of the digit. For example, with rate of convergence O(1/t) every new accuracy digit is 10 times more expensive, in terms of # of steps, than its predecessor.

stèp of the method adds a fixed number rof accuracy digits (for q not too close to 0, $r \approx 1 - q$);

♦ With *superlinear* convergence, every subsequent accuracy digit eventually becomes cheaper than its predecessor – the price of accuracy digit goes to 0 as the position of the digit grows. Equivalently, every additional step adds more and more accuracy digits.

 \diamond With convergence of order p > 1, the price of accuracy digit not only goes to 0 as the position k of the digit grows, but does it rapidly enough – in a geometric progression. Equivalently, eventually every additional step of the method *multiplies by* p the number of accuracy digits. ♣ With the traditional approach, the convergence properties of a method are the better the higher is the "rank" of the method in the above classification. Given a family of problems, traditionally it is thought that linearly converging on every problem of the family method is faster than a sublinearly converging, superlinearly converging method is faster than a linearly converging one, etc.

A <u>Note</u>: Usually we are able to *prove existence* of parameters C and q quantifying linear convergence:

 $\operatorname{Error}_P(t) \leq Cq^t$

or convergence of order p > 1:

 $\operatorname{Error}_P(t+1) \leq C(\operatorname{Error}_P(t))^p$,

but are unable to find numerical values of these parameters – they may depend on "unobservable" characteristics of a particular problem we are solving. As a result, traditional "quantification" of convergence properties is *qualitative* and *asymptotical*.

Solvable Case of MP – Convex Programming

♣ We have seen that as applied to general *MP programs*, optimization methods have a number of severe *theoretical* limitations, including the following major ones:

♦Unless exhaustive search (completely unrealistic in high-dimensional optimization) is used, there are no guarantees of approaching global solution

♦Quantification of convergence properties is of asymptotical and qualitative character. As a result, the most natural questions like:

We should solve problems of such and such structure with such and such sizes and the data varying in such and such ranges. How many steps of method X are sufficient to solve problems within such and such accuracy?

usually do not admit theoretically valid answers. ♣ In spite of their *theoretical* limitations, *in reality* traditional MP algorithms allow to solve many, if not all, MP problems of realworld origin, including those with many thousands variables and constraints.

♣ Moreover, there exists a "solvable case" when practical efficiency admits solid theoretical guarantees – the case of Convex Programming.

• Here is a typical "Convex Programming" result:

Assume we are solving a Convex Programming program

 $\begin{aligned} \text{Opt} &= \min_x \left\{ f(x) : g_j(x) \leq 0, \, j \leq m, |x_i| \leq 1, i \leq n \right\}. \end{aligned}$ where the objective and the constraints are normalized by the requirement

 $|x_i| \leq 1, i \leq n \Rightarrow |f(x)| \leq 1, |g_j(x)| \leq 1, j \leq m$

Given $\epsilon \in (0, 1)$, one can find an ϵ -solution x^{ϵ} to the problem:

 $\underbrace{|x_i^{\epsilon}| \leq 1}_{\forall i \leq n} \And \underbrace{g_j(x^{\epsilon}) \leq \epsilon}_{\forall j \leq m} \And f(x^{\epsilon}) - \mathsf{Opt} < \epsilon$

in no more than

$$2n^2 \ln\left(\frac{2n}{\epsilon}\right)$$

steps, with a single computation of the values and the first order derivatives of $f, g_1, ..., g_m$ at a point and 100(m+n)n additional arithmetic operations per step.
Line Search

Line Search is a common name for techniques for *one-dimensional* "simply constrained" optimization, specifically, for problems

 $\min_{x} \left\{ f(x) : a \le x \le b \right\},\$

where [a, b] is a given segment on the axis (sometimes, we shall allow for $b = +\infty$), and f is a function which is at least once continuously differentiable on (a, b) and is continuous at the segment [a, b] (on the ray $[a, \infty)$, if $b = \infty$).

♣ Line search is used, as a subroutine, in many algorithms for multi-dimensional optimization.

$$\min_{a \le x \le b} f(x) \tag{P}$$

A Zero-order line search. In zero-order line search one uses the values of the objective f in (P) and does not use its derivatives.

♠ To ensure well-posedness of the problem, assume that the objective is *unimodal*, that is, possesses a unique local minimizer x_* on [a, b].

Equivalently: There exists a unique point $x_* \in [a, b]$ such that f(x) strictly decreases on $[a, x_*]$ and strictly increases on $[x_*, b]$:



<u>Main observation</u>: Let f be unimodal on [a,b], and assume we know f(x'), f(x'') for certain x', x'' with

a < x' < x'' < b.

 \oint If $f(x'') \ge f(x')$, then f(x) > f(x'') for x > x'', so that the minimizer belongs to [a, x'']:



 \diamond Similarly, if f(x'') < f(x'), then f(x > f(x')when x < x', so that the minimizer belongs to [x', b].

♠ In both cases, two computations of f at x', x'' allow to reduce the initial "search domain" with a smaller one ([a, x''] or [x', b]). ♣ Choosing x', x'' so that they split $[a_0, b_0] = [a, b]$ into three equal segments, computing f(x'), f(x'') and comparing them to each other, we can build a new segment $[a_1, b_1] \subset [a_0, b_0]$ such that

 \diamond the new segment is a *localizer* – it contains the solution x_*

 \diamond the length of the new localizer is 2/3 of the length of the initial localizer $[a_0, b_0] = [a, b]$.

♠ On the new localizer, same as on the original one, the objective is unimodal, and we can iterate our construction.

♠ In $N \ge 1$ steps (2N computations of f), we shall reduce the size of localizer by factor $(2/3)^N$, that is, we get *linearly converging*, in terms of the argument, algorithm with the convergence ratio

 $q = \sqrt{2/3} = 0.8165...$

Can we do better ? - YES!

$$\begin{bmatrix} a_{t-1}, b_{t-1} \\ x'_t < x''_t \end{bmatrix} \Rightarrow f(x'_t), f(x''_t) \Rightarrow \begin{cases} [a_t, b_t] = [a_{t-1}, x''_t] \\ [a_t, b_t] = [x'_t, b_{t-1}] \end{cases}$$

A Observe that one of two points at which we compute f at a step becomes the endpoint of the new localizer, while the other one is an interior point of this localizer, and therefore we can use it as the one of two points where f should be computed at the next step!

With this approach, only the very first step costs 2 function evaluations, while the subsequent steps cost just 1 evaluation each! Let us implement the idea in such a way that all search points will divide respective localizers in a fixed proportion:

$$x'-a = b - x'' = \theta(b-a)$$

The proportion is given by the equation

$$\theta \equiv \frac{x'-a}{b-a} = \frac{x''-x'}{b-x'} \equiv \frac{1-2\theta}{1-\theta} \Rightarrow \theta = \frac{3-\sqrt{5}}{2}.$$

♣ We have arrived at golden search, where the search points x_{t-1} , x_t of step t are placed in the current localizer $[a_{t-1}, b_{t-1}]$ according to

$$\frac{x'-a}{b-a} = \frac{b-x''}{b-a} = \frac{3-\sqrt{5}}{2}$$

In this method, a step reduces the error (the length of localizer) by factor $1 - \frac{3-\sqrt{5}}{2} = \frac{\sqrt{5}-1}{2}$. The convergence ratio is about

$$\frac{\sqrt{5}-1}{2}\approx 0.6180...$$

 $\min_{x} \left\{ f(x) : a \le x \le b \right\},\$

First order line search: Bisection. Assume that f is differentiable on (a, b) and *strictly unimodal*, that is, it is unimodal, $x_* \in (a, b)$ and f'(x) < 0 for $a < x < x_*$, f'(x) > 0 for $x_* < x < b$.

Let both f and f' be available. In this case the method of choice in *Bisection*.

♠ <u>Main observation</u>: Given $x_1 \in [a,b] \equiv [a_0,b_0]$, let us compute $f'(x_1)$.

 \oint If $f'(x_1) > 0$, then, from strict unimodality, $f(x) > f(x_1)$ to the right of x_1 , thus, x_* belongs to $[a, x_1]$:



 \diamond Similarly, if $f'(x_1) \leq 0$, then $f(x) > f(x_1)$ for $x < x_1$, and x_* belongs to $[a, x_1]$.

♠ In both cases, we can replace the original localizer $[a, b] = [a_0, b_0]$ with a smaller localizer $[a_1, b_1]$ and then iterate the process.

In Bisection, the point x_t where at step $t f'(x_t)$ is computed, is the midpoint of $[a_{t-1}, b_{t-1}]$, so that every step reduces localizer's length by factor 2.

Clearly, Bisection converges linearly in terms of argument with convergence ratio 0.5:

 $a_t - x_* \le 2^{-t}(b_0 - a_0).$

A Many algorithms for multi-dimensional minimization which use Line Search as a sub-routine, in the following way:

 \diamondsuit given current iterate $x_t \in \mathbb{R}^n$, the algorithm defines a search direction $d_t \in \mathbb{R}^n$ which is a direction of decrease of f:

$$d_t^T \nabla f(x_t) < 0.$$

Then Line Search is invoked to minimize the one-dimensional function

$$\phi(s) = f(x_t + \gamma d_t)$$

over $\gamma \geq 0$; the resulting $\gamma = \gamma^t$ defines the stepsize along the direction d_t , so that the new iterate of the outer algorithm is

$$x_{t+1} = x_t + \gamma^t d_t.$$

♠ In many situations of this type, there is no necessity in exact minimization in γ ; an "essential" reduction in ϕ is sufficient. Standard way to define (and to achieve) "essential reduction" is given by Armijo's rule:

Let $\phi(\gamma)$ be continuously differentiable function of $\gamma \ge 0$ such that $\phi'(0) > 0$, and let $\epsilon \in (0,1), \eta > 1$ be parameters (popular choice is $\epsilon = 0.2$ and $\eta = 2$ or $\eta = 10$). We say that a stepsize $\gamma > 0$ is *appropriate*, if

$$\phi(\gamma) \le \phi(0) + \epsilon \gamma \phi'(0), \qquad (*)$$

and is *nearly maximal*, if η times larger step is *not* appropriate:

$$\phi(\eta\gamma) > \phi(0) + \epsilon \eta \gamma \phi'(0). \qquad (**)$$

A stepsize $\gamma > 0$ passes Armijo test (reduces ϕ "essentially"), if its is both appropriate and nearly maximal.

• Fact: Assume that ϕ is bounded below on the ray $\gamma > 0$. Then a stepsize passing Armijo rule does exist and can be found efficiently. Armijo-acceptable step $\gamma > 0$:

 $\phi(\gamma) \le \phi(0) + \epsilon \gamma \phi'(0)$ (*) $\phi(\eta \gamma) > \phi(0) + \epsilon \eta \gamma \phi'(0)$ (**)

Algorithm for finding Armijo-acceptable step: <u>Start:</u> Choose $\gamma_0 > 0$ and check whether it passes (*). If YES, go to Branch A, otherwise go to Branch B.

Branch A: γ_0 satisfies (*). Testing subsequently the values $\eta\gamma_0$, $\eta^2\gamma_0$, $\eta^3\gamma_0$,... of γ , stop when the current value for the first time violates (*); the preceding value of γ passes the Armijo test.

Branch B: γ_0 does not satisfy (*). Testing subsequently the values $\eta^{-1}\gamma_0$, $\eta^{-2}\gamma_0$, $\eta^{-3}\gamma_0$,... of γ , stop when the current value for the first time satisfies (*); this value of γ passes the Armijo test. ♣ <u>Validation of the algorithm</u>: It is clear that *if the algorithm terminates*, then the result indeed passes the Armijo test. Thus, all we need to verify is that the algorithm eventually terminates.

♦Branch A clearly is finite: here we test the inequality

$$\phi(\gamma) > \phi(0) + \epsilon \gamma \phi'(0)$$

along the sequence $\gamma_i = \eta^i \gamma_0 \rightarrow \infty$, and terminate when this inequality is satisfied for the first time. Since $\phi'(0) < 0$ and ϕ is below bounded, this indeed will eventually happen. \diamond Branch B clearly is finite: here we test the inequality

$$\phi(\gamma) \le \phi(0) + \epsilon \gamma \phi'(0)$$
 (*)

along a sequence $\gamma_i = \eta^{-i}\gamma_0 \rightarrow +0$ of values of γ and terminate when this inequality is satisfied for the first time. Since $\epsilon \in (0, 1)$ and $\phi'(0) < 0$, this inequality is satisfied for all small enough positive values of γ , since

$$\phi(\gamma) = \phi(0) + \gamma \Big[\phi'(0) + \underbrace{R(\gamma)}_{\to 0, \gamma \to +0} \Big].$$

For large *i*, γ_i definitely will be "small enough", thus, Branch B is finite.

Methods for Unconstrained Minimization

Unconstrained minimization problem is

$$f_* = \min_x f(x),$$

where f well-defined and continuously differentiable on the entire \mathbb{R}^n .

<u>Note:</u> Most of the constructions to be presented can be straightforwardly extended onto "essentially unconstrained case" where f is continuously differentiable on an <u>open</u> domain D in \mathbb{R}^n and is such that the level sets $\{x \in U : f(x) \leq a\}$ are closed.

$$f_* = \min_x f(x) \tag{P}$$

Gradient Descent

Gradient Descent is the simplest first order method for unconstrained minimization.

<u>The idea:</u> Let x be a current iterate which is not a critical point of f: $f'(x) \neq 0$. We have

$$f(x+th) = f(x) + th^T f'(x) + t ||h||_2 R_x(th)$$

[R_x(s) \rightarrow 0 as s \rightarrow 0]

Since $f'(x) \neq 0$, the unit antigradient direction $g = -f'(x)/||f'(x)||_2$ is a direction of decrease of f:

$$\frac{d}{dt}\Big|_{t=0} f(x+tg) = g^T f'(x) = -\|f'(x)\|_2$$

so that shift $x \mapsto x + tg$ along the direction glocally decreases f "at the rate" $||f'(x)||_2$. \land <u>Note:</u> As far as local rate of decrease is concerned, g is the best possible direction of decrease: for any other unit direction h, we have

$$\frac{d}{dt}\Big|_{t=0}f(x+th) = h^T f'(x) > -\|f'(x)\|_2.$$

In generic Gradient Descent, we update the current iterate x by a step from x in the antigradient direction which reduces the objective:

$$x_t = x_{t-1} - \gamma_t f'(x_{t-1}),$$

where γ_t are positive stepsizes such that

$$f'(x_{t-1}) \neq 0 \Rightarrow f(x_t) < f(x_{t-1}).$$

Standard implementations:
 Steepest GD:

$$\gamma_t = \underset{\gamma \ge 0}{\operatorname{argmin}} f(x_{t-1} - \gamma f'(x_{t-1}))$$

(slight idealization, except for the case of quadratic f)

Armijo GD: $\gamma_t > 0$ is such that

 $f(x_{t-1} - \gamma_t f'(x_{t-1}) \le f(x_{t-1}) - \epsilon \gamma_t \|f'(x_{t-1})\|_2^2$ $f(x_{t-1} - \eta \gamma_t f'(x_{t-1}) > f(x_{t-1}) - \epsilon \eta \gamma_t \|f'(x_{t-1})\|_2^2$

(implementable, provided that $f'(x_{t-1}) \neq 0$ and $f(x_{t-1} - \gamma f'(x_{t-1}))$ is below bounded when $\gamma \geq 0$) <u>Note:</u> By construction, GD is unable to leave a critical point:

$$f'(x_{t-1}) = 0 \Rightarrow x_t = x_{t-1}.$$

Global Convergence Theorem: Assume that the level set of f corresponding to the starting point x_0 :

 $G = \{x : f(x) \le f(x_0)\}$

is compact, and f is continuously differentiable in a neighbourhood of G. Then for both SGD and AGD:

 \diamond the trajectory x_0, x_1, \dots of the method, started at x_0 , is well-defined and never leaves G (and thus is bounded);

♦ the method is monotone:

 $f(x_0) \ge f(x_1) \ge \dots$

and inequalities are strict, unless method reaches a critical point x_t , so that $x_t = x_{t+1} = x_{t+2} = \dots$

 \diamond Every limiting point of the trajectory is a critical point of f.

Sketch of the proof: 1^0 . If $f'(x_0) = 0$, the method never leaves x_0 , and the statements are evident. Now assume that $f'(x_0) \neq 0$. Then the function $\phi_0(\gamma) = f(x_0 - \gamma f'(x_0))$ is below bounded, and the set $\{\gamma \ge 0 : \phi_0(\gamma) \le \phi_0(0)\}$ is compact along with G, so that $\phi_0(\gamma)$ achieves its minimum on the ray $\gamma \ge 0$, and $\phi'_0(0) < 0$. It follows that the first step of GD is well-defined and $f(x_1) < f(x_0)$. The set $\{x : f(x) \le f(x_1)\}$ is a closed subset of G and thus is compact, and we can repeat our reasoning with x_1 in the role of x_0 , etc. We conclude that the trajectory is well-defined, never leaves G and the objective is strictly decreased, unless a critical point is reached.

 2^{0} . "all limiting points of the trajectory are critical points of f":

<u>Fact</u>: Let $x \in G$ and $f'(x) \neq 0$. Then there exists $\epsilon > 0$ and a neighbourhood U of x such that for every $x' \in U$ the step $x' \to x'_+$ of the method from x' reduces f by at least ϵ .

Given Fact, let x be a limiting point of $\{x_i\}$; assume that $f'(x) \neq 0$, and let us lead this assumption to contradiction. By Fact, there exists a neighbourhood U of x such that

$$x_i \in U \Rightarrow f(x_{i+1}) \leq f(x_i) - \epsilon.$$

Since the trajectory visits U infinitely many times and the method is monotone, we conclude that $f(x_i) \rightarrow -\infty$, $i \rightarrow \infty$, which is impossible, since G is compact, so that f is below bounded on G.

Limiting points of Gradient Descent

Lunder assumptions of Global Convergence Theorem, limiting points of GD exist, and all of them are critical points of f. What kind of limiting points could they be?

A <u>nondegenerate maximizer</u> of f cannot be a limiting point of GD, unless the method is started at this maximizer.

A saddle point of f is "highly unlikely" candidate to the role of a limiting point. <u>Practical</u> experience says that limiting points are local minimizers of f.

A <u>nondegenerate global minimizer</u> x_* of f, if any, as an "attraction point" of GD: when starting close enough to this minimizer, the method converges to x_* . In general, we cannot guarantee more than convergence to the set of critical points of f. A natural error measure associated with this set is

$$\delta^2(x) = \|f'(x)\|_2^2.$$

♠ Definition: Let U be an open subset of \mathbb{R}^n , $L \ge 0$ and f be a function defined on U. We say that f is $C^{1,1}(L)$ on U, if f is continuously differentiable in U with locally Lipschitz continuous, with constant L, gradient:

 $[x,y] \in U \Rightarrow ||f'(x) - f'(y)||_2 \leq L||x - y||_2.$ We say that f is $C^{1,1}(L)$ on a set $Q \subset \mathbb{R}^n$, if there exists an open set $U \supset Q$ such that fis $C^{1,1}(L)$ on U.

<u>Note</u>: Assume that f is twice continuously differentiable on U. Then f is $C^{1,1}(L)$ on U iff the norm of the Hessian of f does not exceed L:

 $\forall (x \in U, d \in \mathbb{R}^n) : |d^T f''(x)d| \le L ||d||_2^2.$

<u>Theorem.</u> In addition to assumptions of Global Convergence Theorem, assume that f is $C^{1,1}(L)$ on $G = \{x : f(x) \le f(x_0)\}$. Then \Diamond For SGD, one has

$$\min_{0 \le \tau \le t} \delta^2(x_\tau) \le \frac{2[f(x_0) - f_*]L}{t+1}, t = 0, 1, 2, \dots$$

 \Diamond For AGD, one has

 $\min_{0 \le \tau \le t} \delta^2(x_{\tau}) \le \frac{\eta}{2\epsilon(1-\epsilon)} \cdot \frac{[f(x_0) - f_*]L}{t+1}, \ t = 0, 1, 2, \dots$

Lemma. For $x \in G$, $0 \le s \le 2/L$ one has

$$x - sf'(x) \in G$$
(1)
$$f(x - sf'(x)) \leq f(x) - \delta^2(x)s + \frac{L\delta^2(x)}{2}s^2,$$
(2)

There is nothing to prove when $g \equiv -f'(x) = 0$. Let $g \neq 0$, $s_* = \max\{s \ge 0 : x + sg \in G\}$, $\delta^2 = \delta^2(x) = g^T g$. The function

$$\phi(s) = f(x - sf'(x)) : [0, s_*] \to \mathbb{R}$$

is continuously differentiable and satisfies

(a)
$$\phi'(0) = -g^T g \equiv -\delta^2$$
; (b) $\phi(s_*) = f(x_0)$
(c) $|\phi'(s) - \phi'(0)| = |g^T [f'(x + sg) - f'(x)]| \le Ls\delta^2$
Therefore $\phi(s) \le \phi(0) - \delta^2 s + \frac{L\delta^2}{2}s^2$ (*)

which is (2). Indeed, setting

$$\theta(s) = \phi(s) - [\phi(0) - \delta^2 s + \frac{L\delta^2}{2}s^2],$$

we have

$$\theta(0) = 0, \theta'(s) = \phi'(s) - \phi'(0) - Ls\delta^2 \underbrace{\leq}_{\text{by }(c)} 0.$$

By (*) and (b), we have

$$f(x_0) \le \phi(0) - \delta^2 s_* + \frac{L\delta^2}{2} s_*^2 \le f(x_0) - \delta^2 s_* + \frac{L\delta^2}{2} s_*^2$$

 $\Rightarrow s_* \ge 2/L$

Lemma \Rightarrow **Theorem:** <u>SGD:</u> By Lemma, we have

$$f(x_{t}) - f(x_{t+1}) = f(x_{t}) - \min_{\gamma \ge 0} f(x_{t} - \gamma f'(x_{t}))$$

$$\ge f(x_{t}) - \min_{0 \le s \le 2/L} \left[f(x_{t}) - \delta^{2}(x_{t})s + \frac{L\delta^{2}(x_{t})}{2}s^{2} \right]$$

$$= \frac{\delta^{2}(x_{t})}{2L}$$

$$\Rightarrow f(x_{0}) - f_{*} \ge \sum_{\tau=0}^{t} \left[f(x_{\tau}) - f(x_{\tau+1}) \right] \ge \sum_{\tau=0}^{t} \frac{\delta^{2}(x_{\tau})}{2L}$$

$$\ge (t+1) \min_{0 \le \tau \le t} \delta^{2}(x_{\tau})$$

$$\Rightarrow \min_{0 \le \tau \le t} \delta^{2}(x_{\tau}) \le \frac{2L(f(x_{0}) - f_{*})}{t+1}$$

<u>AGD:</u> <u>Claim:</u> $\gamma_{t+1} > \frac{2(1-\epsilon)}{L\eta}$. Indeed, otherwise by Lemma

$$f(x_{t} - \gamma_{t}\eta f'(x_{t}))$$

$$\leq f(x_{t}) - \gamma_{t+1}\eta\delta^{2}(x_{t}) + \frac{L\delta^{2}(x_{t})}{2}\eta^{2}\gamma_{t+1}^{2}$$

$$= f(x_{t}) - \left[1 - \frac{L}{2}\eta\gamma_{t+1}\right]\eta\gamma_{t+1}\delta^{2}(x_{t})$$

$$\stackrel{\geq \epsilon}{\leq} f(x_{t}) - \epsilon\eta\gamma_{t+1}\delta^{2}(x_{t})$$

which is impossible.

We have seen that $\gamma_{t+1} > \frac{2(1-\epsilon)}{L\eta}$. By Armijo rule,

 $f(x_t) - f(x_{t+1}) \ge \epsilon \gamma_{t+1} \delta^2(x_t) \ge \frac{2\epsilon(1-\epsilon)}{L\eta} \delta^2(x_t);$ the rest of the proof is as for SGD. & <u>Convex case</u>. In addition to assumptions of Global Convergence Theorem, assume that fis convex.

All critical points of a convex function are its global minimizers

 \Rightarrow In Convex case, SGD and AGD converge to the set of global minimizers of $f: f(x_t) \rightarrow f_*$ as $t \rightarrow \infty$, and all limiting points of the trajectory are global minimizers of f.

