

# Minimax Bayes, asymptotic minimax and sparse wavelet priors

Iain M. Johnstone  
Department of Statistics  
Stanford University

December 1992  
Revised April 1993

## Abstract

Pinsker(1980) gave a precise asymptotic evaluation of the minimax mean squared error of estimation of a signal in Gaussian noise when the signal is known *a priori* to lie in a compact ellipsoid in Hilbert space. This ‘Minimax Bayes’ method can be applied to a variety of global non-parametric estimation settings with parameter spaces far from ellipsoidal. For example it leads to a theory of exact asymptotic minimax estimation over norm balls in Besov and Triebel spaces using simple coordinatewise estimators and wavelet bases.

This paper outlines some features of the method common to several applications. In particular, we derive new results on the exact asymptotic minimax risk over weak  $\ell_p$ - balls in  $\mathcal{R}^n$  as  $n \rightarrow \infty$ , and also for a class of ‘local’ estimators on the Triebel scale.

By its very nature, the method reveals the structure of asymptotically least favorable distributions. Thus we may simulate ‘least favorable’ sample paths. We illustrate this for estimation of a signal in Gaussian white noise over norm balls in certain Besov spaces. In wavelet bases, when  $p < 2$ , the least favorable priors are sparse, and the resulting sample paths strikingly different from those observed in Pinsker’s ellipsoidal setting ( $p = 2$ ).

**Key Words.** Minimax Decision theory. Minimax Bayes estimation. Besov, Hölder, Sobolev, Triebel Spaces. Nonlinear Estimation. White Noise Model. Nonparametric regression. Orthonormal Bases of Compactly Supported Wavelets. Threshold rules.

**Acknowledgements.** I am grateful for many conversations with David Donoho and Carl Taswell, and to a referee for helpful comments. This work was supported in part by NSF grants DMS 84-51750, 9209130, and NIH PHS grant GM21215-12.

# 1 Introduction

The minimax approach to the theory of non-parametric estimation of a function  $\theta$  known to lie in a fixed set  $\Theta$  aims in part to quantify the effect of the constraints defining  $\Theta$  on the possible quality of estimation of  $\theta$ . Asymptotic approximations in the small noise or large sample limit are often necessary (e.g. Chentsov (1972), Farrell (1972), Ibragimov and Hasminskii (1981, 1982) and Stone(1982)).

Often these approximations yield information only on the *rate* of estimation possible, and so it was remarkable when Pinsker (1980) was able to identify the exact *constant* in the asymptotic minimax risk of estimation of a signal belonging to an ellipsoid in Hilbert space when observed in Gaussian noise. Subsequently other applications to ellipsoidal constraint sets were given, for example, by Efroimovich and Pinsker (1981, 1982), Nussbaum (1985), Johnstone and Silverman (1990) and Golubev and Nussbaum (1990).

This paper describes an extension of Pinsker's method that we have found useful in deriving asymptotic minimax risks in a number of distinctly non-ellipsoidal settings suggested by the use of wavelet bases (cf. Donoho and Johnstone (1990, 1992a)).

In addition to outlining the method in skeletal form, this paper has three objectives. The first is to use it to give an exact asymptotic evaluation of the minimax mean squared error of estimation over weak  $\ell_p$ - balls in  $\mathcal{R}^n$  as the radius  $r_n$  and noise level  $\epsilon_n$  vary with  $n \rightarrow \infty$ . The results are relevant for describing the best attainable spatial adaptation by non-parametric regression estimators over function classes for which approximation at a given rate is possible. Secondly, we use the method to give a new asymptotic minimax result for a class of 'local' estimators over spaces in the Triebel scale, which includes the classical Sobolev spaces.

Finally, we present some graphs of the sample paths from numerical approximations to the least favorable distributions in a white noise estimation problem when  $\Theta$  is a Besov space  $\Theta_{p,q}^\sigma(C)$ . Donoho and Johnstone (1992a) determined the asymptotic minimax risk for such  $\Theta$  using the approach outlined here. The graphs illustrate how  $p$  acts as an important shape parameter modifying smoothness and support the heuristic that when  $p < 2$ , the least favorable  $\theta$  correspond to relatively sparse signals (at least when viewed in the wavelet domain).

## 2 Minimax risk over ellipsoids

We begin by recalling a special case of Pinsker's result. Consider a homoscedastic Gaussian sequence model

$$y_i = \theta_i + \varepsilon z_i, \quad z_i \stackrel{i.i.d.}{\sim} N(0, 1) \quad i = 1, 2, \dots \quad (1)$$

where it is desired to estimate  $\theta = (\theta_i)$  using squared error loss  $\|\hat{\theta} - \theta\|^2 = \sum_i (\hat{\theta}_i - \theta_i)^2$ . We assume that  $\theta = (\theta_i)$  belongs to the ellipsoid

$$\Theta = \{\theta : \sum a_i^2 \theta_i^2 \leq C^2\}, \quad a_i \nearrow \infty.$$

The model (1) is, for example, the Fourier coefficient form of the usual signal in white Gaussian noise model  $Y(t) = \int_0^t f(s)ds + \varepsilon W(t)$ , where  $W(t)$  is standard Brownian motion.

Pinsker gave an exact evaluation of the asymptotic minimax risk for estimation of  $\theta$  as  $\varepsilon \rightarrow 0$  :

$$R^* = R(\Theta, \varepsilon) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} E \|\hat{\theta} - \theta\|^2 \sim \sup \left\{ \sum_i \frac{\varepsilon^2 t_i^2}{\varepsilon^2 + t_i^2} : \Sigma a_i^2 t_i^2 \leq C^2 \right\} \quad (2)$$

For example, in the common case where  $\Theta$  represents a bound on the mean square  $\sigma^{th}$  derivative, if we put  $a_i = i^2$  and  $r = 2\sigma/(2\sigma + 1)$ , then (2) gives  $R^* \sim \beta_r C^{2(1-r)} \varepsilon^{2r}$  as  $\varepsilon \rightarrow 0$ , where  $\beta_r = (2r/2 - r)^r (1 - r)^{(r-1)}$ . An important building block in Pinsker's proof is the *univariate* Bayes risk for estimation of  $\theta_1$  from  $y_1 \sim N(\theta_1, \varepsilon)$  when  $\theta_1 \sim F(d\theta_1)$  : we write

$$b(F, \varepsilon) = \inf_{\delta} E_F E_{\theta_1} [\delta(y_1) - \theta_1]^2.$$

Of course, if  $F = \Phi_t$ , a centered Gaussian of variance  $t^2$ , then  $b(F, \varepsilon) = \varepsilon^2 t^2 / (\varepsilon^2 + t^2)$ . Thus, an equivalent form of (2) is

$$R^* \sim \sup \left\{ \sum_i b(\Phi_{t_i}, \varepsilon) : t \in \Theta \right\}. \quad (3)$$

The right side is an instance of the type of minimax Bayes problem that we shall study. The infinite dimensional Gaussian priors  $\pi$  formed by making  $\theta_i$  independently distributed as  $\Phi_{t_i}$  are not supported on  $\Theta$ , but Pinsker shows that it is possible to choose  $\pi_\varepsilon$  approaching the supremum in (3) as  $\varepsilon \rightarrow 0$  in such a way that  $\pi_\varepsilon(\Theta) \nearrow 1$ . Thus Gaussian priors are asymptotically least favorable, and the corresponding (linear!) estimators are asymptotically minimax over the *ellipsoid*  $\Theta$ .

### 3 The Minimax Bayes method

In this section, we highlight elements of Pinsker's approach that seem useful in a variety of situations. For simplicity, assume a sequence space version of the Gaussian white noise model, indexed by  $n = 1, 2, \dots$  :

$$\begin{aligned} y_i &= \theta_i + \varepsilon_n z_i, \quad z_i \stackrel{iid}{\sim} N(0, 1) \\ \theta &\in \Theta_n. \end{aligned} \quad (4)$$

The frequentist minimax risk of a class of estimators  $\hat{\Theta}_n$  over  $\Theta_n$  is

$$R_n^* = R(\Theta_n, \varepsilon_n) = \inf_{\hat{\theta} \in \hat{\Theta}_n} \sup_{\theta \in \Theta_n} E_\theta \|\hat{\theta} - \theta\|^2. \quad (5)$$

We consider only squared error loss and exact Gaussian error models in this paper, but this is not essential to the method.

1°. *Minimax Bayes Problem.* Let  $\mathcal{M}_n$  be a collection of prior probability measures on sequence space: in general  $\mathcal{M}_n$  will *not* be supported on  $\Theta_n$ . An important idea behind the method is that a judiciously chosen relaxation of the constraints defining  $\Theta_n$  may be

easier to evaluate, and yet asymptotically equivalent. Thus, consider the Bayes minimax risk

$$B_n^* = B_n(\mathcal{M}_n, \varepsilon) = \inf_{\hat{\theta} \in \hat{\Theta}_n} \sup_{\pi \in \mathcal{M}_n} E_\pi E_\theta \|\hat{\theta} - \theta\|^2. \quad (6)$$

Assume that (i)  $\mathcal{M}_n$  contains  $\delta_\theta$ , point mass one at  $\theta$ , for all  $\theta \in \Theta_n$ , (ii)  $\mathcal{M}_n$  is convex and (iii) that  $\hat{\Theta}_n$  is convex. Assumption (i) entails that  $R_n^* \leq B_n^*$ ; while (ii) and (iii) permit use of the minimax theorem (e.g. Sion (1958, Theorem 4.2'), LeCam (1986, p. 16)). Thus,

$$R_n^* \leq B_n^* = \sup_{\pi \in \mathcal{M}_n} B_n(\pi).$$

where  $B_n(\pi)$  denotes the Bayes risk

$$B_n(\pi) = \inf_{\hat{\theta} \in \hat{\Theta}_n} E_\pi E_\theta \|\hat{\theta} - \theta\|^2.$$

2°. *Reduction of the Minimax Bayes Problem.* In certain cases, there will be a subclass  $\bar{\mathcal{M}}_n \subset \mathcal{M}_n$  which is as difficult as  $\mathcal{M}_n$ , in the sense that

$$\sup_{\bar{\mathcal{M}}_n} B_n(\pi) = \sup_{\mathcal{M}_n} B_n(\pi). \quad (7)$$

while having simpler structure (priors with independent or even i.i.d co-ordinates for example). The space  $\bar{\mathcal{M}}_n$  may allow determination of  $B_n^*$  up to a single constant, for example by symmetry or renormalization arguments. A common strategy for establishing (7) is to construct a mapping for given  $\pi \in \mathcal{M}_n$  to an element  $\bar{\pi} \in \bar{\mathcal{M}}_n$  that is at least as difficult

$$B_n(\bar{\pi}) \geq B_n(\pi).$$

3°. *Asymptotic equivalence of  $R_n^* \leq B_n^*$ .* A standard way to obtain a lower bound for  $R_n^*$  is through the minimax theorem and the collection  $\mathcal{L}_n$  of prior distributions supported on  $\Theta_n$ :

$$R_n^* = \sup_{\pi \in \mathcal{L}_n} B(\pi).$$

A heuristic approach to verifying that  $R_n^* \sim B_n^*$  as  $n \rightarrow \infty$  is to choose  $\pi_n \in \bar{\mathcal{M}}_n$  such that  $B(\pi_n) \sim B_n^*$ , and to form  $\nu_n = \pi_n(\cdot | \Theta_n) \in \mathcal{L}_n$ . The hope is that the structure of  $\bar{\mathcal{M}}_n$  will force  $\pi_n$  to concentrate *asymptotically* on  $\Theta_n$ ; specifically  $\pi_n(\Theta_n) \rightarrow 1$ , and  $B(\pi_n) \sim B(\nu_n)$ . In this case the chain of relations

$$B(\nu_n) \leq R_n^* \leq B_n^* \sim B(\pi_n) \sim B(\nu_n) \quad (8)$$

establishes that  $R_n^* \sim B_n^*$ .

In the final section, we present some extra detail on one approach to making the heuristic ideas of 3° rigorous when there is in fact a *family* of minimax problems linked by a scale parameter  $C$ .

