

Wavelet Based Nonparametric Bayes Methods

Wavelets are the building blocks of wavelet transformations the same way that the functions e^{inx} are the building blocks of the ordinary Fourier transformation. But in contrast to sines and cosines, wavelets can be (or almost can be) supported on an arbitrarily small closed interval. This makes wavelets a very powerful tool in dealing with phenomena which change rapidly in time.

Statistical wavelet modeling and computational research has, in recent years, become a burgeoning area in both theoretical and applied statistics, and is beginning to impact developments in statistical methodology and in various applied scientific fields. Wavelet ideas are developing in statistics in areas such as regression, density and function estimation, factor analysis, modeling and forecasting in time series analysis, spatial statistics, with ranges of application areas in science and engineering. The emerging interests in Bayesian statistical modeling and wavelets is generating exciting new directions for the interface of two research areas, with significant potential for future impact on applied work.

Many nonparametric procedures are in fact infinitely parametric. An example is the orthogonal series regression or density estimator. In order to estimate such functions, shrinkage, tapering or truncation of parameter estimators from an infinite class is necessary (Chencov's orthogonal series density estimators, Stein-type estimation, and so on.).

Wavelet shrinkage is a simple and yet powerful tool in nonparametric statistical modeling. It can be described as a three step procedure. Data are transformed into a set of wavelet coefficients, a shrinkage of the coefficients is performed, and then the shrunken wavelet coefficients are transformed back in the domain of the original data.

Wavelet domain is good modeling environment; several supporting arguments are listed below.

Discrete wavelet transformations tend to "disbalance" the data. Even though the transformations preserve the ℓ_2 norm of the data, the "energy" (an engineering term for the ℓ_2 norm) of the transformed data concentrates in only a few wavelet coefficients. That narrows the class of plausible models and facilitates the thresholding. Mallat (1989) gives an interesting discussion on modeling from the signal-processing point of view. The disbalancing property also yields a variety of criteria for the best basis selection. Standard references are Coifman and Wickerhauser (1992), Donoho (1994) and Wickerhauser (1994).

Wavelets, as building blocks in modeling, are localized well in both time and scale (frequency). Signals with rapid local changes (signals with discontinuities, cusps, sharp spikes, etc.) can be well represented with only a few wavelet coefficients. This is, in general, not true for other standard orthonormal bases which may need many "compensating" coefficients to describe discontinuity artifacts and to suppress Gibbs' effects.

Heisenberg's principle states that in modeling time-frequency phenomena one can not be precise in the time domain and in the frequency domain simultaneously. Wavelets automatically trade-off the time-frequency precision by their innate nature. The parsimony of wavelet transformations can be attributed to the ability of wavelets to handle limitations in the Heisenberg principle in a data-dependent manner.

Also, there is theoretical and empirical evidence that wavelet transformations tend to simplify the dependence structure in the original signal. It is even possible, for any given stationary dependence in input signal, to construct a biorthogonal wavelet basis such that the corresponding wavelet coefficients become uncorrelated (a wavelet counterpart of Karhunen-Loève transformation). For a discussion and examples see Walter (1994).

These arguments identify wavelet bases as suitable tools for effective statistical modeling. More favorable arguments can be given: computational speed of the wavelet transformation, simple descriptors of self-similarity and so on.

The benefits of shrinkage estimation in statistics were first explored in the mid-50s by C. Stein. In the 70s and 80s, many statisticians were interested in statistical benefits of classical and Bayesian shrinkage estimators. Since then, shrinkage estimation has been an active field in statistical research. Rapid increase of computing power and the use of Markov Chain Monte Carlo methods make shrinkage estimation a burgeoning research field in Bayesian statistics today.

Since wavelets are unconditional bases¹ for many important function spaces, by inspecting only the magnitudes of wavelet coefficients, one can determine whether the decomposed function belongs to a particular space or not. The rate of decay of wavelet coefficients by increasing levels of details depends on the global (and local) smoothness of the decomposed function. As an illustration we provide Meyer's result that characterizes membership in the Hölder smoothness space² $\mathcal{H}^\alpha(\mathbb{R})$ of degree α by the rate of decay of wavelet coefficients.

Result (*Y. Meyer, 1988*) f is $\mathcal{H}^\alpha(\mathbb{R})$ iff

$$\begin{aligned} |\langle f, \phi_{0k} \rangle| &\leq C, \quad \forall k \in \mathbb{Z}; \\ |\langle f, \psi_{jk} \rangle| &\leq C 2^{-j(\alpha+1/2)}, \quad \forall j \geq 0, \forall k \in \mathbb{Z}, \end{aligned}$$

where $\{\phi_{0,k}, \psi_{j,k}, j \geq 0, k \in \mathbb{Z}\}$ is a fixed orthonormal wavelet basis and $\langle f, \phi_{0k} \rangle$ and $\langle f, \psi_{jk} \rangle$ are wavelet coefficients.

Many other smoothness spaces (Sobolev, Besov) are characterized similarly. Result of Meyer-type are important in understanding the modeling of smoothness in the wavelet domain. It also provides a convenient tool for a Bayesian to incorporate prior information on smoothness.

Discrete Wavelet Transformations

Prior to describing a formal set-up for Bayesian wavelet shrinkage we provide a brief review of discrete wavelet transformation and the traditional wavelet shrinkage.

Basics on wavelets can be found in many texts, monographs and papers at many different levels of exposition. The interested reader should consult monographs by Daubechies (1992), Walter (1994), Ogden (1996), among others. For the sake of self-containedness of this chapter, we provide a brief overview of the discrete wavelet transformation (DWT).

Discrete Wavelet Transformation

Let \underline{y} be a data vector of dimension (size) n . For simplicity we choose n to be a power of 2, say 2^J . We also assume that measurements \underline{y} belong to an interval and consider periodized wavelet bases. Generalizations to different sample sizes and general wavelet and wavelet-like transformations are straightforward.

Suppose that the vector \underline{y} is wavelet-transformed to a vector \underline{d} . The transformation is linear and orthogonal and can be described by an orthogonal matrix W of dimension $n \times n$. In practice one performs the DWT without exhibiting the matrix W explicitly but by using fast filtering algorithms. The filtering procedure is based on so called quadrature mirror filters which are uniquely determined by the wavelet of choice. The

¹A basis $\{f_n, n \in \mathbb{Z}\}$ is unconditional for the space \mathcal{F} if from $\sum a_i f_i \in \mathcal{F}$ follows $\sum s_i a_i f_i \in \mathcal{F}$, where the sequence $\{s_n\}$ is an arbitrary sequence of +1 and -1.

²The Hölder space $\mathcal{H}^\alpha(\mathbb{R})$ is a generalization of the space of functions that are n times continuously differentiable. It is defined as follows:

$$\begin{aligned} \mathcal{H}^\alpha(\mathbb{R}) &= \{f \in L^\infty(\mathbb{R}); \sup_{x,h} \frac{|f(x+h) - f(x)|}{|h|^\alpha} < \infty, \quad 0 < \alpha < 1\} \\ \mathcal{H}^\alpha(\mathbb{R}) &= \{f \in L^\infty(\mathbb{R}) \cup C^n(\mathbb{R}); f^{(n)} \in \mathcal{H}^{\alpha'}(\mathbb{R}), \quad \alpha = n + \alpha', 0 < \alpha' < 1\} \end{aligned}$$

wavelet decomposition of the vector y can be written as a vector \underline{d}

$$\underline{d} = (H^l y, GH^{l-1} y, \dots, GH^2 y, GH y, G y), \quad (1)$$

\underline{d} has the same length as y and where l is any fixed number between 1 and $J = \log_2 n$. The operators G and H are defined coordinate-wise via

$$(Ha)_k = \sum_m h_{m-2k} a_m, \quad \text{and} \quad (Ga)_k = \sum_m g_{m-2k} a_m,$$

where g and h are high- and low-pass filters corresponding to the decomposing wavelet. Components of g and h are connected via *quadrature mirror* relation $g_n = (-1)^n h_{1-n}$.

The elements of \underline{d} are called ‘‘wavelet coefficients’’. The sub-vectors given in (1) correspond to different levels in a tree-like decomposition where the wavelet coefficients are organized level-wise. For instance, the vector $G y$ contains $n/2$ coefficients representing the level of the finest detail.