♠ In Convex $C^{1,1}(L)$ case, one can quantify the global rate of convergence in terms of the residual $f(x_t) - f_*$:

<u>Theorem.</u> Assume that the set $G = \{x : f(x) \le f(x_0)\}$ is convex compact, f is convex on G and $C^{1,1}(L)$ on this set. Consider AGD, and let $\epsilon \ge 0.5$. Then the trajectory of the method converges to a global minimizer x_* of f, and

$$f(x_t) - f_* \le \frac{\eta L \|x_0 - x_*\|_2^2}{4(1 - \epsilon)t}, \ t = 1, 2, \dots$$

 $\begin{array}{l} \clubsuit \ \underline{\text{Definition:}} \ \text{Let} \ M \ \text{be a convex set in} \ \mathbb{R}^n \\ \text{and} \ 0 < \ell \leq L < \infty. \ \text{A function} \ f \ \text{is called} \\ \underline{\text{strongly convex}}, \ \text{with parameters} \ \ell, L, \ \text{on} \ M, \\ \text{if} \end{array}$

 $[x - y]^T [f'(x) - f'(y)] \ge \ell ||x - y||_2^2. \quad (*)$ The ratio $Q_f = L/\ell$ is called <u>condition number</u> of f.

 \bigstar <u>Comment</u>: If f is $C^{1,1}(L)$ on a convex set M, then

 $x, y \in M \Rightarrow |f(y) - [f(x) + (y - x)^T f'(x)]| \le \frac{L}{2} ||x - y||_2^2.$

If f satisfies (*) on a convex set M, then

 $\forall x, y \in M : f(y) \ge f(x) + (y - x)^T f'(x) + \frac{\ell}{2} ||y - x||_2^2.$

In particular, f is convex on M. \Rightarrow A strongly convex, with parameters ℓ, L , function f on a convex set M satisfies the relation

$$\forall x, y \in M : f(x) + (y - x)^T f'(x) + \frac{\ell}{2} ||y - x||_2^2 \\ \leq f(y) \leq f(x) + (y - x)^T f'(x) + \frac{L}{2} ||y - x||_2^2$$

<u>Note</u>: Assume that f is twice continuously differentiable in a neighbourhood of a convex set M. Then f is (ℓ, L) -strongly convex on M iff for all $x \in M$ and all $d \in \mathbb{R}^n$ one has

$$\begin{split} \ell \|d\|_2^2 &\leq d^T f''(x) d \leq L \|d\|_2^2 \\ & \updownarrow \\ \lambda_{\min}(f''(x)) \geq \ell, \ \lambda_{\max}(f''(x)) \leq L. \end{split}$$

In particular,

♠ A quadratic function

$$f(x) = \frac{1}{2}x^T A x - b^T x + c$$

with positive definite symmetric matrix A is strongly convex with the parameters $\ell = \lambda_{\min}(A)$, $L = \lambda_{\max}(A)$ on the entire space.

GD in strongly convex case.

<u>Theorem.</u> In the strongly convex case, AGD exhibits linear global rate of convergence. Specifically, let the set $G = \{x : f(x) \leq f(x_0)\}$ be closed and convex and f be strongly convex, with parameters ℓ, L , on Q. Then $x \diamondsuit G$ is compact, and the global minimizer x_* of f exists and is unique;

AGD with $\epsilon \geq 1/2$ converges linearly to x_* :

$$||x_t - x_*||_2 \le \theta^t ||x_0 - x_*||_2$$

$$\theta = \sqrt{\frac{Q_f - (2 - \epsilon^{-1})(1 - \epsilon)\eta^{-1}}{Q_f + (\epsilon^{-1} - 1)\eta^{-1}}} = 1 - O(Q_f^{-1}).$$

Besides this,

$$f(x_t) - f_* \le \theta^{2t} Q_f[f(x_0) - f_*].$$

SGD in Strongly convex quadratic case. Assume that $f(x) = \frac{1}{2}x^TAx - b^Tx + c$ is a strongly convex quadratic function: $A = A^T \succ 0$. In this case, SGD becomes implementable and is given by the recurrence

$$g_t = f'(x_t) = Ax_t - b$$

$$\gamma_{t+1} = \frac{g_t^T g_t}{g_t^T A g_t}$$

$$x_{t+1} = x_t - \gamma_{t+1} g_t$$

and guarantees that

$$\underbrace{f(x_{t+1}) - f_*}_{E_{t+1}} \le \left[1 - \frac{(g_t^T g_t)^2}{[g_t^T A g_t][g_t^T A^{-1} g_t]}\right] E_t \le \left(\frac{Q_f - 1}{Q_f + 1}\right)^2 E_t$$

whence

$$f(x_t) - f_* \le \left(\frac{Q_f - 1}{Q_f + 1}\right)^{2t} [f(x_0) - f_*], t = 1, 2, \dots$$

<u>Note:</u> If we know that SGD converges to a *nondegenerate* local minimizer x_* of f, <u>then</u>, under mild regularity assumptions, the *asymptotical* behaviour of the method will be as if f were the strongly convex quadratic form

$$f(x) = \text{const} + \frac{1}{2}(x - x_*)^T f''(x_*)(x - x_*).$$





SGD as applied to quadratic form with $Q_f = 1000$ $\Diamond f(x_0) = 2069.4$, $f(x_{999}) = 0.0232$

Summary on Gradient Descent:

♦ Under mild regularity and boundedness assumptions, both SGD and AGD converge the set of critical points of the objective.

In the case of $C^{1,1}(L)$ -smooth objective, the methods exhibit non-asymptotical O(1/t)-rate of convergence in terms of the error measure $\delta^2(x) = \|f'(x)\|_2^2$.

♦Under the same regularity assumptions, in <u>Convex</u> case the methods converge to the set of glomal minimizers of the objective.

In convex $C^{1,1}(L)$ -case, AGD exhibits nonasymptotical O(1/t) rate of convergence in terms of the residual in the objective $f(x) - f_*$ \Diamond In *Strongly convex case*, AGD exhibits nonasymptotical linear convergence in both the residual in terms of the objective $f(x) - f_*$ and the distance in the argument $||x - x_*||_2$. The convergence ratio is $1 - O(1/Q_f)$, where Q_f is the all condition number of the objective. In other words, to get extra accuracy digit, it takes $O(Q_f)$ steps.

Good news on GD:

♠ Simplicity

Reasonable global convergence properties under mild assumptions on the function to be minimized.

Drawbacks of GD:

 \diamond You are solving the problem $\min_{x} f(x)$ by GD, starting with $x_0 = 0$, Your first search point will be

$$x_1 = -\gamma_1 f'(0).$$

 \diamond I solve the same problem, but in new variables y: x = Ay. My problem is $\min_{y} g(y)$, g(y) = f(Ax), and start with $y_0 = 0$. My first search point will be

$$y_1 = -\hat{\gamma}_1 g'(0) = -\hat{\gamma}_1 A^T f'(0).$$

In x-variables, my search point will be

$$\hat{x}_1 = Ay_1 = -\hat{\gamma}_1 A A^T f'(0)$$

If AA^T is not proportional to the unit matrix, my search point will, in general, be different from yours!
• "Frame-dependence" is common drawback of nearly all *first order* optimization methods, and this is what makes their rate of convergence, even under the most favourable case of strongly convex objective, sensitive to the condition number of the problem. GD is "hyper-sensitive" to the condition number: When minimizing strongly convex function f, the convergence ratio of GD is $1 - O(1/Q_f^{1/2})$, while for better methods it is $1 - O(1/Q_f^{1/2})$.

Consider unconstrained problem

$\min_{x} f(x)$

with *twice* continuously differentiable objective. Assuming second order information available, we approximate f around a current iterate x by the second order Taylor expansion:

$$f(y) \approx f(x) + (y-x)^T f'(x) + \frac{(y-x)^T f''(x)(y-x)}{2}$$

In the Newton method, the new iterate is the minimizer of this quadratic approximation. *If exists*, the minimizer is given by

$$\nabla_{y}[f(x) + (y - x)^{T} f'(x) + \frac{(y - x)^{T} f''(x)(y - x)}{2}] = 0$$

$$\Leftrightarrow f''(x)(y - x) = -f'(x)$$

$$\Leftrightarrow y = x - [f''(x)]^{-1} f'(x)$$

We have arrived at the **Basic Newton method**

$$x_{t+1} = x_t - [f''(x_t)]^{-1} f'(x_t)$$

(step t is undefined when the matrix $f''(x_t)$ is singular).

$$x_{t+1} = x_t - [f''(x_t)]^{-1}f'(x_t)$$

♠ <u>Alternative motivation</u>: We seek for a solution to the Fermat equation

$$f'(x) = 0;$$

given current approximate x_t to the solution, we linearize the left hand side around x_t , thus arriving at the linearized Fermat equation

$$f'(x_t) + f''(x_t)[x - x_t] = 0$$

and take the solution to this equation, that is, $x_t - [f''(x_t)]^{-1}f'(x_t)$, as our new iterate.

$$x_{t+1} = x_t - [f''(x_t)]^{-1} f'(x_t)$$
 (Nwt)

Theorem on Local Quadratic Convergence: Let x_* be a nondegenerate local minimizer of f, so that $f''(x_*) \succ 0$, and let f be three times continuously differentiable in a neighbourhood of x_* . Then the recurrence (Nwt), started close enough to x_* , is well-defined and converges to x_* quadratically.

Proof: 1⁰. Let U be a ball centered at x_* where the third derivatives of f are bounded. For $y \in U$ one has

$$\begin{aligned} \|\nabla f(y) - \nabla^2 f(y)(x_* - y)\|_2 \\ &\equiv \|\nabla f(x_*) - \nabla f(y) - \nabla^2 f(y)(x_* - y)\|_2 \\ &\leq \beta_1 \|y - x_*\|_2^2 \end{aligned}$$
(1)

2⁰. Since f''(x) is continuous at $x = x_*$ and $f''(x_*)$ is nonsingular, there exists a ball $U' \subset U$ centered at x_* such that

$$y \in U' \Rightarrow \|[f''(y)]^{-1}\| \le \beta_2.$$
(2)

<u>Situation</u>: There exists a r > 0 and positive constants β_1, β_2 such that

$$\begin{array}{rcl} \|y - x_*\| < r & \Rightarrow \\ (a) & \|\nabla f(y) - \nabla^2 f(y)(x_* - y)\|_2 & \leq & \beta_1 \|y - x_*\|_2^2 \\ (b) & & \|[f''(y)]^{-1}\| & \leq & \beta_2 \end{array}$$

3⁰. Let an iterate x_t of the method be close to x_* :

$$x_t \in V = \{x : \|x - x_*\|_{\mathfrak{O}} \le \rho \equiv \min[\frac{1}{2\beta_1\beta_2}, r]\}.$$

We have

$$\begin{aligned} \|x_{t+1} - x_*\| &= \|x_t - x_* - [f''(x_t)]^{-1} f'(x_t)\|_2 \\ &= \|\left[[f''(x_t)]^{-1} [f''(x_t)(x_t - x_*) - f'(x_t)] \right] \|_2 \\ &\leq \beta_1 \beta_2 \|x_t - x_*\|_2^2 \leq 0.5 \|x_t - x_*\|_2 \end{aligned}$$

We conclude that the method remains welldefined after step t, and converges to x_* quadratically. A remarkable property of Newton method is affine invariance ("frame independence"): Let x = Ay + b be invertible affine change of variables. Then

 $\begin{array}{ccc} f(x) & \Leftrightarrow & g(y) = f(Ay+b) \\ \bar{x} = A\bar{y} + b & \Leftrightarrow & \bar{y} \end{array}$

$$\begin{aligned} \bar{y}_{+} &= \bar{y} - [g''(\bar{y})]^{-1}g'(\bar{y}) \\ &= \bar{y} - [A^{T}f''(\bar{x})A]^{-1}[A^{T}f'(\bar{x})] \\ &= \bar{y} - A^{-1}[f''(\bar{x})]^{-1}f'(\bar{x}) \\ &\Rightarrow A\bar{y}_{+} + b = [A\bar{y} + b] - [f''(\bar{x})]^{-1}f'(\bar{x}) \\ &= \bar{x} - [f''(\bar{x})]^{-1}f'(\bar{x}) \end{aligned}$$

Difficulties with Basic Newton method.
The Basic Newton method

$$x_{t+1} = x_t - [f''(x_t)]^{-1} f'(x_t),$$

started close enough to nondegenerate local minimizer x_* of f, converges to x_* quadratically. However,

 \diamond Even for a nice strongly convex f, the method, started not too close to the (unique) local=global minimizer of f, may diverge:

$$f(x) = \sqrt{1 + x^2} \Rightarrow x_{t+1} = -x_t^3.$$

⇒ when $|x_0| < 1$, the method converges quadratically (even at order 3) to $x_* = 0$; when $|x_0| > 1$, the method rapidly diverges... ♦ When f is not strongly convex, the Newton direction

$$-[f''(x)]^{-1}f'(x)$$

can be undefined or fail to be a direction of decrease of f...

As a result of these drawbacks, one needs to modify the Basic Newton method in order to ensure global convergence. Modifications include:

♦Incorporating line search

 \diamond Correcting Newton direction when it is undefined or is not a direction of decrease of f.

♣ Incorporating linesearch: Assume that the level set $G = \{x : f(x) \le f(x_0)\}$ is closed and convex, and f is strongly convex on G. Then for $x \in G$ the Newton direction

$$e(x) = -[f''(x)]^{-1}f'(x)$$

is a direction of decrease of f, except for the case when x is a critical point (or, which is the same in the strongly convex case, global minimizer) of f:

$$f'(x) \neq 0 \Rightarrow$$

$$e^{T}(x)f'(x) = -[f'(x)]^{T} \underbrace{[f''(x)]^{-1}}_{\succ 0} f'(x) < 0.$$

In Line Search version of Newton method, one uses e(x) as a search direction rather than the displacement:

$$x_{t+1} = x_t + \gamma_{t+1} e(x_t) = x_t - \gamma_{t+1} [f''(x_t)]^{-1} f'(x_t),$$

where $\gamma_{t+1} > 0$ is the stepsize given by exact
minimization of f in the Newton direction or
by Armijo linesearch.

<u>Theorem</u>: Let the Level set $G = \{x : f(x) \le f(x_0)\}$ be convex and compact, and f be strongly convex on G. Then Newton method with the Steepest Descent or with the Armijo linesearch converges to the unique global minimizer of f.

With proper implementation of the linesearch, convergence is quadratic.

Newton method: Summary

♦ Good news: Quadratic asymptotical convergence, provided we manage to bring the trajectory close to a nondegenerate local minimizer

♦ Bad news:

 relatively high computational cost, coming from the necessity to compute and to invert the Hessian matrix

— necessity to "cure" the method in the non-strongly-convex case, where the Newton direction can be undefined or fail to be a direction of decrease...

Modifications of the Newton method

Modifications of the Newton method are aimed at overcoming its shortcomings (difficulties with nonconvex objectives, relatively high computational cost) while preserving its major advantage – rapid asymptotical convergence. There are three major groups of modifications:

Modified Newton methods based on secondorder information

Modifications based on first order information:

— conjugate gradient methods

— quasi-Newton methods

♠ All modifications of Newton method exploit a natural Variable Metric idea.

When speaking about GD, it was mentioned that the method

$$x_{t+1} = x_t - \gamma_{t+1} \underbrace{BB^T}_{A^{-1} \succ 0} f'(x_t) \qquad (*)$$

with nonsingular matrix B has the same "right to exist" as the Gradient Descent

$$x_{t+1} = x_t - \gamma_{t+1} f'(x_t);$$

the former method is nothing but the GD as applied to

$$g(y) = f(By).$$

$$x_{t+1} = x_t - \gamma_{t+1} A^{-1} f'(x_t) \qquad (*)$$

Equivalently: Let A be a positive definite symmetric matrix. We have exactly the same reason to measure the "local directional rate of decrease" of f by the quantity

$$\frac{d^T f'(x)}{\sqrt{d^T d}} \tag{a}$$

as by the quantity

$$\frac{d^T f'(x)}{\sqrt{d^T A d}} \tag{b}$$

 \diamond When choosing, as the current search direction, the direction of steepest decrease in terms of (*a*), we get the anti-gradient direction -f'(x) and arrive at GD.

 \diamond When choosing, as the current search direction, the direction of steepest decrease in terms of (*b*), we get the "scaled anti-gradient direction" $-A^{-1}f'(x)$ and arrive at "scaled" GD (*).

We have motivated the scaled GD

$$x_{t+1} = x_t - \gamma_{t+1} A^{-1} f'(x_t) \qquad (*)$$

Why not to take one step ahead by consider a generic Variable Metric algorithm

$$x_{t+1} = x_t - \gamma_{t+1} A_{t+1}^{-1} f'(x_t) \qquad (VM)$$

with "scaling matrix" $A_{t+1} \succ 0$ varying from step to step?

• <u>Note</u>: When $A_{t+1} \equiv I$, (VM) becomes the generic Gradient Descent;

When f is strongly convex and $A_{t+1} = f''(x_t)$, (VM) becomes the generic Newton method...

♠ <u>Note</u>: When x_t is *not* a critical point of f, the search direction $d_{t+1} = -A_{t+1}^{-1}f'(x_t)$ is a direction of decrease of f:

$$d_{t+1}^T f'(x_t) = -[f'(x_t)]^T A_{t+1}^{-1} f'(x_t) < 0.$$

Thus, we have no conceptual difficulties with *monotone* linesearch versions of (VM)...

$$x_{t+1} = x_t - \gamma_{t+1} A_{t+1}^{-1} f'(x_t) \qquad (VM)$$

A It turns out that Variable Metric methods possess good global convergence properties: <u>Theorem</u>: Let the level set $G = \{x : f(x) \le f(x_0)\}$ be closed and bounded, and let f be twice continuously differentiable in a neighbourhood of G.

Assume, further, that the policy of updating the matrices A_t ensures their *uniform positive definiteness and boundedness*:

$\exists 0 < \ell \leq L < \infty : \ell I \preceq A_t \preceq L I \ \forall t.$

Then for both the Steepest Descent and the Armijo versions of (VM) started at x_0 , the trajectory is well-defined, belongs to G (and thus is bounded), and f strictly decreases along the trajectory unless a critical point of f is reached. Moreover, all limiting points of the trajectory are critical points of f.

♣ Implementation via Spectral Decomposition: ♦ Given x_t , compute $H_t = f''(x_t)$ and then find spectral decomposition of H_t :

$$H_t = V_t \mathsf{Diag}\{\lambda_1, ..., \lambda_n\} V_t^T$$

 \diamond Given once for ever chosen tolerance $\delta > 0$, set

$$\widehat{\lambda}_i = \max[\lambda_i, \delta]$$

and

$$A_{t+1} = V_t \mathsf{Diag}\{\hat{\lambda}_1, ..., \hat{\lambda}_n\} V_t^T$$

<u>Note</u>: The construction ensures uniform positive definiteness and boundedness of $\{A_t\}_t$, provided the level set $G = \{x : f(x) \le f(x_0)\}$ is compact and f is twice continuously differentiable in a neighbourhood of G.



$$A_{t+1} = \epsilon_t I + H_t,$$

where $\epsilon_t \geq 0$ is chosen to ensure that $A_{t+1} \succeq \delta I$ with once for ever chosen $\delta > 0$.

 $\diamond \epsilon_t$ is found by Bisection as applied to the problem

$$\min\left\{\epsilon : \epsilon \ge 0, H_t + \epsilon I \succeq \delta I\right\}$$

Objection requires to check whether the condition

$$H_t + \epsilon I \succ \delta I \Leftrightarrow H_t + (\epsilon - \delta) I \succ 0$$

holds true for a given value of ϵ , and the underlying test comes from Choleski decomposition.

\clubsuit <u>Choleski Decomposition</u>. By Linear Algebra, a symmetric matrix P is \succ 0 iff

$$P = DD^T \tag{(*)}$$

with lower triangular nonsingular matrix D. When Choleski Decomposition (*) exists, it can be found by a simple algorithm as follows:

 \Diamond Representation (*) means that

$$p_{ij} = d_i d_j^T,$$

where

$$d_i = (d_{i1}, d_{i2}, ..., d_{ii}, 0, 0, 0, 0, ..., 0)$$

$$d_j = (d_{j1}, d_{j2}, ..., d_{ji}, ..., d_{jj}, 0, ..., 0)$$

are the rows of D.

 \diamond In particular, $p_{i1} = d_{11}d_{i1}$, and we can set $d_{11} = \sqrt{p_{11}}$, $d_{i1} = p_{i1}/d_{11}$, thus specifying the first column of D.

 \diamond Further, $p_{22} = d_{21}^2 + d_{22}^2$, whence $d_{22} = \sqrt{p_{22} - d_{21}^2}$. After we know d_{22} , we can find all remaining entries in the second column of D from the relations

$$p_{i2} = d_{i1}d_{21} + d_{i2}d_{22} \Rightarrow d_{i2} = \frac{p_{i2} - d_{i1}d_{21}}{d_{22}}, i > 2.$$

 \diamond We proceed in this way: after the first (k-1) columns in D are found, we fill the k-th column according to

$$d_{kk} = \sqrt{p_{kk} - d_{k1}^2 - d_{k2}^2 - \dots - d_{k,k-1}^2}$$

$$d_{ik} = \frac{p_{ik} - d_{i1}d_{k1} - \dots - d_{i,k-1}d_{k,k-1}}{d_{kk}}, i > k.$$

 \blacklozenge The outlined process either results in the required D, or terminates when you cannot carry out current pivot, that is, when

$$p_{kk} - d_{k1}^2 - d_{k2}^2 - \ldots - d_{k,k-1}^2 \le 0$$

This "bad termination" indicates that P is not positive definite.

The outlined *Choleski Algorithm* allows to find the Choleski decomposition, if any, in $\approx \frac{n^3}{6}$ a.o. It is used routinely to solve linear systems

$$Px = p \tag{S}$$

with $P \succ 0$. To solve the system, one first computes the Choleski decomposition

$$P = DD^T$$

and then solves (S) by two *back-substitutions*

$$b \mapsto y : Dy = b, \ y \mapsto x : D^T x = y,$$

that is, by solving two triangular systems of equations (which takes just $O(n^2)$ a.o.

Another application of the algorithm (e.g., in Levenberg-Marquardt method) is to check positive definiteness of a symmetric matrix. <u>Note:</u> The Levenberg-Marquardt method produces uniformly positive definite bounded sequence $\{A_t\}$, provided that the set $G = \{x : f(x) \le f(x_0)\}$ is compact and f is twice continuously differentiable in a neighbourhood of G. A The "most practical" implementation of Modified Newton Method is based on running the Choleski decomposition as applied to $H_t = f''(x_t)$. When in course of this process the current pivot (that is, specifying d_{kk}) becomes impossible or results in $d_{kk} < \delta$, one increases the corresponding diagonal entry in H_t until the condition $d_{kk} = \delta$ is met.

With this approach, one finds a diagonal correction of H_t which makes the matrix "well positive definite" and ensures uniform positive definiteness and boundedness of the resulting sequence $\{A_t\}$, provided that the set $G = \{x : f(x) \le f(x_0)\}$ is compact and f is twice continuously differentiable in a neighbourhood of G.