*Remark.* Minimax Bayes problems related to (6) have often been considered in the literature (see for example Hodges and Lehmann (1952), Bickel (1983), Morris (1983) and the references therein). Recently, Feldman (1991) considered a particular minimax Bayes problem of the kind also reviewed in Section 4.1 below. The focus here is on their use as a device in evaluating the frequentist minimax risk over  $\Theta$  in (5).

## 4 $\ell_p$ -balls in $\mathbb{R}^n$ , strong and weak

For this section, we consider the  $n$ -dimensional version of model (4), namely  $y \sim N_n(\theta, \varepsilon^2 I)$ . A simple example of the use of the Minimax Bayes technique lies in the asymptotic evaluation of minimax risk over balls in  $\ell_p$  norms in  $\mathbb{R}^n$ . These also form a building block for the main results on Besov balls in wavelet bases discussed in Donoho and Johnstone (1992a).

We then turn to new results on minimax risk over Marcinkiewicz, or weak  $\ell_p$  balls. The sequence space weak- $\ell_p$  is relevant here as a representation of an approximation space, namely the collection of all functions on  $[0, 1]$  that can be approximated in  $L^2[0, 1]$  norm at rate  $N^{-\sigma}$ ,  $\sigma = 1/p - 1/2$ , using any one of a number of non-linear approximation methods: largest wavelet coefficients, piecewise polynomials or dyadic splines with variable breakpoints, and rational functions. See for example DeVore (1989). The full range  $0 < p < 2$  is thus of interest. Minimax risks over weak  $\ell_p$  balls are in this way connected to the study of best attainable spatial adaptation by variable bandwidth estimators (to be described in detail elsewhere.) We concentrate here on finite dimensional analogues in  $\mathbb{R}^n$  (which are compact!), as  $n$  increases and such that the normalised radius  $\eta_n = n^{-1/p} r_n / \varepsilon_n \rightarrow 0$ . Although we allow radius  $r_n$  and noise level  $\varepsilon_n$  to vary freely with  $n$ , the calibration of most interest in function estimation is  $r_n = r, \varepsilon_n = \sigma n^{-1/2}$ , so that  $\eta_n = \sigma n^{1/2-1/p} \rightarrow 0$  when  $p < 2$ .

### 4.1 Strong $\ell_p$

We discuss this only in outline: to emphasise the steps listed in Section 3, and to collect consequences for use later in the paper. Fuller details may be found in [DJ 90]. Let

$$\Theta_{n,p}(r) = \{\theta \in \mathbb{R}^n : \sum_1^n |\theta_i|^p \leq r^p\}.$$

A convex set of measures containing  $\Theta_{n,p}(r)$  is

$$\mathcal{M}_n = \{\pi(d\theta) : E_\pi \sum_1^n |\theta_i|^p \leq r^p\},$$

and a symmetric subset that is as difficult as  $\mathcal{M}_n$  is composed of priors that make the co-ordinates i.i.d. :

$$\begin{aligned} \bar{\mathcal{M}}_n &= \{\pi \in \mathcal{M}_n : \pi = G^n \text{ is i.i.d.}\} \\ &= \{\pi = G^n : E_G |\theta_1|^p \leq n^{-1} r^p\} \end{aligned}$$

For a fixed prior  $\pi$  in  $\mathcal{M}$  with marginals  $\pi_i$ , form the average marginal  $G = n^{-1} \sum_1^n \pi_i$ . The additive structure of the loss function together with concavity of Bayes risk shows that  $\bar{\pi} = G^n$  is harder than  $\pi : B(\bar{\pi}) \geq B(\pi)$ .

The symmetry of  $\bar{\mathcal{M}}$  makes it possible to express  $B^*$  in terms of a univariate minimax problem. Let  $b(F) = \inf_d \{E_F E_\theta (d(x) - \xi)^2\}$  denote the Bayes risk of prior  $F(d\xi)$  for estimating  $\xi$  from  $x \sim N(\xi, 1)$ . Write  $\mathcal{F}_p(\eta)$  for the class of probability measures  $F$  on  $\mathbb{R}$

with absolute  $p$ th moment bounded by  $\eta : E_F|\xi|^p \leq \eta^p$ . Then the corresponding univariate Bayes risk is

$$\rho_p(\eta) = \sup_{F \in \mathcal{F}_p(\eta)} b(F). \quad (9)$$

This can be evaluated explicitly for  $p = 2$ , and for  $p < 2$  we use a numerical approximation in Section 6.

The multivariate minimax Bayes risk is obtained from  $\rho$  by independence and rescaling:

$$B^* = n\varepsilon^2 \rho_p(\eta_n), \quad \eta_n^p = n^{-1}(r/\varepsilon)^p.$$

For asymptotic equivalence in the case  $\eta_n \rightarrow \eta$ , one chooses a prior  $F(d\xi)$  in (9) that is near optimal for  $\eta(1 - \delta)$  and sets  $\pi_n(d\theta) = F^n(\varepsilon^{-1}d\xi)$ . The law of large numbers ensures that  $n^{-1} \sum_1^n |\theta_i/\varepsilon|^p \xrightarrow{P} E_F|\xi|^p \leq \eta^p(1 - \delta)^p$ , so that

$$\pi_n\{\theta \in \Theta_{n,p}(r_n)\} \rightarrow 1.$$

The remaining details of the argument are completed as described in Section 7.

*Properties of  $\rho_p(\eta)$  and approximations.* When  $p < 2$ , as  $\eta \rightarrow 0$ ,

$$\rho_p(\eta) \sim \eta^p (2 \log \eta^{-p})^{1-p/2},$$

and corresponding asymptotically least favorable priors have the symmetric (and sparse!) three-point form (c.f. Bickel, 1983)

$$F = (1 - \epsilon)\delta_0 + \epsilon/2(\delta_\mu + \delta_{-\mu}).$$

Here we assume that a sequence  $a = a(\eta) \rightarrow \infty$  is given and that  $\mu(\eta)$  and  $\epsilon(\eta)$  are then chosen to be solutions of the equations

$$\epsilon\mu^p = \eta^p \quad (10)$$

$$\phi(\mu + a) = \epsilon\phi(a). \quad (11)$$

The corresponding Bayes estimators  $d_F(x) = E_F[\xi|x]$  approximate  $d(x) = \text{sign}(x)\mu I\{|x| > \mu + a\}$  as  $\eta \rightarrow 0$ .

A simple family of non-linear estimators with attractive risk properties is given by (soft) thresholding

$$d_\lambda(x) = d(x; \lambda) = \text{sign } x(|x| - \lambda)_+.$$

The threshold minimax risk

$$\bar{\rho}_p(t) = \inf_\lambda \sup_F \left\{ E_F E_\xi (d_\lambda(x) - \xi)^2 : E_F |\xi|^p \leq t^p \right\}$$

is not much worse than the unrestricted risk:

$$\sup_{0 < t < \infty} \frac{\bar{\rho}_p(t)}{\rho_p(t)} = \Lambda(p) < \infty \quad (12)$$

and, for example,  $\Lambda(1) \leq 1.6$ . In addition if we set  $\lambda(\eta) = \sqrt{2 \log \eta^{-p}}$ , then the soft thresholds  $d_{\lambda(\eta)}$  are asymptotically minimax over  $\mathcal{F}_p(\eta)$  as  $\eta \rightarrow 0$ .

## 4.2 Weak $\ell_p$

Again suppose  $y \sim N(0, 1)$ . Consider the Marcinkiewicz, or weak  $\ell_p$  ball

$$\Theta^* = \Theta_{n,p}^*(r_n) = \{\theta : k^{1/p}|\theta|_{(k)} \leq r_n, \quad k = 1, \dots, n\}.$$

with minimax risk  $R_n^* = R(\Theta_{n,p}^*(r_n), \varepsilon_n)$  given by (5). A weak  $\ell_p$  ball contains the corresponding strong  $\ell_p$  ball, and by contrast it is not convex for *any*  $p < \infty$ . Again, let  $\eta_n = n^{-1/p}(r_n/\varepsilon_n)$ .

Let  $\mathcal{F}_p^*(\eta)$  denote the class of probability measures  $F(d\xi)$  on  $\mathbb{R}$  whose survivor functions  $\tilde{F}(t) = F\{\xi : |\xi| > t\}$  satisfy  $\tilde{F}(\eta t) \leq t^{-p}$  for  $t \geq 0$ . Equivalently,  $|\xi|$  is stochastically smaller than  $\eta X$ , where  $X \sim \text{Pareto}(p)$ . The minimax Bayes risk over  $\mathcal{F}_p^*(\eta)$  is

$$\rho_p^*(\eta) = \sup_{F \in \mathcal{F}_p^*(\eta)} b(F).$$

**Theorem 1** *If either (i)  $p \geq 2$ , or (ii)  $0 < p < 2$ , and*

$$(\varepsilon_n/r_n)^2 \log n (\varepsilon_n/r_n)^p = o((\log n)^{-6/p}) \quad (13)$$

*then*

$$R_n^* \sim n \varepsilon_n^2 \rho_p^*(\eta_n). \quad (14)$$

*Of particular interest is case (ii) with  $\eta_n \rightarrow 0$ . In this case*

$$R_n^* \sim \gamma_p n \varepsilon_n^2 \eta_n^p (2 \log \eta_n^{-p})^{1-p/2}, \quad \gamma_p = 2/(2-p), \quad (15)$$

*and the soft thresholding rules  $\hat{\theta}_n$  with  $\hat{\theta}_{n,i} = \varepsilon d(\cdot; \varepsilon \lambda_n)$  and  $\lambda_n = \sqrt{2 \log \eta_n^{-p}}$  are asymptotically minimax. An asymptotically least favorable sequence of distributions is obtained by setting  $\theta_i = \varepsilon W_i$  where the  $W_i$  are i.i.d. and  $W_1$  is defined in terms of a  $\text{Pareto}(p)$  variable  $X$  by  $W_1 = \min(\eta_n X, \mu_n)$ . Here  $\mu_n$  and implicitly  $\varepsilon_n$  is defined from  $\eta_n$  and  $a(\eta_n) \rightarrow \infty$  via equations (10) and (11).*

**Remarks 1.** The minimax risk for a weak  $\ell_p$  - ball is asymptotically larger than the risk for the corresponding strong  $\ell_p$  - ball of the same radius by the factor  $\gamma_p = 2/(2-p)$ . The asymptotically least favorable priors differ in that the atom of mass  $(1-\varepsilon)$  at 0 for the  $\ell_p$  - ball case is smeared over the interval  $[\eta_n, \mu_n]$  (and its reflection) according to a scaled Pareto distribution which is the extremal member of the family  $\mathcal{F}_p^*(\eta)$ . It would be interesting to explore the extension of these results to the Lorentz spaces  $\ell_{p,q}$  (e.g. Peetre (1976), Bergh and Löfström (1976) ), which are intermediate between  $\ell_p = \ell_{p,p}$  and weak  $\ell_p = \ell_{p,\infty}$ .

2. Somewhat more specific statements than (14) can be made in cases (i) and (ii) according as  $\eta_n$  converges to  $\infty$ , or to  $\eta \in (0, \infty)$  or 0. The results are entirely analogous to Theorem 5 of [DJ 90].

Following the method of Section 3, we introduce a bounding Bayes minimax problem. Let  $t_{kn} = (n/k)^{1/p}$ ,  $k = 1, \dots, n$ , and define the convex set

$$\mathcal{M}_n = \left\{ \pi(d\theta) : E_\pi n^{-1} \sum_{i=1}^n I\{|\theta_i| > \varepsilon_n \eta_n t_{kn}\} < t_{kn}^{-p}, \quad 1 \leq k \leq n \right\}.$$

Since  $\varepsilon_n \eta_n t_{kn} = r_n k^{-1/p}$ , point masses at  $\theta \in \Theta_{n,p}^*(r)$  automatically belong to  $\mathcal{M}_n$  and so  $R_n^* \leq B(\mathcal{M}_n, \varepsilon)$ .