When $l = J$ in (1), the vectors $GH^{J-1} y = \{d_{00}\}$ and $H^J y = \{c_{00}\}$ contain a single coefficient each and represent the coarsest level of detail and the smooth part in wavelet decomposition. In general, the level j in wavelet decomposition of y contains 2^j elements, and is represented as

$$GH^{n-j-1} y = (d_{j,0}, d_{j,1}, \dots, d_{j,2^j-1}). \quad (2)$$

Wavelet Shrinkage

Wavelet shrinkage methods are now widely recognized as a useful tool in nonparametric function estimation and signal recovery. The simplest nonlinear wavelet shrinkage technique is thresholding. The coordinates of \underline{d} are replaced by 0 if they are smaller in absolute value than a fixed threshold λ .

The two most common thresholding policies are **hard** and **soft** with corresponding transformations given by:

$$\begin{aligned} \theta^h(d, \lambda) &= d \mathbf{1}(|d| > \lambda), \\ \theta^s(d, \lambda) &= (d - \text{sign}(d)\lambda) \mathbf{1}(|d| > \lambda), \end{aligned}$$

where $\mathbf{1}(A)$ is the indicator of A .

Graphs of hard- and soft-thresholding rules are given in Figure 1.

Bayes and Wavelets

Bayesian approaches to choosing the shrinkage method are less ad-hoc than earlier proposals, and have been shown to be effective. It is known that, in general, Bayes rules are ‘‘shrinkers’’ and that their shape in many cases has a desirable property for wavelet shrinkage: it can heavily shrink small arguments and only slightly large arguments. If we use Bayes models for the wavelet coefficients the resulting optimal actions can be very close to thresholding.

Suppose that y_i are ‘noisy’ measurements, i.e. the sums of an unknown signal f_i and noise ϵ_i :

$$y_i = f_i + \epsilon_i, \quad i = 1, \dots, n, \quad (3)$$

or in vector notation, $y = \underline{f} + \underline{\epsilon}$. After applying the wavelet transformation W , the resulting vector $\underline{d} = W(y)$ is again the sum of the transformation of the signal $\underline{\theta} = W(\underline{f})$ and the transformation of the noise $\underline{\eta} = W(\underline{\epsilon})$. This is a consequence of the linearity of W . If the components of the noise vector $\underline{\epsilon}$, ϵ_i , are modeled by independent normals with mean 0 and variance σ^2 , then the orthogonality of W implies that the components

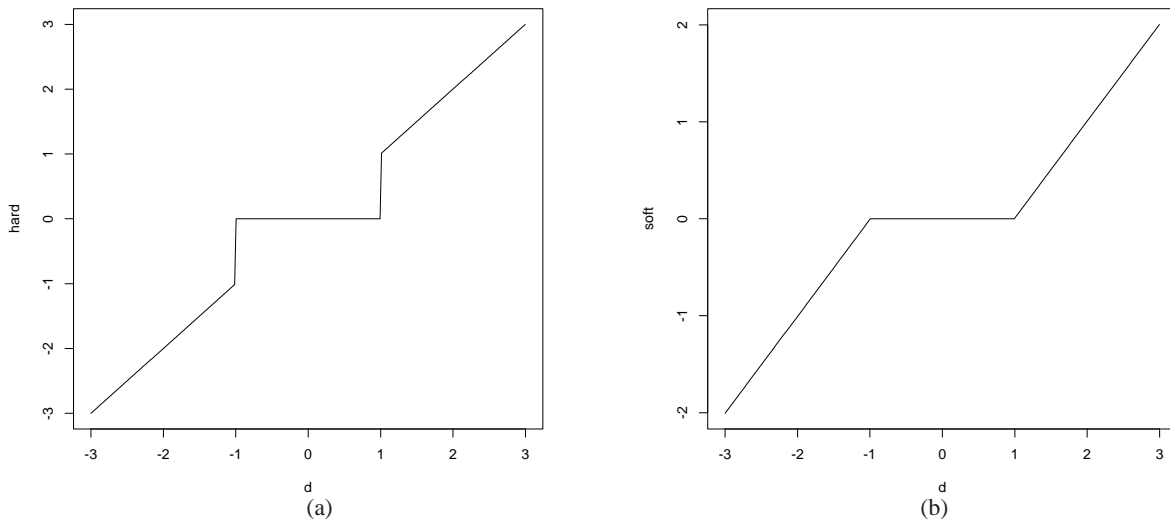


Figure 1: Hard and soft thresholding rules with $\lambda = 1$.

of η have the same distribution. The case of correlated normal noise can be addressed as well, see Johnstone and Silverman (1996).

Thus, instead of estimating the function directly, we estimate its wavelet transformation. Equivalently, we estimate the means θ_i in the model $[d_i] \sim N(\theta_i, \sigma^2)$. The goal of Bayesian wavelet shrinkage is to exhibit such models so that resulting optimal actions, resulting from the Bayesian inference, mimic thresholding rules.

To illustrate Bayesian shrinkage in the wavelet domain we start with an example from Vidakovic (1994).

Example 0.1 Let

$$[d|\theta, \sigma^2] \sim N(\theta, \sigma^2), \quad \sigma^2 \text{ unknown.} \quad (4)$$

be a model for a wavelet coefficient. Because of practical (computational) reasons and noninformative properties³ the prior distribution on σ^2 is chosen to be exponential,

$$[\sigma^2] \sim \mathcal{E}(\mu) \quad (f(\sigma^2|\mu) = \mu e^{-\mu\sigma^2}), \quad (5)$$

since the exponential distribution is the entropy maximizer among all distributions supported on $(0, \infty)$ with a fixed first moment. The marginal model (marginal likelihood) is double exponential,

$$[d|\theta] \sim \mathcal{DE}(\theta, \frac{1}{\sqrt{2\mu}}), \quad (6)$$

with $f(d|\theta) = \frac{1}{2}\sqrt{2\mu}e^{-\sqrt{2\mu}|d-\theta|}$. This follows from the fact that the double exponential distribution is a scale mixture of normals.

The prior on θ is t with location 0, scale τ and n degrees of freedom,

$$[\theta] \sim t_n(0, \tau). \quad (7)$$

³Exponential distribution minimizes the Fisher information in the class of all distributions supported on $[0, \infty)$ with a fixed first moment.

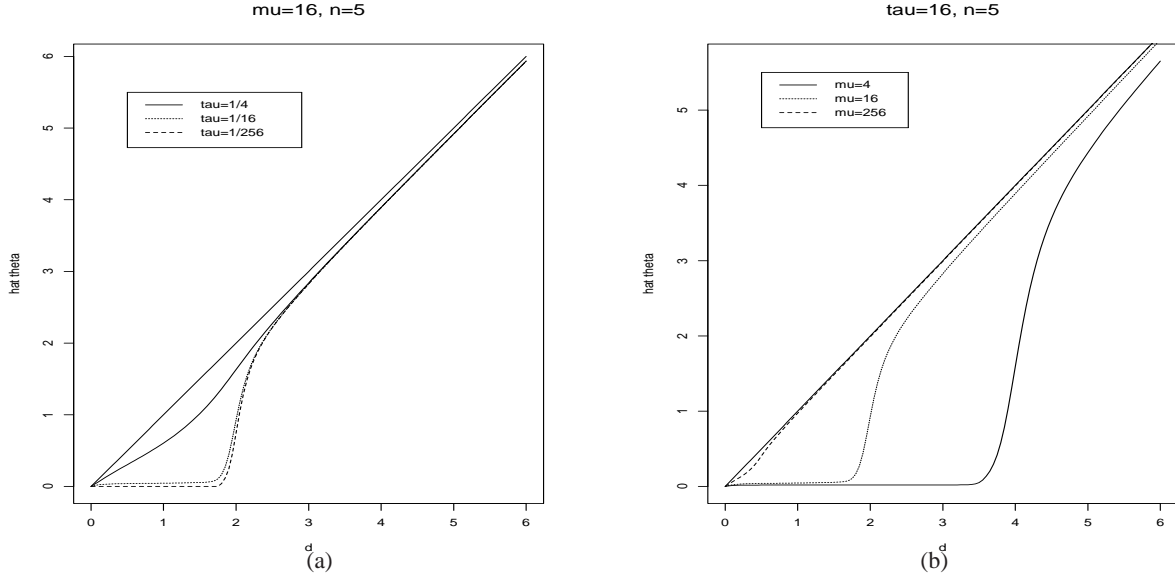


Figure 2: Bayes rules for selected values of τ and μ . The rules are odd and given only for d positive

Several graphs of Bayes rules, with respect to this model, are given in Figure 2. The rules are odd functions and only positive values are plotted. The hyperparameters μ, τ and n can be specified by empirical Bayes arguments.