Conjugate Gradient methods

Consider a problem of minimizing a positive definite quadratic form

$$f(x) = \frac{1}{2}x^T H x - b^T x + c$$

Here is a "conceptual algorithm" for minimizing f, or, which is the same, for solving the system

$$Hx = b$$
:

Given starting point x_0 , let $g_0 = f'(x_0) = Hx_0 - b$, and let

 $E_k = \operatorname{Lin}\{g_0, Hg_0, H^2g_0, \dots, H^{k-1}g_0\},$

and

$$x_k = \underset{x \in x_0 + E_k}{\operatorname{argmin}} f(x).$$

<u>Fact I</u>: Let k_* be the smallest integer k such that $E_{k+1} = E_k$. Then $k_* \leq n$, and x_{k_*} is the unique minimizer of f on \mathbb{R}^n <u>Fact II</u>: One has

$$f(x_k) - \min_x f(x) \le 4 \left[\frac{\sqrt{Q_f} - 1}{\sqrt{Q_f} + 1} \right]^{2k} \left[f(x_0) - \min_x f(x) \right]$$

<u>Fact III</u>: The trajectory $\{x_k\}$ is given by explicit recurrence: \Diamond Initialization: Set

$$d_0 = -g_0 \equiv -f'(x_0) = b - Hx_0;$$

♦ Step *t*: if $g_{t-1} \equiv \nabla f(x_{t-1}) = 0$, terminate, x_{t-1} being the result. Otherwise set

$$\gamma_t = -\frac{g_{t-1}^T d_{t-1}}{d_{t-1}^T H d_{t-1}}$$

$$x_t = x_{t-1} + \gamma_t d_{t-1}$$

$$g_t = f'(x_t) \equiv H x_t - b$$

$$\beta_t = \frac{g_t^T H d_{t-1}}{d_{t-1}^T H d_{t-1}}$$

$$d_t = -g_t + \beta_t d_{t-1}$$

and loop to step t + 1. <u>Note:</u> In the above process, \diamond The gradients $g_0, ..., g_{k_*-1}, g_{k_*} = 0$ are mutually orthogonal \diamond The directions d_k d_k d_k are H orthogonal

 \diamond The directions $d_0, d_1, ..., d_{k_*-1}$ are *H*-orthogonal:

$$i \neq j \Rightarrow d_i^T H d_j = 0$$

♦One has

$$\gamma_t = \operatorname{argmin}_{\gamma} f(x_{t-1} + \gamma d_{t-1})$$

$$\beta_t = \frac{g_t^T g_t}{g_{t-1}^T g_{t-1}}$$

Conjugate Gradient method as applied to a strongly convex quadratic form *f* can be viewed as an iterative algorithm for solving the linear system

Hx = b.

As compared to "direct solvers", like Choleski Decomposition or Gauss elimination, the advantages of CG are:

 \diamond Ability, in the case of exact arithmetic, to find solution in at most *n* steps, with a single matrix-vector multiplication and O(n) additional operations per step.

 \Rightarrow The cost of finding the solution is at most O(n)L, where L is the arithmetic price of matrix-vector multiplication.

<u>Note</u>: When *H* is sparse, $L \ll n^2$, and the price of the solution becomes much smaller than the price $O(n^3)$ for the direct LA methods.

 \Diamond In principle, there is no necessity to assemble H – all we need is the possibility to multiply by H

♦ The non-asymptotic error bound

$$f(x_k) - \min_x f(x) \le 4 \left[\frac{\sqrt{Q_f} - 1}{\sqrt{Q_f} + 1} \right]^{2k} \left[f(x_0) - \min_x f(x) \right]$$

indicates rate of convergence completely independent of the dimension and depending only on the condition number of H.

♠ <u>Illustrations:</u>

 \diamond System 1000 × 1000, $Q_f = 1.e2$:

Itr	$f - f_*$	$ x - x_* _2$
1	2.297e + 003	2.353e + 001
11	1.707e + 001	4.265e + 000
21	3.624e - 001	6.167e - 001
31	6.319 <i>e</i> – 003	8.028e - 002
41	1.150e - 004	1.076e - 002
51	2.016e - 006	1.434e - 003
61	3.178e - 008	1.776e - 004
71	5.946e - 010	2.468e - 005
81	9.668 <i>e</i> - 012	3.096 <i>e</i> – 006
91	1.692e - 013	4.028e - 007
94	4.507e - 014	2.062e - 007

\diamond System 1000 × 1000, $Q_f = 1.e4$:

Itr	$f - f_*$	$ x - x_* _2$
1	1.471e + 005	2.850e + 001
51	1.542e + 002	1.048e + 001
101	1.924e + 001	4.344e + 000
151	2.267e + 000	1.477e + 000
201	2.248e - 001	4.658e - 001
251	2.874e - 002	1.779e - 001
301	3.480 <i>e</i> - 003	6.103 <i>e</i> – 002
351	4.154e - 004	2.054e - 002
401	4.785e - 005	6.846 <i>e</i> – 003
451	4.863 <i>e</i> - 006	2.136e - 003
501	4.537e - 007	6.413 <i>e</i> - 004
551	4.776 <i>e</i> – 008	2.109e - 004
601	4.954e - 009	7.105e - 005
651	5.666e - 010	2.420e - 005
701	6.208e - 011	8.144e - 006
751	7.162e - 012	2.707e - 006
801	7.850 <i>e</i> - 013	8.901e - 007
851	8.076e - 014	2.745e - 007
901	7.436 <i>e</i> – 015	8.559e - 008
902	7.152e - 015	8.412e - 008

 \diamondsuit System 1000 \times 1000, $Q_f = 1.e6$:

Itr	$f - f_*$	$ x - x_* _2$
1	9.916e + 006	2.849e + 001
1000	7.190e + 000	2.683e + 000
2000	4.839 <i>e</i> - 002	2.207e - 001
3000	4.091e - 004	1.999e - 002
4000	2.593 <i>e</i> - 006	1.602e - 003
5000	1.526e - 008	1.160e - 004
6000	1.159e - 010	1.102e - 005
7000	6.022e - 013	7.883 <i>e</i> - 007
8000	3.386 <i>e</i> - 015	5.595 <i>e</i> - 008
8103	1.923e - 015	4.236 <i>e</i> - 008

\diamondsuit System 1000 \times 1000, $Q_f = 1.e12$:

Itr	$f - f_*$	$ x - x_* _2$
1	5.117e + 012	3.078e + 001
1000	1.114e + 007	2.223e + 001
2000	2.658e + 006	2.056e + 001
3000	1.043e + 006	1.964e + 001
4000	5.497e + 005	1.899e + 001
5000	3.444e + 005	1.851e + 001
6000	2.343e + 005	1.808e + 001
7000	1.760e + 005	1.775e + 001
8000	1.346e + 005	1.741e + 001
9000	1.045e + 005	1.709e + 001
10000	8.226 <i>e</i> + 004	1.679e + 001

Non-Quadratic Extensions: CG in the form

$$d_{0} = -g_{0} = -f'(x_{0})$$

$$\gamma_{t} = \operatorname*{argmin}_{\gamma} f(x_{t-1} + \gamma d_{t-1})$$

$$x_{t} = x_{t-1} + \gamma_{t} d_{t-1}$$

$$g_{t} = f'(x_{t})$$

$$\beta_{t} = \frac{g_{t}^{T} g_{t}}{g_{t-1}^{T} g_{t-1}}$$

$$d_{t} = -g_{t} + \beta_{t} d_{t-1}$$

can be applied to *whatever* function *f*, not necessarily quadratic one (Fletcher-Reevs CG), and similarly for another equivalent *in the quadratic case* form:

$$d_{0} = -g_{0} = -f'(x_{0})$$

$$\gamma_{t} = \operatorname*{argmin}_{\gamma} f(x_{t-1} + \gamma d_{t-1})$$

$$x_{t} = x_{t-1} + \gamma_{t} d_{t-1}$$

$$g_{t} = f'(x_{t})$$

$$\beta_{t} = \frac{(g_{t} - g_{t-1})^{T} g_{t}}{g_{t-1}^{T} g_{t-1}}$$

$$d_{t} = -g_{t} + \beta_{t} d_{t-1}$$

(Polak-Ribiere CG).

A Being equivalent in the quadratic case, these (and other) forms of CG become different in the non-quadratic case! Non-quadratic extensions of CG can be used with and without *restarts*.

 \Diamond In quadratic case CG, modulo rounding errors, terminate in at most n steps with exact solutions. In non-quadratic case this is not so.

 \Diamond In non-quadratic CG with restarts, execution is split in *n*-step *cycles*, and cycle t + 1 starts from the last iterate x^t of the previous cycle as from the starting point.

In contrast to this, with no restarts the recurrence like

$$d_{0} = -g_{0} = -f'(x_{0})$$

$$\gamma_{t} = \operatorname*{argmin}_{\gamma} f(x_{t-1} + \gamma d_{t-1})$$

$$x_{t} = x_{t-1} + \gamma_{t} d_{t-1}$$

$$g_{t} = f'(x_{t})$$

$$\beta_{t} = \frac{(g_{t} - g_{t-1})^{T} g_{t}}{g_{t-1}^{T} g_{t-1}}$$

$$d_{t} = -g_{t} + \beta_{t} d_{t-1}$$

is never "refreshed".

<u>Theorem</u>: Let the level set $\{x : f(x) \le f(x_0)\}$ of f be compact and f be twice continuously differentiable in a neighbourhood of G. When minimizing f by Fletcher-Reevs or Polak-Ribiere Conjugate Gradients with exact linesearch and restarts,

 \diamond the trajectory is well-defined and bounded $\diamond f$ never increases

 \diamond all limiting points of the sequence x^t of concluding iterates of the subsequent cycles are critical points of f.

 \Diamond If, in addition, x^t converge to a nondegenerate local minimizer x_* of f and f is 3 times continuously differentiable around x_* , then x^t converge to x_* quadratically.

Quasi-Newton Methods

Quasi-Newton methods are variable metric methods of the generic form

$$x_{t+1} = x_t - \gamma_{t+1} \underbrace{S_{t+1}}_{=A_{t+1}^{-1}} f'(x_t)$$

where $S_{t+1} \succ 0$ and γ_{t+1} is given by line-search.

♠ In contrast to Modified Newton methods, in Quasi-Newton algorithms one operates directly on matrix S_{t+1} , with the ultimate goal to ensure, under favourable circumstances, that

$$S_{t+1} - [f''(x_t)]^{-1} \to 0, \ t \to \infty.$$
 (*)

♠ In order to achieve (*), in Quasi-Newton methods one updates S_t into S_{t+1} in a way which ensures that

 $\diamond S_{t+1}$ is $\succ 0$ $\diamond S_{t+1}(g_t - g_{t-1}) = x_t - x_{t-1}$, where $g_{\tau} = f'(x_{\tau})$.

Generic Quasi-Newton method:

Initialization: Choose somehow starting point x_0 , matrix $S_1 \succ 0$, compute $g_0 = f'(x_0)$. Step t: given x_{t-1} , $g_{t-1} = f'(x_{t-1})$ and $S_t \succ 0$, terminate when $g_{t-1} = 0$, otherwise \diamond Set $d_t = -S_t g_{t-1}$ and perform exact line search from x_{t-1} in the direction d_t , thus getting new iterate

$$x_t = x_{t-1} + \gamma_t d_t;$$

 \diamond compute $g_t = f'(x_t)$ and set

$$p_t = x_t - x_{t-1}, q_t = g_t - g_{t-1};$$

 \diamond update S_t into positive definite symmetric matrix S_{t+1} in such a way that

$$S_{t+1}q_t = p_t$$

and loop.

 \blacklozenge Requirements on the updating rule $S_t \mapsto S_{t+1}$:

♦In order for d_{t+1} to be direction of decrease of f, the rule should ensure $S_{t+1} \succeq 0$ ♦In the case of strongly convex quadratic f, the rule should ensure that $S_t - [f''(\cdot)]^{-1} \to 0$ as $t \to \infty$.
Davidon-Fletcher-Powell method:

$$S_{t+1} = S_t + \frac{1}{p_t^T q_t} p_t p_t^T - \frac{1}{q_t^T S_t q_t} S_t q_t q_t^T S_t.$$

♠ The Davidon-Fletcher-Powell method, as applied to a strongly convex quadratic form, finds exact solution in no more than *n* steps. The trajectory generated by the method initialized with $S_1 = I$ is exactly the one of the Conjugate Gradient method, so that the DFP (Davidon-Fletcher-Powell) method with the indicated initialization is a Conjugate Gradient method.

The Broyden family.

Broyden-Fletcher-Goldfarb-Shanno updating formula:

$$S_{t+1}^{BFGS} = S_t + \frac{1 + q_t^T S_t q_t}{(p_t^T q_t)^2} p_t p_t^T - \frac{1}{p_t^T q_t} \left[p_t q_t^T S_t + S_t q_t p_t^T \right]$$

can be combined with the Davidon-Fletcher-Powell formula

$$S_{t+1}^{DFP} = S_t + \frac{1}{q_t^T p_t} p_t p_t^T - \frac{1}{q_t^T S_t q_t} S_t q_t q_t^T S_t.$$

to yield a single-parametric *Broyden* family of updating formulas

$$S_{t+1}^{\phi} = (1 - \phi)S_{t+1}^{DFP} + \phi S_{t+1}^{BFGS}$$

where $\phi \in [0, 1]$ is parameter.

Facts:

 \diamond As applied to a strongly convex quadratic form f, the Broyden method minimizes the form exactly in no more than n steps, n being the dimension of the design vector. If S₀ is proportional to the unit matrix, then the trajectory of the method on f is exactly the one of the Conjugate Gradient method.

 \diamond all Broyden methods, independently of the choice of the parameter ϕ , being started from the same pair (x_0, S_1) and equipped with the same exact line search and applied to the same problem, generate the same sequence of iterates (although not the same sequence of matrices S_t !).

A Broyden methods are thought to be the most efficient in practice versions of the Conjugate Gradient and quasi-Newton methods, with the pure BFGS method ($\phi = 1$) seemingly being the best.

Convergence of Quasi-Newton methods

A Global convergence of Quasi-Newton methods without restarts is proved only for certain versions of the methods and only under strong assumptions on f.

• For methods with restarts, where the updating formulas are "refreshed" every m steps by setting $S = S_0$, one can easily prove that under our standard assumption that the level set $G = \{x : f(x) \le f(x_0)\}$ is compact and f is continuously differentiable in a neighbourhood of G, the trajectory of starting points of the cycles is bounded, and all its limiting points are critical points of f.

Local convergence:

 \diamond For scheme with restarts, one can prove that if m = n and $S_0 = I$, then the trajectory of starting points x^t of cycles, <u>if</u> it converges to a nondegenerate local minimizer x_* of fsuch that f is 3 times continuously differentiable around x_* , converges to x_* quadratically.

 $fightharpoints \sum_{i=1}^{n} [Powell, 1976].$ Consider the BFGS method without restarts and assume that the method converges to a nondegenerate local minimizer x^* of a three times continuously differentiable function f. Then the method converges to x^* superlinearly.

Convex Programming

A Convex Programming program is an optimization program of the form

$$\min_{X} f(x) \tag{P}$$

where

- $X \subset \mathbb{R}^n$ is a convex compact set, $int X \neq \emptyset$;
- f is a continuous convex function on \mathbb{R}^n .

 \blacklozenge Convexity of a set $X \subset \mathbb{R}^n$ means that whenever X contains a pair of points x, y, it contains the entire segment linking x and y:

 $x, y \in X \Rightarrow \lambda x + (1 - \lambda)y \in X \ \forall \lambda \in [0, 1].$

• Convexity of a function $f : X \to \mathbb{R}$ defined on a convex set X means that on every segment in X, the function is below the corresponding secant:

 $x, y \in X, \lambda \in [0, 1] \Rightarrow$ $f(\lambda x + (1 - \lambda)y) \le \lambda f(x) + (1 - \lambda)f(y).$ Equivalently: The epigraph of f - the set $\mathsf{Epi}\{f\} = \{(x, t) : x \in X, t \ge f(x)\} \subset \mathbb{R}^{n+1}$

- is convex.

$$\min_{X} f(x) \tag{P}$$

Assume that our "environment" when solving convex program (P) is as follows:

1. We have access to a Separation Oracle Sep(X) for X - a routine which, given on input a point $x \in \mathbb{R}^n$, reports whether $x \in X$, and in the case of $x \notin X$, returns a separator -a vector $e \neq 0$ such that

$$e^T x \ge \sup_{y \in X} e^T y$$

• <u>Note</u>: When X is convex and $x \notin intX$, a separator does exist.

$$\min_{X} f(x) \tag{P}$$

2. We have access to a *First Order* oracle which, given on input a point $x \in X$, returns the value f(x) and a *subgradient* f'(x) of f at x.

• A subgradient of f at $x \in X$ is a vector g such that

$$f(y) \ge f(x) + g^T(y - x) \ \forall y \in X.$$

Assuming $f: X \to \mathbb{R}$ convex,

- a subgradient exists at every point $x \in intX$;
- if f is differentiable at $x \in X$, the gradient $\nabla f(x)$ is a subgradient of f at x (and this is the only subgradient of f at x, provided $x \in intX$);

• if f is Lipschitz continuous on X, subgradient of f exists at every point $x \in X$.

$$\min_{X} f(x) \tag{P}$$

3. We are given two positive reals R, r such that for some (unknown) c one has

 $\{x : \|x - c\| \le r\} \subset X \subset \{x : \|x\|_2 \le R\}.$

Example: Consider an optimization program

$$\min_{x \in X} f(x) \equiv \max_{1 \le \ell \le L} f_{\ell}(x)$$
$$X = \{ x \in \mathbb{R}^n : g(x) = \max_{\ell=1,\dots,L} g_{\ell}(x) \le 0 \}$$

where all f_{ℓ} , g_{ℓ} are convex differentiable functions on \mathbb{R}^n . In this case

To build Separation and First Order oracles is the same as to build a routine which is capable to compute, at a given x, all $f_{\ell}(x)$, $\nabla f_{\ell}(x)$, $g_{\ell}(x)$, $\nabla g_{\ell}(x)$.

Indeed, for a differentiable convex function h, its gradient is a subgradient as well:

 $h(y) \ge h(x) + (y - x)^T \nabla h(x) \quad \forall x, y,$ and for the maximum $h = \max_{\ell=1,...,L}$ of convex differentiable functions <u>a</u> subgradient at x is given by

$$\nabla h_{\ell(x)}(x), \ \ell(x) : \max_{\ell=1,\dots,L} h_{\ell}(x) = h_{\ell(x)}(x).$$

Thus,

• If we know how to compute f_{ℓ} , ∇f_{ℓ} , we have in our disposal the First Order oracle for f: we can take, as a subgradient of f at x, the gradient of (any) function f_{ℓ} which is the largest at x.

• If we know how to compute g_{ℓ} , ∇g_{ℓ} , we automatically have in our disposal a Separation oracle for X. Indeed, $x \in X$ iff $g(x) = \max_{\ell=1,...,L} g_{\ell}(x) \leq 0$. Therefore

• In order to check whether $x \in X$, it suffices to compute g(x) and to check whether $g(x) \leq 0$;

• In order to separate x from X when $x \notin X$, it suffices

a) to find $\ell = \ell_*$ such that $g_{\ell}(x) > 0$

b) to set $e = \nabla g_{\ell_*}(x)$. Indeed, we have

 $y \in X \implies [g_{\ell_*}(x) + (y - x)^T e \leq] g_{\ell_*}(y) \leq 0$ $\implies e^T y \leq e^T x - g_{\ell_*}(x) \leq e^T x.$ <u>Theorem</u> ["Polynomial solvability" of convex programs] In the outlined "working environment", for every given $\epsilon > 0$ it is possible to find an ϵ -solution to (P) – a point $x_{\epsilon} \in X$ with

$$f(x_{\epsilon}) \le \min_{x \in X} f(x) + \epsilon$$

- in no more than $N(\epsilon)$ subsequent calls to the Separation and the First Order oracles plus no more than $O(1)n^2N(\epsilon)$ arithmetic operations to process the answers of the oracles, with

$$N(\epsilon) = O(1)n^2 \ln\left(2 + \frac{\operatorname{Var}_R(f)R}{\epsilon \cdot r}\right).$$

Here

$$\operatorname{Var}_{R}(f) = \max_{\|x\|_{2} \le R} f(x) - \min_{\|x\|_{2} \le R} f(x).$$

• We are about to build a "good" solution method for a convex program

$$\min_{x \in X} f(x) \tag{P}$$

• $X \subset \mathbb{R}^n$ is a <u>closed and bounded</u> convex set with a nonempty interior equipped with Separation oracle,

• $f: X \to \mathbb{R}$ is convex and continuous function represented by a First Order oracle,

• We are given R > r > 0:

 $\{x : ||x - \bar{x}|| \le r\} \subset X \subset \{x : ||x|| \le R\}.$

Fo get an idea, consider the one-dimensional case. Here a good solution method is <u>Bisection</u>. When solving a problem

 $\min_{x} \{f(x) : x \in X = [a, b] \subset [-R, R]\},\$ by Bisection, we recursively update <u>localizers</u> $\Delta_t = [a_t, b_t] \text{ of the optimal set } X_{\text{opt}}: X_{\text{opt}} \subset \Delta_t.$

- Initialization: Set $\Delta_0 = [-R, R] [\supset X_{opt}]$
- <u>Step t</u>: Given $\Delta_{t-1} \supset X_{opt}$ let c_t be the midpoint of Δ_{t-1} . Calling Separation and First Order oracle at e_t , we always can replace Δ_{t-1} by twice smaller localizer Δ_t .

"Given $\Delta_{t-1} \supset X_{opt}$ let c_t be the midpoint of Δ_{t-1} . Calling Separation and First Order oracle at e_t , we always can replace Δ_{t-1} by twice smaller localizer Δ_t :"



$\min_{x \in X} f(x) \tag{P}$

♣ In the multi-dimensional case, one can use the *Ellipsoid method* – a simple generalization of Bisection with ellipsoids playing the role of segments.

Cutting Plane Scheme



 \clubsuit We build a sequence of *localizers* G_t such that

- $G_t \supset X_{\mathsf{opt}}$
- G_t decrease as t grows



Straightforward implementation: Centers of Gravity Method

•
$$G_0 = X$$

•
$$x_t = \frac{1}{\operatorname{Vol}(G_{t-1})} \int_{G_{t-1}} x dx$$

<u>Theorem:</u> For the Center of Gravity method, one has

$$Vol(G_t) \leq \left(1 - \left(\frac{n}{n+1}\right)^n\right) Vol(G_{t-1}) \\ \leq \left(1 - \frac{1}{e}\right) Vol(G_{t-1}) \\ < 0.6322 Vol(G_{t-1}).$$
(*)

As a result,

$$\min_{t \leq T} f(x_t) - f_* \leq (1 - 1/e)^{T/n} \operatorname{Var}_X(f)$$

$$Vol(G_t) \leq \left(1 - \left(\frac{n}{n+1}\right)^n\right) Vol(G_{t-1}) \\ \leq \left(1 - \frac{1}{e}\right) Vol(G_{t-1}) \\ < 0.6322 Vol(G_{t-1}).$$
(*)

Reason for (*): Brunn-Minkowski Symmeterization Principle:

Let Y be a convex compact set in \mathbb{R}^n , e be a unit direction and Z be "equi-crosssectional" to X body symmetric w.r.t. e, so that

- Z is symmetric w.r.t. the axis e
- for every hyperplane $H = \{x : e^T x = const\}$, one has

 $\operatorname{Vol}_{n-1}(X \cap H) = \operatorname{Vol}_{n-1}(Z \cap H)$

Then Z is a convex compact set. Equivalently: Let U, V be convex compact nonempty sets in \mathbb{R}^n . Then

 $\operatorname{Vol}^{1/n}(U+V) \ge \operatorname{Vol}^{1/n}(U) + \operatorname{Vol}^{1/n}(V).$ In fact, convexity of U, V is redundant!

Why progress in volumes ensures progress in accuracy?

Let $Vol(G_t) \leq \kappa Vol(G_{t-1})$ with $\kappa < 1$. We claim that

$$\min_{t \le T} f(x_t) - f_* \le \kappa^{T/n} \operatorname{Var}_X(f).$$

Indeed, let $\nu \in (\kappa^{T/n}, 1)$, let $x_* = \underset{X}{\operatorname{argmin}} f$ and let $X_{\nu} = x_* + \nu(X - x_*)$. Then $\operatorname{Vol}(X_{\nu}) = \nu^n \operatorname{Vol}(X) > \kappa^T \operatorname{Vol}(X) >$ $\operatorname{Vol}(G_T)$, whence

$$X_{\nu} \setminus G_T \neq \emptyset.$$

We see that there exists $z \in X$ such that

$$y \equiv x_* + \nu(z - x_*) \not\in G_T.$$

Thus, $y \in G_{t_*-1} \setminus G_{t_*}$ for certain $t_* \leq T$, that is,

$$(y - x_{t_*})^T f'(x_{t_*}) > 0 \Rightarrow f(x_{t_*}) < f(y).$$

On the other hand,

$$f(y) = f(x_* + \nu(z - x_*))$$

$$\leq f(x_*) + \nu(f(z) - f(x_*))$$

$$\leq f(x_*) + \nu \operatorname{Var}_X(f).$$

We see that

 $\min_{t \leq T} f(x_t) \leq f(x_{t_*}) < f(y) \leq f(x_*) + \nu \operatorname{Var}_X(f).$

Passing to limits as $\nu \to \kappa^{T/n} + 0,$ we arrive at

$$\min_{t \leq T} f(x_t) \leq f(x_*) + \kappa^{T/n} \operatorname{Var}_X(f).$$

$$\min_{x \in X} f(x) \tag{P}$$

Ellipsoid method – the idea. Assume we already have an *n*-dimensional ellipsoid

$$E = \{x = c + Bu : u^T u \le 1\}$$
$$[B \in \mathbb{R}^{n \times n}, \mathsf{Det}(B) \neq 0]$$

which covers the optimal set X_{opt} .