To obtain a related set based on i.i.d. priors we write  $S_\varepsilon F$  for the measure  $F$  scaled by  $\varepsilon$  (that is,  $S_\varepsilon F(A) = F(\varepsilon^{-1}A)$ ) and define

$$\bar{\mathcal{M}}_n = \{\pi = (S_\varepsilon F)^n : F \in \mathcal{F}_{p,n}^*(\eta_n)\}$$

where

$$\mathcal{F}_{p,n}^*(\eta) = \{F : \tilde{F}(\eta t) \leq t^{-p} + n^{-1}, \quad 1 \leq t \leq n^{1/p}\}.$$

Good behavior at  $n$  points qualifies for membership in  $\mathcal{F}_{p,n}^*(\eta)$ :

$$\tilde{F}(\gamma t_{kn}) \leq t_{kn}^{-p}, \quad 1 \leq k \leq n \Rightarrow \tilde{F} \in \mathcal{F}_{p,n}^*(\eta). \quad (16)$$

For a given  $\pi \in \bar{\mathcal{M}}_n$ , the i.i.d. measure  $\bar{\pi} = (av\varepsilon(\pi_i))^n$  is less favorable:  $B(\bar{\pi}) \geq B(\pi)$ . Further,  $\bar{\pi} \in \bar{\mathcal{M}}_n$  from (16) and because membership in  $\mathcal{M}_n$  depends only on the average of the *marginal* distributions of  $\pi$ . Consequently  $B(\mathcal{M}_n, \varepsilon) = B(\bar{\mathcal{M}}_n, \varepsilon)$ .

Using the symmetry of this decision problem and of  $\mathcal{M}_n$ , we obtain

$$B(\bar{\mathcal{M}}_n, \varepsilon) = n\varepsilon^2 \rho_p^*(\eta_n). \quad (17)$$

Now define an (asymptotically least favorable) distribution  $F_\eta \in \mathcal{F}_p^*(\eta)$  as follows. Let  $\varepsilon = \varepsilon(\eta)$  and  $\mu = \mu(\eta)$  be the solutions to (10) and (11) for a sequence  $a = a(\eta) \nearrow \infty$  slowly enough that  $a = o(\mu)$ . Define a probability measure on  $\mathbb{R}_+$  by

$$F_+(d\xi) = p\eta^p \xi^{-1-p} I\{\eta \leq \xi \leq \mu\} d\xi + \eta^p \mu^{-p} \delta_{\{\mu\}}, \quad (18)$$

and let  $F(d\xi) = \frac{1}{2}F_+(d\xi) + \frac{1}{2}F_-(d\xi)$ , where  $F_-$  is the reflection of  $F_+$  about 0. Alternatively, we may say that  $F_+$  is the distribution of  $\eta$ Pareto( $p$ ) with all mass located beyond  $\mu$  (totalling  $\varepsilon$ , by (10)) lumped together at  $\mu$ . It is easily verified that

$$\int_0^\infty x^{2k} F_\eta(dx) = \frac{2k}{2k-p} \eta^p \mu^{2k-p}. \quad (19)$$

**Theorem 2** *Suppose  $0 < p < 2$ . As  $\eta_n \rightarrow 0$*

$$\rho_p^*(\eta_n) = \sup_{\mathcal{F}_p^*(\eta_n)} b(F) \sim \sup_{\mathcal{F}_{p,n}^*(\eta_n)} b(F) \sim \gamma_p \eta_n^p (2 \log \eta_n^{-p})^{1-p/2}. \quad (20)$$

(Proofs are collected in the final subsection.) Having thus established the asymptotics of  $B_n^*$ , there remains the somewhat tedious task of using Step 3 to verify asymptotic equivalence of  $R_n^*$  and  $B_n^*$ . Let us note here only that we may define  $\pi_n$  by fixing  $\delta > 0$ ,  $\bar{\eta}_n = (1 - \delta)\eta_n$  and making  $\theta_i$  i.i.d.  $\varepsilon X_i$  with  $X_i$  drawn from  $F_{\bar{\eta}_n}$ . This yields the necessary asymptotic support property:

**Lemma 1** *If  $\eta_n^p \mu^{-p}(\eta_n) \gg n^{-1} \log^3 n$ , then  $\pi_n \{\theta \in \Theta_{n,p}^*(r_n)\} \rightarrow 1$ .*



## 5 Minimax risks for Besov and Tiebel balls

We begin this section by reviewing some of the results of Donoho and Johnstone (1992a) as an example of a somewhat more elaborate use of the Minimax Bayes strategy. This lays the groundwork for generating the pictures of sample paths from near least favorable priors in the next section. We then give a separate application of the approach for the Triebel scale.

Consider now an *infinite*-dimensional Gaussian sequence estimation problem  $y_I = \theta_I + \varepsilon z_I$ , where the index  $I$  represents a pair  $(j, k)$  with  $j \in \mathcal{N}$  and  $k \in \{0, 1, \dots, 2^j - 1\}$ , and the  $z_I$  are i.i.d.  $N(0, 1)$ . This may be thought of as a wavelet-coefficient form of the standard signal in Gaussian noise regression model

$$Y(t) = \int_0^t f(s) ds + \varepsilon W(t) \quad 0 \leq t \leq 1, \quad (21)$$

where  $W(t)$  is a standard Brownian motion. Thus, for example,  $\theta_{jk} = \int_0^1 f \psi_{jk}$  measures the frequency content of  $f$  at frequencies near  $2^j$  at locations near the subinterval  $I_{jk} = [2^{-j}k, 2^{-j}(k+1)]$ .

The Besov and Triebel-Lizorkin scales of function spaces on  $[0, 1]$  can be defined through convergence conditions on the components of  $\theta$ . Define normalized indicator functions  $\chi_{jk}(s) = 2^{j/2} I_{jk}(s) = 2^{j/2} I\{s \in I_{jk}\}$  and

$$f(t, j) = \sum_k 2^{j\sigma} \theta_{jk} \chi_{jk}(t).$$

Besov spaces require the  $L_p(dt)$  norms of  $f(\cdot, j)$  to belong to  $\ell_q$ ; equivalently

$$\|\theta\|_{b_{p,q}^\sigma}^q = \sum_j 2^{sjq} \left( \sum_k |\theta_{jk}|^p \right)^{q/p}, \quad s = \sigma + 1/2 - 1/p.$$

The Triebel spaces on the other hand ask for the  $\ell_q$  norms of  $f(t, \cdot)$  to belong to  $L_p(dt)$ ; namely

$$\|\theta\|_{f_{p,q}^\sigma}^p = \int_0^1 \left( \sum_j 2^{sjq} |\theta_{jk}|^q I_{jk} \right)^{p/q}$$

for further details, see for example, Frazier, Jawerth and Weiss (1991). These scales coincide if  $p = q < \infty$  and contain the Sobolev ( $f_{p,2}^\sigma$ ) and Hölder ( $b_{\infty,\infty}^\sigma$ ) norms as well as others, such as the Bump Algebra ( $b_{1,1}^1$ ), and (by bracketing between  $b_{1,1}^1$  and  $b_{1,\infty}^1$ ) Total Variation, which capture other forms of prior information of scientific relevance.

### 5.1 Besov Balls

The minimax Bayes approach enables exact evaluations of asymptotic minimax risk  $R(\Theta, \varepsilon) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} E \|\hat{\theta} - \theta\|^2$  as  $\varepsilon \rightarrow 0$  over Besov balls

$$\Theta = \Theta_{p,q}^\sigma(C) = \{\theta : \|\theta\|_{b_{p,q}^\sigma} \leq C\}$$

for  $q \geq p > 0$  and  $\sigma > (1/p - 1/2)_+$ . In particular, we recover (when  $p = q = 2$ ) the homoscedastic case of Pinsker's theorem presented in Section 2.

For the minimax Bayes risk  $B(\mathcal{M}, \varepsilon)$  we take

$$\mathcal{M} = \mathcal{M}_B(\Theta_{p,q}^\sigma(C)) = \{\pi(d\theta) : \sum_j 2^{sjq} (\sum_k E_\pi |\theta_{jk}|^p)^{q/p} \leq C^q\} \quad (22)$$

which is convex when  $q \geq p$ .

Define the simpler subclass of measures  $\bar{\mathcal{M}}$  as those distributions in  $\mathcal{M}$  that make  $\theta = (\theta_{jk})$  independent *between* resolution levels, and furthermore, i.i.d. *within* levels. For any given  $\pi \in \mathcal{M}$ , a less favorable  $\bar{\pi}$  (in the series that  $B(\bar{\pi}) \geq B(\pi)$ ) is constructed by making  $\theta_{jk}$  i.i.d. within levels, distributed as  $\bar{\pi}_j = \text{ave}_k(\pi_{jk})$ , where  $\pi_{jk}$  denotes the  $(jk)^{\text{th}}$  univariate marginal of  $\pi$ . Since membership in  $\mathcal{M}$  is determined only by average properties of the univariate marginals of  $\pi$ , it follows that  $\bar{\pi} \in \mathcal{M}$  also, and hence  $B(\mathcal{M}, \varepsilon) = B(\bar{\mathcal{M}}, \varepsilon)$ .

Priors in  $\bar{\mathcal{M}}$  have Bayes rules that operate co-ordinatewise: they depend only on the level of the component:  $d_{\pi,I}(y) = d_I(y_I)$  when  $I \in \mathcal{I}_j$ . Here  $d_I$  is the univariate Bayes rule corresponding to the  $I^{\text{th}}$  marginal of  $\pi$ . Thus the overall risk  $E\|\hat{\theta}_\pi - \theta\|^2$  becomes a sum of univariate Bayes risks for priors  $\pi$ ; with multiplicity  $2^j$  at the  $j^{\text{th}}$  level. The moment condition for  $\bar{\mathcal{M}}$  becomes a Besov norm condition on the  $p^{\text{th}}$  moments of the univariate  $\pi_j$  and a renormalization argument leads to

$$B(\bar{\mathcal{M}}, \varepsilon) \sim \varepsilon^{2r} C^{2(1-r)} \gamma(C\varepsilon^{-1}) \quad \text{as } \varepsilon \rightarrow 0, \quad r = 2\sigma/(2\sigma + 1). \quad (23)$$

where  $\gamma = \gamma(\cdot, \sigma, p, q)$  is a positive continuous function of period 1, and

$$\gamma(b) = \sup \left\{ \sum_{j=-\infty}^{\infty} 2^j \rho_p(t_j) : \sum_{j=-\infty}^{\infty} 2^{j\beta q} t_j^q \leq b^q \right\}. \quad (24)$$

Here  $\beta = \sigma + 1/2$ , and  $\rho_p(t)$  is the univariate Bayes minimax risk defined at (9): namely estimate  $\xi$  from  $x \sim N(\xi, 1)$ , when  $\xi$  is distributed as  $F$  constrained only by the moment condition  $E_F |\xi|^p \leq t^p$ .

Asymptotic equivalence of  $R_n^*$  and  $B_n^*$  is verified from (40) and (41) after using (23) to establish (38). We refer to Donoho and Johnstone (1992a) for details.

## 5.2 Local estimators and the Triebel scale

Exact asymptotic minimax evaluations are not known for the Triebel scale or for the Besov scale with  $p < q$ . We outline here an exact asymptotic minimax result for a restricted class of estimators in the Triebel scale:

$$\Theta = \Phi_{p,q}^\sigma(C) = \{\theta : \|\theta\|_{f_{p,q}^\sigma} \leq C\}$$

Further details are given in Section 7.3.

We consider a class of *local* estimators defined in terms of the natural ordering on indices  $I$ : if  $I = (jk)$ ,  $I' = (j'k')$  and  $j > j'$ , then say that  $I > I'$  if the interval  $I_{jk} \subset I_{j'k'}$ . Let  $A(I) = \{I' : I' < I\}$  denote the ‘ancestors’ of  $I$ , and let  $F(I) = \{I\} \cup A(I)$ . The notation  $y_F$  stands for  $(y_I : I \in F)$ . Then the class of local estimators will be

$$\hat{\Theta}_L = \{\hat{\theta}(y) : \hat{\theta}_I(y) = \hat{\theta}_I(y_{F(I)}) \forall I, y\}$$

Thus, the  $I^{th}$  component of  $\hat{\theta}$  depends only on observations  $y_{I'}$  associated with intervals  $I'$  containing  $I$ . This allows for estimators that ‘borrow strength’ from information in data at coarser resolution levels, but in a manner localized through the containment constraint.

In general the minimax risk

$$R_L^* = \inf_{\hat{\theta}_L} \sup_{\theta} E \|\hat{\theta} - \theta\|^2 \geq R^*, \quad (25)$$

although we have seen that in the Besov case with  $p \leq q$  they are asymptotically equivalent.