It is possible to get an analytic expression for the Bayes rule under the model described, when the prior on θ is general but symmetric, i.e. $\pi(\theta) = \pi(-\theta)$. The Bayes rule with respect to the squared error loss is:

$$\delta(d) = d - \frac{\Pi_1'(c) - \Pi_2'(c)}{\Pi_1(c) + \Pi_2(c)}, \quad (8)$$

where Π_1 and Π_2 are the one-sided Laplace transforms of functions $\pi(\theta + d)$ and $\pi(\theta - d)$, $\theta \in (0, \infty)$, and $c = \sqrt{2\mu}$.

Regression Problems

Somewhat simplistic model in Example 0.1 can be replaced with more powerful and coherent Bayesian models that exploit a hierarchical structure.

In the context of wavelet regression we will discuss two approaches in more detail. The first one is Adaptive Bayesian Wavelet Shrinkage (ABWS) proposed by Chipman *et al.* (1997). The approach they take is based on the stochastic search variable selection (SSVS) model proposed by George and McCulloch, with the assumption that σ is known.

Chipman *et al.*(1997) start with the model

$$[d|\theta, \sigma^2] \sim N(\theta, \sigma^2). \quad (9)$$

The prior on θ is defined as a mixture of two normals (Figure 3, Left)

$$[\theta|\gamma_j] \sim \gamma_j N(0, (c_j \tau_j)^2) + (1 - \gamma_j) N(0, \tau_j^2), \quad (10)$$

where

$$[\gamma_j] \sim \text{Bernoulli}(p_j). \quad (11)$$

Because the hyperparameters $p_j, c_j,$ and τ_j depend on the level j to which the corresponding θ (or d) belongs, and can be level-wise different, the method is adaptive.

The Bayes rule under squared error loss for θ (from the level j) has an explicit form,

$$\delta(d) = [P(\gamma_j = 1|d)\frac{(c_j\tau_j)^2}{\sigma^2 + (c_j\tau_j)^2} + P(\gamma_j = 0|d)\frac{\tau_j^2}{\sigma^2 + \tau_j^2}] d, \quad (12)$$

where

$$P(\gamma_j = 1|d) = \frac{p_j\pi(d|\gamma_j = 1)}{(1 - p_j)\pi(d|\gamma_j = 0)}$$

and

$$\pi(d|\gamma_j = 1) \sim N(0, \sigma^2 + (c_j\tau_j)^2) \text{ and } \pi(d|\gamma_j = 0) \sim N(0, \sigma^2 + \tau_j^2).$$

The shrinkage rule (12, Figure 3, Right) can be viewed as smooth interpolation between two lines of slope $\tau_j^2/(\sigma^2 + \tau_j^2)$ and $(c_j\tau_j)^2/(\sigma^2 + (c_j\tau_j)^2)$. The authors give very sophisticated empirical Bayes argument for tuning the hyperparameters level-wise. They provide full posterior analysis on the coefficients and function values. The simulations presented in [5] show that ABWS method is superior than the VisuShrink and SureShrink over the standard Donoho-Johnstone test functions.

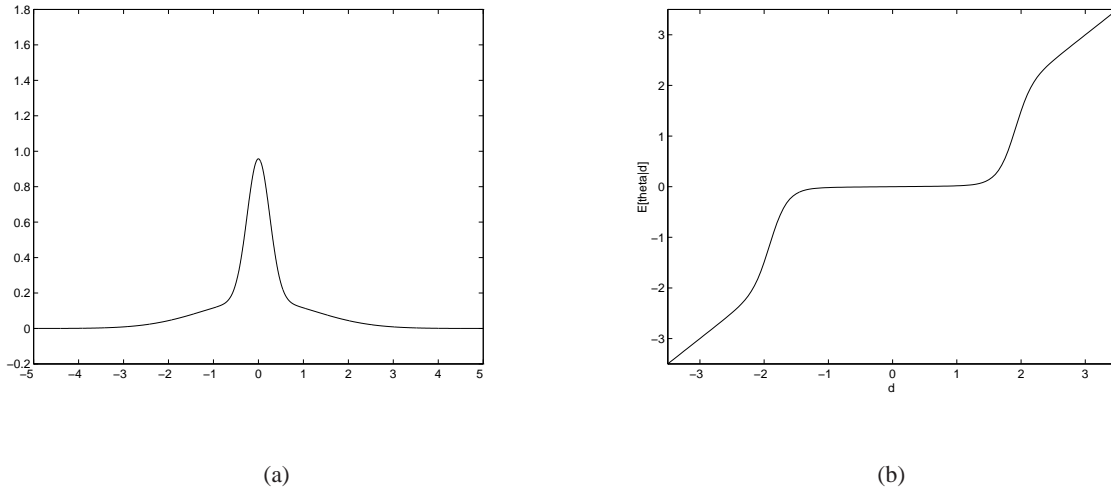


Figure 3: Left: Prior on θ as a mixture of two normal distributions with different variances; Right: Shrinkage rule in [5].

The approach used by Clyde, Parmigiani and Vidakovic (1998) is based on a limiting form of the conjugate SSVS prior in George and McCulloch (1994). Clyde *et al.* (1996) consider a prior for θ which is a mixture of a point mass at 0 if the variable is excluded from the wavelet regression and a normal distribution if it is included,

$$[\theta|\gamma_j, \sigma^2] \sim N(0, (1 - \gamma_j) + \gamma_j c_j \sigma^2). \quad (13)$$

The γ_j are indicator variables that specify which basis element or column of W should be selected. As before, the subscript j points to the level to which θ belongs. The set of all possible vectors $\underline{\gamma}$'s will be referred to as the subset space. The prior distribution for σ^2 is inverse χ^2 i.e.

$$[\lambda\nu/\sigma^2] \sim \chi_\nu^2,$$

where λ and ν are fixed hyperparameters and the γ_j 's are independently distributed as Bernoulli (p_j) random variables.

The posterior mean of $\underline{\theta}|\underline{\gamma}$ is

$$E(\underline{\theta}|\underline{d}, \underline{\gamma}) = \Gamma(I_n + C^{-1})^{-1}\underline{d} \quad (14)$$

where Γ and C are diagonal matrices with γ_{jk} and c_{jk} respectively on the diagonal and 0 elsewhere. For a particular subset determined by the ones in $\underline{\gamma}$ (14) corresponds to thresholding with linear shrinkage.

The posterior mean is obtained by averaging over all models. Model averaging leads to multiple shrinkage estimator of $\underline{\theta}$:

$$E(\underline{\theta}|\underline{d}) = \sum_{\underline{\gamma}} \pi(\underline{\gamma}|\underline{d})\Gamma(I_n + C^{-1})^{-1}\underline{d}, \quad (15)$$

where $\pi(\underline{\gamma}|\underline{d})$ is the posterior probability of a particular subset $\underline{\gamma}$.

An additional nonlinear shrinkage of the coefficients to 0 results from the uncertainty in which subsets should be selected.

Calculating the posterior probabilities of $\underline{\gamma}$ and the mixture estimates for the posterior mean of $\underline{\theta}$ above involve sums over all 2^n values of $\underline{\gamma}$. The calculational complexity of the mixing is prohibitive even for problems of moderate size, and either approximations or stochastic methods for selecting subsets $\underline{\gamma}$ possessing high posterior probability must be used.

In the orthogonal case, Clyde, DeSimone, and Parmigiani (1996) obtain an approximation to the posterior probability of $\underline{\gamma}$ which is adapted to the wavelet setting in [8]. The approximation can be achieved by either conditioning on σ (plug-in approach) or by assuming independence of the elements in $\underline{\gamma}$.