In order to replace E with a smaller ellipsoid containing X_{opt} , we

1) Call Sep(X) to check whether $c \in X$

1.a) If Sep(X) says that $c \notin X$ and returns a separator e:

$$e \neq 0, e^T c \ge \sup_{y \in X} e^T y.$$

we may be sure that

 $X_{\text{opt}} = X_{\text{opt}} \bigcap E \subset \widehat{E} = \{x \in E : e^T(x-c) \le 0\}.$

1.b) If Sep(X) says that $c \in X$, we call $\mathcal{O}(f)$ to compute f(c) and f'(c), so that

$$f(y) \ge f(c) + (y - c)^T f'(c) \quad \forall y \quad (*)$$

If f'(c) = 0, c is optimal for (P) by (*), otherwise (*) says that

$$X_{\text{opt}} = X_{\text{opt}} \bigcap E \subset \widehat{E} = \{x \in E : \underbrace{[f'(c)]}_{e}^{T}(x-c) \le 0\}$$

• Thus, we either terminate with optimal solution, or get a "half-ellipsoid"

$$\widehat{E} = \{ x \in E : e^T x \le e^T c \} \quad [e \neq 0]$$

containing X_{opt} .

1) Given an ellipsoid

$$E = \{x = c + Bu : u^T u \le 1\}$$
$$[B \in \mathbb{R}^{n \times n}, \text{Det}(B) \neq 0]$$

containing the optimal set of (P) and calling the Separation and the First Order oracles at the center c of the ellipsoid, we pass from Eto the half-ellipsoid

$$\widehat{E} = \{ x \in E : e^T x \le e^T c \} \quad [e \neq 0]$$

containing X_{opt} .

2) It turns out that the half-ellipsoid \widehat{E} can be included into a new ellipsoid E^+ of n-dimensional volume less than the one of E:

$$E^{+} = \{x = c^{+} + B^{+}u : u^{T}u \leq 1\},\$$

$$c^{+} = c - \frac{1}{n+1}Bp,\$$

$$B^{+} = B\left(\frac{n}{\sqrt{n^{2}-1}}(I_{n} - pp^{T}) + \frac{n}{n+1}pp^{T}\right)\$$

$$= \frac{n}{\sqrt{n^{2}-1}}B + \left(\frac{n}{n+1} - \frac{n}{\sqrt{n^{2}-1}}\right)(Bp)p^{T},\$$

$$p = \frac{B^{T}e}{\sqrt{e^{T}BB^{T}e}}\$$

$$/ol(E^{+}) = \left(\frac{n}{\sqrt{n^{2}-1}}\right)^{n-1}\frac{n}{n+1}Vol(E)\$$

$$\leq \exp\{-1/(2n)\}Vol(E)$$

♠ In the Ellipsoid method, we iterate the above construction, starting with $E_0 = \{x : \|x\|_2 \le R\} \supset X$, thus getting a sequence of ellipsoids $E_{t+1} = (E_t)^+$ with volumes "rapidly" converging to 0, all of them containing X_{opt} .

Verification of the fact that

$$\begin{bmatrix} E = \{x = c + Bu : u^{T}u \leq 1\} \\ \hat{E} = \{x \in E : e^{T}x \leq e^{T}c\} \\ [e \neq 0] \end{bmatrix}$$

$$\stackrel{\Downarrow}{} \\ \hat{E} \subset E^{+} = \{x = c^{+} + B^{+}u : u^{T}u \leq 1\} \\ c^{+} = c - \frac{1}{n+1}Bp,$$

$$B^{+} = B\left(\frac{n}{\sqrt{n^{2}-1}}(I_{n} - pp^{T}) + \frac{n}{n+1}pp^{T}\right) \\ = \frac{n}{\sqrt{n^{2}-1}}B + \left(\frac{n}{n+1} - \frac{n}{\sqrt{n^{2}-1}}\right)(Bp)p^{T},$$

$$p = \frac{B^{T}e}{\sqrt{e^{T}BB^{T}e}}$$

is easy:

• E is the image of the unit ball under a oneto-one affine mapping

• therefore \widehat{E} is the image of a half-ball $\widehat{W}=\{u~:~u^Tu~\leq~\mathbf{1}, p^Tu~\leq~\mathbf{0}\}$ – under the same mapping

• Ratio of volumes remains invariant under affine mappings, and therefore all we need is to find a "small" ellipsoid containing half-ball \widehat{W} .



Since \widehat{W} is "highly symmetric", to find the smallest possible ellipsoid containing \widehat{W} is a simple exercise in elementary Calculus.

and

$$\rho_0 = R, \ L_0 = 0$$

$$\begin{bmatrix} \rho_t : \operatorname{Vol}(E_t) = \operatorname{Vol}\left\{\{x : \|x\|_2 \le \rho_t\} \\ L_t : L_t \le \operatorname{Var}_R(f) = \max_{E_0} f - \min_{E_0} f \end{bmatrix}$$

<u>Step t, t = 1, 2, ...</u> At the beginning of step t, we have the data c_{t-1}, B_{t-1} of the previous ellipsoid

 $E_{t-1} = \{x = c_{t-1} + B_{t-1}u : u^T u \le 1\}$ $[c_{t-1} \in \mathbb{R}^n, B_{t-1} \in \mathbb{R}^{n \times n}, \text{Det}B_{t-1} \neq 0]$ along with the quantities $L_{t-1} \ge 0$ and ρ_{t-1} .
• At step t,

1) We call Sep(X), c_{t-1} being the input. If $x_{t-1} \notin X$ ("non-productive step"), Sep(X) returns a separator

$$e \neq 0$$
: $e^T c_{t-1} \ge \sup_{y \in X} e^T y.$

In this case, we set

$$e_t = e, \ L_t = L_{t-1}$$

and go to 3). If $c_{t-1} \in X$ ("productive step"), we go to 2).

2) We call $\mathcal{O}(f)$, c_{t-1} being the input, to get $f(c_{t-1})$, $e = f'(c_{t-1})$. If e = 0, we terminate and claim that c_{t-1} is an optimal solution to (P).

If $e \neq 0$, we set

$$e_t = e$$
,

compute the quantity

$$\ell_t = \max_{y \in E_0} [e_t^T y - e_t^T c_{t-1}] = R ||e_t||_2 - e_t^T c_{t-1},$$

update L by setting

$$L_t = \max\{L_{t-1}, \ell_t\}$$

and go to 3). <u>Note</u>: Since $f(y) - f(c_{t-1}) \ge e_t^T(y - c_{t-1})$, and $c_{t-1} \in X \subset E_0$, we have $\ell_t \le \max_{E_0} f(y) - f(c_{t-1}) \le \operatorname{Var}_R(f)$, whence $L_t \le \operatorname{Var}_R(f)$. 3) We set

$$\widehat{E}_t = \{ x \in E_{t-1} : e_t^T x \le e_t^T c_{t-1} \}$$

and define the new ellipsoid

$$E_t = \{x = c_t + B_t u : u^T u \le 1\}$$

by setting

$$p_{t} = \frac{B_{t-1}^{T} e_{t}}{\sqrt{e_{t}^{T} B_{t-1} B_{t-1}^{T} e_{t}}}$$

$$c_{t} = c_{t-1} - \frac{1}{n+1} B_{t-1} p_{t},$$

$$B_{t} = \frac{n}{\sqrt{n^{2}-1}} B_{t-1} + \left(\frac{n}{n+1} - \frac{n}{\sqrt{n^{2}-1}}\right) (B_{t-1} p_{t}) p_{t}^{T}.$$
We also set

$$\rho_t = |\mathsf{Det}B_t|^{1/n} = \left(\frac{n}{\sqrt{n^2 - 1}}\right)^{\frac{n-1}{n}} \left(\frac{n}{n+1}\right)^{\frac{1}{n}} \rho_{t-1}$$

and go to 4).

4) [Termination test]. We check whether the inequality

$$\frac{\rho_t}{r} < \frac{\epsilon}{L_t + \epsilon} \tag{(*)}$$

(r > 0 is a given in advance radius of Euclidean ball contained in X).

• If (*) holds, we terminate and output the best (with the smallest value of f) of the "search points" $c_{\tau-1}$ associated with productive steps $\tau \leq t$.

• If (*) does not hold, we pass to step t + 1.

Numerical example:

 $f(x) = \frac{1}{2}(1.4435x_1 + 0.6232x_2 - 7.9574)^2$ $+5(-0.3509x_1 + 0.7990x_2 + 2.8778)^4$ $-1 \le x_1, x_2 \le 1,$

• $X_{opt} = \{(1, -1)\}, Opt = 70.030152768...$

Here are the best objective values of feasible solutions found in course of the first t steps, t = 1, ..., 256:

t	best value	t	best value
1	374.61091739	16	76.838253451
2	216.53084103		
3	146.74723394	32	70.901344815
4	112.42945457		
5	93.84206347	64	70.031633483
6	82.90928589		
7	82.90928589	128	70.030154192
8	82.90928589		
		256	70.030152768
• The initial phase of the process looks as follows:



Ellipses E_{t-1} and search points c_{t-1} , t = 1, 2, 3, 4Arrows: gradients of the objective f(x)Unmarked segments: tangents to the level lines of <u>Theorem</u> [Complexity of the Ellipsoid Algorithm] Let the Ellipsoid Algorithm be applied to convex program

$$\min_{x \in X} f(x) \tag{P}$$

such that X contains a Euclidean ball of a given radius r > 0 and is contained in the ball $E_0 = \{ ||x||_2 \le R \}$ of a given radius R. For every input accuracy $\epsilon > 0$, the Ellipsoid method terminates after no more than

$$= \operatorname{Ceil}\left(2n^{2}\left[\ln\left(\frac{R}{r}\right) + \ln\left(\frac{\epsilon + \operatorname{Var}_{R}(f)}{\epsilon}\right)\right]\right) + 1$$

steps, where

$$\operatorname{Var}_R(f) = \max_{E_0} f - \min_{E_0} f,$$

Vol is the *n*-dimensional volume and Ceil(a) is the smallest integer $\geq a$.

Moreover, the result \hat{x} generated by the method is a feasible ϵ -solution to (P):

$$\widehat{x} \in X \text{ and } f(x) - \min_{X} f \leq \epsilon.$$

<u>Proof.</u> Let t be such that the method does not terminate before this step and does <u>not</u> terminate at this step because of

$$c_{t-1} \in X, f'(c_{t-1}) = 0.$$

Then for $1 \leq \tau \leq t$ one has

(a)
$$E_0 \supset X;$$

(b) $E_{\tau} \supset \hat{E}_{\tau} = \left\{ x \in E_{\tau-1} : e_{\tau}^T x \leq e_{\tau}^T c_{\tau-1} \right\},$
(c) $\operatorname{Vol}(E_{\tau}) = \rho_{\tau}^n \operatorname{Vol}(E_0)$
 $= \left(\frac{n}{\sqrt{n^2 - 1}} \right)^{n-1} \frac{n}{n+1} \operatorname{Vol}(E_{\tau-1})$
 $\leq \exp\{-1/(2n)\} \operatorname{Vol}(E_{\tau-1}).$

By (c) we have

$$\rho_{\tau} \leq \exp\{-\tau/(2n^2)\}R, \ \tau = 1, ..., t.$$

1^0 . We claim that

(!) If the Ellipsoid Algorithm terminates at certain step t, then the result \hat{x} is well-defined and is a feasible ϵ -solution to (P). (!) "If the Ellipsoids method terminates at certain step t, <u>then</u> the result \hat{x} is well-defined and is a feasible ϵ -solution to (P)."

Indeed, there are two possible reasons for termination at step t:

• We arrive at the situation where

$$c_{t-1} \in X \text{ and } f'(c_{t-1}) = 0.$$
 (A)

In this case

 $\forall y : f(y) \ge f(c_{t-1}) + (y - c_{t-1})^T f'(c_{t-1}) = f(c_{t-1}),$ and since $c_{t-1} \in X$, c_{t-1} is optimal solution to $\min_X f$.

We arrive at the situation where

$$\frac{\rho_t}{r} < \frac{\epsilon}{L_t + \epsilon} \tag{B}$$

Let us prove that (!) is valid in the case of (B).

 2^{0} . Let the method terminate at step t according to

$$\frac{\rho_t}{r} < \frac{\epsilon}{L_t + \epsilon} \equiv \nu \tag{B}$$

There exists ν' : $\frac{\rho_t}{r} < \nu' < \nu$ [\leq 1]. (P) is solvable (f is continuous on X, X is compact). Let $x_* \in X_{\text{opt}}$, and let

$$X^{+} = x_{*} + \nu'(X - x_{*}) = \{x = (1 - \nu')x_{*} + \nu'z : z \in X\}$$

Note that

$$Vol(X^{+}) = (\nu')^{n} Vol(X) \geq (\nu')^{n} r^{n} Vol(\{x : ||x||_{2} \le 1\}) > \rho_{t}^{n} Vol(\{x : ||x||_{2} \le 1\}) = Vol(E_{t}),$$

and consequently

$$\begin{array}{l} X_+ \setminus E_t \neq \emptyset \\ \downarrow \\ \exists y = (1 - \nu') x_* + \nu' z \notin E_t \qquad [z \in X] \end{array}$$



Since $y \in X \subset E_0$ and $y \notin E_t$, y was "cut off":

$$e_{\tau}^T y > e_{\tau}^T c_{\tau-1} \tag{+}$$

at certain step $\tau \leq t$.

• Exists
$$\tau \leq t$$
:

$$e_{\tau}^T y > e_{\tau}^T c_{\tau-1} \tag{+}$$

Observe that the step τ is productive, since otherwise e_{τ} separates $c_{\tau-1}$ and X, while $y \in X$. Consequently, $e_{\tau} = f'(c_{\tau-1})$.

• <u>Situation:</u> We are solving problem

$$\min_{x \in X} f(x) \tag{P}$$

and the Ellipsoid Algorithm terminates at step t.

• Target: To prove that the result is an ϵ -solution to (P)

• <u>Current state of the proof</u>: We are exploring the case

$$\exists y = (1 - \nu')x_* + \nu'z \quad \left[x_* \in X_{\text{opt}}, z \in X, \nu' < \frac{\epsilon}{L_t + \epsilon}\right]$$

and a productive step $\tau \leq t$ such that

$$(y - c_{\tau-1})^T f'(c_{\tau-1}) > 0$$
 (+)

$$(y - c_{\tau-1})^T f'(c_{\tau-1}) > 0 \qquad (+)$$
• We have
$$\begin{array}{c|c} f(x_*) \geq f(c_{\tau-1}) + (x_* - c_{\tau-1})^T f'(c_{\tau-1}) & \times & (1 - \nu') \\ + & & + \\ L_{\tau} \geq (z - c_{\tau-1})^T f'(c_{\tau-1}) & & \times & \nu' \end{array}$$

$$(1 - \nu')f(x_*) + \nu'L_{\tau} \geq (1 - \nu')f(c_{\tau-1}) \\ + ([(1 - \nu')x_* + \nu'z] - c_{\tau-1})^T f'(c_{\tau-1}) \\ = (1 - \nu')f(c_{\tau-1}) + (y - c_{\tau-1})^T f'(c_{\tau-1}) \\ \geq (1 - \nu')f(c_{\tau-1}) & [by (+)] \end{array}$$

$$\begin{array}{c} \downarrow \\ f(c_{\tau-1}) \leq f(x_*) + \frac{\nu'L_{\tau}}{1 - \nu'} \\ \leq f(x_*) + \frac{\nu'L_{t}}{1 - \nu'} \\ [since L_{\tau} \leq L_t \text{ in view of } \tau \leq t] \\ \leq f(x_*) + \epsilon \\ [by definition of \nu \text{ and since } \nu' < \nu] \\ = Opt(C) + \epsilon. \end{array}$$

• <u>Conclusion</u>: In course of running the method, a feasible solution $c_{\tau-1}$ with $f(c_{\tau-1}) \leq \text{Opt}(C) + \epsilon$ was found. • <u>Situation</u>: We are solving problem

$$\min_{x \in X} f(x) \tag{P}$$

and the Ellipsoid Algorithm terminates at step t.

• Target: To prove that the result is an ϵ -solution to (P)

• <u>Current state of the proof</u>: It was shown that in course of running the method, a feasible solution $c_{\tau-1}$ with $f(c_{\tau-1}) \leq \operatorname{Opt}(C) + \epsilon$ was found.

• By construction, the result of the Ellipsoid Algorithm is the best – with the smallest value of f – of feasible solutions c_{s-1} generated in course of t steps, so that if the method terminates at a step t, then the result is an ϵ -solution to (P).

 3^{0} . It remains to prove that the method does terminate in course of the first

$$\underbrace{\operatorname{Ceil}\left(2n^{2}\left[\ln\left(\frac{R}{r}\right) + \ln\left(\frac{\epsilon + \operatorname{Var}_{R}(f)}{\epsilon}\right)\right]\right) + 1}_{N_{(P)}(\epsilon)}$$

steps.

We have seen that $L_{\tau} \leq \operatorname{Var}_{R}(f)$ for all τ and that

$$\rho_{\tau} \le \exp\{-\tau/(2n^2)\}R$$

for all τ .

It follows that the premise in termination rule

$$\frac{\rho_t}{r} < \underbrace{\frac{\epsilon}{L_t + \epsilon}}_{\geq \frac{\epsilon}{\operatorname{Var}_R(f) + \epsilon}} \Rightarrow \text{ Termination}$$

indeed is satisfied in course of the first $N_{(\mathsf{P})}(\epsilon)$ steps.

Traditional methods for general constrained problems

$$\min_{x} \left\{ f(x): \begin{array}{l} g_{j}(x) \leq 0, \ j = 1, ..., m \\ h_{i}(x) = 0, \ i = 1, ..., k \end{array} \right\} \quad (P)$$

can be partitioned into

◇Primal methods, where one mimics unconstrained approach, travelling along the feasible set in a way which ensures progress in objective at every step

♦Penalty/Barrier methods, which reduce constrained minimization to solving a sequence of essentially unconstrained problems

 \diamond Lagrange Multiplier methods, where one focuses on dual problem associated with (*P*). A posteriori the Lagrange multiplier methods, similarly to the penalty/barrier ones, reduce (*P*) to a sequence of unconstrained problems, but in a "smart" manner different from the penalty/barrier scheme

 \diamond Sequential Quadratic Programming methods, where one directly solves the KKT system associated with (P) by a kind of Newton method.

Penalty/Barrier Methods

Penalty Scheme, Equality Constrainsts. Consider equality constrained problem

 $\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$ (P)

and let us "approximate" it by unconstrained problem

$$\min_{x} f_{\rho}(x) = f(x) + \underbrace{\frac{\rho}{2} \sum_{i=1}^{k} h_{i}^{2}(x)}_{\substack{\text{penalty} \\ \text{term}}} \qquad (P[\rho])$$

 $\rho > 0$ is penalty parameter. <u>Note:</u> (A) On the feasible set, the penalty term vanishes, thus $f_{\rho} \equiv f$;

(B) When ρ is large and x is infeasible, $f_{\rho}(x)$ is large:

$$\lim_{\rho \to \infty} f_{\rho}(x) = \begin{cases} f(x), & x \text{ is feasible} \\ +\infty, & \text{otherwise} \end{cases}$$

 \Rightarrow It is natural to expect that solution of $(P[\rho])$ approaches, as $\rho \to \infty$, the optimal set of (P).

Penalty Scheme, General Constraints. In the case of general constrained problem

$$\min_{x} \left\{ f(x): \begin{array}{l} h_{i}(x) = 0, \ i = 1, ..., k \\ g_{j} \leq 0, \ j = 1, ..., m \end{array} \right\}, \quad (P)$$

the same idea of penalizing the constraint violations results in approximating (P) by unconstrained problem

$$\min_{x} f_{\rho}(x) = f(x) + \underbrace{\frac{\rho}{2} \left[\sum_{i=1}^{k} h_{i}^{2}(x) + \sum_{j=1}^{m} [g_{j}(x)^{+}]^{2} \right]}_{\substack{\text{penalty}\\ \text{term}}}$$

where

$$g_j^+(x) = \max[g_j(x), 0]$$

and $\rho > 0$ is penalty parameter. Here again

$$\lim_{
ho
ightarrow\infty}f_
ho(x)=\left\{egin{array}{cc} f(x), & x ext{ is feasible} \ +\infty, & ext{otherwise} \end{array}
ight.$$

and we again may expect that the solutions of $(P[\rho])$ approach, as $\rho \to \infty$, the optimal set of (P).

Barrier scheme normally is used for inequality constrained problems

 $\min_{x} \left\{ f(x) : g_j(x) \le 0, \, j = 1, ..., m \right\}$ (P)

satisfying "Slater condition": the feasible set

$$G = \left\{ x : g_j(x) \le 0, \, j \le m \right\}$$

of (P) possesses a nonempty interior $\inf G$ which is dense in G, and $g_j(x) < 0$ for $x \in \inf G$.

♠ Given (P), one builds a *barrier* (≡interior penalty) for G – a function F which is welldefined and smooth on intG and blows up to +∞ along every sequence of points $x_i \in intG$ converging to a boundary point of G:

 $x_{i} \in \operatorname{int} G, \lim_{i \to \infty} x_{i} = x \notin \operatorname{int} G \Rightarrow F(x_{i}) \to \infty, i \to \infty.$ $\underbrace{\mathsf{Examples:}}_{\diamondsuit \mathsf{Log-barrier}} F(x) = -\sum_{j} \ln(-g_{j}(x))$ $\diamondsuit \mathsf{Carrol Barrier} F(x) = -\sum_{j} \frac{1}{g_{j}(x)}$

 $\min_{x} \{f(x) : g_j(x) \le 0, j = 1, ..., m\}$ (P) After interior penalty F for the feasible domain of (P) is chosen, the problem is approx-

imated by the "essentially unconstrained" problem

$$\min_{x \in \text{int}_G} F_{\rho}(x) = f(x) + \frac{1}{\rho} F(x) \qquad (P[\rho])$$

When *penalty parameter* ρ is large, the function F_{ρ} is close to f everywhere in G, except for a thin stripe around the boundary.

 \Rightarrow It is natural to expect that solutions of $(P[\rho])$ approach the optimal set of (P) as $\rho \to \infty$,

Investigating Penalty Scheme

Let us focus on equality constrained problem

$$\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$$
 (P)

and associated penalized problems

$$\min_{x} f_{\rho}(x) = f(x) + \frac{\rho}{2} ||h(x)||_{2}^{2} \qquad (P[\rho])$$

(results for general case are similar).

Questions of interest:

 \diamond Whether indeed unconstrained minimizers of the penalized objective f_{ρ} converge, as $\rho \rightarrow \infty$, to the optimal set of (P)?