For a minimax Bayes problem, we take a slightly different choice than (22); namely

$$\mathcal{M} = \mathcal{M}_T(\Phi_{p,q}^\sigma(C)) = \{\pi : E_\pi \|\theta\|_{f_{p,q}^\sigma}^p \leq C^p\}.$$

This is trivially convex, but no longer depends only on the univariate marginals of  $\pi$ ; but rather on the joint distribution of wavelet coefficients at all levels which are associated with intervals  $I_{jk}$  containing  $t$ .

Let  $\kappa(d\xi)$  be a probability measure on sequences  $(\xi_0, \xi_1, \dots)$ . The notation  $\xi_{(j)}$  denotes initial segments  $(\xi_0, \xi_1, \dots, \xi_j)$ . Call a probability measure  $\pi(d\theta)$  *subordinated to  $\kappa$*  (written  $\pi = S\kappa$ ) if  $\pi$  may be constructed from  $\kappa$  via the following recipe. Let  $\mathcal{F}_{j-1} = \{\theta_{j'k'}, k' \in \mathcal{I}_{j'}, j' < j\}$ , and require that

- (i)  $(\theta_I, I \in \mathcal{I}_j)$  are conditionally independent given  $\mathcal{F}_{j-1}$ ,
- (ii)  $\mathcal{L}_\pi(\theta_{jk} | \mathcal{F}_{j-1}) = \mathcal{L}_\pi(\theta_{jk} | \theta_{A(jk)})$ ,
- (iii)  $\mathcal{L}_\pi(\theta_{jk} | \theta_{A(jk)} = t_{(j-1)}) = \mathcal{L}_\kappa(\xi_j | \xi_{(j-1)} = t_{(j-1)})$ .

We define

$$\bar{\mathcal{M}} = \{\pi \in \mathcal{M} : \pi \text{ is subordinated to some } \kappa\}.$$

A key property of subordinated priors  $\pi$  is that the sequence  $(\theta_{j,[2^j r]})_{j=0}^\infty$  is distributed as  $\kappa$  for every  $r \in [0, 1]$ . Conversely, given a prior  $\pi \in \mathcal{M}$ , we construct  $\kappa$  and then a subordinated  $\bar{\pi} \in \bar{\mathcal{M}}$  by introducing  $R \sim U(0, 1)$  independent of  $\theta$  and setting

$$\xi_j(\theta, r) = \theta_{j,[2^j r]}.$$

This produces  $\kappa$  and  $\bar{\pi}$  is obtained from recipe (i) – (iii) above.

Consider a (*doubly* infinite) auxiliary estimation problem with data  $v_i = \xi_i + \varepsilon z_i$ ,  $i \in \mathcal{Z}$ , and it is desired to estimate  $\xi_j$  using estimators  $\hat{\xi}_j \in \hat{\Xi}_j$  having the property that  $\hat{\xi}_j(y)$  depends only on  $y_{(j)} = (\dots, y_{j-1}, y_j)$ . Denote the corresponding Bayes risk

$$b_j(\kappa, \varepsilon) = \inf_{\hat{\xi}_j \in \hat{\Xi}_j} E_\kappa E_\theta [\hat{\xi}_j - \theta_j]^2.$$

Here of course  $\kappa(d\xi)$  is a probability distribution on bilateral sequences  $(\xi_j, j \in \mathcal{Z})$ .

One shows that  $B_L(\mathcal{M}, \varepsilon) = B_L(\bar{\mathcal{M}}, \varepsilon)$  by using a version of this estimation problem using data only for  $i \in \mathcal{N}$ . The restriction to local estimators in (25) is critical here, as it forces the Bayes estimator for  $\bar{\pi}$  to be local.

Proceeding now in analogy with the Besov case, a renormalization argument shows that

$$B_L(\bar{\mathcal{M}}, \varepsilon) \sim \varepsilon^{2r} C^{2(1-r)} \gamma_L(C\varepsilon^{-1}) \quad \text{as } \varepsilon \rightarrow 0 \quad (26)$$

where

$$\gamma_L(b) = \sup_{-\infty}^{\infty} \left\{ \sum_{-\infty}^{\infty} 2^j b_j(\kappa, \varepsilon) : \kappa \in \mathcal{K}(C) \right\} \quad (27)$$

$$\mathcal{K}(C) = \left\{ \kappa(d\xi) : E_{\kappa} \left( \sum_j 2^{sjq} |\xi_j|^q \right)^{p/q} \leq C^p \right\}. \quad (28)$$

We describe the main idea in the asymptotic equivalence argument that  $R_L^*(\varepsilon) \sim B_L(\mathcal{M}, \varepsilon)$  as  $\varepsilon \rightarrow 0$ . Construct a near optimal solution  $\kappa$  to (27) with  $\xi_j \equiv 0$  *w.p.1* unless  $|j| \leq J$ , for  $J$  sufficiently large. Restrict  $\varepsilon$  to values  $\varepsilon_h = 2^{-sh}$  for  $h$  sufficiently large. For  $h > J$ , construct  $2^{h-J}$  independent realizations  $\xi^l$  from the dyadic prior subordinated to  $\kappa$  (shifted so as to start at 0 rather than  $-J$ ). ‘Attach’ each realization to the subtree of indices rooted at  $(jk) = (h - J, l)$ , for  $l = 0, 1, \dots, 2^{h-J} - 1$ , and multiply by  $\varepsilon_h$  to give a realization from  $\pi_h$ . Under this distribution  $\|\theta\|_{f_{p,q}^{\sigma}}^p = \text{Ave}(V_l, l = 0, 1, \dots, 2^{h-J} - 1)$  where  $(V_l)$  are i.i.d. with  $EV_1 < C^p$ , from the construction of  $\kappa$ . Thus condition (40) will follow from the law of large numbers.

Thus  $R_L^*(\varepsilon)$  also satisfies the exact asymptotic relation (26). In general, approximate numerical evaluation of  $\gamma_L(b)$  will be more difficult than for  $\gamma(b)$ , because of the multivariate dependencies involved. Instead, we return in the next section to the Besov case, and numerical approximation of  $\gamma(b)$ .

## 6 Sample Paths from Priors

As mentioned in the previous section, the dyadic sequence space model is a representation in suitable wavelet bases of the classical signal in noise regression model (21). The minimax Bayes approach together with renormalization shows that an asymptotically least favorable prior distribution may be built from the solution to the optimization problem (24). In this section we attempt this numerically, in order to exhibit the variety in appearance of sample paths from these distributions as the parameters of the Besov space vary. In particular the *separate* variation in sparsity and smoothness supports the argument that these spaces capture genuinely different and scientifically relevant forms of prior information about the unknown function.

The basic building block is the univariate problem of estimating  $\xi$  from  $x \sim N(\xi, 1)$  with  $\xi$  distributed as  $F$  constrained only by the moment condition  $E_F|\xi|^p \leq t^p$ , as studied in Section 4. Denote by  $F^{p,t}(d\xi)$  a least favorable prior for this setting. Then a least favorable prior for

$$B^* = B^*(\Theta_{p,q}^{\sigma}, \varepsilon) = \inf_{\hat{\theta}} \sup_{\mu \in \mathcal{M}} E_{\mu} E_{\theta} \|\hat{\theta} - \theta\|^2 \quad (29)$$

(when  $\varepsilon^2$  is restricted to a subsequence  $\varepsilon_h^2 = 2^{-\beta h}$ ) may be constructed as

$$\theta_{\bar{j}k} = \varepsilon_h X_{j,k} \quad j = \bar{j} - h, \quad (30)$$

where  $(X_j) = (X_{j,1}, \dots, X_{j,2^j})$  is a vector distributed i.i.d. as  $F^{p,t_j}$ , the vectors  $\{(X_j), j = -h, -h+1, \dots\}$  are independent and the vector  $(t_{-h}, \dots, t_j, \dots)$  is a solution to optimization problem (24).

Sample path realizations from this least favorable prior may be constructed by treating  $(\theta_I)$  as wavelet coefficients of a random (periodic) function on  $[0, 1]$ :

$$X(t) = \sum_{\bar{j}=0}^{\infty} \sum_{k=1}^{2^{\bar{j}}} \theta_{\bar{j}k} \psi_{\bar{j}k}(t).$$

Here  $\psi_{\bar{j}k}(t) = 2^{\bar{j}/2} \psi(2^{\bar{j}}t - k)$  is a suitable orthonormal wavelet basis of regularity at least  $\sigma$ . In fact we will confine attention to *periodic* functions on  $[0, 1]$  and work in effect with *periodic* wavelets (cf. for example, Daubechies (1992, Sec. 9.3)). This computational simplification affects only a fixed number of wavelet coefficients at each resolution level and does not alter the qualitative phenomena we wish to present.

We carry out approximate constructions of sample paths for a small selection of values of  $p = q$  and  $\sigma$ . For  $p \neq 2$ , neither the minimax Bayes risk  $\rho_p(t)$  nor the least favorable distribution  $F^{p,t}$  is available explicitly, and so we consider a modified risk  $\rho_{\lambda,p}(t)$  for soft threshold rules for which a simpler, though still numerical, evaluation is possible.

We begin, however, with Pinsker's case,  $p = q = 2$ , in which  $\rho_2(t) = t^2/(1+t^2)$  and the least favorable distribution  $F^{2,t} = N(0, t^2)$  is Gaussian. In this case, we need only solve for the least favorable sequence  $(t_j)$  in

$$\sup \left\{ \sum_{j=-\infty}^{\infty} 2^j t_j^2 / (1 + t_j^2) : \sum_{j=-\infty}^{\infty} 2^{(2\sigma+1)j} t_j^2 = 1 \right\} \quad (31)$$

where we have made the approximation of replacing  $h$  by  $-\infty$  in the summations (c.f. [DJ 1992a]). As in Pinsker (1980), the solution has the form

$$t_j^2 = (\beta_0 2^{\sigma(j_0-j)} - 1)_+,$$

where  $\beta_0(\sigma)$  and integer  $j_0(\sigma)$  are determined to satisfy the equality constraint in (31). In fact, set  $a = 2^\sigma - 1$ ,  $b = 1 - 2^{-\sigma-1}$  and  $c = 1 - 2^{-2\sigma-1}$ ; then

$$j_0 = \lceil \frac{1}{2\sigma+1} \log_2 \frac{bc}{a} \rceil, \quad \beta_0 = b(2^{-j_0(2\sigma+1)} + c^{-1}).$$

For the cases  $\sigma = 1$  and  $\sigma = 2$ , Table 1 shows the values of  $t_j^2$  and the components of (31). In particular, there is no contribution for  $j > 0$ .

We have chosen the nominal noise level  $\epsilon_h = 2^{-\beta h/2}$  to be similar in the two cases. For  $\sigma = 1$ , we set  $h = 8$  and  $\epsilon_h = 2^{-6}$  and for  $\sigma = 2$ , we take  $h = 5$  and so  $h\beta = 12.5$  and  $\epsilon_h = 2^{-6.25}$ . Figure 1 shows sample paths from these least favorable distributions, namely

$$X(i/N) = \sum_{\bar{j}=0}^{11} \sum_{k=1}^{2^{\bar{j}}} \theta_{\bar{j}k} W_{\bar{j}k}(i/N) = (W_N \theta)_i \quad i = 1, \dots, N \quad (32)$$

where  $\tilde{\theta}_{\bar{j}k} = X_{\bar{j}-6,k}$ ,  $k = 1, \dots, 2^{\bar{j}}$  are i.i.d.  $\sim N(0, t_{\bar{j}-6}^2)$ .

Here  $W_N$  denotes a periodized form of the discrete wavelet transform of length  $N = 2^{12} = 4096$ . This is derived from a periodized form of the usual cascade algorithm and a pair of finite filter sequences, as described, for example, in Daubechies (1992, Ch 5). We

used the  $N = 8$  instance of the “closest to linear phase” Daubechies wavelet (coefficients listed in Table 6.3 in Daubechies, 1992), and the specific MATLAB code is documented in [DJT, 1992c].