The approximate model probabilities, for the conditional case, are functions of the data through the regression sum of squares and are given by:

$$\begin{aligned} \pi(\underline{\gamma}|\underline{d}) &\approx \tilde{\pi}(\underline{\gamma}|y) = \prod_{j,k} \rho_{jk}^{\gamma_{jk}} (1 - \rho_{jk})^{1-\gamma_{jk}} \\ \rho_{jk}(\underline{d}, \sigma) &= \frac{a_{jk}(\underline{d}, \sigma)}{1 + a_{jk}(\underline{d}, \sigma)} \end{aligned} \quad (16)$$

where

$$\begin{aligned} a_{jk}(\underline{d}, \sigma) &= \frac{p_{jk}}{1 - p_{jk}} (1 + c_{jk})^{-1/2} \cdot \exp \left\{ \frac{1}{2} \frac{S_{jk}^2}{\sigma^2} \right\} \\ S_{jk}^2 &= d_{jk}^2 / (1 + c_{jk}^{-1}). \end{aligned}$$

The p_{jk} can be used to obtain a direct approximation to the multiple shrinkage Bayes rule. The independence assumption leads to more involved formulas. Thus, the posterior mean for θ_{jk} is approximately

$$\rho_{jk} (1 + c_{jk}^{-1})^{-1} d_{jk}. \quad (17)$$

Equation (17) can be viewed as a level dependent wavelet shrinkage rule, generating a variety of nonlinear rules. Depending on the choice of prior hyperparameters shrinkage may be monotonic, if there are no level dependent hyperparameters, or non-monotonic; see Figure 4, Left. Authors report good MSE performance of approximation rules.

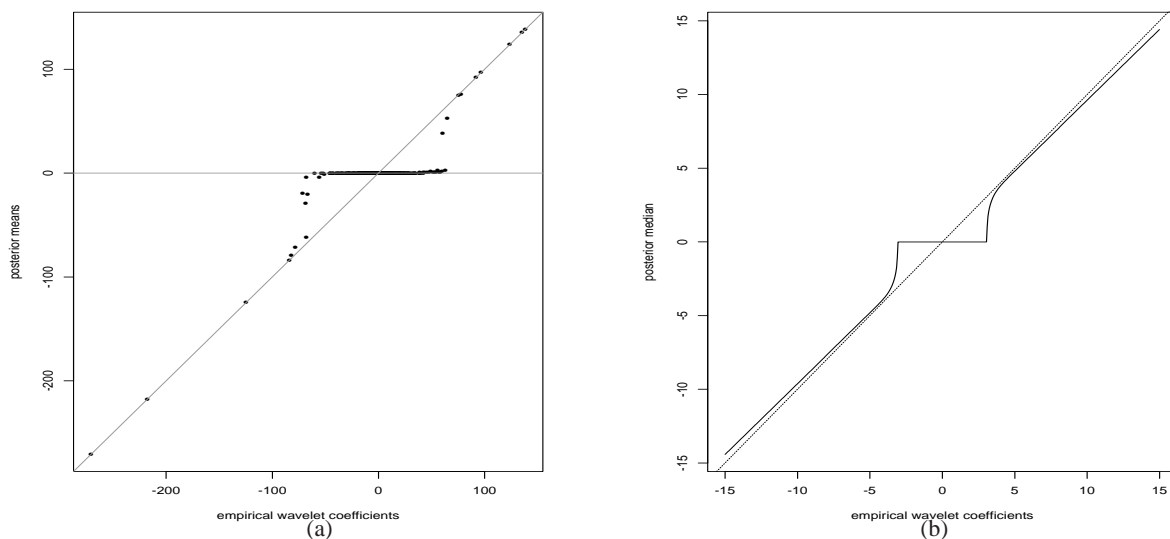


Figure 4: Left: Shrinkage rule from [8] based on independence approximation (17); Right: Posterior median thresholding rule from [1].

Bayesian Thresholding Rules

Bayes rules under the squared error loss and regular models are never thresholding rules. We discuss two possible approaches for obtaining *bona fide* thresholding rules in a Bayesian manner. The first one is via hypothesis testing while the second one uses weighted absolute error loss.

Donoho and Johnstone (1994, 1995) gave a heuristic for the selection of the universal threshold via rejection regions of suitable hypotheses tests. Testing a precise hypothesis in Bayesian fashion requires a prior which has a point mass component. A method based on Bayes factors is discussed first. For details see [39].

Let

$$[d|\theta] \sim f(d|\theta).$$

After observing the coefficient d , the hypothesis $H_0 : \theta = 0$, versus $H_1 : \theta \neq 0$ is tested. If the hypothesis H_0 is rejected, θ is estimated by d . Let

$$[\theta] \sim \pi(\theta) = \pi_0 \delta_0 + \pi_1 \xi(\theta), \tag{18}$$

where $\pi_0 + \pi_1 = 1$, δ_0 is a point mass at 0, and $\xi(\theta)$ is a prior that describes distribution of θ when H_0 is false.

The resulting Bayesian procedure is:

$$\hat{\theta} = d \mathbf{1}(P(H_0|d) < \frac{1}{2}), \tag{19}$$

where

$$P(H_0|d) = (1 + \frac{\pi_1}{\pi_0} \frac{1}{B})^{-1}, \tag{20}$$

is the posterior probability of H_0 hypothesis, and $B = \frac{f(d|0)}{\int_{\theta \neq 0} f(d|\theta)\xi(\theta)d\theta}$ is the Bayes factor in favor of H_0 . The following result holds

For instance, let, as in Example 0.1,

$$d|\theta \sim \mathcal{DE}(\theta, \frac{1}{\sqrt{2\mu}})$$

and

$$\pi(\theta) = \pi_0\delta_0 + \pi_1\xi(\theta),$$

then d will be thresholded if

$$\frac{\pi_0 e^{-c|d|}}{\pi_0 e^{-c|d|} + \pi_1(\Pi_1(c) + \Pi_2(c))} \geq \frac{1}{2}, \quad (21)$$

where Π_1 and Π_2 are one-sided Laplace transformations of $\xi(\theta - d)$ and $\xi(\theta + d)$.

Abramovich *et al.* (1998) use weighted absolute error loss and show that for a prior on θ

$$[\theta] \sim \pi_j N(0, \tau_j^2) + (1 - \pi_j)\delta(0) \quad (22)$$

and normal $N(\theta, \sigma^2)$ likelihood, the posterior median is

$$\text{Med}(\theta|d) = \text{sign}(d) \max(0, \zeta). \quad (23)$$

Here

$$\begin{aligned} \zeta &= \frac{\tau_j^2}{\sigma^2 + \tau_j^2} |d| - \frac{\tau_j \sigma}{\sqrt{\sigma^2 + \tau_j^2}} \Phi^{-1}\left(\frac{1 + \min(\omega, 1)}{2}\right), \text{ and} \\ \omega &= \frac{1 - \pi_j}{\pi_j} \frac{\sqrt{\tau_j^2 + \sigma^2}}{\sigma} \exp\left\{-\frac{\tau_j^2 d^2}{2\sigma^2(\tau_j^2 + \sigma^2)}\right\}. \end{aligned}$$

The index j , as before, points to the level containing θ (or d) facilitating adaptivity. The plot of the thresholding function (23) is given in Figure 4, Right.

Abramovich *et al.* (1998) assume

$$\tau_j^2 = C_1 2^{-\alpha j} \text{ and } \pi_j = \min(1, C_2 2^{-\beta j}) \quad (24)$$

where C_1, C_2, α , and β are non-negative hyperparameters. The hyperparameters α and β are determined from the assumption that the function to be estimated belongs to a particular Besov space, while C_1 and C_2 are determined in an empirical Bayes fashion.

The authors compare their BayesThresh rule (??) with several existing methods: Cross-Validation, False Discovery Rate, VisuShrink and GlobalSure and report very good MSE performance.

Density Estimation Problem

Delyon and Juditsky (1993), Donoho *et al.* (1995), Vannucci (1995), Pinheiro and Vidakovic (1997), Walter and Shen (1997), among others, applied wavelets in density estimation from a classical and data analytic perspective.

Chencov (1962) proposed projection type density estimators in terms of an arbitrary orthogonal basis. In the case of wavelet basis Chencov's estimator has the form

$$\hat{f}(x) = \sum_k c_{j_0 k} \phi_{j_0 k}(x) + \sum_{j_0 \leq j \leq j_1} \sum_k d_{jk} \psi_{jk}(x), \quad (25)$$

where the coefficients c_{jk} and d_{jk} constituting the vector \underline{d} are defined via standard empirical counterparts of $\langle f, \phi_{jk} \rangle$ and $\langle f, \psi_{jk} \rangle$. Let X_1, \dots, X_n be a random sample from f . Then

$$c_{jk} = \frac{1}{n} \sum_{i=1}^n \phi_{jk}(X_i), \quad \text{and} \quad d_{jk} = \frac{1}{n} \sum_{i=1}^n \psi_{jk}(X_i). \quad (26)$$

The local nature of wavelet functions makes the wavelet estimator superior to projection estimators that use classical orthonormal bases (Fourier, Hermite, etc.).