♦ What are our possibilities to minimize the penalized objective?

$$\min_{x} \{f(x) : h_{i}(x) = 0, i = 1, ..., k\} \quad (P) \\
\Downarrow \\
\min_{x} f_{\rho}(x) = f(x) + \frac{\rho}{2} \|h(x)\|_{2}^{2} \quad (P[\rho])$$

Simple fact: Let (P) be feasible, the objective and the constraints in (P) be continuous and let f possess bounded level sets $\{x : f(x) \leq a\}$. Let, further X_* be the set of global solutions to (P). Then X_* is nonempty, approximations problems $(P[\rho])$ are solvable, and their global solutions approach X_* as $\rho \to \infty$:

$$\forall \epsilon > 0 \exists \rho(\epsilon) : \rho \ge \rho(\epsilon), x_*(\rho) \text{ solves } (P[\rho]) \\\Rightarrow \mathsf{dist}(x_*(\rho), X_*) \equiv \min_{x_* \in X_*} \|x_*(\rho) - x_*\|_2 \le \epsilon$$

Proof. 1⁰. By assumption, the feasible set of (P) is nonempty and closed, f is continuous and $f(x) \to \infty$ as $||x||_2 \to \infty$. It follows that f attains its minimum on the feasible set, and the set X_* of global minimizers of f on the feasible set is bounded and closed.

$$\min_{x} \{f(x) : h_{i}(x) = 0, i = 1, ..., k\}$$
(P)
$$\underset{x}{\Downarrow} \prod_{x} f_{\rho}(x) = f(x) + \frac{\rho}{2} ||h(x)||_{2}^{2}$$
(P[\rho])

2⁰. The objective in $(P[\rho])$ is continuous and goes to $+\infty$ as $||x||_2 \to \infty$; consequently, $(P[\rho])$ is solvable.

$$\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$$
 (P)

$$\min_{x} f_{\rho}(x) = f(x) + \frac{\rho}{2} ||h(x)||_{2}^{2} \qquad (P[\rho])$$

3⁰. It remains to prove that, for every $\epsilon > 0$, the solutions of $(P[\rho])$ with large enough value of ρ belong to ϵ -neighbourhood of X_* . Assume, on the contrary, that for certain $\epsilon > 0$ there exists a sequence $\rho_i \to \infty$ such that an optimal solution x_i to $(P[\rho_i])$ is at the distance $> \epsilon$ from X_* , and let us lead this assumption to contradiction.

 \diamond Let f_* be the optimal value of (P). We clearly have

$$f(x_i) \le f_{\rho_i}(x_i) \le f_*, \tag{1}$$

whence $\{x_i\}$ is bounded. Passing to a subsequence, we may assume that $x_i \rightarrow \overline{x}$ as $i \rightarrow \infty$.

$$\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$$
 (P)

$$\min_{x} f_{\rho}(x) = f(x) + \frac{\rho}{2} ||h(x)||_{2}^{2} \qquad (P[\rho])$$

$$x_{i} \in \operatorname{Argmin}_{x} f_{\rho_{i}}(x), x_{i} \to \overline{x} \notin X_{*}$$

$$\Rightarrow f(x_{i}) \leq f_{\rho_{i}}(x_{i}) \leq f_{*}$$
(1)

 \diamond We claim that $\bar{x} \in X_*$, which gives the desired contradiction. Indeed,

— \bar{x} is feasible, since otherwise

$$\lim_{i \to \infty} \underbrace{ [f(x_i) + \frac{\rho_i}{2} ||h(x_i)||_2^2]}_{\substack{f_{\rho_i}(x_i) \\ \\ = f(\bar{x}) + \lim_{i \to \infty} \frac{\rho_i}{2} \underbrace{||h(x_i)||_2^2}_{\rightarrow ||h(\bar{x})||_2^2 > 0} = +\infty,$$

in contradiction to (1); — $f(\bar{x}) = \lim_{i \to \infty} f(x_i) \le f_*$ by (1); since \bar{x} is feasible for (P), we conclude that $\bar{x} \in X_*$. A Shortcoming of Simple Fact: In non-convex case, we cannot find/approximate global minimizers of the penalized objective, so that Simple Fact is "unsubstantial"...

$$\min_{x} \{f(x) : h_{i}(x) = 0, i = 1, ..., k\} \quad (P) \\
\Downarrow \\
\min_{x} f_{\rho}(x) = f(x) + \frac{\rho}{2} ||h(x)||_{2}^{2} \quad (P[\rho])$$

<u>Theorem.</u> Let x^* be a <u>nondegenerate</u> locally optimal solution to (P), i.e., a feasible solution such that

 $\Diamond f, h_i$ are twice continuously differentiable in a neighbourhood of $x_*,$

 \diamond the gradients of the constraints taken at x_* are linearly independent,

 \diamond at x_* , the Second Order Sufficient Optimality condition is satisfied.

There exists a neighbourhood V of x^* and $\bar{\rho}>0$ such that

 \Diamond for every $\rho \geq \overline{\rho}$, f_{ρ} possesses in V exactly one critical point $x_*(\rho)$;

 $\Diamond x_*(\rho)$ is a nondegenerate local minimizer of f_ρ and a minimizer of f_ρ in V;

 $\Diamond x_*(\rho) \to x_* \text{ as } \rho \to \infty.$

In addition,

• The local "penalized optimal value"

$$f_{\rho}(x_*(\rho)) = \min_{x \in V} f_{\rho}(x)$$

is nondecreasing in ρ

Indeed, $f_{\rho}(\cdot) = f(\cdot) + \frac{\rho}{2} ||h(\cdot)||_{2}^{2}$ grow with ρ • The constraint violation $||h(x_{*}(\rho))||_{2}$ monotonically goes to 0 as $\rho \to \infty$ Indeed, let $\rho'' > \rho'$, and let $x' = x_{*}(\rho')$, $x'' = x_{*}(\rho'')$. Then

$$\begin{aligned} f(x') + \frac{\rho''}{2} \|h(x')\|_{2}^{2} &\geq f(x'') + \frac{\rho''}{2} \|h(x'')\|_{2}^{2} \\ f(x'') + \frac{\rho'}{2} \|h(x'')\|_{2}^{2} &\geq f(x') + \frac{\rho'}{2} \|h(x')\|_{2}^{2} \\ &\Rightarrow f(x') + f(x'') + \frac{\rho''}{2} \|h(x')\|_{2}^{2} + \frac{\rho'}{2} \|h(x'')\|_{2}^{2} \\ &\geq f(x') + f(x'') + \frac{\rho''}{2} \|h(x'')\|_{2}^{2} + \frac{\rho'}{2} \|h(x')\|_{2}^{2} \\ &\Rightarrow \frac{\rho'' - \rho'}{2} \|h(x')\|_{2}^{2} &\geq \frac{\rho'' - \rho'}{2} \|h(x'')\|_{2}^{2} \end{aligned}$$

• The true value of the objective $f(x_*(\rho))$ at $x_*(\rho)$ is nondecreasing in ρ

• The quantities $\rho h_i(x_*(\rho))$ converge to optimal Lagrange multipliers of (P) at x_* Indeed,

$$0 = f'_{\rho}(x_*(\rho)) = f'(x_{\rho}) + \sum_i (\rho h_i(x_*(\rho))) h'_i(x_*(\rho)).$$

Solving penalized problem

$$\begin{split} \min_{x} f_{\rho}(x) &\equiv f(x) + \frac{\rho}{2} \|h(x)\|_{2}^{2} \qquad (P[\rho]) \\ \diamondsuit In \ principle, \ \text{one can solve} \ (P[\rho]) \ \text{by what-} \\ \text{ever method for unconstrained minimization.} \\ \diamondsuit \underline{\text{However:}} \ The \ conditioning \ of \ f \ deterio- \\ rates \ as \ \rho \to \infty. \\ \text{Indeed, as } \rho \to \infty, \ \text{we have} \end{split}$$

$$d^{T} f_{\rho}''(\underbrace{x_{*}(\rho)}_{x})d = d^{T} \left[f''(x) + \sum_{i} \rho h_{i}(x)h_{i}''(x) \right] d$$
$$+ \underbrace{\rho \sum_{i}^{\neg \nabla_{x}^{2}L(x_{*},\mu^{*})}}_{\substack{\rho \sum_{i} (d^{T} h_{i}'(x))^{2}}} \right]$$
$$+ \underbrace{\rho \sum_{i}^{\neg \infty, \rho \to \infty}}_{\text{except for } d^{T} h'(x_{*}) = 0}$$

 \Rightarrow slowing down the convergence and/or severe numerical difficulties when working with large penalties...

Barrier Methods

$$\min_{x} \{ f(x) : x \in G \equiv \{ x : g_j(x) \le 0, \, j = 1, ..., m \} \}$$
(P)

$$\min_{x} F_{\rho}(x) \equiv \overset{\vee}{f}(x) + \frac{1}{\rho}F(x) \qquad (P[\rho])$$

F is interior penalty for G = cl(intG):

Ш

 $\diamondsuit F$ is smooth on intG

 $\Diamond F$ tends to ∞ along every sequence $x_i \in$ int*G* converging to a boundary point of *G*. <u>Theorem.</u> Assume that $G = \operatorname{cl}(\operatorname{int} G)$ is bounded and f, g_j are continuous on *G*. Then the set X_* of optimal solutions to (*P*) and the set $X_*(\rho)$ of optimal solutions to (*P*[ρ]) are nonempty, and the second set converges to the first one as $\rho \to \infty$: for every $\epsilon > 0$, there exists $\rho = \rho(\epsilon)$ such that

$$\rho \ge \rho(\epsilon), x_*(\rho) \in X_*(\rho) \Rightarrow \operatorname{dist}(x_*(\rho), X_*) \le \epsilon.$$

In the case of <u>convex</u> program

$$\min_{x \in G} f(x) \tag{P}$$

with closed and bounded convex G and convex objective f, the domain G can be in many ways equipped with a twice continuously differentiable *strongly convex* penalty F(x).

 \blacklozenge Assuming f twice continuously differentiable on intG, the aggregate

$$F_{\rho}(x) = \rho f(x) + F(x)$$

is strongly convex on intG and therefore attains its minimum at a single point

$$x_*(\rho) = \underset{x \in \text{int}G}{\operatorname{argmin}} F_{\rho}(x).$$

• It is easily seen that the path $x_*(\rho)$ is continuously differentiable and converges, as $\rho \to \infty$, to the optimal set of (P).

$$\min_{\substack{x \in G \\ \psi}} f(x) \qquad (P)$$

$$\lim_{\substack{x \in G \\ x \in \text{int}G}} F_{\rho}(x) = \rho f(x) + F(x) \qquad (P[\rho])$$

$$\lim_{\substack{x \in \text{int}G \\ x \in \text{int}G}} F_{\rho}(x) \underset{\rho \to \infty}{\to} \operatorname{Argmin} f$$

for the inclassical path-following scheme (Fiacco and McCormic, 1967), one traces the path $x_*(\rho)$ as $\rho \to \infty$ according to the following generic scheme:

 \bigcirc Given $(x_i \in \text{int}G, \rho_i > 0)$ with x_i close to $x_*(\rho_i)$,

— update ρ_i into a larger value ρ_{i+1} of the penalty

— minimize $F_{\rho_{i+1}}(\cdot)$, x_i being the starting point, until a new iterate x_{i+1} close to

$$x_*(\rho_{i+1}) = \underset{x \in \text{int}G}{\operatorname{argmin}} F_{\rho_{i+1}}(x)$$

is built, and loop.

• To update a tight approximation x_i of argmin $F_{\rho_i}(x)$ into a tight approximation x_{i+1} of argmin $F_{\rho_i}(x)$, one can apply to $F_{\rho_{i+1}}(\cdot)$ a method for "essentially unconstrained" minimization, preferably, the Newton method • When Newton method is used, one can try to increase penalty at a "safe" rate, keeping x_i in the domain of quadratic convergence of the Newton method as applied to $F_{\rho_{i+1}}(\cdot)$ and thus making use of fast local convergence of the method.

Questions: • How to choose *F*?

- How to measure closeness to the path?
- How to ensure "safe" penalty updating without slowing the method down?

<u>Note</u>: As $\rho \to \infty$, the condition number of $F_{\rho}''(x_*(\rho))$ may blow up to ∞ , which, according to the traditional theory of the Newton method, makes the problems of updating x_i into x_{i+1} more and more difficult. Thus, slowing down seems to be unavoidable...

♣ In late 80's, it was discovered that the classical path-following scheme, associated with properly chosen barriers, admits "safe" implementation without slowing down. This discovery led to invention of Polynomial Time Interior Point methods for convex programs.

A Majority of Polynomial Time Interior Point methods heavily exploit the classical pathfollowing scheme; the novelty is in what are the underlying barriers – these are specific *self-concordant* functions especially well suited for Newton minimization. Let G be a closed convex domain with nonempty interior which does not contain lines. A 3 times continuously differentiable convex function

F(x): int $G \to \mathbb{R}$

is called self-concordant, if $\Diamond F$ is an interior penalty for G:

 $x_i \in \text{int}G, x_i \to x \in \partial G \Rightarrow F(x_i) \to \infty$

 $\Diamond F$ satisfies the relation

$$\left|\frac{d^3}{dt^3}\Big|_{t=0}F(x+th)\right| \le 2\left(\frac{d^2}{dt^2}\Big|_{t=0}F(x+th)\right)^{3/2}$$

• Let $\vartheta \ge 1$. F is called ϑ -self-concordant barrier for G, if, in addition to being selfconcordant on G, F satisfies the relation

$$\left|\frac{d}{dt}\right|_{t=0} F(x+th) \right| \le \sqrt{\vartheta} \left(\frac{d^2}{dt^2}\Big|_{t=0} F(x+th)\right)^{1/2}$$

 ϑ is called the *parameter* of s.-c.b. *F*.

Every convex program

$\min_{x\in G} f(x)$

can be converted into a convex program with *linear* objective, namely,

 $\min_{t,x} \left\{ t : x \in G, f(x) \le t \right\}.$

Assuming that this transformation has been done at the very beginning, we can w.l.o.g. focus on convex program with *linear* objective

$$\min_{x \in G} c^T x \tag{P}$$

$$\min_{x \in G} c^T x \tag{P}$$

Assume that G is a closed and bounded convex set with a nonempty interior, and let F be a ϑ -s.c.b. barrier for G.

 \Diamond <u>Fact I:</u> *F* is strongly convex on int*G*: *F*"(*x*) ≻ 0 for all *x* ∈ int*G*. Consequently,

$$F_{\rho}(x) \equiv \rho c^T x + F(x)$$

also is strongly convex on intG. In particular, the quantity

$$\lambda(x, F_{\rho}) = \left([F_{\rho}'(x)]^T \Big[\underbrace{[F_{\rho}''(x)]}_{=F''(x)} \Big]^{-1} F_{\rho}'(x) \right)^{1/2}$$

called the Newton decrement of F_{ρ} at x is well-defined for all $x \in \operatorname{int} G$ and all $\rho > 0$. Note: • $\frac{1}{2}\lambda^2(x, F_{\rho}) = F_{\rho}(x) - \min_y \left[F_{\rho}(x) + (y - x)^T F_{\rho}'(x) + \frac{1}{2}(y - x)^T F_{\rho}''(x)(y - x)\right]$ • $\lambda(x, F_{\rho}) \ge 0$ and $\lambda(x, F_{\rho}) = 0$ iff $x = x_*(\rho)$, so that the Newton decrement can be viewed as a "proximity measure" – a kind of distance from x to $x_*(\rho)$.

$$c_* = \min_{x \in G} c^T x \tag{P}$$

Fact II: Let (P) be solved via the classical penalty scheme implemented as follows:

 \diamond The barrier underlying the scheme is a ϑ -s.-c.b. F for G;

 \diamond "Closeness" of x and $x_*(\rho)$ is specified by the relation $\lambda(x, F_{\rho}) \leq 0.1$;

 \diamond The penalty update is $\rho_{i+1} = \left(1 + \frac{\gamma}{\sqrt{\vartheta}}\right)\rho_i$, where $\gamma > 0$ is a parameter;

 \diamond To update x_i into x_{i+1} , we apply to $F_{\rho_{i+1}}$ the Damped Newton method started as x_i :

$$x \mapsto x - \frac{1}{1 + \lambda(x, F_{\rho_{i+1}})} [F_{\rho_{i+1}}''(x)]^{-1} F_{\rho_{i+1}}'(x)$$

• The method is well-defined, and the number of damped Newton steps in updating $x_i \mapsto x_{i+1}$ depends solely on γ (and is as small as 1 for $\gamma = 0.1$)

• One has $c^T x_i - c_* \leq \frac{2\vartheta}{\rho_i}$

 \Rightarrow With the outlined method, it takes $O(\sqrt{\vartheta})$ Newton steps to reduce inaccuracy $c^T x - c_*$ by absolute constant factor! Fact III: • Every convex domain $G \subset \mathbb{R}^n$ admits O(n)-s.-c.b.

• For typical feasible domains arising in Convex Programming, one can point out explicit "computable" s.-c.b.'s. For example,

 \blacklozenge Let G be given by m convex quadratic constraints:

$$G = \{x : \underbrace{x^T A_j^T A_j x + 2b_j^T x + c_j}_{g_j(x)} \le 0, \ 1 \le j \le m\}$$

satisfying the Slater condition. Then the logarithmic barrier

$$F(x) = -\sum_{j=1}^{m} \ln(-g_j(x))$$

is m-s.-c.b. for G.

♠ Let G be given by Linear Matrix Inequality

$$G = \{x : \underbrace{A_0 + x_1 A_1 + \dots + x_n A_n}_{\mathcal{A}(x):m \times m} \succeq 0\}$$

satisfying the Slater condition: $\mathcal{A}(\bar{x}) \succ 0$ for some \bar{x} . Then the log-det barrier

$$F(x) = -\ln \mathsf{Det}(\mathcal{A}(x))$$

is m-s.-c.b. for G.

Consider an LP

$$\min_{z} \left\{ c^{T} z : A z - b \ge \mathbf{0} \right\}$$
(P)

with $m \times n$ matrix A, Null $(A) = \{0\}$, along with the dual problem

$$\max_{y} \left\{ b^{T} y : A^{T} y = c, \ y \ge 0 \right\}$$
 (D)

and assume that both problems are strictly feasible:

$$\exists \overline{z} : A\overline{z} - b > 0 \& \exists y > 0 : A^T y = c$$

<u>Note:</u> Passing from z to "primal slack" x = Az - b, we can rewrite (P) as

$$\min_{x} \left\{ e^{T}x : x \ge 0, x \in L = \operatorname{Im} A - b \right\} \qquad (P')$$

where e is a vector satisfying $A^T e = c$, so that

$$e^T x = e^T (Az-b) = (A^T e)^T z - \text{const} = c^T z - \text{const}$$

• Let $\Phi(x) = -\sum_{i=1}^{m} \ln x_i$. Equipping the domain of (P) with *m*-s.c.b. $F(z) = \Phi(Az - b)$, consider

$$z_*(\rho) = \arg\min_{z} [\rho c^T z + F(z)]$$

=
$$\arg\min_{z} [\rho e^T (Az - b) + \Phi(Az - b)]$$

<u>Observation</u>: The point $x_*(\rho) = Az_*(\rho) - b$ minimizes $\rho e^T x + \Phi(x)$ over the feasible set of (P'):

 $x > 0, x + b \in \text{Im}A, \rho e + \Phi'(x) \in (\text{Im}A)^{\perp}.$ $\Rightarrow y = y_*(\rho) = -\rho^{-1}\Phi'(x_*(\rho)) \text{ satisfies}$

 $y > 0, y - e \in (\text{Im}A)^{\perp}, -\rho b + \Phi'(y) \in \text{Im}A$ i.e., the point $y_*(\rho)$ minimizes $-\rho b^T y + \Phi(y)$ over the feasible set of (D).
♣ We arrive at a nice symmetric picture: ♣ The primal central path $x \equiv x_*(\rho)$ which minimizes the primal aggregate

$$\rho c^T x + \Phi(x) \qquad [\Phi(x) = -\sum_i \ln x_i]$$

over the primal feasible set is given by

$$x > 0, x + b \in \operatorname{Im} A, \rho c + \Phi'(x) \in (\operatorname{Im} A)^{\perp}$$

♣ The dual central path $y \equiv y_*(\rho)$ which minimizes the dual aggregate

$$-\rho b^T y + \Phi(y)$$
 $[\Phi(y) = -\sum_i \ln y_i]$

over the dual feasible set is given by

$$y > 0, y - e \in (\operatorname{Im} A)^{\perp}, -\rho b + \Phi'(y) \in \operatorname{Im} A$$

$$y = -\rho^{-1}\Phi'(x) \Leftrightarrow x = -\Phi'(y_*(\rho)) \Leftrightarrow x_i y_i = \frac{1}{\rho} \forall i.$$

 $\Rightarrow \mathsf{DualityGap}(x, y) = x^T y = \begin{bmatrix} c^T x - \mathsf{Opt}(P) \end{bmatrix} \\ + [\mathsf{Opt}(D) - b^T y]$

on the path is equal to $m\rho^{-1}$.

The paths are linked by

♣ <u>Generic Primal-Dual IPM for LP</u>: ♦Given current iterate — primal-dual strictly feasible pair x^i, y^i and value ρ_i of penalty, update it into new iterate $x^{i+1}, y^{i+1}, \rho_{i+1}$ by ♦Updating $\rho_i \mapsto \rho_{i+1} \ge \rho_i$

♦ Applying a Newton step to the system

$$x > 0, \ x + b \in \operatorname{Im}A; \ y > 0, \ y - e \in (\operatorname{Im}A)^{\perp}$$

$$\operatorname{Diag}\{x\}y = \frac{1}{\rho_{+}} \underbrace{(1, ..., 1)^{T}}_{e}$$

defining the primal-dual central path:

$$\begin{split} x^{i+1} &= x^i + \Delta x, \ y^{i+1} = y^i + \Delta y \\ \text{where } \Delta x, \Delta y \text{ solve the linear system} \\ \Delta x \in \text{Im}A, \ \Delta y \in (\text{Im}A)^{\perp}, \\ \text{Diag}\{x^i\}\Delta y + \text{Diag}\{y^i\}\Delta x = \frac{e}{\rho_{i+1}} - \text{Diag}\{x^i\}y^i \end{split}$$

The classical path-following scheme as applied to (P) and the *m*-s.c.b. $F(z) = \Phi(Az-b)$ allows to trace the path $z_*(\rho)$ (and thus $x_*(\rho) = Az_*(\rho) - b$).

More advanced *primal-dual* path-following methods *simultaneously* trace the primal and the dual central paths, which results in algorithmic schemes with better practical performance than the one of the "purely primal" scheme. A Both approaches, with proper implementation, result in the best known so far theoretical complexity bounds for LP. According to these bounds, the "arithmetic cost" of generating ϵ -solution to a primal-dual pair of strictly feasible LP's with $m \times n$ matrix Ais

 $O(1)mn^2 \ln\left(\frac{mn\Theta}{\epsilon}\right)$

operations, where O(1) is an absolute constant and Θ is a data-dependent constant.

♣ In practice, properly implemented primaldual methods by far outperform the purely primal ones and solve in few tens of Newton iterations real-world LPs with tens and hundreds of thousands of variables and constraints.

Augmented Lagrangian methods

$$\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$$
 (P)

♣ Shortcoming of penalty scheme: in order to solve (P) to high accuracy, one should work with large values of penalty, which makes the penalized objective

$$f_{\rho}(x) = f(x) + \frac{\rho}{2} ||h(x)||_2^2$$

difficult to minimize.

Augmented Lagrangian methods use the penalty mechanism in a "smart way", which allows to avoid the necessity to work with very large values of ρ .

Ingredient I: Local Lagrange Duality

$$\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$$
 (P)

Let x_* be a nondegenerate local solution to (P), so that there exists λ such that

(a)
$$\nabla_x L(x_*, \lambda^*) = 0$$

(b) $d^T \nabla_x^2 L(x_*, \lambda^*) d > 0 \ \forall 0 \neq d \in T_{x_*}$
 $\begin{bmatrix} L(x, \lambda) = f(x) + \sum_i \lambda_i h_i(x) \\ T_{x^*} = \{d : d^T h'_i(x) = 0, i = 1, ..., k\} \end{bmatrix}$

Assume for the time being that instead of (b), a stronger condition hods true:

(!) the matrix $\nabla_x^2 L(x_*, \lambda^*)$ is positive definite on the entire space

\clubsuit Under assumption (!), x_* is a nondegenerate unconstrained local minimizer of the smooth function

$L(\cdot,\lambda^*)$

and as such can be found by methods for unconstrained minimization.

 $\min_{x} \{f(x) : h_i(x) = 0, i = 1, ..., k\}$ (P)

• Intermediate Summary: If (a) we are clever enough to guess the vector λ^* of Lagrange multipliers, (b) we are lucky to have $\nabla_x^2 L(x_*, \lambda^*) \succ 0$, then x_* can be found by unconstrained optimization technique. ♠ How to become smart when being lucky: Local Lagrange Duality.