For  $p < 2$ , as previously noted, we consider an approximating minimax Bayes problem obtained by restricting estimators to soft threshold rules  $d_\lambda(x) = \text{sign } x(|x| - \lambda)_+$ . As noted at (12), the threshold minimax risk is not much worse than the unrestricted risk. The corresponding bound holds for Besov balls also (DJ 92):

$$B_\lambda^*(\varepsilon, \Theta) = \inf_{(\lambda_I)} \sup_{\mu \in \mathcal{M}(\Theta)} E_\mu E_\theta \|\hat{\theta}^\lambda - \theta\|^2 \leq \Lambda(p) B^*(\varepsilon, \Theta).$$

We shall therefore use the least favorable distributions for  $B_\lambda^*(\varepsilon, \Theta)$  as an approximation to those for  $B^*(\varepsilon, \Theta)$ , reassured by the bound (12) and the fact that the ratio  $\rho_p(t)/\bar{\rho}_p(t)$  approaches 1 as  $t$  approaches 0 or  $\infty$ . Thus, for extreme values of  $t$ , the least favorable priors for thresholding are in fact nearly unrestricted least favorable. The renormalization argument of [DJ 1992a] shows that when  $\varepsilon_h^2 = 2^{-\beta h}$ ,

$$B_\lambda^*(\varepsilon_h, \Theta_{pq}^\sigma) = (\varepsilon^2)^{2\sigma/2\sigma+1} \sup \left\{ \sum_{j=-h}^{\infty} 2^j \bar{\rho}_p(t_j) : \sum_{j=-h}^{\infty} 2^{j\beta q} t_j^q \leq c^q \right\} \quad (33)$$

So  $B_\lambda^*(\varepsilon_h, \Theta_{pq}^\sigma)$  is given by (24) with  $\rho_p(t_j)$  replaced by  $\bar{\rho}_p(t_j)$ . An evaluation of  $\bar{\rho}_p(t_j)$  is obtained by writing

$$\bar{\rho}_p(t) = \inf_\lambda \rho_p(\lambda, t) \quad (34)$$

$$\rho_p(\lambda, t) = \sup_F \left\{ \int r(\lambda, \xi) dF(\xi) : F \in \mathcal{F}_p(\eta) \right\} \quad (35)$$

Here  $r(\lambda, \xi) = E_\xi(d_\lambda(x) - \xi)^2$  is the univariate threshold risk function. It is shown in [DJ 90] that

$$\rho_p(\lambda, t) = \begin{cases} (1 - \varepsilon_0)r(\lambda, 0) + \varepsilon_0 r(\lambda, \xi_0) & \xi_0 \geq t \\ r(\lambda, t) & \xi_0 < t \end{cases} \quad (36)$$

where  $\xi_0 = \xi_0(\lambda, p)$  maximizes  $\xi \rightarrow (1 - t^p \xi^{-p})r(\lambda, 0) + t^p \xi^{-p} r(\lambda, \xi)$  and  $\varepsilon_0 = \varepsilon_0(\lambda, p, t) = t^p \xi_0^{-p}$ . The corresponding least favorable priors for (35) are given by

$$F_{\varepsilon_0, \xi_0} = \begin{cases} (1 - \varepsilon_0) + \varepsilon_0/2(\delta_{\xi_0} + \delta_{-\xi_0}) & \xi_0 \geq t \\ \frac{1}{2}(\delta_t + \delta_{-t}) & \xi_0 \leq t \end{cases}.$$

The function  $\bar{\rho}_p(t)$  is therefore evaluated by minimizing (36), and from the minimizing  $\lambda = \lambda(t)$  the least favorable prior is the two or three point distribution corresponding to  $\varepsilon_0(\lambda(t), p, t)$  and  $\xi_0(\lambda(t), p)$ . Numerical values for  $\rho_p(t)$  for  $p = 1$  and  $.5$  and  $t$  in the logarithmically spaced grid  $\log_{10} t = -5(.1)1$  are given in Table 1.

This table was used in conjunction with a constrained optimization program (`constr` in MATLAB) to obtain solutions to (33) in which the index  $j$  in both objective and constraint sums was restricted to a range  $[j_1, j_2]$ .

The function  $\bar{\rho}_p(t)$  is strictly concave in  $\tau = t^p$ , and when  $p = q$ , the constraint is linear in  $\tau$ , so (33) and its finite sum approximants have unique optima. For  $p = q = 1, \sigma = 1$  and

$p = q = .5, \sigma = 2$ , optima were computed numerically by decreasing  $j_1$  and increasing  $j_2$  until either the contribution or total risk became small or the probability  $\varepsilon_j$  of the non-zero atoms  $\pm\xi_j$  became small or both. The adequacy of the approximation is affected in part by the range of  $t$  values used in the table look up, and perhaps also by the number of variables in the constrained optimisation. The results shown in Figure 4 are reasonably satisfactory for  $p = q = 1$ , while for  $p = q = 1/2$  extra levels with  $j > 4$  appear to be needed to give the total minimax risk  $\Sigma 2^j \rho(t_j)$  [For this an extension of the table of  $\bar{\rho}_p(t)$  to  $t < 10^{-5}$  would be required]. Note also that to adequately represent a sample path including coefficients drawn from the least favorable distribution for these levels we would have to use more than the  $N = 4096$  points employed in our pictures, shown in Figure 2.

Comparison of the  $\sigma = 1$  plots (corresponding to an MSE of order  $(\epsilon^2)^{2/3}$ ) shows that the case  $p = q = 1$  already exhibits a degree of sparsity in the near least favorable path. The contrast is clearer in the examples chosen for smoothness  $\sigma = 2$  (corresponding to MSE of order  $(\epsilon^2)^{4/5}$ ): here the case  $p = q = .5$  shows considerable sparsity by comparison to the ellipsoid case.

Also shown in Figure 3 are plots of the wavelet coefficients associated with each of the sample paths. The norm constraint is expressed very differently according as  $p = 2$  or  $p < 2$ . For the ellipsoid case, at higher levels, essentially all wavelet coefficients are present (with Gaussian distribution of size), but with decreasing magnitude. For  $p < 2$ , however, the coefficients are increasingly rare (relative to the  $2^{\bar{j}}$  possible) as the levels  $\bar{j}$  increase, but the size of the non-zero coefficients does not change so rapidly between levels.

In summary, this small subset of cases already illustrates the variety of forms of prior information captured by spaces in the Besov scale. This variety adds a degree of suppleness to minimax analysis of estimators and makes the Besov and Triebel scales attractive for studying adaptivity properties of estimators.

## 7 Proofs and Details

### 7.1 Asymptotic equivalence of $R_n^*$ and $B_n^*$ .

We suppose that there is in fact a *family* of minimax problems linked by a scale parameter  $C$ : let  $R(\varepsilon_n, C)$  denote the (frequentist) minimax risk over  $C\Theta_n$ , and  $B(\varepsilon_n, C)$  the corresponding minimax Bayes risk.

It is easily shown that the following two conditions entail the equivalence  $R(\varepsilon_n, C) \sim B(\varepsilon_n, C)$ :

$$\forall \gamma < 1, R(\varepsilon_n, C) \geq \gamma B(\varepsilon_n, \gamma C)(1 + o(1)), \quad (37)$$

and

$$\liminf_{\varepsilon_n \rightarrow 0} \frac{B(\varepsilon_n, \gamma C)}{B(\varepsilon_n, C)} \rightarrow 1 \text{ as } \gamma \uparrow 1. \quad (38)$$

The second condition is often easily checked from the asymptotic evaluation of  $B(\varepsilon_n, C)$ . Some general comments can be made on the first condition which are often useful:

Let  $\pi_n$  be any prior distribution with  $\pi_n(C\Theta_n) > 0$  and set  $\nu_n = \pi_n(\cdot | C\Theta_n)$ , and let  $\hat{\theta}_{\nu_n}$  be the Bayes estimator of  $\theta$  for the conditioned prior  $\nu_n$ . Then from the definitions

$$B(\pi_n) \leq E_{\pi_n} \left\{ \|\hat{\theta}_{\nu_n} - \theta\|^2 | C\Theta_n \right\} \pi_n(C\Theta_n) + E_{\pi_n} \left\{ \|\hat{\theta}_{\nu_n} - \theta\|^2, C\Theta_n^c \right\}$$

$$\leq B(\nu_n)\pi_n(C\Theta_n) + 2E_{\pi_n} \left\{ \|\theta_{\nu_n}\|^2 + \|\theta\|^2, C\Theta_N^\varepsilon \right\}.$$

Denote by  $\Delta_n$  the second term in the final bound above. For fixed  $\gamma \in (0, 1)$ , choose  $\pi_n \in \bar{\mathcal{M}}_n$  so that

$$B(\pi_n) \geq \gamma B(\varepsilon_n, \gamma C) \quad (39)$$

From (37) and since  $\nu_n$  is supported on  $C\Theta_n$ ,

$$\gamma B(\varepsilon_n, \gamma C) \leq R(\varepsilon_n, \gamma C)\pi_n(C\Theta_n) + \Delta_n.$$

In summary, condition (37) (and hence asymptotic equivalence) will follow if priors  $\pi_n \in \bar{\mathcal{M}}_n$  can be chosen satisfying (39) and

$$\pi_n(C\Theta_n) \rightarrow 1, \quad (40)$$

$$\Delta_n = o(B(\varepsilon_n, C)) \text{ as } n \rightarrow \infty. \quad (41)$$

## 7.2 Proofs for weak $\ell_p$

We first note two immediate consequences of the equations (10) and (11) defining  $\epsilon(\eta)$  and  $\mu(\eta)$ :

$$\epsilon\mu^2 = \eta^p \mu^{2-p}, \quad \mu^2 \sim 2 \log \eta^{-p}. \quad (42)$$

### 7.2.1 Proof of Theorem 2

*Upper bound.* Denote the mean squared error of a univariate soft threshold estimate  $d_\lambda$  with risk function  $r(\lambda, \xi) = E(d_\lambda(X) - \xi)^2$ . Some simple properties of  $\xi \rightarrow r(\lambda, \xi)$  will be used here without proof – extra details are in [DJ, 1992b].

The minimax risk over  $\mathcal{F}_{p,n}^*(\eta)$  and  $\mathcal{F}_p^*(\eta)$  is bounded above by

$$\sup \left\{ \int r(\lambda, \xi) F(d\xi) : F \in \mathcal{F}_{p,n}^*(\eta) \right\}. \quad (43)$$

Since  $\xi \rightarrow r(\lambda, \xi)$  is symmetric about 0 and monotone increasing in  $\xi$  to a limit  $1 + \lambda^2$  at  $\xi = \infty$ , the least favorable prior in the optimisation (43) is bounded above by a measure with density equal to that of  $\eta$ .Pareto( $p$ ) for  $\eta \leq \xi \leq \eta n^{1/p}$  and with point mass  $2n^{-1}$  at  $\xi = \infty$ . Thus

$$\rho^* \leq \int_\eta^\infty r(\lambda, \xi) dF_{\eta+}(\xi) + 2(1 + \lambda^2)n^{-1}. \quad (44)$$

Let us set  $\lambda = \lambda_\eta = \sqrt{2 \log \eta^{-p}}$ . From (42) we have  $\lambda \sim \mu$  as  $\eta \rightarrow 0$ , and from our assumption that  $\eta_n^p / \mu_n^p \gg n^{-1} \log^3 n$ , we conclude that

$$\lambda^2/n = o(\eta^p \lambda^{2-p}).$$

Apply integration by parts in (44) and note that  $(\partial/\partial\xi)r(\lambda, \xi) = 2\xi[\Phi(\lambda - \xi) - \Phi(-\lambda - \xi)]$ . Thus

$$B^* \leq r(\lambda, \xi) + 2\eta^p I(\lambda, \eta) + o(\eta^p \lambda^{2-p}), \quad \text{and} \quad (45)$$

$$I(\lambda, \eta) = \int_\eta^\infty [\Phi(\lambda - \xi) - \Phi(-\lambda - \xi)] \xi^{1-p} d\xi. \quad (46)$$



Set  $\lambda_{\pm} = \lambda \pm \sqrt{2 \log \lambda}$ : asymptotics of  $I(\lambda(\eta), \eta)$  follow easily from

$$(1 - \lambda^{-1}) \int_{\eta}^{\lambda^{-}} \xi^{1-p} d\xi \leq I(\lambda, \eta) \leq \int_{\eta}^{\lambda^{+}} \xi^{1-p} d\xi + o(\lambda^{2-p}). \quad (47)$$

Study of the risk function  $\xi \rightarrow r(\lambda, \xi)$  shows that

$$r(\lambda, \eta) \leq r(\lambda, 0) + \eta^2 \leq c\eta^p \lambda^{-p} \phi(0) + \eta^2 = o(\eta^p). \quad (48)$$

Finally, inserting (47) and (48) into (45) yields  $B^* \leq \gamma_p \eta^p \lambda_{\eta}^p (1 + o(1))$  as required.