Brunk (1978) first proposed the shrinkage of coefficients in the projection density estimate based on linear Bayesian estimation. When an \mathbb{L}_2 density f has the Fourier expansion $f(x) = \sum_n a_n \xi(x)$ in some orthonormal basis $\xi_n(x)$, $n \in \mathbb{N}$, Brunk (see also Wahba, 1981) assumed normal $\mathcal{N}(0, b_n)$ and independent priors on the coefficients a_n . The suggested choice for b_n was C/m^n , where C is a constant and $m > 1$.

Brunk's linear Bayes shrinkage paradigm, in the wavelet setup, was explored by Vidakovic and Müller (1995). They assumed

$$[\underline{d} | \underline{\theta}, \sigma^2] \sim \mathcal{MVN}(\underline{\theta}, \sigma^2 I). \quad (27)$$

The rationale behind the model in (27) is that both c_{jk} and d_{jk} are averages of n iid random variables, see (26). Given $\underline{\theta}$, the components of \underline{d} are independent. The model in (27) was completed by specifying the prior distribution on the location $\underline{\theta}$ and scale σ .

Let

$$[\underline{\theta}, \sigma^2] \sim \mathcal{NIG}(\alpha, \delta, 0, \Sigma), \quad (28)$$

where $\mathcal{NIG}(\alpha, \delta, m, \Sigma)$ is the normal-inverse-gamma distribution with density function $f(\underline{\theta}, \sigma^2 | \alpha, \delta, m, \Sigma) = C \cdot (\sigma^2)^{\frac{\delta+p+2}{2}} \exp[-\{(\underline{\theta} - \underline{m})' \Sigma^{-1} (\underline{\theta} - \underline{m}) + \alpha\} / (2\sigma^2)]$.

Bayesian updating is straightforward because of the conjugate structure. The choice of a prior covariance matrix Σ corresponds to the choice of variances b_n in the traditional (Fourier) problems discussed in [2] and [41]. Authors consider a more general case in which the wavelet coefficients are not independent. More precisely, the coefficients in different levels are assumed independent, but the coefficients in the same level are assumed correlated. If the coefficients between levels j_0 and j_1 are considered in the estimator (25), covariance structure can be described by the following covariance matrix:

$$\Sigma = \lambda_{j_0, \phi} \Sigma_{j_0, \phi} \oplus \bigoplus_{j=j_0}^{j_1} \lambda_{j, \psi} \Sigma_{j, \psi}, \quad (29)$$

where \oplus is direct sum operation. Σ in (29), is a block diagonal matrix in which the block sub-matrices $\Sigma_{j_0, \phi}, \Sigma_{j, \psi}, j = j_0, \dots, j_1$ describe correlations within their corresponding levels.

The posterior for $[\theta, \sigma^2]$ is again the normal-inverse-gamma distribution

$$[\theta, \sigma^2 | \underline{d}] \sim \mathcal{NIG}(\alpha^*, \delta^*, m^*, \Sigma^*),$$

with

$$\begin{aligned} \Sigma^* &= (I + \Sigma^{-1})^{-1} \\ \hat{\theta} &= \Sigma^* \underline{d} \quad (= E(\theta | \underline{d})) \\ \alpha^* &= \alpha + \|\underline{d}\| + (\theta^*)' (\Sigma^*)^{-1} \theta^* \\ \delta^* &= \delta + n. \end{aligned}$$

The Bayes estimator of θ is $\hat{\theta} = \Sigma^* \underline{d}$. Naturally, we estimate the density in (25) by replacing c_{jk} and d_{jk} from \underline{d} by their affine Bayes estimators.

Example 0.2 In practical implementations of the method the matrix Σ was block diagonal matrix consisting of Laurent submatrices $\Sigma_{j_0, \phi}, \Sigma_{j_0, \psi}, \dots, \Sigma_{j_1, \psi}$ of appropriate dimensions with entries $\sigma_{i,j} = \rho^{|i-j|}$, $|\rho| < 1$.

Figure 5 shows shrinkage of empirical wavelet coefficients \underline{d} for the *galaxy velocities* data set (Roeder, 1990). There were 176 empirical coefficients between levels $j_0 = 0$ and $j_1 = 6$. The following values for λ s and ρ were adopted: $\lambda_{0, \phi} = 1000$, $\lambda_{0, \psi} = 1000$, $\lambda_{1, \psi} = 100$, $\lambda_{2, \psi} = 10$, $\lambda_{3, \psi} = 1$, $\lambda_{4, \psi} = 0.1$, $\lambda_{5, \psi} = 0.01$, $\lambda_{6, \psi} = 0.001$, and $\rho = 0.8$.

Building on (27)-(29), Huerta (1997) proposed the following model

$$[\underline{d} | \theta, \sigma^2] \sim MVN(\theta, \sigma^2 I), \quad [\sigma^2] \sim IG(\alpha_1, \delta_1),$$

$$[\theta | \tau^2] \sim MVN(\mathbf{0}, \tau^2 \Sigma), \quad [\tau^2] \sim IG(\alpha_2, \delta_2).$$

The closed form for $E(\theta | \underline{d})$ in the above model is not possible, however the full conditionals can be specified. The shrinkage is performed by Gibbs sampling. The prior covariance matrix Σ was defined as in (29). Obtained simulation results results are promising.

Reduction of number of hyperparameters in (27)-(29) and imposing a prior structure on remaining hyperparameters in a hierarchical fashion was suggested by Vannucci and Corradi (1996). They applied the shrinkage on regression problem call their method BayesShrink. Assuming that θ corresponds to an autoregressive process in the time domain the authors demonstrate that the matrix Σ depends on only two hyperparameters, λ and ρ . The parameter ρ is the ‘‘autocovariance index’’ and λ is the precision parameter. The covariance matrix $\Sigma(\lambda, \rho) = \lambda \Sigma(\rho)$ has an interesting ‘‘finger-like’’ structure.

Authors in [36] suggest

$$[\lambda] \sim IG(p/2, q/2), \text{ and}$$

$$[\rho] \propto (C - \rho)^{r_1 - 1} (C + \rho)^{r_2 - 1}, \quad |\rho| < C,$$

and simulate from the posterior distribution of θ using MCMC method.

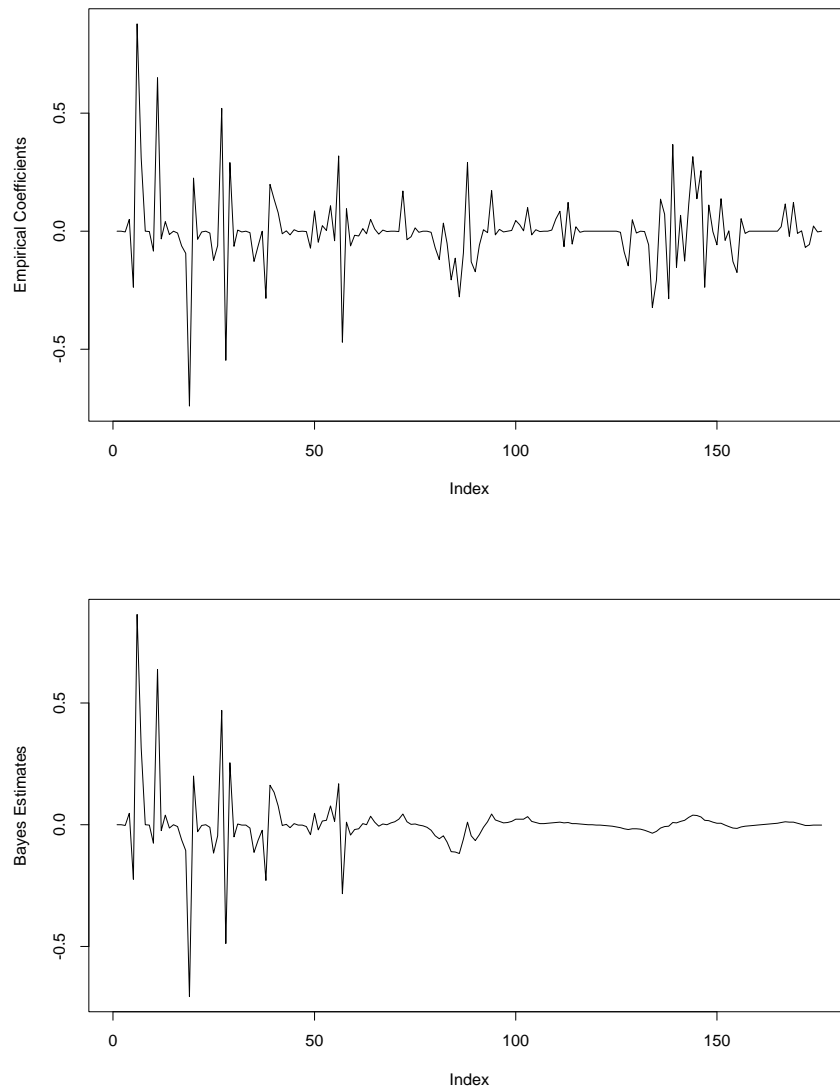


Figure 5: Empirical coefficients (26) for the `galaxy` data and their affine Bayes shrinkage. The coefficients are arranged as a concatenation of levels in the decomposition, starting with scaling coefficients and ending with coefficients of fine detail. Notice that the Bayes estimate preserves coarse structure and heavily shrinks detail coefficients that cause overfitting.