Situation: x_* is a nondegenerate local solution to

 $\min_{x} \{f(x) : h_i(x) = 0, i = 1, ..., k\}$ (P) and we are lucky:

 $\exists \lambda^* : \nabla_x L(x_*, \lambda^*) = 0, \ \nabla_x^2 L(x_*, \lambda^*) \succ 0 \quad (!)$

<u>Fact:</u> Under assumption (!), there exist convex neighbourhood V of x_* and convex neighbourhood Λ of λ^* such that

(i) For every $\lambda \in \Lambda$, function $L(x,\lambda)$ is strongly convex in $x \in V$ and possesses uniquely defined critical point $x_*(\lambda)$ in Vwhich is continuously differentiable in $\lambda \in \Lambda$. $x_*(\lambda)$ is a nondegenerate local minimizer of $L(\cdot, \lambda)$;

(ii) The function

$$\underline{L}(\lambda) = L(x_*(\lambda), \lambda) = \min_{x \in V} L(x, \lambda)$$

is C²-smooth and concave in Λ ,

$$\underline{L}'(\lambda) = h(x_*(\lambda)),$$

and λ_* is a nondegenerate maximizer of $\underline{L}(\lambda)$ on Λ .

$$\min_{x} \{f(x) : h_i(x) = 0, i = 1, ..., k\} \quad (P)$$

$$\Rightarrow \quad L(x, \lambda) = f(x) + \sum_i \lambda_i h_i(x)$$

Situation: $\nabla_x L(x_*, \lambda^*) = 0, \ \nabla_x^2 L(x_*, \lambda^*) \succ 0$

$$\lambda^* = \operatorname{argmax} \underline{L}(\lambda) = \min_{x \in V} L(x, \lambda)$$

 $x_* = \operatorname{argmin}_{x \in V} L(x, \lambda)$

⇒ We can solve (P) by maximizing $\underline{L}(\lambda)$ over $\lambda \in \Lambda$ by a first order method for unconstrained minimization.

The first order information on $\underline{L}(\lambda)$ required by the method can be obtained by solving auxiliary unconstrained problems

$$x_*(\lambda) = \underset{x \in V}{\operatorname{argmin}} L(x, \lambda)$$

via

$$\underline{L}(\lambda) = L(x_*(\lambda), \lambda)$$

$$\underline{L}'(\lambda) = h(x_*(\lambda))$$

<u>Note:</u> In this scheme, there are no "large parameters"! However: How to ensure luck? How to ensure luck: convexification by penalization

Observe that the problem of interest

$$\min_{x} \{f(x) : h_i(x) = 0, i = 1, ..., k\}$$
(P)

for every $\rho \ge 0$ is exactly equivalent to

$$\min_{x} \left\{ f_{\rho}(x) = f(x) + \frac{\rho}{2} \|h(x)\|_{2}^{2} : \begin{array}{c} h_{i}(x) = 0, \\ i \le k \end{array} \right\} \\ (P_{\rho})$$

It turns out that

(!) If x_* is a nondegenerate locally optimal solution of (P) and ρ is large enough, then x_* is a locally optimal and "lucky" solution to (P_{ρ}) .

⇒ We can solve (P) by applying the outlined "primal-dual" scheme to (P_ρ), provided that ρ is appropriately large!

<u>Note</u>: Although in our new scheme we do have penalty parameter which should be "large enough", we still have an advantage over the straightforward penalty scheme: in the latter, ρ should go to ∞ as $O(1/\epsilon)$ as required inaccuracy ϵ of solving (P) goes to 0, while in our new scheme <u>a single</u> "large enough" value of ρ will do!

$$\min_{x} \{ f(x) : h_i(x) = 0, i = 1, ..., k \}$$
(P)

 $\min_{x} \left\{ f_{\rho}(x) = f(x) + \frac{\rho}{2} \|h(x)\|_{2}^{2} : \begin{array}{c} h_{i}(x) = 0, \\ i \leq k \end{array} \right\} \quad (P_{\rho})$ <u>Justifying the claim:</u> Let

$$L_{\rho}(x,\lambda) = f(x) + \frac{\rho}{2} ||h(x)||_{2}^{2} + \sum_{i} \lambda_{i} h_{i}(x)$$

be the Lagrange function of (P_{ρ}) ; the Lagrange function of (P) is then $L_o(x, \lambda)$. Given nondegenerate locally optimal solution x_* to (P), let λ^* be the corresponding Lagrange multipliers. We have

$$\nabla_{x}L_{\rho}(x_{*},\lambda^{*}) = \nabla_{x}L_{0}(x_{*},\lambda^{*}) + \rho \sum_{i} h_{i}(x_{*})h'_{i}(x_{*}) \\
= \nabla_{x}L_{0}(x_{*},\lambda^{*}) = 0 \\
\nabla_{x}^{2}L_{\rho}(x_{*}\lambda^{*}) = \nabla_{x}^{2}L(x_{*},\lambda^{*}) + \rho \sum_{i} h_{i}(x_{*})h''_{i}(x_{*}) \\
+ \rho \sum_{i} h'_{i}(x_{*})[h'_{i}(x_{*})]^{T} \\
= \nabla_{x}^{2}L_{0}(x_{*},\rho^{*}) + \rho H^{T}H, \\
H = \begin{bmatrix} [h'_{1}(x_{*})]^{T} \\ \cdots \\ [h'_{k}(x_{*})]^{T} \end{bmatrix}$$

$$\nabla_x^2 L_\rho(x_*\lambda^*) = \nabla_x^2 L_0(x_*,\rho^*) + \rho H^T H$$
$$H = \begin{bmatrix} [h'_1(x_*)]^T \\ \cdots \\ [h'_k(x_*)]^T \end{bmatrix}$$

Directions d orthogonal to $h'_i(x_*)$, i = 1, ..., k, are exactly the directions d such that Hd = 0. Thus,

 \diamond For all $\rho \ge 0$, at x_* the Second Order sufficient optimality condition for (P_{ρ}) holds true:

$$Hd = 0, d \neq 0 \Rightarrow d^T \nabla_x^2 L_\rho(x_*, \lambda^*) d > 0$$

All we need in order to prove that x^* is a "lucky" solution for large ρ , is the following Linear Algebra fact:

Let Q be a symmetric $n \times n$ matrix, and Hbe a $k \times n$ matrix. Assume that Q is positive definite on the null space of H:

$$d \neq 0, Hd = 0 \Rightarrow d^T Qd > 0.$$

Then for all large enough values of ρ the matrix $Q + \rho H^T H$ is positive definite.

Let Q be a symmetric $n \times n$ matrix, and H be a $k \times n$ matrix. Assume that Q is positive definite on the null space of H:

$$d \neq 0, Hd = 0 \Rightarrow d^T Qd > 0.$$

Then for all large enough values of ρ the matrix $Q + \rho H^T H$ is positive definite.

Proof: Assume, on the contrary, that there exists a sequence $\rho_i \rightarrow \infty$ and d_i , $||d_i||_2 = 1$:

$$d_i^T [Q + \rho_i H^T H] d_i \le 0 \ \forall i.$$

Passing to a subsequence, we may assume that $d_i \to d$, $i \to \infty$. Let $d_i = h_i + h_i^{\perp}$ be the decomposition of d_i into the sum of its projections onto Null(H) and [Null(H)]^{\perp}, and similarly $d = h + h^{\perp}$. Then

$$d_i^T H^T H d_i = \|Hd_i\|_2^2 = \|Hh_i^{\perp}\|_2^2 \to \|Hh^{\perp}\|_2^2 \Rightarrow$$

$$0 \ge d_i^T [Q + \rho_i H^T H] d_i = \underbrace{d_i^T Q d_i}_{\rightarrow d^T Q d} + \rho_i \underbrace{\|Hh_i^{\perp}\|_2^2}_{\rightarrow \|Hh^{\perp}\|_2^2} (*)$$

If $h^{\perp} \neq 0$, then $||Hh^{\perp}||_2 > 0$, and the right hand side in (*) tends to $+\infty$ as $i \to \infty$, which is impossible. Thus, $h^{\perp} = 0$. But then $0 \neq d \in \text{Null}(H)$ and therefore $d^T Q d > 0$, so that the right hand side in (*) is positive for large *i*, which again is impossible. Putting things together: Augmented Lagrangian Scheme

Generic Augmented Lagrangian Scheme: For a given value of ρ , solve the dual problem

$$\max_{\lambda} \underline{L}_{\rho}(\lambda)$$

$$\left[\underline{L}_{\rho}(\lambda) = \min_{x} L_{\rho}(x, \lambda)\right]$$
(D)

by a first order method for unconstrained minimization, getting the first order information for (D) from solving the auxiliary problems

$$x_{\rho}(\lambda) = \underset{x}{\operatorname{argmin}} L_{\rho}(x, \lambda) \qquad (P^{\lambda})$$

via the relations

$$\underline{L}_{\rho}(\lambda) = L_{\rho}(x_{\rho}(\lambda), \lambda)$$

$$\underline{L}'_{\rho}(\lambda) = h(x_{\rho}(\lambda))$$

<u>Note:</u> If ρ is large enough and the optimizations in (P^{λ}) and in (D) and are restricted to appropriate convex neighbourhoods of nondegenerate local solution x_* to (P_{ρ}) and the corresponding vector λ^* of Lagrange multipliers, respectively, then

— the objective in (D) is concave and C^2 , and λ^* is a nondegenerate solution to (D)— the objectives in (P^{λ}) are convex and C^2 , and $x_*(\lambda) = \underset{x}{\operatorname{argmin}} L_{\rho}(x, \lambda)$ are nondegenerate local solutions to (P^{λ})

— as the "master method" working on (D) converges to λ^* , the corresponding primal iterates $x_*(\lambda)$ converge to x_* .

Implementation issues: Solving auxiliary problems

$$x_{
ho}(\lambda) = \operatorname{argmin}_{x} L_{
ho}(x, \lambda)$$
 (P^{λ})

— the best choices are Newton method with linesearch or Modified Newton method, provided that the second order information is available; otherwise, one can use Quasi-Newton methods, Conjugate Gradients, etc. ♦ Solving the master problem

$$\max_{\lambda} \left\{ \underline{L}_{\rho}(\lambda) \equiv \min_{x} L_{\rho}(x,\lambda) \right\}$$
 (D)

Surprisingly, the method of choice here is the simplest gradient ascent method with constant step:

$$\lambda^t = \lambda^{t-1} + \rho \underline{L}'_{\rho}(\lambda^{t-1}) = \lambda^{t-1} + \rho h(x^{t-1}),$$

where x^{t-1} is (approximate) minimizer of $L_{\rho}(x, \lambda^{t-1})$ in x. <u>Motivation:</u> We have

$$\begin{array}{rcl} 0 &\approx & \nabla_x L_{\rho}(x^{t-1}, \lambda^{t-1}) \\ &= & f'(x^{t-1}) + \sum\limits_i [\lambda_i^{t-1} + \rho h_i(x^{t-1})] h_i'(x^{t-1}) \end{array}$$

which resembles the KKT condition

$$0 = f'(x_*) + \sum_i \lambda_i^* h'_i(x_*).$$

$$\max_{\lambda} \left\{ \underline{L}_{\rho}(\lambda) \equiv \min_{x} L_{\rho}(x,\lambda) \right\}$$
 (D)

$$\Rightarrow \begin{cases} \lambda^t = \lambda^{t-1} + \rho h(x^{t-1}) \\ x^{t-1} = \operatorname{argmin}_x L_\rho(x, \lambda^{t-1}) \end{cases} (*)$$

Justification: Direct computation shows that

$$\Psi_{\rho} \equiv \nabla_{\lambda}^{2} \underline{L}_{\rho}(\lambda^{*}) = -H[Q + \rho H^{T}H]^{-1}H^{T},$$

$$Q = \nabla_{x}^{2} L_{0}(x_{*}, \lambda^{*})$$

$$H = \begin{bmatrix} [h'_{1}(x_{*})]^{T} \\ \cdots \\ [H'_{k}(x_{*})]^{T} \end{bmatrix}$$

whence $-\rho\Psi_{\rho} \rightarrow I$ as $\rho \rightarrow \infty$.

Consequently, when ρ is large enough and the starting point λ_0 in (*) is close enough to λ^* , (*) ensures linear convergence of λ^t to λ^* with the ratio tending to 0 as $\rho \to +\infty$. Indeed, asymptotically the behaviour of (*) is as if $\underline{L}_{\rho}(\lambda)$ were quadratic function

$$\operatorname{const} - \frac{1}{2} (\lambda - \lambda^*)^T \Psi_{\rho} (\lambda - \lambda_*),$$

and for this model recurrence (*) becomes

$$\lambda^{t} - \lambda^{*} = \underbrace{(I + \rho \Psi_{\rho})}_{\to 0, \rho \to \infty} (\lambda^{t-1} - \lambda^{*}).$$

Adjusting penalty parameter:

$$\Rightarrow \begin{cases} \lambda^t = \lambda^{t-1} + \rho h(x^{t-1}) \\ x^{t-1} = \operatorname{argmin}_x L_\rho(x, \lambda^{t-1}) \end{cases} (*)$$

When ρ is "large enough", so that (*) converges linearly with reasonable convergence ratio, $\|\underline{L}'_{\rho}(\lambda^t)\|_2 = \|h(x^t)\|_2$ should go to 0 linearly with essentially the same ratio.

⇒ We can use progress in $||h(\cdot)||_2$ to control ρ , e.g., as follows: when $||h(x^t)||_2 \leq 0.25 ||h(x^{t-1})||_2$, we keep the current value of ρ intact, otherwise we increase penalty by factor 10 and recompute x^t with the new value of ρ .

Incorporating Inequality Constraints

Given a general-type constrained problem

$$\min_{x} \left\{ f(x): \begin{array}{l} h_i = 0, i \leq m \\ g_j(x) \leq 0, j \leq m \end{array} \right\}$$

we can transform it equivalently into the equality constrained problem

$$\min_{x,s} \left\{ f(x): \begin{array}{l} h_i(x) = 0, i \le m \\ g_j(x) + s_j^2 = 0, j \le k \end{array} \right\}$$

and apply the Augmented Lagrangian scheme to the reformulated problem, thus arriving at Augmented Lagrangian

$$L_{\rho}(x,s;\lambda,\mu) = f(x) + \sum_{i} \lambda_{i} h_{i}(x) + \sum_{j} \mu_{j} [g_{j}(x) + s_{j}^{2}] + \frac{\rho}{2} \left[\sum_{i} h_{i}^{2}(x) + \sum_{j} [g_{j}(x) + s_{j}^{2}]^{2} \right]$$

The corresponding dual problem is

$$\max_{\lambda,\mu} \left\{ \underline{L}_{\rho}(\lambda,\mu) = \min_{x,s} L_{\rho}(x,s;\mu,\lambda) \right\} \qquad (D)$$

$$L_{\rho}(x,s;\lambda,\mu) = f(x) + \sum_{i} \lambda_{i}h_{i}(x) + \sum_{j} \mu_{j}[g_{j}(x) + s_{j}^{2}] + \frac{\rho}{2} \left[\sum_{i} h_{i}^{2}(x) + \sum_{j} [g_{j}(x) + s_{j}^{2}]^{2} \right] \downarrow \max_{\lambda,\mu} \left\{ \underline{L}_{\rho}(\lambda,\mu) \equiv \min_{x,s} L_{\rho}(x,s;\mu,\lambda) \right\}$$

We can carry out the minimization in s analytically, arriving at

$$\underline{L}_{\rho}(\lambda,\mu) = \min_{x} \left\{ f(x) + \frac{\rho}{2} \sum_{j=1}^{k} \left(g_{j}(x) + \frac{\mu_{j}}{\rho} \right)_{+}^{2} \right. \\ \left. + \sum_{i=1}^{m} \lambda_{i} h_{i}(x) + \frac{\rho}{2} \sum_{i=1}^{m} h_{i}(x)^{2} \right\} \\ \left. - \sum_{j=1}^{k} \frac{\mu_{j}^{2}}{2\rho} \right\}$$

where $a_+ = \max[0, a]$.

⇒ The auxiliary problems arising in the Augmented Lagrangian Scheme are problems in the initial design variables!

$$\min_{x} \left\{ f(x): \begin{array}{l} h_{i}(x) = 0, \ i \leq k \\ g_{j}(x) \leq 0, \ j \leq m \end{array} \right\} \quad (P) \\
\downarrow \\
\min_{x,s} \left\{ f(x): \begin{array}{l} h_{i}(x) = 0, \ i \leq k \\ g_{j}(x) + s_{j}^{2} = 0, \ j \leq m \end{array} \right\} \quad (P')$$

♣ Theoretical analysis of Augmented Lagrangian scheme for problems with equality constraints was based on assumption that we are trying to approximate nondegenerate local solution. Is it true that when applying reducing the inequality constrained problem to an equality constrained one, we preserve nondegeneracy of the local solution? Yes!

<u>Theorem.</u> Let x_* be a nondegenerate local solution to (P). Then the point

$$(x_*, s^*)$$
: $s_j^* = \sqrt{-g_j(x_*)}, j = 1, ..., m$

is a <u>nondegenerate</u> local solution to (P').

Convex case: Augmented Lagrangians

Consider a <u>convex</u> optimization problem

$$\min_{x} \left\{ f(x) : g_j(x) \le 0, \, j = 1, ..., m \right\}$$
(P)

(*f*, g_j are convex and C² on \mathbb{R}^n). <u>Assumption</u>: (*P*) is solvable and satisfies the Slater condition:

$$\exists \bar{x} : g_j(\bar{x}) < 0 \ j = 1, ..., m$$

♠ In the convex situation, the previous local considerations can be globalized due to the Lagrange Duality Theorem.

$$\min_{x} \left\{ f(x) : g_j(x) \le 0, \, j = 1, ..., m \right\}$$
(P)

<u>Theorem</u>: Let (P) be convex, solvable and satisfy the Slater condition. Then the <u>dual</u> problem

$$\max_{\lambda \ge 0} \underline{L}(\lambda) \equiv \min_{x} \left[f(x) + \sum_{j} \lambda_{j} g_{j}(x) \right] \quad (D)$$
$$\underbrace{L(x,\lambda)}$$

possess the following properties:

 \diamond dual objective <u>L</u> is concave

 $\Diamond(D)$ is solvable

 \Diamond for every optimal solution λ^* of (D), all optimal solutions of (P) are contained in the set Argmin_x $L(x, \lambda^*)$.

Implications:

 \diamond Sometimes we can build (*D*) explicitly (e.g., in Linear, Linearly Constrained Quadratic and Geometric Programming). In these cases, we may gain a lot by solving (*D*) and then recovering solutions to (*P*) from solution to (*D*).

$$\min_{x} \left\{ f(x) : g_{j}(x) \leq 0, \ j = 1, ..., m \right\} \quad (P)$$

$$\max_{\lambda \geq 0} \underline{L}(\lambda) \equiv \min_{x} \left[f(x) + \sum_{j} \lambda_{j} g_{j}(x) \right] \quad (D)$$

$$\underline{L}(x,\lambda)$$

 \Diamond In the general case one can solve (D) numerically by a first order method, thus reducing a problem with general convex constraints to one with simple linear constraints. To solve (D) numerically, we should be able to compute the first order information for <u>L</u>. This can be done via solving the auxiliary problems

$$x_* = x_*(\lambda) = \min_x L(x,\lambda)$$
 (P_{\lambda})

due to

$$\underline{L}(\lambda) = L(x_*(\lambda), \lambda)$$

$$\underline{L}'(\lambda) = g(x_*(\lambda))$$

<u>Note:</u> (P_{λ}) is a convex unconstrained program with smooth objective!

$$\min_{x} \left\{ f(x) : g_{j}(x) \leq 0, \ j = 1, ..., m \right\} \quad (P)$$

$$\max_{\lambda \geq 0} \underline{L}(\lambda) \equiv \min_{x} \left[f(x) + \sum_{j} \lambda_{j} g_{j}(x) \right] \quad (D)$$

$$\underline{L}(x,\lambda)$$

Potential difficulties:

 $\Diamond \underline{L}(\cdot)$ can be $-\infty$ at some points; how to solve (D)?

 \Diamond after λ^* is found, how to recover optimal solution to (*P*)? In may happen that the set Argmin_x $L(x, \lambda^*)$ is much wider than the optimal set of (*P*)!

Example: LP. (P) : $\min_{x} \{c^T x : Ax - b \leq 0\}$. Here

$$\underline{L}(\lambda) = \min_{x} \left[c^{T}x + (A^{T}\lambda)^{T}x - b^{T}\lambda \right]$$
$$= \begin{cases} -b^{T}\lambda, & A^{T}\lambda + c = 0\\ -\infty, & \text{otherwise} \end{cases}$$

— how to solve (D) ???

At the same time, for every λ the function $L(x,\lambda)$ is linear in x; thus, $\underset{x}{\operatorname{Argmin}} L(x,\lambda)$ is either \emptyset , or \mathbb{R}^n – how to recover x_* ???

Observation: Both outlined difficulties come from possible non-existence/non-uniqueness of solutions to the auxiliary problems

$$\min_{x} L(x,\lambda) \equiv \min_{x} [f(x) + \sum_{j} \lambda_{j} g_{j}(x)] \quad (P_{\lambda})$$

Indeed, if solution $x_*(\lambda)$ to (P_{λ}) exists and is unique and continuous in λ on certain set Λ , then $\underline{L}(\lambda)$ is finite and continuously differentiable on Λ due to

$$\underline{L}(\lambda) = L(x_*(\lambda), \lambda)$$

$$\underline{L}'(\lambda) = g(x_*(\lambda))$$

Besides this, if $\lambda^* \in \Lambda$, then there is no problem with recovering optimal solution to (*P*) from λ_* .

Example: Assume that the function

$$r(x) = f(x) + \sum_{j=1}^{k} g_j(x)$$

is locally strongly convex $(r''(x) \succ 0 \forall x)$ and is such that

$$r(x)/\|x\|_2 \to \infty, \ \|x\|_2 \to \infty.$$

Then $x_*(\lambda)$ exists, is unique and is continuous in λ on the set $\Lambda = \{\lambda > 0\}$.

In Augmented Lagrangian scheme, we ensure local strong convexity of

 $r(\cdot) = f(x) + \text{sum of constraints}$

by passing from the original problem

 $\min_{x} \left\{ f(x) : g_j(x) \le 0, \ j = 1, ..., m \right\}$ (P) to the *equivalent* problem

 $\min_{x} \left\{ f(x) : \theta_{j}(g_{j}(x)) \leq 0, \ j = 1, ..., m \right\} \quad (P')$ where $\theta_{j}(\cdot)$ are *increasing strongly convex* smooth functions satisfying the normalization

 $\theta_j(0) = 0, \ \theta'_j(0) = 1.$

$$\min_{x} \left\{ f(x) : g_{j}(x) \leq 0, \ j = 1, ..., m \right\} \quad (P) \\
\downarrow \\
\min_{x} \left\{ f(x) : \theta_{j}(g_{j}(x)) \leq 0, \ j = 1, ..., m \right\} \quad (P') \\
\left[\theta_{j}(0) = 0, \ \theta_{j}'(0) = 1 \right]$$

Facts:

 $\Diamond(P')$ is convex and equivalent to (P) \Diamond optimal Lagrange multipliers for (P) and (P') are the same:

 $\nabla_x [f(x) + \sum_j \lambda_j^* g_j(x)] = 0 & \lambda_j^* g_j(x) = 0 \forall j$ $\widehat{\nabla}_x [f(x) + \sum_j \lambda_j^* \theta_j(g_j(x))] = 0 & \lambda_j^* g_j(x) = 0 \forall j$

♦under mild regularity assumptions,

$$r(x) = f(x) + \sum_{j} \theta_{j}(g_{j}(x))$$

is locally strongly convex and $r(x)/||x||_2 \to \infty$ as $||x||_2 \to \infty$.

$$\min_{x} \left\{ f(x) : g_{j}(x) \leq 0, \ j = 1, ..., m \right\} \quad (P) \\
\downarrow \\
\min_{x} \left\{ f(x) : \theta_{j}(g_{j}(x)) \leq 0, \ j = 1, ..., m \right\} \quad (P') \\
\left[\theta_{j}(0) = 0, \ \theta_{j}'(0) = 1 \right]$$

\clubsuit With the outlined scheme, one passes from the classical Lagrange function of (P)

$$L(x,\lambda) = f(x) + \sum_{j} \lambda_{j} g_{j}(x)$$

to the augmented Lagrange function

$$\widetilde{L}(x,\lambda) = f(x) + \sum_{j} \lambda_{j} \theta_{j}(g_{j}(x))$$

of the problem, which yields the dual problem

$$\max_{\lambda \ge 0} \underline{\widetilde{L}}(\lambda) \equiv \max_{\lambda \ge 0} \min_{x} \widetilde{L}(x, \lambda)$$

better suited for numerical solution and recovering a solution to (P) than the usual Lagrange dual of (P).