*Remark.* The upper bound may also be derived using the ‘oracle inequality’ of [DJ 1992d].

*Lower Bound.* We study  $F_{\eta}$  and the corresponding Bayes estimator  $d_{\eta}$  and proceed to bound the Bayes risk

$$b(F_{\eta}) = \int_{\eta}^{\mu} r(d_{\eta}, \xi) p\eta^p \xi^{-1-p} d\xi + \epsilon r(d_{\eta}, \mu). \quad (49)$$

Let us assume for now the following important and rather remarkable property of the Winsorized Pareto prior  $F_{\eta}$ :

$$d_{\eta}(\mu + a/2) \leq 1, \quad \eta < \eta_0. \quad (50)$$

The identity  $d'_{\eta}(x) = \text{Var}_F(\xi|x) \geq 0$  along with  $d_{\eta}(0) = 0$  guarantees that  $0 \leq d_{\eta}(x) \leq 1$  on the event  $A_{\eta} = \{x : 0 \leq x \leq \mu + a/2\}$ . This implies  $0 \leq E^A d_{\eta} \leq 1$ , where  $E^A$  denotes expectation conditional on  $A_{\eta}$ . Bounding mean squared error by squared bias yields

$$r(d_{\eta}, \xi) \geq P_{\xi}(A_{\eta}) E^A (d_{\eta} - \xi)^2 \geq s_{\eta} (\xi - 1)^2 \quad (51)$$

where  $s_{\eta} = \inf \{P_{\xi}(A_{\eta}) : a/2 \leq \xi \leq \mu\}$  increases to 1 as  $\eta$  decreases. Now substitute (51) into (49), recall that  $a = o(\mu)$ , and conclude that

$$b(F_{\eta}) \geq s_{\eta} p / (2 - p) [\eta^p \mu^{2-p} + o(\eta^p \mu^{2-p})] + s_{\eta} \epsilon \mu^2 + o(\epsilon \mu^2).$$

Appeal now to (42) to derive

$$\begin{aligned} b(F_{\eta}) &\geq \gamma_p \eta^p \mu^{2-p} (1 + o(1)) \\ &\sim \gamma_p \eta^p (2 \log \eta^{-p})^{2-p} (1 + o(1)). \end{aligned}$$

Derivation of (50). For notational convenience, set  $x = x_{\eta} = \mu + a/2$ . Expressing  $d_{\eta}$  explicitly as a posterior mean and neglecting terms gives the bound

$$d_{\eta}(x) \leq \frac{\phi(x - \mu) \mu^{1-p} + \int_{\eta}^{\mu} \phi(x - \xi) \xi^{-p} d\xi}{p \int_{\eta}^{\mu} \phi(x - \xi) \xi^{-p-1} d\xi}. \quad (52)$$

Define

$$I_r(\eta) = \int_{\eta}^{\mu} \phi(x - \xi) \xi^{-r} d\xi = \int_{\eta}^{\mu} e^{h(\xi)} d\xi / \sqrt{2\pi},$$

where  $h(\xi) = -1/2(x - \xi)^2 - r \log \xi$ . Over the interval  $[\eta, \mu(\eta)]$  it attains its only maxima at the endpoints and there is a single minimum at the smaller root  $\mu_- = \mu_-(\eta) \sim r/x_\eta$  of the quadratic equation obtained from  $h'(\xi) = 0$ . Decompose  $[\eta, \mu]$  into low and high subintervals  $[\eta, \mu_-]$  and  $[\mu_-, \mu]$ , and set  $I_r = I_{L,r} + I_{H,r}$ . The asymptotic behaviour of these integrals is given as  $\eta \rightarrow 0$  by

$$I_{L,r}(\eta) \sim \phi(x_\eta) \int_{\mu_-}^{\mu} \eta \xi^{-r} d\xi, \quad (53)$$

$$I_{H,r}(\eta) \sim \tilde{\Phi}(a/2) \mu_\eta^{-r}. \quad (54)$$

To verify, for example, (54), regard the ratio of the left side over the right as the expectation of  $\Lambda_\eta = e^{x\xi - \xi^2/2}$  for  $\xi$  distributed as a random variable with density proportional to  $\xi^{-r}$  on  $[\eta, \mu_-]$ . Then  $\Lambda_\eta$  is uniformly bounded by the constant  $e^r$  on  $[\eta, \mu_-]$  and converges to one in probability as  $\eta \rightarrow 0$ .

We return to bounding the behavior of  $d_\eta(x)$ . The first ratio in (52) is asymptotically bounded as  $\eta \rightarrow 0$  via (54) by

$$\frac{\phi(x_\eta - \mu) \mu^{1-p}}{\phi(x_\eta) \eta^{-p}} = \mu e^{-\mu a/2} \rightarrow 0.$$

[The equality is obtained after eliminating  $\epsilon$  from equations (10) and (11).] This same argument shows that the component  $I_{L,p+1}$  dominates in the denominator integral in (52), and also dominates  $I_{H,p}$ . Thus the second ratio in (52) is asymptotically bounded by

$$\frac{I_{L,p}}{p I_{L,p+1}} \sim \eta^p \int_{\eta}^{\mu} \xi^{-p} d\xi \rightarrow 0 \quad \text{as } \eta \rightarrow 0.$$

This completes the verification of (50).

### 7.2.2 Lemma 1

**Proof of Lemma** The distribution  $\pi_n$  sets  $\theta_i = \varepsilon_n W_i$ , where  $W_i$  are distributed i.i.d. as  $F_{\bar{\eta}}$ , defined at (18), where  $\bar{\eta} = (1 - \delta)\eta_n$ . The event  $\{\theta \in \Theta^*\}$  is equivalent to  $A_n = \{|W|_{(k)} \leq t_{kn} \eta_n \quad \forall k\}$ , but because of the support bound on  $\mathcal{L}(W)$ , it is enough to consider those  $k$  for which  $t_{kn} \eta_n \leq \mu(\bar{\eta})$ . Expressing  $A_n$  in terms of the empirical distribution  $G_n$  of  $|W_i|$ , we find

$$A_n \supset \{(1 - G_n)(t\eta_n) \leq t^{-p} - n^{-1} \text{ for all } t\eta_n \leq \mu(\bar{\eta})\}$$

For values of  $t$  such that  $t\eta_n \leq \mu(\bar{\eta})$ , the distribution of  $W$  matches that of  $U^{-1/p}$ , where  $U \sim \text{Uniform}[0,1]$ . Consequently

$$\begin{aligned} G_n(t\eta_n) &= n^{-1} \#\{i : \bar{\eta} W_i \leq \eta_n t\} \\ &= n^{-1} \#\{i : U_i \geq (1 - \delta)^p u\}, \quad u = t^{-p}. \end{aligned}$$

Define  $v = (1 - \delta)^p u$  and  $F_{U,n}$  for the empirical distribution of  $n$  uniform observations. Then

$$\begin{aligned} P(A_n) &\geq P\{F_{U,n}(v) \leq (1 - \delta)^{-p} v - n^{-1} \forall v \geq (\bar{\eta}/\mu(\bar{\eta}))^p\} \\ &\geq P\{F_{U,n}(v) \leq \gamma^r v \forall v \geq n^{-1} \log^3 n\} \end{aligned}$$

for some  $\gamma = \gamma(\delta) > 0$  from the assumption on  $\eta_n/\mu(\eta_n)$ . This last probability converges to one, as follows from an empirical (uniform) process convergence result for a non-uniform metric:

$$\sup_{0 < t < 1} \left| \frac{\sqrt{n}(F_{U,n}(v) - v) - W^o(v)}{v^{1/2} \log v^{-1}} \right| \xrightarrow{P} 0, \quad n \rightarrow \infty.$$

(c.f. Shorack and Wellner, p 140, Theorem 1). Here  $W^o(t)$  is a standard Brownian bridge on a common probability space with  $(U_i)$ .

### 7.2.3 Completion of Theorem 1

First note that condition (38) follows easily from (17) and (20). It remains to verify conditions (40) and (41), and we note that it suffices to take  $C = 1$  (otherwise simply redefine  $\eta_n$  by multiplication by  $C$ ). Condition (40) follows from (13), (42) and Lemma 1.

We now show that  $\pi_n$  satisfies (41). Let  $A_n = \{\theta \in \Theta_{n,p}^*(r_n)\}$  and  $W_n = \|\theta\|^2 + E[\|\theta\|^2 | y]$  and observe that  $\Delta_n \leq E(W_n, A_n^c)$ . Now using (19), we have

$$EW_n = 2n\varepsilon^2 EX_1^2 = \frac{4}{2-p} n\varepsilon^2 \eta_n^p \mu_n^{2-p} \sim CB_n.$$

and so it suffices to show that

$$E\left(\frac{W_n}{EW_n}, A_n^c\right) \leq \frac{\sqrt{\text{Var } W_n}}{EW_n} + P(A_n^c) \rightarrow 0.$$

By properties of conditional expectation and (19)

$$\text{Var } W_n \leq 4 \text{Var } \|\theta\|^2 \leq 4n\varepsilon EX_1^4 = Cn\varepsilon^4 \eta_n^p \mu_n^{4-p}.$$

Thus  $\sqrt{\text{Var } W_n}/EW_n = (n\eta_n^p \mu_n^{-p})^{-1/2} \rightarrow 0$  by (13). This completes the proof of (14).

## 7.3 Minimax Risk over Triebel bodies, continued

Write the sequence space white noise model in dyadically indexed form

$$y_I = \theta_I + \varepsilon z_I \quad z_I \stackrel{iid}{\sim} N(0, 1)$$

where  $I \in \mathcal{I} = \{(j, k), j \geq 0, k = 0, 1, \dots, 2^j - 1\}$ . This may be thought of as the expression of the signal-in-Gaussian-white noise model  $dY(t) = f(t)dt + \varepsilon dW(t), t \in [0, 1]$ , in an orthonormal basis of wavelets  $(\psi_I)$  for  $L^2[0, 1]$  such as constructed by Meyer (1991).

We assume that  $\theta$  is constrained to lie in a ball in a Tiebel-Lizorkin norm:

$$\Phi_{p,q}^s(C) = \{\theta : \|\theta\|_{f_{p,q}^s}^p \leq C^p\}$$

where

$$\|\theta\|_{f_{p,q}^s}^p = \int_0^1 \left( \sum_j 2^{sjq} \sum_{k=0}^{2^j-1} |\theta_{jk}|^q \chi_{jk}(r) \right)^{p/q} dr.$$

Here  $\chi_{jk}(r) = I\{k2^{-j} \leq r < (k+1)2^{-j}\}$ . For further information and background we refer to (DJ 1992a).

1. *Bayes-Minimax problem.* Let  $\mathcal{M} = \mathcal{M}_{p,q}^s(c)$  consist of all probability measures  $\pi(d\theta)$  such that  $E_\pi \|\theta\|_{f_{p,q}^s}^p \leq C^p$ . This set is certainly convex and contains point masses in  $\Phi_{p,q}^s(C)$ , but allows arbitrary dependencies among the co-ordinates  $(\theta_I)$ .

2. *Reduction of Bayes-Minimax problem.* We use the partial order derived from the obvious tree structure on  $\mathcal{I}$  in which a ‘parent’ node  $I = (jk)$  gives rise to two daughter nodes at level  $j+1$ . Formally, for  $j > j', I = (j,k) < I' = (j',k')$  iff  $[k2^{j'-j}] = k'$ . Define the ‘ancestor’ set  $A(I)$  of  $I$  to be the (linearly ordered) set of  $I' > I$ . The notation  $\theta_F$  denotes the collection  $(\theta_I; I \in F)$ .