Full Bayesian Model

Müller and Vidakovic (1994) parameterize unknown density $f(x)$ by a wavelet series on its logarithm, and propose a prior model which explicitly defines geometrically decreasing prior probabilities for non-zero wavelet coefficients at higher levels of detail.

The unknown probability density function $f(\cdot)$ is modeled by:

$$\log f(x) = \sum_{k \in Z} \xi_{j_0 k} \phi_{j_0, k}(x) + \sum_{j \geq j_0, k \in Z} s_{jk} \theta_{jk} \psi_{jk}(x) - \log K, \quad (30)$$

where $K = \int f(x) dx$ is the normalization constant and $s_{jk} \in \{0, 1\}$ is an indicator variable that performs model induced thresholding. A prior on $[\theta_{jk} | s_{jk}]$ which assigns considerable prior probability mass at 0 was used.

The dependence of $f(x)$ on the vector $\theta = (\xi_{j_0, k}, s_{jk}, \theta_{jk}, j = j_0, \dots, j_1, k \in Z)$ of wavelet coefficients and indicators is expressed by $f(x) = p(x|\theta)$. The sample $X = \{X_1, \dots, X_n\}$ defines a likelihood function $p(X|\theta) = \prod_{i=1}^n p(X_i|\theta)$.

The model is completed by a prior probability distribution for θ . Without loss of generality $j_0 = 0$ can be assumed. Also, any particular application will determine a maximum level of detail j_1 .

$$\begin{aligned} \xi_{0k} &\sim N(0, \tau r_0), \\ \theta_{jk} | s_{jk} = 1 &\sim N(0, \tau r_j), \quad r_j = 2^{-j}, \\ \theta_{jk} | s_{jk} = 0 &\sim h(\theta_{jk}), \\ s_{jk} &\sim \text{Bernoulli}(\alpha^j), \\ \alpha &\sim \text{Beta}(a, b), \\ 1/\tau &\sim \text{Ga}(a_\tau, b_\tau). \end{aligned} \quad (31)$$

The wavelet coefficients θ_{jk} are non-zero with geometrically decreasing probabilities. Given that a coefficient is non-zero, it is generated from a normal distribution. When a coefficient is not included in (30), that is, $s_{jk} = 0$, a ‘‘pseudo-prior’’ $h(\theta_{jk})$ is assumed. The parameter vector θ is augmented in order to include all model parameters, i.e. $\theta = (\theta_{jk}, \xi_{jk}, s_{jk}, \alpha, \tau)$.

The scale factor r_j contributes to the adaptivity of the method. Wavelet shrinkage is controlled by both: the factor r_j and geometrically decreasing prior probabilities for non-zero coefficient, α^j .

The conditional prior $p(\theta_{jk} | s_{jk} = 0) = h(\theta_{jk})$ in model (31) is a pseudo-prior as discussed in Carlin and Chib (1995). The choice of $h(\cdot)$ has no bearing on the inference about $f(x)$. In fact, the model could be alternatively formulated by dropping θ_{jk} under $s_{jk} = 0$. However, this would lead to a parameter space of varying dimension. Carlin and Chib (1995) argue that the pseudo-prior $h(\theta_{jk})$ should be chosen to produce values for θ_{jk} which are consistent with the data. The normal distribution $p(\theta_{jk} | s_{jk} = 0) = N(\hat{\theta}_{jk}, \sigma_{jk})$, where $\hat{\theta}_{jk}$ is some rough preliminary estimate of θ_{jk} was proposed.

The particular MCMC simulation scheme used to estimate model (30) and (31) is described. Starting with some initial values for $\theta_{jk}, j = 0, \dots, j_1, \xi_{00}, \alpha$ and τ , the following Markov chain was implemented by Peter Müller:

- 1. For each $j = 0, \dots, j_1$ and $k = 1, \dots, 2^j - 1$ go over the steps 2 and 3.
- 2. Updating s_{jk} . Let θ_0 and θ_1 indicate the current parameter vector θ with s_{jk} replaced by 0 and 1, respectively. Compute $p_0 = p(y|\theta_0) \cdot (1 - \alpha^j) h(\theta_{jk})$ and $p_1 = p(y|\theta_0) \cdot \alpha^j p(\theta_{jk} | s_{jk} = 1)$. With probability $p_1 / (p_0 + p_1)$ set $s_{jk} = 1$, else $s_{jk} = 0$.

- 3a. Updating θ_{jk} . If $s_{jk} = 1$, generate $\tilde{\theta}_{jk} \sim g(\tilde{\theta}_{jk}|\theta_{jk})$. Use, for example, $g(\tilde{\theta}_{jk}|\theta_{jk}) = N(\theta_{jk}, 0.25\sigma_{jk})$, where σ_{jk} is some rough estimate of the posterior standard deviation of θ_{jk} . We will discuss alternative choices for the probing distribution $g(\cdot)$ below.

Compute

$$a(\theta_{jk}, \tilde{\theta}_{jk}) = \min \left[1, \frac{p(y|\tilde{\theta})p(\tilde{\theta}_{jk})}{p(y|\theta)p(\theta_{jk})} \right],$$

where $\tilde{\theta}$ is the parameter vector θ with θ_{jk} replaced by $\tilde{\theta}_{jk}$, and $p(\theta_{jk})$ is the p.d.f. of the normal prior distribution given in (31).

With probability $a(\theta_{jk}, \tilde{\theta}_{jk})$ replace θ_{jk} by $\tilde{\theta}_{jk}$; else keep θ_{jk} unchanged.

- 3b. If $s_{jk} = 0$, generate θ_{jk} from the full conditional posterior $p(\theta_{jk} | \dots, X) = p(\theta_{jk} | s_{jk} = 0) = h(\theta_{jk})$.
- 4. Update ξ_{00} . Generate $\tilde{\xi}_{00} \sim g(\tilde{\xi}_{00}|\xi_{00})$. Use, for example, $g(\tilde{\xi}_{00}|xi_{00}) = N(xi_{00}, 0.25\rho_{00})$, where ρ_{00} is some rough estimate of the posterior standard deviation of xi_{00} . Analogously to step 3a, compute an acceptance probability a and replace ξ_{00} with probability a .
- 5. Update α . Generate $\tilde{\alpha} \sim g_{\alpha}(\tilde{\alpha}|\alpha)$ and compute

$$a(\alpha, \tilde{\alpha}) = \min \left[1, \frac{\prod_{jk} \tilde{\alpha}^{j s_{jk}} (1 - \tilde{\alpha}^j)^{s_{jk}}}{\prod_{jk} \alpha^{j s_{jk}} (1 - \alpha^j)^{s_{jk}}} \right].$$

With probability $a(\alpha, \tilde{\alpha})$ replace α by $\tilde{\alpha}$, else keep α unchanged. See below for comments about $g_{\alpha}(\cdot)$.

- 6. Update τ . Resample τ from the complete inverse Gamma conditional posterior.
- 7. Iterate over steps 1 through 6 until the chain is judged to have practically converged.

The algorithm implements a Metropolis chain changing one parameter at a time in the parameter vector. See, for example, Tierney (1994) for a description and discussion of Metropolis chains for posterior exploration. For a practical implementation, g should be chosen such that the acceptance probabilities a are neither close to zero, nor close to one. In the implementations $g(\tilde{\theta}_{jk}|\theta_{jk}) = N(\theta_{jk}, 0.25\sigma_{jk})$ with $\sigma_{jk} = 2^{-j}$, was used.

Example 0.3 Described wavelet based density estimation model is illustrated on the galaxy data set (Roeder, 1990). The data is rescaled to the interval $[0, 1]$. The hyperparameters were fixed as $a = 10$, $b = 10$ and $a_{\tau} = b_{\tau} = 1$. The $Beta(10, 10)$ prior distribution on α is reasonably non-informative compared to the likelihood based on $n = 82$ observations.