Further flexibility is added by penalty mechanism:

$$\widetilde{L}(x,\lambda) \Rightarrow f(x) + \sum_{j} \lambda_{j} \rho^{-1} \theta_{j}(\rho g_{j}(x))$$

equivalent to "rescaling"

$$\theta_j(s) \Rightarrow \theta_j^{(\rho)}(s) = \rho^{-1} \theta_j(\rho s).$$

The larger is ρ , the faster is convergence of the first order methods as applied to (\widetilde{D}) and the more difficult become the auxiliary problems

$$\min_{x} \left[f(x) + \sum_{j} \lambda_{j} \rho^{-1} \theta_{j}(\rho g_{j}(x)) \right]$$

Sequential Quadratic Programming

SQP is thought of to be the most efficient technique for solving general-type optimization problems with smooth objective and constraints.

SQP methods directly solve the KKT system of the problem by a Newton-type iterative process. • Consider an equality constrained problem $\min_{x} \left\{ f(x) : h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0 \right\} \quad (P)$ $\Rightarrow L(x, \lambda) = f(x) + h^{T}(x)\lambda$

The KKT system of the problem is

$$\nabla_{x}L(x,\lambda) \equiv f'(x) + [h'(x)]^{T}\lambda = 0$$

$$\nabla_{\lambda}L(x,\lambda) \equiv h(x) = 0$$

(KKT)

Every locally optimal solution x_* of (P) which is regular (that is, the gradients $\{h'_i(x_*)\}_{i=1}^k$ are linearly independent) can be extended by properly chosen $\lambda = \lambda^*$ to a solution of (KKT).

(KKT) is a system of nonlinear equations with n + k equations and n + k unknowns. We can try to solve this system by Newton method.

Newton method for solving nonlinear systems of equations

\clubsuit To solve a system of N nonlinear equations with N unknowns

 $P(u) \equiv (p_1(u), ..., p_N(u))^T = 0,$

with C^1 real-valued functions p_i , we act as follows:

Given current iterate \bar{u} , we linearize the system at the iterate, thus arriving at the linearized system

$$P(\bar{u}) + P'(\bar{u})(u - \bar{u}) \\ \equiv \begin{bmatrix} p_1(\bar{u}) + [p'_1(\bar{u})]^T(u - \bar{u}) \\ \vdots \\ p_N(\bar{u}) + [p'_N(\bar{u})]^T(u - \bar{u}) \end{bmatrix} = 0.$$

Assuming the $N \times N$ matrix $P'(\bar{u})$ nonsingular, we solve the linearized system, thus getting the new iterate

$$\bar{u}^{+} = \bar{u} \underbrace{-[P'(\bar{u})]^{-1}P(\bar{u})}_{\text{Newton}};$$
displacement

$$\bar{u} \mapsto \bar{u}^+ = \bar{u} - [P'(\bar{u})]^{-1} P(\bar{u}) \qquad (N)$$

<u>Note:</u> The Basic Newton method for unconstrained minimization is nothing but the outlined process as applied to the Fermat equation

$$P(x) \equiv \nabla f(x) = 0.$$

Same as in the optimization case, the Newton method possesses fast local convergence: <u>Theorem.</u> Let $u_* \in \mathbb{R}^N$ be a solution to the square system of nonlinear equations

$$P(u) = 0$$

with components of P being C^1 in a neighbourhood of u_* . Assuming that u_* is nondegenerate (i.e., $Det(P'(u_*)) \neq 0$), the Newton method (N), started close enough to u_* , converges to u_* superlinearly.

If, in addition, the components of P are C^2 in a neighbourhood of u_* , then the above convergence is quadratic.

Applying the outlined scheme to the KKT system

$$\nabla_{x}L(x,\lambda) \equiv f'(x) + [h'(x)]^{T}\lambda = 0$$

$$\nabla_{\lambda}L(x,\lambda) \equiv h(x) = 0$$

(KKT)

we should answer first of all the following crucial question:

(?) When a KKT point (x_*, λ^*) is a nondegenerate solution to (KKT)?

Let us set

$$P(x,\lambda) = \nabla_{x,\lambda} L(x,\lambda)$$

=
$$\begin{bmatrix} \nabla_x L(x,\lambda) \equiv f'(x) + [h'(x)]^T \lambda \\ \nabla_\lambda L(x,\lambda) \equiv h(x) \end{bmatrix}$$

Note that

$$P'(x,\lambda) = \begin{bmatrix} \nabla_x^2 L(x,\lambda) & [h'(x)]^T \\ h'(x) & 0 \end{bmatrix}$$
$$\min_{x} \left\{ f(x) : h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0 \right\} (P)$$

$$\Rightarrow \quad L(x, \lambda) = f(x) + h^{T}(x)\lambda$$

$$\Rightarrow \quad P(x, \lambda) = \nabla_{x,\lambda}L(x, \lambda)$$

$$= \left[\begin{array}{c} \nabla_{x}L(x, \lambda) \equiv f'(x) + [h'(x)]^{T}\lambda \\ \nabla_{\lambda}L(x, \lambda) \equiv h(x) \end{array} \right]$$

$$\Rightarrow \quad P'(x, \lambda) = \left[\begin{array}{c} \nabla_{x}^{2}L(x, \lambda) & [h'(x)]^{T} \\ h'(x) & 0 \end{array} \right]$$

<u>Theorem.</u> Let x_* be a nondegenerate local solution to (P) and λ^* be the corresponding vector of Lagrange multipliers. Then (x_*, λ^*) is a nondegenerate solution to the KKT system

 $P(x,\lambda)=0,$

that is, the matrix $P' \equiv P'(x_*, \lambda^*)$ is nonsingular.

Proof. Setting $Q = \nabla_x^2 L(x_*, \lambda^*)$, $H = \nabla h(x_*)$, we have

$$P' = \left[\begin{array}{cc} Q & H^T \\ H & 0 \end{array} \right]$$

$$Q = \nabla_x^2 L(x_*, \lambda^*), \ H = \nabla h(x_*),$$
$$P' = \begin{bmatrix} Q & H^T \\ H & 0 \end{bmatrix}$$

<u>We know</u> that $d \neq 0, Hd = 0 \Rightarrow d^T Qd > 0$ and that the rows of H are linearly independent. We should prove that if

$$0 = P' \begin{bmatrix} d \\ g \end{bmatrix} \equiv \begin{bmatrix} Qd + H^Tg \\ Hd \end{bmatrix},$$

then d = 0, g = 0. We have Hd = 0 and

$$0 = Qd + H^T g \Rightarrow d^T Qd + (Hd)^T g = d^T Qd,$$

which, as we know, is possible iff d = 0. We now have $H^Tg = Qd + H^Tg = 0$; since the rows of H are linearly independent, it follows that g = 0.

Structure and interpretation of the Newton displacement

In our case the Newton system

$$P'(u)\Delta = -P(u)$$
 $[\Delta = u^+ - u]$

becomes

$$\begin{aligned} [\nabla_x^2 L(\bar{x}, \bar{\lambda})] \Delta x + [\nabla h(\bar{x})]^T \Delta \lambda &= -f'(\bar{x}) \\ & -[h'(\bar{x})]^T \lambda \\ [h'(\bar{x})] \Delta \lambda &= -h(\bar{x}) \end{aligned}$$

where $(\bar{x}, \bar{\lambda})$ is the current iterate. Passing to the variables Δx , $\lambda^+ = \bar{\lambda} + \Delta \lambda$, the system becomes

$$\begin{bmatrix} \nabla_x^2 L(\bar{x}, \bar{\lambda}) \end{bmatrix} \Delta x + \begin{bmatrix} h'(\bar{x}) \end{bmatrix}^T \lambda^+ = -f'(\bar{x}) \\ h'(\bar{x}) \Delta x = -h(\bar{x}) \end{bmatrix}$$

$$\begin{bmatrix} \nabla_x^2 L(\bar{x}, \bar{\lambda}) \end{bmatrix} \Delta x + \begin{bmatrix} h'(\bar{x}) \end{bmatrix}^T \lambda^+ = -f'(\bar{x}) \\ h'(\bar{x}) \Delta x = -h(\bar{x}) \end{bmatrix}$$

Interpretation.

Assume for a moment that we know the optimal Lagrange multipliers λ^* and the tangent plane T to the feasible surface at x_* . Since $\nabla_x^2 L(x_*, \lambda^*)$ is positive definite on T, and $\nabla_x L(x_*, \lambda^*)$ is orthogonal to T, x_* is a nondegenerate local minimizer of $L(x, \lambda^*)$ over $x \in T$, and we could find x_* by applying the Newton minimization method to the function $L(x, \lambda^*)$ restricted onto T:

$$\bar{x} \in T \mapsto \bar{x} + \underset{\bar{x} + \Delta x \in T}{\operatorname{argmin}} \left[L(\bar{x}, \lambda^*) + \Delta x^T \nabla_x L(\bar{x}, \lambda^*) + \frac{1}{2} \Delta x^T \nabla_x^2 L(\bar{x}, \lambda^*) \Delta x \right]$$

In reality we do not know neither λ^* , nor T, only current approximations \overline{x} , $\overline{\lambda}$ of x_* and λ^* . We can use these approximations to *approximate* the outlined scheme:

• Given \bar{x} , we approximate T by the plane

$$\overline{T} = \{y = \overline{x} + \Delta x : [h'(\overline{x})]\Delta x + h(\overline{x}) = 0\}$$

• We apply the outlined step with λ^* replaced with $\bar{\lambda}$ and T replaced with \bar{T} :

$$\bar{x} \in T \mapsto \bar{x} + \underset{\bar{x} + \Delta x \in \bar{T}}{\operatorname{argmin}} \left[L(\bar{x}, \bar{\lambda}) + \Delta x^T \nabla_x L(\bar{x}, \bar{\lambda}) + \frac{1}{2} \Delta x^T \nabla_x^2 L(\bar{x}, \bar{\lambda}) \Delta x \right]$$

Note: Step can be simplified to

$$\bar{x} \in T \mapsto \bar{x} + \underset{\bar{x} + \Delta x \in \bar{T}}{\operatorname{argmin}} \left[f(\bar{x}) + \Delta x^T f'(\bar{x}) + \frac{1}{2} \Delta x^T \nabla_x^2 L(\bar{x}, \bar{\lambda}) \Delta x \right]$$

due to the fact that for $\bar{x}+\Delta x\in\bar{T}$ one has

$$\Delta x^T \nabla_x L(\bar{x}, \bar{\lambda}) = \Delta x^T f'(\bar{x}) + \bar{\lambda}^T [h'(\bar{x})] \Delta x$$

= $\Delta x^T f'(\bar{x}) - \bar{\lambda}^T h(\bar{x})$

& We have arrived at the following scheme: Given approximations $(\bar{x}, \bar{\lambda})$ to a nondegenerate KKT point x_*, λ^* of equality constrained problem

$$\min_{x} \left\{ f(x) : h(x) \equiv (h_1(x), ..., h_k(x))^T = 0 \right\}$$
(P)

solve the auxiliary quadratic program

$$\begin{split} \min_{\Delta x} \left\{ f(\bar{x}) + \Delta x^T f'(\bar{x}) + \frac{1}{2} \Delta x^T \nabla_x^2 L(\bar{x}, \bar{\lambda}) \Delta x : \\ h(\bar{x}) + h'(\bar{x}) \Delta x = 0 \right\} \end{split}$$

$$(QP)$$

and replace \bar{x} with $\bar{x} + \Delta x_*$.

<u>Note</u>: (QP) is a nice Linear Algebra problem, provided that $\nabla^2 L(\bar{x}, \bar{\lambda})$ is positive definite on the feasible plane $\bar{T} = \{\Delta x : h(\bar{x}) + h'(\bar{x})\Delta x = 0\}$ (which indeed is the case when $(\bar{x}, \bar{\lambda})$ is close enough to $(x_*\lambda^*)$).

$$\min_{x} \left\{ f(x) : h(x) \equiv (h_1(x), ..., h_k(x))^T = 0 \right\}$$
(P)

Step of the Newton method as applied to the KKT system of (P):

$$(\bar{x},\bar{\lambda}) \mapsto (\bar{x}^{+} = \bar{x} + \Delta x, \lambda^{+}) :$$

$$[\nabla_{x}^{2}L(\bar{x},\bar{\lambda})]\Delta x + [h'(\bar{x})]^{T}\lambda^{+} = -f'(\bar{x})$$

$$h'(\bar{x})\Delta x = -h(\bar{x})$$
(N)

$$\min_{\Delta x} \left\{ f(\bar{x}) + \Delta x^T f'(\bar{x}) + \frac{1}{2} \Delta x^T \nabla_x^2 L(\bar{x}, \bar{\lambda}) \Delta x : h(\bar{x}) + h'(\bar{x}) \Delta x = 0 \right\}$$

$$(QP)$$

<u>Crucial observation</u>: Let the Newton system underlying (N) be a system with nonsingular matrix. Then the Newton displacement Δx given by (N) is the unique KKT point of the quadratic program (QP), and λ^+ is the corresponding vector of Lagrange multipliers.

$$\begin{bmatrix} \nabla_x^2 L(\bar{x},\bar{\lambda}) \end{bmatrix} \Delta x + \begin{bmatrix} h'(\bar{x}) \end{bmatrix}^T \lambda^+ &= -f'(\bar{x}) \\ h'(\bar{x}) \Delta x &= -h(\bar{x}) \end{bmatrix}$$
(N)
$$\min_{\Delta x} \left\{ f(\bar{x}) + \Delta x^T f'(\bar{x}) \\ + \frac{1}{2} \Delta x^T \nabla_x^2 L(\bar{x},\bar{\lambda}) \Delta x : h'(\bar{x}) \Delta x = -h(\bar{x}) \right\}$$
(QP)

Proof of Critical Observation: Let z be a KKT points of (QP), and μ be the corresponding vector of Lagrange multipliers. The KKT system for (QP) reads

$$f'(\bar{x}) + \nabla_x^2 L(\bar{x}, \bar{\lambda})z + [h'(\bar{x})]^T \mu = 0$$

$$h'(\bar{x})z = -h(\bar{x})$$

which are exactly the equations in (N). Since the matrix of system (N) is nonsingular, we have $z = \Delta x$ and $\mu = \lambda^+$.

$$\min_{x} \left\{ f(x) : h(x) \equiv (h_1(x), ..., h_k(x))^T = 0 \right\}$$
(P)

Free Newton method as applied to the KKT system of (*P*) works as follows: Given current iterate $(\bar{x}, \bar{\lambda})$, we linearize the constraints, thus getting "approximate feasible set"

$$\bar{T} = \{\bar{x} + \Delta x : h'(\bar{x})\Delta x = -h(\bar{x})\},\$$

and minimize over this set the quadratic function

$$f(\bar{x}) + (x - \bar{x})^T f'(\bar{x}) + \frac{1}{2} (x - \bar{x})^T \nabla_x^2 L(\bar{x}, \bar{\lambda}) (x - \bar{x}).$$

The solution of the resulting quadratic problem with linear equality constraints is the new x-iterate, and the vector of Lagrange multipliers associated with this solution is the new λ -iterate.

<u>Note:</u> The quadratic part in the auxiliary quadratic objective comes from the Lagrange function of (P), and not from the objective of (P)!

General constrained case

"Optimization-based" interpretation of the Newton method as applied to the KKT system of equality constrained problem can be extended onto the case of general constrained problem

$$\min_{x} \begin{cases} f(x): & h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0 \\ g(x) = (g_{1}(x), ..., g_{m}(x))^{T} \leq 0 \end{cases} \\ (P)$$

and results in the Basic SQP scheme: Given current approximations x_t , λ_t , $\mu_t \ge 0$ to a nondegenerate local solution x_* of (P) and corresponding optimal Lagrange multipliers λ^* , μ^* , we solve auxiliary linearly constrained quadratic problem

$$\begin{split} \min_{\Delta x} \left\{ f(x_t) + \Delta x^T f'(x_t) \\ + \frac{1}{2} \Delta x^T \nabla_x^2 L(x_t; \lambda_t, \mu_t) \Delta x : \begin{array}{l} h'(x_t) \Delta t &= -h(x_t) \\ g'(x_t) \Delta x &\leq -g(x_t) \end{array} \right\} \\ L(x; \lambda, \mu) &= f(x) + h^T(x) \lambda + g^T(x) \mu \\ (QP_t) \end{split}$$

set $x_{t+1} = x_t + \Delta x_*$ and define λ_{t+1} , μ_{t+1} as the optimal Lagrange multipliers of (QP_t).

$$\min_{x} \left\{ f(x): \begin{array}{l} h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0\\ g(x) = (g_{1}(x), ..., g_{m}(x))^{T} \leq 0 \end{array} \right\}$$
(P)

<u>Theorem.</u> Let $(x_*; \lambda^*, \mu^*)$ be a nondegenerate locally optimal solution to (P) and the corresponding optimal Lagrange multipliers. The Basic SQP method, started close enough to $(x_*; \lambda^*, \mu^*)$, and restricted to work with appropriately small Δx , is well defined and converges to $(x_*; \lambda^*, \mu^*)$ quadratically.

Difficulty: From the "global" viewpoint, the auxiliary quadratic problem to be solved may be bad (e.g., infeasible or below unbounded). In the <u>equality constrained</u> case, this never happens when we are close to the nondegenerate local solution; in the general case, bad things may happen even close to a nondegenerate local solution. • <u>Cure</u>: replace the matrix $\nabla_x^2 L(x_t; \lambda^t, \mu^t)$ when it is not positive definite on the entire space by a positive definite matrix B_t , thus arriving at the method where the auxiliary quadratic problem is

$$\min_{\Delta x} \left\{ f(x_t) + \Delta x^T f'(x_t) + \frac{1}{2} \Delta x^T B_t \Delta x : \begin{array}{l} h'(x_t) \Delta t = -h(x_t) \\ g'(x_t) \Delta x \leq -g(x_t) \end{array} \right\} \\ (QP_t)$$

With this modification, the auxiliary problems are convex and solvable with unique optimal (provided that they are feasible, which indeed is the case when x_t is close to a nondegenerate solution to (P)).

Ensuring global convergence

"Cured" Basic SQP scheme possesses nice local convergence properties; however, it in general is not globally converging.

Indeed, in the simplest unconstrained case SQP becomes the basic/modified Newton method, which is not necessarily globally converging, unless linesearch is incorporated.

♠ To ensure global convergence of SQP, we incorporate linesearch. In the scheme with linesearch, the solution $(\Delta x, \lambda^+, \mu^+)$ to the auxiliary quadratic problem

$$\min_{\Delta x} \left\{ f(x_t) + \Delta x^T f'(x_t) + \frac{1}{2} \Delta x^T B_t \Delta x : \begin{array}{l} h'(x_t) \Delta t = -h(x_t) \\ g'(x_t) \Delta x \leq -g(x_t) \end{array} \right\}$$
(QP_t)

is used as *search direction* rather than as a new iterate. The new iterate is

$$x_{t+1} = x_t + \gamma_{t+1} \Delta x$$

$$\lambda_{t+1} = \lambda_t + \gamma_{t+1} (\lambda^+ - \lambda_t)$$

$$\mu_{t+1} = \mu_t + \gamma_{t+1} (\mu^+ - \mu_t)$$

where $\gamma_{t+1} > 0$ is the stepsize given by line-search.

Question: What should be minimized by the linesearch?

In the constrained case, the auxiliary objective to be minimized by the linesearch cannot be chosen as the objective of the problem of interest. In the case of SQP, a good auxiliary objective ("merit function") is

$$M(x) = f(x) + \theta \left[\sum_{i=1}^{m} |h_i(x)| + \sum_{j=1}^{k} g_j^+(x) \right]$$

where $\theta > 0$ is parameter.

Fact: Let x_t be current iterate, B_t be a positive definite matrix used in the auxiliary quadratic problem, Δx be a solution to this problem and $\lambda \equiv \lambda_{t+1}$, $\mu \equiv \mu_{t+1}$ be the corresponding Lagrange multipliers. Assume that θ is large enough:

$$\theta \geq \max\{|\lambda_1|, ..., |\lambda_k|, \mu_1, \mu_2, ..., \mu_m\}$$

Then either $\Delta x = 0$, and then x_t is a KKT point of the original problem, or $\Delta x \neq 0$, and then Δx is a direction of decrease of $M(\cdot)$, that is,

$$M(x + \gamma \Delta x) < M(x)$$

for all small enough $\gamma > 0$.

SQP Algorithm with Merit Function

Generic SQP algorithm with merit function is as follows:

 \triangle <u>Initialization</u>: Choose $\theta_1 > 0$ and starting point x_1

 \diamond Step t: Given current iterate x_t ,

— choose a matrix $B_t \succ 0$ and form and solve auxiliary problem

$$\min_{\Delta x} \left\{ f(x_t) + \Delta x^T f'(x_t) + \frac{1}{2} \Delta x^T B_t \Delta x : \begin{array}{l} h'(x_t) \Delta t = -h(x_t) \\ g'(x_t) \Delta x \leq -g(x_t) \end{array} \right\}$$
(QP_t)

thus getting the optimal Δx along with associated Lagrange multipliers λ, μ . — if $\Delta x = 0$, terminate: x_t is a KKT point of the original problem, otherwise proceed as follows:

- check whether

$$\theta_t \geq \overline{\theta}_t \equiv \max\{|\lambda_1|, ..., |\lambda_k|, \mu_1, ..., \mu_m\}.$$

if it is the case, set $\theta_{t+1} = \theta_t$, otherwise set

$$\theta_{t+1} = \max[\overline{\theta}_t, 2\theta_t];$$

— Find the new iterate

$$x_{t+1} = x_t + \gamma_{t+1} \Delta x$$

by linesearch aimed to minimize the merit function

$$M_{t+1}(x) = f(x) + \theta_{t+1} \left[\sum_{i=1}^{m} |h_i(x)| + \sum_{j=1}^{k} g_j^+(x) \right]$$

on the search ray $\{x_t + \gamma \Delta x \mid \gamma \ge 0\}$. Replace t with t + 1 and loop.

$$\min_{x} \left\{ f(x): \begin{array}{l} h(x) = (h_{1}(x), ..., h_{k}(x))^{T} = 0\\ g(x) = (g_{1}(x), ..., g_{m}(x))^{T} \leq 0 \end{array} \right\}$$
(P)

<u>Theorem:</u> Let general constrained problem be solved by SQP algorithm with merit function. Assume that

• there exists a compact $\Omega \subset \mathbf{R}^n$ such that for $x \in \Omega$ the solution set D(x) of the system of linear inequality constraints

 $S(x): \quad h'(x)\Delta x = -h(x), g'(x)\Delta x \leq -g(x)$

with unknowns Δx is nonempty, and each vector $\Delta x \in D(x)$ is a regular solution of system S(x);

• the trajectory $\{x_t\}$ of the algorithm belongs to Ω and is infinite (i.e., the method does not terminate with exact KKT point);

• the matrices B_t used in the method are uniformly bounded and uniformly positive definite: $cI \leq B_t \leq CI$ for all t, with some $0 < c \leq C < \infty$.

Then all accumulation points of the trajectory of the method are KKT points of (P).

Separation Theorems and Statistical Estimation

Consider the Linear Regression problem as follows:

Problem I. Given indirect noisy observations

 $y = Ax + \sigma\xi$

of "signal" (vector of parameters) xknown to belong to a given set $X \subset \mathbb{R}^n$, infer from the observations some information on x.

- A: given $m \times n$ matrix
- $\xi \sim \mathcal{N}(0, I_m)$: standard Gaussian noise
- $\sigma > 0$: known intensity of the noise

Basic assumption: X is a closed and bounded *convex* set.

♠ Example: Imaging

• A (discretized) 1D/2D/3D image is a function on a *n*-element grid of "pixels" – small 1D/2D/3D boxes. The value of this function at a pixel *j* is the intensity (brightness, "blackness", etc.) of the image in the pixel. \Rightarrow An image is a vector $x \in \mathbb{R}^n$, where x_j is the intensity of the "physical image" in *j*-th pixel.

• With no noise, the output of a typical scanner depends linearly on the input image x, and in many cases the noise η enters the observations z additively: $z = Bx + \eta$.