Let  $\kappa(d\psi)$  be a probability measure on sequences  $(\psi_0, \psi_1, \psi_2, \dots)$ . The notation  $\psi_{(j)}$  denotes initial segments  $(\psi_0, \dots, \psi_j)$ . We will call a probability distribution  $\pi(d\theta)$  *subordinated* to  $\kappa$ , (written  $\pi = S\kappa$ ), if  $\pi$  may be constructed via the following recipe:

- (i)  $(\theta_{jk}, k \in \mathcal{I}_j)$  are conditionally independent given  $\mathcal{F}_{j-1} = (\theta_{j',k'}, k' \in I_{j'}, j' < j)$  and
- (ii)  $\mathcal{L}_\pi(\theta_{jk} | \mathcal{F}_{j-1}) = \mathcal{L}_\pi(\theta_{jk} | \theta_{A(jk)})$ ,
- (iii)  $\mathcal{L}_\pi(\theta_{jk} | \theta_{A(jk)} = t_{(j-1)}) = \mathcal{L}_k(\psi_j | \psi_{(j-1)} = t_{(j-1)})$ .

Finally, define  $\bar{\mathcal{M}} = \{\pi \in \mathcal{M} : \pi \text{ is subordinated to some } \kappa \text{ on } \mathbb{R}^N\}$ . A key property of subordinated priors  $\pi$  is that the sequence  $(\theta_{j,[2^j r]})_{j=0}^\infty$  is distributed as  $\kappa$  for every  $r \in [0, 1]$ . In particular, the function

$$\begin{aligned} F(r) &= \sum_j 2^{sjq} \sum_k |\theta_{jk}|^q \chi_{jk}(r) \\ &= \sum_j 2^{sjq} |\theta_{j,[2^j r]}|^q \end{aligned}$$

is identically distributed for each  $r \in [0, 1]$ , and

$$E_\pi \|\theta\|_{f_{p,q}^s}^p = E_\pi \int_0^1 \{F(r)\}^{p/q} dr \tag{55}$$

$$= E_k \left( \sum_j 2^{sjq} |\psi_j|^q \right)^{p/q} \tag{56}$$

Given a prior  $\pi \in \mathcal{M}$ , we construct a subordinated  $\bar{\pi} \in \bar{\mathcal{M}}$  by the following recipe: define  $\kappa(d\psi)$  by introducing  $R \sim U(0, 1)$  independent of  $\theta$  and setting

$$\psi_j(\theta, r) = \theta_{j,[2^j r]}. \tag{57}$$

Now set  $\bar{\pi} = S\kappa$ , i.e. construct  $\bar{\pi}$  through the recipe (i)-(iii) set out earlier.

That  $\bar{\pi} \in \bar{\mathcal{M}}$  follows from (55) and (57);

$$\begin{aligned} E_{\bar{\pi}} \|\theta\|_{f_{p,q}^s}^p &= E_{\kappa} \left( \sum_j 2^{sjq} |\psi_j|^q \right)^{p/q} \\ &= E_{\pi} \int_0^1 \left( \sum_j 2^{sjq} |\theta_{j,[2^j r]}|^q \right)^{p/q} dr \\ &= E_{\pi} \|\theta\|_{f_{p,q}^s}^p \leq C^p. \end{aligned}$$

To show that  $B(\bar{\pi}) \geq B(\pi)$ , we introduce an auxiliary estimation problem in which the data are  $v_i = \psi_i + \varepsilon z_i$ ,  $z_i \stackrel{iid}{\sim} N(0, 1)$ ,  $i \in \mathcal{N}$  and  $\psi_{(j)} = (\psi_0, \psi_1, \dots, \psi_j)$  has prior distribution  $\sigma$ . The Bayes risk for estimation of  $\psi_j$  is denoted

$$b_j^{(0)}(\sigma) = \inf_{\hat{\psi}(v_{(j)})} E_{\sigma} E_{\psi_{(j)}} [\hat{\psi}(v_{(j)}) - \psi_j]^2.$$

Below we use the notation  $F(I) = I \cup A(I)$  and  $\pi_F$  for the marginal distribution of  $(\theta_I)_{I \in F}$ . Using the posterior mean form of the Bayes risk and noting that ignoring variables  $y_{I'}$  for  $I' \notin F(I)$  can only increase risk, we find

$$\begin{aligned} B(\pi) &= \sum_I E_{\pi} E_{\theta} [E(\theta_I | y) - \theta_I]^2 \\ &\leq \sum_I E_{\pi_{F(I)}} E_{\theta_{F(I)}} [E(\theta_I | y_{F(I)}) - \theta_I]^2 \\ &= \sum_{jk} b_j^{(0)}(\pi_{F(jk)}) \\ &\leq \sum_j 2^j b_j^{(0)}(\text{ave}_k \pi_{F(jk)}). \end{aligned}$$

where the last inequality uses concavity of Bayes risk.

On the other hand, since  $\bar{\pi}_{F(jk)}$  are identical for each  $k$  (for fixed  $j$ ), the right side equals exactly  $B^0(\bar{\pi})$ .

We use a renormalization argument to determine the  $\varepsilon$ -dependence of  $B^0(\bar{\mathcal{M}}, \varepsilon)$ . To this end, introduce a two sided version of the previous auxiliary estimation problem:  $v_i = \psi_i + \varepsilon z_i$ ,  $z_i \stackrel{iid}{\sim} N(0, 1)$  but now with  $i \in \mathcal{Z}$ . Let  $b_j(\kappa, \varepsilon)$  denote the Bayes risk for estimation of  $\psi_j$  given prior  $\kappa(d\psi)$  and data  $(v_i, i \leq j)$ . Let  $\mathcal{K}^0$  and  $\mathcal{K}$  respectively denote the collection of probability measures  $\kappa(d\psi)$  on  $\mathbb{R}^{\mathcal{N}}$  and  $\mathbb{R}^{\mathcal{Z}}$  satisfying

$$J_{spq}^p(\kappa) = \int \left( \sum_{-\infty}^{\infty} 2^{sjq} |\psi_j|^q \right)^{p/q} d\kappa(\psi) \leq C^p.$$

We identify  $\mathcal{K}^0$  as a subset of  $\mathcal{K}$  in the obvious way by setting  $\psi_j \equiv 0$  for  $j < 0$ . In this way  $b_j^{(0)}(\kappa, \varepsilon) = b_j(\kappa, \varepsilon)$ . Set now

$$J_{\varepsilon}(\kappa) = \sum_{-\infty}^{\infty} 2^j b_j(\kappa, \varepsilon);$$

we have  $B^0(\bar{\mathcal{M}}, \varepsilon) = \text{val}(P_{\varepsilon, C})$ , where  $P_{\varepsilon, C}$  denotes the optimization problem  $\sup\{J_\varepsilon(\kappa) : \kappa \in \mathcal{K}^0\}$ . Let  $Q_{\varepsilon, C}$  denote the corresponding two-sided optimization over  $\kappa \in \mathcal{K} : \sup\{J_\varepsilon(\kappa) : \kappa \text{ on } \mathbb{R}^Z \text{ satisfies } J_{spq}(\kappa) \leq C\}$ .

This is asymptotically equivalent to  $P_{\varepsilon, C}$  because

$$\text{val}(Q_{\varepsilon, C}) - \varepsilon^2 \leq \text{val}(P_{\varepsilon, C}) \leq \text{val}(Q_{\varepsilon, C}). \quad (58)$$

The first inequality follows because (i) restricting estimators to depend on  $(v_0, \dots, v_j)$  instead of  $(\dots, v_{-1}, v_0, \dots, v_j)$  can only increase Bayes risk, and (ii)  $2^j b(\kappa, \varepsilon) \leq 2^j \varepsilon^2$  which is summed over  $j < 0$ . The second inequality uses the embedding of  $\mathcal{K}^0$  in  $\mathcal{K}$ .

If  $(\psi_j) \sim \kappa$ , Let  $S_{a, h}\kappa$  denote the distribution of the scaled and shifted sequence  $(a\psi_{j-h})$ ,  $a > 0, h \in \mathcal{Z}$ . From the transformation relations

$$J_\varepsilon(S_{a, h}\kappa) = a^{-2} J_{a\varepsilon}(\kappa), \quad J_{spq}(S_{a, h}\kappa) = a2^{hs} J_{spq}(\kappa),$$

it follows that  $v(\varepsilon, C) = \text{val}(Q_{\varepsilon, C})$  satisfies

$$v(\varepsilon, C) = \varepsilon^2 2^h v(1, C\varepsilon^{-1} 2^{-sh}) \quad h \in \mathcal{Z}.$$

Choose  $h \in \mathcal{Z}$  and  $\eta \in [0, 1]$  so that  $s^{-1} \log C\varepsilon^{-1} = h + \eta$ ; i.e.  $C\varepsilon^{-1} = 2^{sh} 2^{s\eta}$ . Then

$$v(\varepsilon, C) = \varepsilon^{2-1/s} C^{1/s} 2^{-\eta} v(1, 2^{s\eta})$$

Combining this with the bounds provided by (58), we have

$$B^0(\bar{\mathcal{M}}, \varepsilon) = \varepsilon^{2r} C^{2(1-r)} \gamma(C\varepsilon^{-1})(1 + o(1)) \quad (59)$$

where  $\gamma(\cdot) = \gamma(\cdot; s, p, q)$  is a continuous periodic function of  $\eta$  and hence  $C\varepsilon^{-1}$ .

*Asymptotic Equivalence.* We use the conditions developed in Section 7.1. The expansion (59) for  $B^0(\bar{\mathcal{M}}, \varepsilon)$ , together with continuity of the periodic function  $\gamma$ , establishes (38) and combined with monotonicity in  $\varepsilon$  of  $R(\varepsilon, C)$ , shows that it suffices to verify (37) on dyadic subsequences of the form  $\varepsilon_h = 2^{-s\eta_0} 2^{-sh}$  as  $h \in \mathcal{N} \nearrow \infty$  with  $\eta_0$  remaining fixed.

Let us construct ‘near-optimal’ priors  $\bar{\pi}_{\varepsilon_h} \in \bar{\mathcal{M}}$  which satisfy conditions (40) and (41). Fix  $\eta$ , and choose  $J, M$  and  $\delta_0$  and a distribution  $\kappa(d\psi)$  in  $\mathcal{K}$  such that

- a)  $\psi_j \equiv 0$  w.p. 1 unless  $|j| \leq J$ ,
- b)  $|\psi_j| \leq M$  w. p. 1. if  $|j| \leq J$ ,
- c)  $\sum_{-J}^J 2^j b_j(\kappa, 1) \geq v(1, C)(1 - \eta)$ ,
- d)  $E_\kappa(\sum_{-J}^J 2^{sjq} |\psi_j|^q)^{p/q} \leq C^p(1 - \delta_0)$ .

For  $h > J$ , construct  $2^{h-J}$  independent realizations  $\xi^{(l)} = (\xi_{j,k}^{(l)}, j, k \in \mathcal{I})$  from the dyadic prior measure subordinated to  $T_{1, J}\kappa$ . The prior  $\pi_{\varepsilon_h}$  is obtained by ‘attaching’ each of these replicas as ‘trees’ rooted at  $(h - J, l), l = 0, \dots, 2^{h-J} - 1$ . Specifically

$$\theta_{jk} := \varepsilon_h \xi_{j-(h-J), k}^{(l)}$$

where  $k'$  is such that  $k = l2^{j-(h-J)} + k'$ .