Initially, all s_{jk} are set to one, and α to its prior mean $\alpha = 0.5$. The first 10 iterations as burn-in period were discarded, then 1000 iterations of steps 1 through 6 were simulated. For each j, k Step 3 was repeated three times. The maximum level of detail selected was $j_1 = 5$.

Figures 6 and 7 describe some aspects of the analysis.

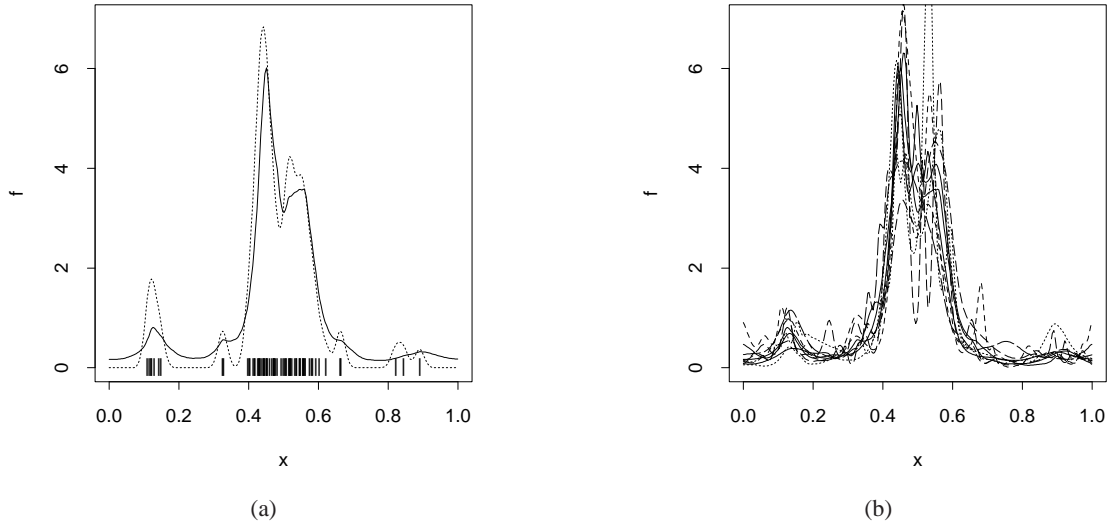


Figure 6: [27] Left: The estimated p.d.f. $\hat{f}(x) = \int p(x|\theta)dp(\theta|X)$. The dotted line plots a conventional kernel density estimate for the same data. Right: The posterior distribution of the unknown density $f(x) = p(x|\theta)$ induced by the posterior distribution $p(\theta|X)$. The lines plot $p(x|\theta_i)$ for ten simulated draws posterior $\theta_i \sim p(\theta|X), i = 1, \dots, 10$.

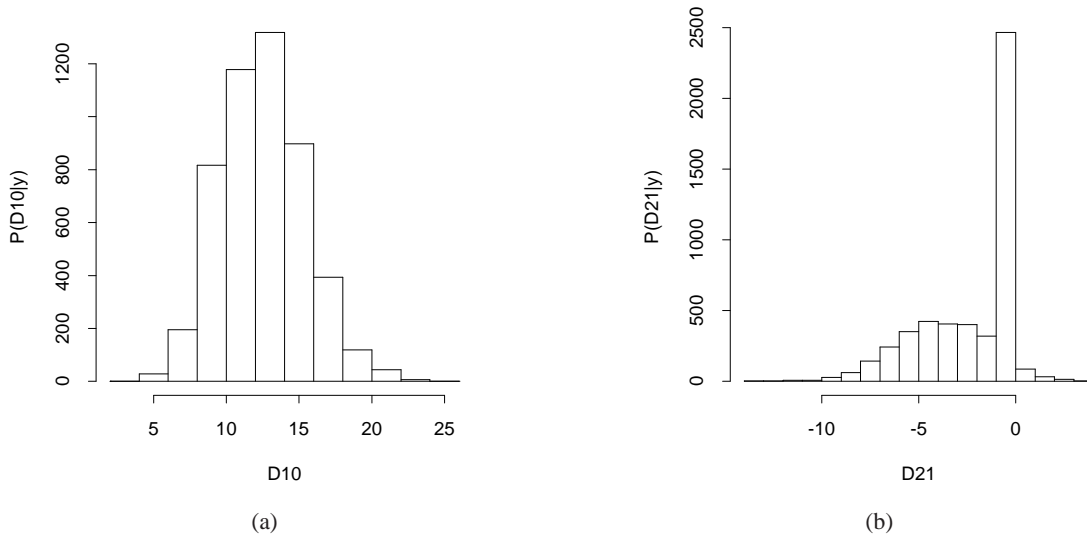


Figure 7: Posterior distributions $p((s_{10}\theta_{10})|X)$ and $p((s_{21}\theta_{21})|X)$. While $s_{10}\theta_{10}$ is non-zero with posterior probability close to one, the posterior distribution $p((s_{21}\theta_{21})|X)$ is a mixture of a point mass at zero and a continuous part.

Other Problems

Lina and MacGibbon (1997) apply Bayesian approach to the wavelet regression with complex valued Daubechies wavelets. To some extent they exploit redundancy in the representation of real signals by the complex wavelet coefficients. Their shrinkage technique is based on the observation that modulus and phase of wavelet coefficients encompass very different information about the signal. Bayesian shrinkage model is constructed for the modulus taking into account the corresponding phase.

Simoncelli and Adelson (1996) discuss Bayes “coring” procedure in the context of image processing. The prior on signal is Mallat’s model, see [25], while the noise is assumed normal. They implement their noise reduction scheme on an oriented multiresolution representation - known as the *steerable pyramid*. They report that Bayesian coring outperforms classical Wiener filtering.

Malfait and Roose (1995) and Malfait *et al.* (1996) employed the theory of Markov random fields in a wavelet denoising of images. The idea is that the thresholding of a coefficient in a 2-D array depends not only on its magnitude but also on the magnitudes of its neighbors. This dependence is modeled by using masks consisting of 0 and 1 and exhibiting properties of a Markov field. Crouse *et al.* (1997) also consider hidden Markov fields in a similar setup. They develop Efficient Expectation Maximization (EEM) algorithm to fit their model.

Pesquet et al (1996) develop a Bayesian-based approach to the best basis problem, while preserving the classical tree search efficiency in wavelet packets and local trigonometric bases. Kohn and Marron (1997) use the model similar to one in [21] but in the context of the best basis selection. Gendron (1997) builds on Bayesian hypothesis testing idea in the wavelet domain. He uses Bayesian paradigm to construct a cost function which points to the best wavelet packet. Using the best basis he develops a filtering procedure for seismic signals by applying posterior expectations when the prior is a mixture of the uniform distribution and a point mass at 0.

Ogden (1996) considers the change-point problem from a Bayesian point.

Ruggeri and Vidakovic (1997) discuss Bayesian decision theoretic thresholding. In the set of all hard thresholding rules they find restricted Bayes rule under variety of models, priors, and loss functions. They identify model-prior pairs that work well (in sense of exhibiting hard thresholding rule that minimizes Bayes risk) and show that in presence of prior information on noise and signal their procedure outperforms SureShrink and VisuShrink methods.

Lu, Huang, and Tung (1997) introduce linear Bayesian wavelet shrinkage in a non-parametric mixed-effect model. Their formulation is conceptually inspired by duality between reproducing kernel Hilbert spaces and random processes as well as on connections between smoothing splines and Bayesian regressions.

The unknown function f in the standard non-parametric regression formulation $y = f(x_i) + \sigma\epsilon_i$, $i = 1, \dots, n$; $0 \leq x \leq 1$; $\sigma > 0$; $Cov(\epsilon_1, \dots, \epsilon_n) = R$; is given a prior of the form, $f(x) = \sum_k \alpha_{Jk} \phi_{Jk}(x) + \delta Z(x)$; $Z(x) \sim \sum_{j \geq J} \sum_k \theta_{jk} \psi_{jk}(x)$ where θ_{jk} are uncorrelated random variables such that $E\theta_{jk} = 0$ and $E\theta_{jk}^2 = \lambda_j$.

The authors propose a linear, empirical Bayes estimator \hat{f} of f that enjoys Gauss-Markov type of optimality. Several non-linear versions of the estimator are proposed, as well.