• When, as it often is the case, the noise is Gaussian, we can pass from the observations $z = Bx + \eta$ to the observations $y \equiv Cz = CBx + C\eta \equiv Ax + \xi$ in such a way that the Gaussian noise ξ is white: $\xi \sim \mathcal{N}(0, I)$.

• In many cases, a priori information, like $0 \le x_j \le L$, or $0 \le x_j$, $\sum_j x_j \le L$, allows to localize the image in a known in advance convex compact set X.

$y = Ax + \sigma \xi$ with $x \in X$

• X, $A \in \mathbb{R}^{m \times n}$, $\sigma > 0$: given • $\xi \sim \mathcal{N}(0, I_m)$

Case I: Hypotheses Testing

(HT): Given two closed convex subsets X_1 , X_2 of X, test the hypothesis $P_1 : x \in X_1$ vs. the alternative $P_2 : x \in X_2$.

• <u>A test</u> for HT is a (measurable) function $\psi(y)$ of observations taking just 2 values 1 and 2. Given observation y, we accept *i*-th hypothesis when $\psi(y) = i$, i = 1, 2.

• Quantifying risk of a test. The most natural ray to quantify the risk of a test ψ , is to look at the *error probabilities*

 $\begin{aligned} \epsilon_1[\psi] &= \sup_{x \in X_1} \operatorname{Prob}\{\psi(Ax + \sigma\xi) \neq 1\} \\ \epsilon_2[\psi] &= \sup_{x \in X_2} \operatorname{Prob}\{\psi(Ax + \sigma\xi) \neq 2\} \end{aligned}$

– the worst-case probabilities to reject a hypothesis when it is true, and at the the test error

$$\epsilon[\phi] = \max[\epsilon_1[\phi], \epsilon_2[\phi]]$$

$y = Ax + \sigma \xi$ with $x \in X$

• $X, A \in \mathbb{R}^{m \times n}, \sigma > 0$: given • $\xi \sim \mathcal{N}(0, I_m)$ • Solving HT. Let $Y_i = AX_i, i = 1, 2$. Y_i are closed convex compact sets along with X_1 and X_2 . We can separate them by the widest possible stripe by solving the optimization problem

Opt =
$$\min_{u,v} \left\{ \frac{1}{2} ||u - v||_2 : u \in Y_1, v \in Y_2 \right\}$$

= $\frac{1}{2} ||u_* - v_*||_2$
[$u_* \in Y_1, v_* \in Y_2$]
• when Opt = 0 ($\Leftrightarrow u_* = v_*$), the sets Y_1 and
 Y_2 cannot be strictly separated
• when Opt > 0, setting $f = \frac{u_* - v_*}{||u - v||_2}$, $e = \frac{u_* + v_*}{2}$,
we have

$$f^{T}[y-e] \left\{ \begin{array}{ll} \geq \mathsf{Opt}, & y \in Y_{1} \\ \leq -\mathsf{Opt}, & y \in Y_{2} \end{array}
ight.$$

<u>Theorem</u> (i) [lower bounds] The test error of every test satisfies the inequality

 $\epsilon[\phi] \geq \mathsf{Erf}(\mathsf{Opt}/\sigma),$

where $\operatorname{Erf}(s) = \int\limits_{s}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\{-t^2/2\} dt$ is the error function.

(ii) [optimal test] Assuming $Y_1 \cap Y_2 = \emptyset$, consider the test

$$\psi_*(y) = \left\{ egin{array}{cc} 1, & f^T(y-e) > 0 \ 2, & f^T(y-e) < 0 \end{array}
ight.$$

Then

 Proof of lower bound: We have $Opt = \frac{1}{2}||u_* - v_*||_2 = \frac{1}{2}||A[x^1 - x^2]||_2$, where $x^i \in X_i$, i = 1, 2. Let us replace the "complex" hypotheses $P_i : x \in X_i$, i = 1, 2, with simple hypotheses $\Pi_i : x = x^i$, i = 1, 2. Every test capable to distinguish between Π_1 and Π_2 with certain test error and error probabilities, distinguishes between Π_1 and Π_2 with the same or smaller test error and error probabilities. Now let us use the following simple result:

Lemma Let ψ be a test for distinguishing between two simple hypotheses Π_1 , Π_2 on the distribution of observation $y \in \mathbb{R}^m$, Π_i stating that the density of the distribution is $p_i(y)$, i = 1, 2. Then

$$2\epsilon[\phi] \ge \int \min[p_1(y), p_2(y)] dy$$

Proof of Lemma. Consider a randomized test ψ which, given an observation y, accepts Π_1 with probability p(y), and Π_2 with probability q(y) = 1 - p(y). Then

$$\epsilon_{1}[\phi] = \int q(y)p_{1}(y)dy$$

$$\epsilon_{2}[\phi] = \int p(y)p_{2}(y)dy$$

$$\Rightarrow 2\epsilon[\phi] \geq \epsilon_{1}[\phi] + \epsilon_{2}[\phi]$$

$$= \int [p(y)p_{2}(y) + (1 - p(y))p_{1}(y)]dy$$

$$\geq \int \min[p_{1}(y), p_{2}(y)]dy.$$

Lemma \Rightarrow lower bound: We are in the situation of $p_i(y) = \left[\frac{1}{\sqrt{2\pi\sigma}}\right]^m \exp\{-\frac{\|y-Ax^i\|_2^2}{2\sigma^2}\},\ i = 1, 2, \text{ whence}$ $2\epsilon[\phi] \ge \int \left[\frac{1}{\sqrt{2\pi\sigma}}\right]^m \exp\{\min\left[-\frac{\|y-u_*\|_2^2}{2\sigma^2}, -\frac{\|y-v_*\|_2^2}{2\sigma^2}\right]\}dy$ [setting $y - \frac{u_* + v_*}{2} = \sigma z, \ h = \frac{u_* - v_*}{2\text{Opt}}, \ \rho = \text{Opt}/\sigma$] $= \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\max[\|u-\rho h\|_2^2, \|u+\rho h\|_2^2]}{2}\}du$ $= 2 \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|u+\rho h\|_2^2}{2}\}du$ $= 2 \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|u+\rho h\|_2^2}{2}\}du$ $= 2 \int \frac{1}{\sqrt{2\pi}}\exp\{-\frac{[s+\rho]^2}{2}\}ds = 2\text{Erf}(\rho).$

Upper bound: Let $x \in X_2$. Then

$$\begin{aligned} \operatorname{Prob}\{\psi_*(Ax + \sigma\xi) &= 1\} \\ &\leq \int \left[\frac{1}{\sqrt{2\pi\sigma}}\right]^m \exp\{-\frac{\|u - Ax\|_2^2}{2\sigma^2}\} du \\ &u: f^T[u - e] \ge 0 \end{aligned} \\ \begin{aligned} \operatorname{[setting } u - Ax &= \sigma z] \\ &\int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|z\|_2^2}{2}\} dz \end{aligned} \\ f^Tz \ge \frac{1}{\sigma} f^T[e - Ax] \\ &\leq \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|z\|_2^2}{2}\} dz \end{aligned} \\ f^Tz \ge \frac{\operatorname{Opt}}{\sigma} \\ &= \int \frac{1}{\sqrt{2\pi}} \exp\{-s^2/2\} = \operatorname{Erf}(\operatorname{Opt}/\sigma) \end{aligned}$$

and similarly for $x \in X_1$ it holds

$$\begin{aligned} & \operatorname{Prob}\{\psi_*(Ax + \sigma\xi) = 2\} \\ & \leq \int \left[\frac{1}{\sqrt{2\pi\sigma}}\right]^m \exp\{-\frac{\|u - Ax\|_2^2}{2\sigma^2}\} du \\ & = \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|z\|_2^2}{2}\} dz \\ & = \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|z\|_2^2}{2}\} dz \\ & \leq \int \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|z\|_2^2}{2}\} dz = \operatorname{Erf}(\operatorname{Opt}/\sigma) \\ & = \int \int_{T_z \leq -\frac{\operatorname{Opt}}{\sigma}} \left[\frac{1}{\sqrt{2\pi}}\right]^m \exp\{-\frac{\|z\|_2^2}{2}\} dz \end{aligned}$$

 $y = Ax + \sigma \xi$ with $x \in X$

• X, $A \in \mathbb{R}^{m \times n}$, $\sigma > 0$: given • $\xi \sim \mathcal{N}(0, I_m)$

(ELF): Given a linear form $g^T z$ of $z \in \mathbb{R}^n$, estimate $g^T x$.

♠ <u>An estimate</u> for ELF is a (measurable) function $\hat{g}(y)$ of observations taking real values; $\hat{g}(y)$ is the estimate of $g^T x$ associated with observation y.

• Quantifying risk of an estimate. Most natural ray to quantify the risk of an estimate $\hat{g}(\cdot)$ at a given $x \in X$ is via the mean squared estimation error $\mathbf{E}\{[\hat{g}(Ax + \sigma\xi) - g^Tx]^2\}$.

• Estimate $\hat{g}(\cdot)$ on the entire X is quantified by its worst-case, over $x \in X$, risk

 $\mathsf{Risk}[\widehat{g}] = \sup_{x \in X} \mathbf{E}\{[\widehat{g}(Ax + \sigma\xi) - g^T x]^2\}$

• It makes sense to compare this risk with the *minimax optimal* risk

 $\mathsf{Risk}_* = \inf_{\widehat{g}(\cdot)} \mathsf{Risk}[\widehat{g}].$

An estimate $\hat{g}(\cdot)$ is called *affine*, if it is *an* affine function of observations:

$$\widehat{g}(y) = h^T y + c.$$

For an affine estimate, its risk at a point is

$$\begin{split} \mathbf{E} \{ [h^T (Ax + \sigma\xi) + c - g^T x]^2 \} \\ = \mathbf{E} \{ \left[\underbrace{[h^T Ax + c - f^T x]}_{\text{bias}} + \underbrace{\sigma h^T \xi}_{\text{stochastic}} \right]^2 \} \\ = \begin{bmatrix} h^T Ax + c - f^T x \end{bmatrix}^2 + \sigma^2 ||h||_2^2 \end{split}$$

whence

$$\operatorname{Risk}[\widehat{g}] = \left[\max_{x \in X} |h^T A x + c - f^T x|\right]^2 + \sigma^2 ||h||_2^2.$$
(*)

Note: When X is "computationally tractable" (e.g., is a polytope), the right hand side in (*) can be efficiently minimized in h, c. Thus, we can "reach efficiently" the optimal estimation risk RiskAff_{*} achievable with affine estimates.

In contrast, in general it is completely unclear how to reach the minimax optimal risk $Risk_*$ – the underlying estimate "exists in the nature", but we have no idea what is it.

Theorem [D. Donoho, 1994] RiskAff_{*} is within absolute constant factor (in fact, 5/4) of Risk_{*}.

<u>Note:</u> When estimating the signal itself rather than an affine form of the signal, the risks of affine estimates can be incomparably worse than the risks of non-affine ones...

Proof of Donoho's Theorem. 1. It is easily seen that we can normalize the situation by assuming $\sigma = 1$ and that RiskAff_{*} > 1; all we need to prove that in this normalized situation Risk_{*} $\geq O(1)$.

2. Replacing \mathbb{R}^n with Aff(X), we may assume that $intX \neq \emptyset$. Let $X_* = X - X$, so that X_* is a convex compact set symmetric w.r.t. $0 \in intX$. Thus, X_* is the unit ball of a norm $p(\cdot)$ on \mathbb{R}^n . Let

$$p_*(u) = \max_z \left\{ u^T z : p(z) \le 1 \right\}$$

$$\equiv \max_z \left\{ u^T x : z \in X_* \right\}$$

$$\equiv \max_{x,x'} \left\{ u^T [x - x'] : x, x' \in X \right\}$$

be the *conjugate* norm.

<u>General fact</u>: For every norm $p(\cdot)$ on \mathbb{R}^n , its twice conjugate norm $(p_*)_*(\cdot)$ is $p(\cdot)$ itself. **Proof.** From definition of the conjugate norm it follows that

$$\forall (u,v) : |u^T v| \leq p(u)p_*(v).$$

Therefore

$$(p_*)_*(u) = \max_{v:p_*(v) \le 1} v^T u \le \max_{v:p_*(v) \le 1} p(u)p_*(v)$$

 $\le p(u).$

It remains to lead to a contradiction the assumption that $(p_*)_*(\bar{u}) < p(\bar{u})$ for certain \bar{u} . Indeed, assume that it is the case. By homogeneity, we may assume that $p(\bar{u}) > 1$, while $(p_*)_*(\bar{u}) < 1$. Since $p(\bar{u}) > 1$, we can strongly separate \bar{u} and the unit ball of $p(\cdot)$: there exists \bar{v} such that

$$\bar{v}^T \bar{u} > \max_{u: p(u) \le 1} \bar{v}^T u = p_*(\bar{v}).$$

But then

$$p_*(\overline{v}) < \overline{v}^T \overline{u} \le p_*(\overline{v})(p_*)_*(\overline{u}),$$

whence

$$(p_*)_*(\bar{u}) \ge 1,$$

which is a desired contradiction.

3. Let us verify that when $\alpha \ge 0, \beta \ge 0$ are such that $\alpha^2 + \beta^2 \le 1$, no $h \in \mathbb{R}^m$ satisfies

 $p_*(A^Th - g) \leq 2\alpha$ and $||h||_2 \leq \beta$.

Indeed, assume such an h exists. Then

$$2\alpha \\ \ge p_*(A^T h - g) = \max_{\substack{x, x' \in X}} [A^T h - g]^T [x - x'] \\ = \max_{x \in X} [A^T h - g]^T x - \min_{x \in X} [A^T h - g]^T x.$$

Setting

$$c = -\frac{1}{2} \left[\max_{x \in X} [A^T h - g]^T x + \min_{x \in X} [A^T h - g]^T x \right],$$

we conclude that

$$\max_{x \in X} |[h^T A x + c] - g^T x| \le \alpha,$$

whence for the affine estimate $\hat{g}(y) = h^T y + c$ one has

$$\mathsf{Risk}[\widehat{g}] \le \alpha^2 + \beta^2 = 1,$$

which is impossible, since $RiskAff_* > 1$.

4. We have seen that when $\alpha^2 + \beta^2 = 1$ and $\alpha, \beta \ge 0$, no *h* satisfies the relations

$$p_*(A^Th - g) \leq 2\alpha$$
 and $||h||_2 \leq \beta$,

meaning that the compact convex sets

 $\{u : p_*(u - g) \le 2\alpha\}, \ \{A^T h : \|h\|_2 \le \beta\}$ can be strongly separated:

$$\exists z: \min_{u:p_*(u-g) \le 2\alpha} u^T z > \max_{\|h\|_2 \le \beta} (A^T h)^T z$$

We have

$$\max_{\|h\|_{2} \le \beta} (A^{T}h)^{T}z = \max_{\|h\|_{2} \le \beta} h^{T}(Az) = \beta \|Az\|_{2}$$

and

$$\min_{\substack{u:p_*(u-g) \le 2\alpha \\ = g^T z - 2\alpha(p_*)_*(z) = g^T z - 2\alpha p(z).}} \min_{\substack{u:p_*(d) \le 2\alpha \\ = g^T z - 2\alpha(p_*)_*(z) = g^T z - 2\alpha p(z).}} [g^T z + d^T z]$$

Thus, $g^T z - 2\alpha p(z) > \beta ||Az||_2$. This inequality remains valid when z is replaced with θz , $\theta >$ 0; thus, we may assume that p(z) = 1, i.e., z = r - s with $r, s \in X$. The bottom line is:

$$\forall (\alpha \ge 0, \beta \ge 0, \alpha^2 + \beta^2 = 1) : \\ \exists r, s \in X : g^T[r-s] \ge 2\alpha + \beta \|A[r-s]\|_2.$$

5. Setting
$$\alpha = \beta = \sqrt{1/2}$$
, we get

$$\exists r, s \in X : g^{T}[r-s] > \sqrt{2} + \frac{1}{\sqrt{2}} \|A[r-s]\|_{2}.$$
(!)

We claim that then

 $\begin{aligned} \exists u, v \in X : g^{T}[u-v] \geq \sqrt{2} \text{ and } \|A[u-v]\|_{2} \leq 2. \\ \text{Indeed, when } \|A[r-s]\|_{2} \leq 2, \text{ we can take} \\ u = r, v = s. \text{ When } \|A[r-s]\|_{2} > 2, \text{ we can} \\ \text{set } u = s + \frac{2}{\|A[r-s]\|_{2}}[r-s], v = s, \text{ so that} \\ \|A[u-v]\|_{2} = \frac{2}{\|A[r-s]\|_{2}}\|A[r-s]\|_{2} = 2 \text{ and} \\ g^{T}[u-v] = \frac{2}{\|A[r-s]\|_{2}}g^{T}[r-s] \\ \geq \frac{2}{\|A[r-s]\|_{2}}\frac{1}{\sqrt{2}}\|A[r-s]\|_{2} = \sqrt{2}. \end{aligned}$

6. Since $||A[u-v]||_2 \le 2$, for every test ψ for distinguishing between the hypotheses P_1 : x = u and P_2 : x = v we have $\epsilon[\psi] \ge \text{Erf}(1)$.

Now let $\hat{g}(\cdot)$ be an estimate of $g^T x$, and let

$$\psi(y) = \begin{cases} 1, \ \hat{g}(y) \ge g^T e, \\ 2, \ \hat{g}(y) < g^T e \end{cases}, e = \frac{u+v}{2}.$$

When x = u and the test accepts P_2 , we have $\hat{g}(y) \leq g^T e$, while

$$g^T x = g^T r = g^T e + \frac{1}{2}g^T [u - v] \ge g^T e + \frac{1}{\sqrt{2}},$$

that is, the squared estimation error is $\geq 1/2$, whence

$$\epsilon_1[\psi] \leq \frac{\operatorname{Risk}[\widehat{g}]}{1/2} = 2\operatorname{Risk}[\widehat{g}].$$

Similarly, when x = v and the test accepts P_1 , we have $\hat{g}(y) \ge g^T e$ and $g^T x = g^T v = g^T e - \frac{1}{2}g^T[u-v] \le g^T e - \frac{1}{\sqrt{2}}$, whence

$$\epsilon_2[\psi] \leq \frac{\operatorname{Risk}[\widehat{g}]}{1/2} = 2\operatorname{Risk}[\widehat{g}].$$

Thus,

 $2\operatorname{Risk}[\widehat{g}] \ge \epsilon[\psi] \ge \operatorname{Erf}(1) \Rightarrow \operatorname{Risk}[\widehat{g}] \ge \operatorname{Erf}(1)/2.$ Since \widehat{g} is an arbitrary estimate, we get

$$\mathsf{Risk}_* \geq O(1).$$

Illustration: Participles detection

<u>Situation:</u> A stream of participles contaminated by "background signal" inputs a detector. The output of the detector *with no output noise* is

$$z(t) = \int_0^{\Delta} D(t-s) \left[\sum_j \mu_j \delta(s-s_j) + w(s) \right] ds$$

where

• μ_j is the energy of *j*-th participle, s_j is the moment when it arrives and $\delta(\cdot)$ is the Dirac δ -function

• $w(\cdot)$ is the "background signal"

• $D(\cdot)$ is the "impulse response" of the detector – its output when the input is $\delta(\cdot)$. We assume D(t) to be supported on $[0, \Delta]$. • After discretization in time with resolution

 $\delta t = \Delta/N$, the model becomes

$y_{\tau} = \sum_{\ell} \left[\mu_{\ell} + w_{\ell} \delta t \right] D_{\tau - \ell} + \sigma \xi_{\tau}$

 $\begin{bmatrix} \mu_{\ell} : \text{ energy of participle arriving at time } \ell \cdot \delta t \\ \{\xi_{\tau} \sim \mathcal{N}(0, 1)\}_{\tau = -\infty}^{\infty} : \text{ independent output noises} \\ w_{\ell} = w(\ell \cdot \delta t), D_r = D(r \cdot \delta t), \sigma > 0 : \text{ noise level} \end{bmatrix}$ $\frac{\text{The goal: To infer from observations information on } \{\mu_{\tau}\}_{\tau = -\infty}^{\infty}.$

$y_{\tau} = \sum_{\ell} \left[\mu_{\ell} + w_{\ell} \delta t \right] D_{\tau - \ell} + \sigma \xi_{\tau}$

Strategy:

• We fix a number K of consecutive "kernel widths" to be considered, thus focusing on recovering μ_{τ} 's for $0 \leq \tau \leq KN - 1$. Observations at these instants are defined, up to observation noise, by two vectors

$$\mu = [\mu_{-N+1}; \mu_{-N+2}; ...; \mu_{KN-1}]$$

$$w = [w_{-N+1}; w_{-N+2}; ...; w_{KN-1}]$$

according to

 $[y_0; ...; y_{KN-1}] \equiv y = A[\mu; w] + \sigma\xi$ $[\xi \sim \mathcal{N}(0, I_{KN})]$

• Assume we know an upper bound $\hat{\mu}$ on participle's energy and an upper bound n on the number of participles which can arrive at instants -N + 1, -N + 2, ..., KN - 1. Then μ belongs to the convex compact set

$$\mathcal{M} = \left\{ \mu \in \mathbb{R}^{(K+1)N-1} : \begin{array}{l} \mu \ge 0, \mu_{\tau} \le \widehat{\mu}, \\ \sum_{\tau} \mu_{\tau} \le n\widehat{\mu} \end{array} \right\}$$
• We model the background signal w(t) as smooth with given parameters of smoothness, say, twice differentiable with $|w(\cdot)| \leq C_0$ and $|w''(\cdot)| \leq C_2$, which translates into the restriction

$$w \in \mathcal{W} = \left\{ \begin{aligned} |w_{\tau}| &\leq C_0, \ -N+1 \leq \tau < KN \\ w &: \ |w_{\tau} - 2w_{\tau+1} + w_{\tau+2}| \leq C_2 \delta t^2, \\ -N+1 \leq \tau < KN-2 \end{aligned} \right\}$$

• Thus, we arrive at the situation where we are given noisy indirect observations

$$y = Ax + \sigma\xi \qquad [x = [\mu; w]]$$

of a signal $x \in X = \mathcal{M} \times \mathcal{W}$ and are interested to recover linear forms

$$g_{\tau}^{T}[\mu;w] = \mu_{\tau}$$

of the signal.

<u>Note</u>: In fact, it suffices to find a single $\tau = \tau_*$ for which μ_{τ} can be well estimated via our observations $y_{-N+1}, ..., y_{KN-1}$. Given the corresponding estimate \hat{g} , we can estimate μ_{τ} for every τ by the quantity

$$\hat{g}(y_{\tau-\tau_*-N+1},...,y_{\tau-\tau_*+KN-1}).$$

Numerical illustration:

Setup:

• Impulse response $D(s) = c [\exp\{as\} - \exp\{bs\}],$ $0 \le s \le 1.$

• Discretization: N = 40 grid points per 1 sec ($\delta t = 1/40$)



• K = 2, i.e., 80 observations and dim $\mu = \dim w = 119$

• Participles: $\mu_{\tau} \in [0, \hat{\mu} = 0.2]$, with at most n = 5 participles per 119 consecutive discrete time instants

- Background signal: $|w(\cdot)| \leq 1$, $|w''(\cdot)| \leq 1$
- Noise intensity: $\sigma = 4.33e-4$

<u>Results:</u>

• The best position for estimating μ_{τ} : $\tau =$ 41. The estimate of μ_{τ} :



Estimation weights h_i

• Risk of the estimate:



• Sample recovery:



• Recovering energy distribution of participles:

In our simulation, the energy distribution of a participle was

 $\mathsf{Prob}\{\mu_{\tau} \leq s\widehat{\mu} \, \middle| \, \mu_{\tau} \neq 0\} = \sqrt{s}, 0 \leq s \leq 1.$

To recover this distribution from observations, we

• computed the 0.99-confidence level ℓ of our estimate $\hat{\mu}_{\tau}$ when μ_{τ} is 0:

 $\ell = bias + Erf(0.01) \cdot \sigma \cdot ||h||_2 = 0.0156$

• in a long simulation run, filtered out all time instants τ with $\hat{\mu}_{\tau} < \ell$ and computed the empirical distribution of the remaining $\hat{\mu}_{\tau}$. The result:



Magenta: true energy distribution Blue: estimated energy distribution