The point of this construction is that

$$\begin{aligned} \|\theta\|_{f_{p,q}^s}^p &= \int_0^1 \left( \sum_{j=h-J}^{h+J} 2^{sjq} |\theta_{j,[2^j t]}|^q \right)^{p/q} dt \\ &= \text{Ave}(V_l; l = 0, \dots, 2^{h-J} - 1), \end{aligned}$$

where  $V_l$  are i.i.d. copies distributed as

$$V_1 = \int_0^1 \left( \sum_{-J}^J 2^{sjq} |\xi_{j,[2^{j+J} t]}|^q \right)^{p/q} dt.$$

Since  $\xi_{j+J,[2^{j+J} t]}, |j| \leq J \stackrel{D}{=} (\psi_j, |j| \leq J)$ ,  $EV_1 = E(\sum_{-J}^J 2^{sjq} |\psi_j|^q)^{p/q} \leq C^p(1 - \delta_0)$ , condition (40) follows from the law of large numbers:

$$\pi_{\varepsilon_h}(\|\theta\|_{f_{p,q}^s} \leq C) \rightarrow 1.$$

To verify condition (41), first note that the bound on the support of  $\psi_j$  ensures that

$$|\theta_{\nu_n, jk}| = |E_{\nu_n}(\theta_{jk}|y)| \leq M\varepsilon_h,$$

so that

$$\|\theta_{\nu_n}\|^2 + \|\theta\|^2 \leq 2M^2 \varepsilon_h^2 \sum_{h-J}^{h+J} 2^j \asymp 2^h \varepsilon_h^2 = \varepsilon_h^{2r}.$$

Consequently, in the notation of section 7.1, since  $B(\varepsilon_h, C) \asymp \varepsilon_h^{2r}$  as  $h \rightarrow \infty$

$$\Delta_h/B(\varepsilon_h, C) \leq c\varepsilon_h^{2r} \pi_h(C\Theta^c)/\varepsilon_h^{2r} \rightarrow 0$$

which establishes condition (41) and hence the required asymptotic equivalence.

## References

- [1] Bergh, J. and Löfström. (1976) *Interpolation Spaces. An Introduction*. New York, Springer-Verlag.
- [2] Bickel, P. J. (1983). Minimax estimation of a normal mean subject to doing well at a point. In *Recent Advances in Statistics* (M. H. Rizvi, J. S. Rustagi, and D. Siegmund, eds.), Academic Press, New York, 511–528.
- [3] Chentsov, N. N. (1972). *Statistical Decision Functions and Optimal Inference*. Translations of Mathematical Monographs, Vol. 53. American Math. Society. Providence, R.I.
- [4] Daubechies, I. (1992) *Ten Lectures on Wavelets*. SIAM. Philadelphia.

- [5] DeVore, R. A. (1989) Degree of Nonlinear Approximation. *in Approximation Theory VI, Volume 1*, C.K. Chui, L.L. Schumaker and J.D. Ward (eds.) pp. 175 – 201. Academic Press.
- [6] Donoho, D. L. and Johnstone, I. M (1990) Minimax risk over  $\ell_p$ -balls. Technical Report, Department of Statistics, University of California, Berkeley.
- [7] Donoho, D. L. and Johnstone, I. M (1992a) Minimax estimation via Wavelet Shrinkage. Technical Report, Department of Statistics, Stanford University.
- [8] Donoho, D. L. and Johnstone, I. M (1992b) Ideal Spatial Adaptation via Wavelet Shrinkage. Technical Report, Department of Statistics, Stanford University.
- [9] Donoho, D. L., Johnstone, I. M and Taswell, C. (1992c) Wavelet software for data analysis using MATLAB. Manuscript.
- [10] Donoho, D. L. and Johnstone, I. M (1992d) Non-classical minimax problems, thresholding, adaptation. Manuscript.
- [11] Efroimovich, S.Y. and Pinsker, M.S. (1981) Estimation of square-integrable [spectral] density based on a sequence of observations. *Problemy Peredatsii Informatsii* **17** 50-68 (in Russian); *Problems of Information Transmission* (1982) 182-196 (in English).
- [12] Efroimovich, S.Y. and Pinsker, M.S. (1982) Estimation of square-integrable probability density of a random variable. *Problemy Peredatsii Informatsii* **18** 19-38 (in Russian); *Problems of Information Transmission* (1983) 175-189 (in English).
- [13] Farrell, R.H. (1972) On the best obtainable asymptotic rates of convergence in estimation of a density function at a point. *Ann. Math. Statist.* **43**. 170–180.
- [14] Feldman, I. (1991) Constrained minimax estimation of the mean of the normal distribution. *Ann. Statist.* **19**. 2259–2265.
- [15] Frazier, M., Jawerth, B. and Weiss G. (1991) *Littlewood-Paley theory and the study of function spaces*. CBMS – Conference Lecture Notes 79, American Mathematical Society, Providence, R.I.
- [16] Golubev, G.K. and Nussbaum, M. (1990). A Risk Bound in Sobolev class regression. *Ann. Statist.* **18** 758–778.
- [17] Hodges, J.L. Jr. and Lehmann, E.L. (1952). The use of previous experience in reaching statistical decisions. *Ann. Math. Statist.* bf 23, 396–407.
- [18] Ibragimov, I.A. and Has'minskii, R.Z. (1981) *Statistical Estimation, Asymptotic Theory*. Translation: Springer-Verlag, New York.
- [19] Ibragimov, I.A. and Has'minskii, R.Z. (1982) Bounds for the risk of nonparametric regression estimates. *Theory Probab. Appl.* **27** 84-99.



- [20] Johnstone, I.M. and Silverman, B.W. (1990) Speed of estimation in positron emission tomography and related inverse problems. *Ann. Statist.* **18** 251–280.
- [21] LeCam, L. (1986) *Asymptotic methods in statistical decision theory*. Springer Verlag, New York.
- [22] Morris, C. (1983) Parametric empirical Bayes inference: Theory and applications. *J. Amer. Statist. Assoc.* **78**, 47–65.
- [23] Nemirovskii, A.S. (1985) Nonparametric estimation of smooth regression functions. *Izv. Akad. Nauk. SSR Tekhn. Kibernet.* **3**, 50-60 (in Russian). *J. Comput. Syst. Sci.* **23**, 6, 1-11, (1986) (in English).
- [24] Nemirovskii, A.S., Polyak, B.T. and Tsybakov, A.B. (1985) Rate of convergence of nonparametric estimates of maximum-likelihood type. *Problems of Information Transmission* **21**, 258-272.
- [25] Nussbaum, M. (1985) Spline smoothing and asymptotic efficiency in  $L_2$ . *Ann. Statist.*, **13**, 984–997.
- [26] Peetre, J. (1976) *New Thoughts on Besov Spaces*. Duke University Mathematics Series. Number 1.
- [27] Pinsker, M.S. (1980) Optimal filtering of square integrable signals in Gaussian white noise. *Problemy Peredatsii Informatsii* **16** 52-68 (in Russian); *Problems of Information Transmission* (1980) 120-133 (in English).
- [28] Shorack, G.R. and Wellner, J.A. (1986) *Empirical Processes with Applications to Statistics*. Wiley, New York.
- [29] Sion, M. On general minimax theorems. *Pacific J. Math.* **8** , 171–176.
- [30] Speckman, P. (1985) Spline smoothing and optimal rates of convergence in nonparametric regression models. *Ann. Statist.*, **13**, 970-983.
- [31] Stone, C. (1982). Optimal global rates of convergence for nonparametric estimators. *Ann. Statist.*, **10**, 1040-1053.

## Captions

**Table 1**

$j$	index in renormalized problem (24)
$\bar{j} = j + h$	resolution level in figures
$t_j$	constraint on $p$ -th moment in $j$ -th level (cf. (24) ). In ellipsoid case, same as standard deviation of level $j$ prior. (standardised with noise level $\epsilon = 1$ ).
$2^{j\beta q} t_j^q$	contribution of this level to constraint (since $c = 1$ , these are percentages.)
$\rho_p(t_j)$	For $p = 2$ , equals ellipsoid Bayes risk, for $p < 2$ , equals thresholding Bayes risk $\bar{\rho}_p$ for the least favorable prior (note, lies in $[0, 1]$ ).
$2^j \rho_p(t_j)$	contribution of this level to total Bayes risk.
$(\epsilon_j, \mu_j)$	for $p < 2$ , size and location (in the standardised case $\epsilon = 1$ ) of non-zero atom in least favorable three point distribution for level $j$ .
$E_j$	equals $2^j \epsilon_j$ , the expected number of non-zero coefficients at the $j$ th level. (compare Figure 3).

**Figure 1:** Sample paths from least favorable priors in the ellipsoidal case, solving (31) for the indicated parameter values. The corresponding noise level is  $\epsilon_h^2 = 2^{-\beta h}$ , which equals either  $2^{-12}$  or  $2^{-12.5}$ . Index  $\sigma$  measures smoothness, corresponding roughly to the number of derivatives in the  $L_2$  sense. Sample path constructed using (32) from wavelet coefficients drawn from Gaussian distributions with standard deviations given in Table 1.  $N = 2^{12} = 4096$  plotting points. Actual (scaled) wavelet coefficients shown in Figure 3. Note that the same pseudo-Gaussian input sequence is used in all plots in Figures 1 and 2 to make the sample paths more readily comparable.

**Figure 2:** Sample paths from near least favorable priors for two parameter configurations corresponding to “sparsity”. Wavelet coefficients drawn from the three point distributions at each level given in Table 1. The sparsity of the signal is apparent for  $p = 1$ , but much more pronounced for  $p = .5$ .

**Figure 3:** Wavelet coefficients generated from parameters in Table 1. At resolution level  $\bar{j}$  there are  $2^{\bar{j}}$  coefficients. Coefficients are scaled for plotting purposes so that the largest coefficient is .9. The parameter  $h$  represents the shift from level  $j$  in the renormalized problem to level  $\bar{j}$  used in producing the plots corresponding to noise level  $\epsilon^2 = 2^{-\beta h}$ .

**Figure 4:** Graphical representation of the parameters associated with the numerical solution of the univariate minimax thresholding problem (34) and (35). Solid line corresponds to  $p = 1$ , dashed line to  $p = .5$  and dotted line to  $p = 2$ . Horizontal axis is  $L_p$  moment constraint  $\eta$ .

$j$	$\bar{j}$	$t_j$	$2^{j\beta q} t_j^q$	$\rho_p(t_j)$	$2^j \rho_p(t_j)$	$\epsilon_j$	$\mu_j$	$E_j$
$p = q = 2, \sigma = 1, \beta = 1.5, h = 8$								
-8	0	20.26	0.000	0.998	0.00			
-7	1	14.31	0.000	0.995	0.01			
-6	2	10.09	0.000	0.990	0.02			
-5	3	7.10	0.001	0.981	0.03			
-4	4	4.97	0.006	0.961	0.06			
-3	5	3.44	0.023	0.922	0.12			
-2	6	2.33	0.085	0.844	0.21			
-1	7	1.49	0.277	0.689	0.34			
0	8	0.78	0.607	0.378	0.38			
$p = q = 2, \sigma = 2, \beta = 2.5, h = 5$								
-5	0	42.66	0.000	0.999	0.03			
-4	1	21.31	0.000	0.998	0.06			
-3	2	10.62	0.003	0.991	0.12			
-2	3	5.24	0.027	0.965	0.24			
-1	4	2.47	0.191	0.859	0.43			
0	5	0.88	0.778	0.438	0.44			
$p = q = 1, \sigma = 1, \beta = 1.5, h = 8$								
-4	4	1.5849	0.025	0.945	0.06	0.9890	1.60	16
-3	5	1.2589	0.056	0.867	0.11	0.7724	1.63	25
-2	6	0.8423	0.105	0.707	0.18	0.4954	1.70	32
-1	7	0.5346	0.189	0.534	0.27	0.2965	1.81	38
0	8	0.2512	0.251	0.314	0.31	0.1254	2.00	32
1	9	0.0794	0.225	0.127	0.25	0.0341	2.33	18
2	10	0.0158	0.127	0.033	0.13	0.0057	2.78	6
3	11	0.0010	0.023	0.003	0.02	0.0003	3.52	.6
$p = q = .5, \sigma = 2, \beta = 2.5, h = 5$								
-4	1	0.1980	0.014	0.599	0.04	0.2891	2.37	.5
-3	2	0.1001	0.023	0.485	0.06	0.2020	2.45	.8
-2	3	0.0500	0.039	0.386	0.10	0.1402	2.55	1.1
-1	4	0.0252	0.067	0.304	0.15	0.0977	2.64	1.6
0	5	0.0079	0.089	0.200	0.20	0.0533	2.80	1.7
1	6	0.0032	0.139	0.147	0.28	0.0329	2.92	2.1
2	7	0.0012	0.198	0.098	0.39	0.0200	3.06	2.5
3	8	0.0002	0.209	0.051	0.41	0.0086	3.27	2.3
4	9	0.0000	0.226	0.027	0.43	0.0038	3.48	1.9