Independently, and by using different techniques, Huang and Cressie (1997) consider the same problem and derive a Bayesian estimate.

References

- [1] ABRAMOVICH, F., SAPATINAS, T. and SILVERMAN, B. W. (1998). Wavelet Thresholding via Bayesian Approach. *Journal of the Royal Statistical Society*, Vol. 60, No.3, 1998.

- [2] BRUNK, H. (1978). Univariate density estimation by orthogonal series, *Biometrika* **65**, 3, 521-528.
- [3] CARLIN, B. and CHIB, S. (1995). Bayesian model choice via Markov chain Monte Carlo. *Journal of the Royal Statistical Society, Series B*, 57, 473-484.
- [4] CHENCOV, N. N. (1962). Evaluation of an unknown distribution density from observations. *Doklady*, **3**, 1559-1562.
- [5] CHIPMAN, H., MCCULLOCH, R. and KOLACZYK, E. (1997). Adaptive Bayesian Wavelet Shrinkage, To appear in the *Journal of the American Statistical Association*.
- [6] CLYDE, M., DESIMONE, H. and PARMIGIANI, G. (1996). Prediction via orthogonalized model mixing. To appear *Journal of the American Statistical Association*
- [7] CLYDE, M., PARMIGIANI, G. and VIDAKOVIC, B. (1996). Bayesian Strategies for Wavelet Analysis, Statistical Computing & Graphics, Newsletter of Statistical Computing and Statistical Graphics sections of ASA, August 1996 (Special issue on Bayesian Function Estimation).
- [8] CLYDE, M., PARMIGIANI, G., and VIDAKOVIC, B. (1995). Multiple Shrinkage and Subset Selection in Wavelets, Discussion Paper **95-37**, ISDS, Duke University. To appear in *Biometrika*.
- [9] COIFMAN, R. and WICKERHAUSER, V. (1992). Entropy based algorithms for best basis selection. *IEEE Trans. Inform. Theory*, 38 (2), 713-718.
- [10] CROUSE, M., NOWAK, R. and BARANIUK, R. (1997). Statistical Signal Processing Using Wavelet-Domain Hidden Markov Models, Proceedings of SPIE, Wavelet Applications in Signal and Image Processing V, vol. 3169, 248-259.
- [11] DAUBECHIES, I. (1992). *Ten Lectures on Wavelets*. Society for Industrial and Applied Mathematics.
- [12] DELYON, B. and JUDITSKY, A. (1993). Wavelet Estimators, Global Error Measures: Revisited. *Publication interne no. 782, IRISA-INRIA*.
- [13] DONOHO, D. (1994). On minimum entropy segmentation, In *Wavelets: Theory, Algorithms, and Applications*, Eds Chui, Montefusco and Puccio, Academic Press.
- [14] DONOHO, D. and JOHNSTONE, I. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika* 81, 425-455.
- [15] DONOHO, D. and JOHNSTONE, I. (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association* 90, 1200-1224.
- [16] DONOHO, D., JOHNSTONE, I., KERKYACHARIAN, G. and PICKARD, D. (1995). Wavelet shrinkage: Asymptopia? *J. R. Statis. Soc.* **57** (2) 301-369.
- [17] GEORGE, E.I. and MCCULLOCH, R. (1994). Approaches to Bayesian Variable Selection. Tech Report Graduate School of Business, University of Chicago.
- [18] HERNÁNDEZ, E. and WEISS, G. (1996). *A First Course on Wavelets*. Boca Raton: CRC Press Inc.
- [19] HUERTA, G. (1997). Bayes Wavelet Shrinkage and Applications to Data Denoising. In electronic proceedings of International Workshop on Wavelets in Statistics, Duke University, 12-13 October 1997.
- [20] JOHNSTONE, I. and SILVERMAN B. W. (1996). Wavelet threshold estimators for data with correlated noise. *Journal of the Royal Statistical Society B*, **59**, 319-351.
- [21] KOHN R, and MARRON J. S. (1997). Bayesian Wavelet Shrinkage, International Workshop on Wavelets in Statistics, Duke University, 12-13 October 1997.
- [22] LINA, J-M. and MacGibbon, B. (1997). Non-Linear Shrinkage Estimation with Complex Daubechies Wavelets, Proceedings of SPIE, Wavelet Applications in Signal and Image Processing V, vol. 3169, 67-79.
- [23] LU, H. H-S., HUANG, S-Y., and TUNG Y-C. (1997). Wavelet Shrinkage for nonparametric Mixed-Effects Models Technical Report Institute of Statistics, National Chiao Tung University Also in: Electronic proceedings of International Workshop on Wavelets in Statistics, Duke University, 12-13 October 1997.

- [24] MALFAIT, M. and ROOSE, D. (1995). Wavelets and Markov Random Fields in a Bayesian Framework, In *Wavelets and Statistics*, Eds: A. Antoniadis and G. Oppenheim. lecture Notes in Statistics, Springer-Verlag 103 225-238.
- [25] MALLAT, S. (1989). A theory for multiresolution signal decomposition: the wavelet representation. *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 11, No. 7, 674-693.
- [26] MEYER, Y. (1992). *Wavelets and Operators*, Cambridge.
- [27] MÜLLER, P. and VIDA KOVIC, B. (1995). Bayesian Inference with Wavelets: Density Estimation. Discussion Paper **95-33**, ISDS, Duke University.
- [28] OGDEN, T. (1996). *Essential Wavelets for Statistical Applications and Data Analysis*, Birkhauser, Boston.
- [29] OGDEN, T. (1996). Wavelets in Bayesian Change-Point Analysis. Technical Report, Department of Statistics, University of South Carolina.
- [30] PESQUET, J., KRIM, H., LEPORINI, D., and HAMMAN, E. (1996). Bayesian approach to best basis selection, In: *IEEE International Conference on Acoustics, Speech, and Signal Processing*, **5**, 2634-2637.
- [31] RIOS, D. and VIDA KOVIC, B. (1997) "Wavelet-Based Random Densities" Discussion Paper **97-05**, ISDS, Duke University.
- [32] ROEDER, K. (1990). Density estimation with confidence sets exemplified by superclusters and voids in the galaxies. *Journal of the American Statistical Association*, 85 617-624.
- [33] RUGGERI, F. and VIDA KOVIC, B. (1995). A Bayesian Decision Theoretic Approach to Wavelet Thresholding. Discussion Paper **95-35**, ISDS, Duke University.
- [34] SIMONCELLI, E. and ADELSON, E. (1996). Noise removal via Bayesian wavelet coring. Presented at: 3rd IEEE International Conference on Image Processing, Lausanne, Switzerland.
- [35] TIERNEY, L. (1994). Markov chains for exploring posterior distributions. *Annals of Statistics*, 22, 1701-1728.
- [36] VANNUCCI, M. and CORRADI, F. (1997). Some Findings on the Covariance Structure of Wavelet Coefficients: Theory and Models in a Bayesian Perspective. Report UKC/IMS **97-05**, Institute of Maths and Stats, University of Kent at Canterbury, UK.
- [37] VANNUCCI, M. (1995). Nonparametric Density Estimation Using Wavelets. Discussion Paper **95-26**, ISDS, Duke University.
- [38] VANNUCCI, M. and VIDA KOVIC, B. (1995). Preventing the Dirac Disaster: Wavelet Based Density Estimation, Discussion Paper **95-27**, ISDS, Duke University.
- [39] VIDA KOVIC, B. (1994). Nonlinear wavelet shrinkage with Bayes rules and Bayes factors, Discussion Paper **94-A-24**, ISDS, Duke University. To appear in *J. Amer. Statist. Assoc.* March 1998 issue.
- [40] VIDA KOVIC, B. and MÜLLER, P. (1995). Wavelet shrinkage with affine Bayes rules with applications. Discussion Paper **95-34**, ISDS, Duke University.
- [41] WAHBA, G. (1981). Data-based optimal smoothing of orthogonal series density estimates. *The Annals of Statistics*, **9**, 146-156.
- [42] WALTER, G.G. (1994). *Wavelets and Others Orthogonal Systems with Applications*. CRC Press, Boca Raton, FL.
- [43] WALTER, G. and SHEN, X. (1997). Continuous Non-negative Wavelets and Their Use in Density Estimation, In electronic proceedings of International Workshop on Wavelets in Statistics, Duke University, 12-13 October 1997,
- [44] WICKERHAUSER, M. V. (1994). *Adapted Wavelet Analysis from Theory to Software*. A.K. Peters, Ltd., Wellesley, MA.