

# On minimax estimation of a sparse normal mean vector

Iain M. Johnstone  
Department of Statistics  
Stanford University

*July 1991*  
*Revised May 1993*

## Abstract

Mallows has conjectured that among distributions which are Gaussian but for occasional contamination by additive noise, the one having least Fisher information has (two-sided) geometric contamination. A very similar problem arises in estimation of a non-negative vector parameter in Gaussian white noise when it is known also that most, i.e.  $(1 - \epsilon)$ , components are zero.

We provide a partial asymptotic expansion of the minimax risk as  $\epsilon \rightarrow 0$ . While the conjecture seems unlikely to be exactly true for finite  $\epsilon$ , we verify it asymptotically up to the accuracy of the expansion. Numerical work suggests the expansion is accurate for  $\epsilon$  as large as .05. The best  $\ell_1$ -estimation rule is first but not second order minimax. The results bear on an earlier study of maximum entropy estimation and various questions in robustness and function estimation using wavelet bases.

*Key Words and Phrases* Fisher information, minimax decision theory, least favorable prior, nearly black object, robustness, white noise model,

*AMS 1980 (1985 Revision) Subject Classifications* Primary 62C20, Secondary 62C10 62G05

# 1 Introduction

In many estimation settings, there is definite prior information concerning the values of a parameter vector  $\theta$ . One may have bounds on the individual components  $\theta_i$  — “all  $\theta_i$  lie between 0 and 1,” or on particular functionals of the whole vector— “the squared length of  $\theta$  is at most  $c$ ,” or “most of the  $\theta_i$  are zero.” Many estimation methods have been developed to capitalize on such prior information, either explicitly in the form of constraints on an optimization procedure (e.g. positivity-constrained least squares) or implicitly in the sense that the estimator performs well “if and only if”  $\theta$  satisfies the prior constraints. An example of the latter are maximum entropy regularization estimates in the case of “nearly black”  $\theta$  (e.g. Frieden (1972), Gull and Daniel (1978)).

How does one compare the performance of various possible estimators when such prior information is present? One common, admittedly conservative, approach is the worst-case analysis: given some error measure, compute the maximum expected error over the restricted parameter space, and then seek the estimator that minimizes this maximum risk. The resulting best or minimax risk provides (i) a benchmark against which to measure other estimators, and (ii) a measure of the value of the prior information (by comparison with the minimax risk computed ignoring the prior information).

This paper presents an asymptotic evaluation of this minimax risk in a simple, but hopefully representative, context: estimation of non-negative signals that are mostly zero, such as spectra, star maps and the like. Our model is idealized in certain significant ways: a signal-plus-noise model is adopted, and the noise is assumed to be i.i.d. Gaussian. These assumptions permit a more detailed analysis, but at the price of a degree of non-robustness to departures from the model.

To state the problem precisely, let  $X \sim N(\theta, 1)$ . Denote the mean squared error, or risk, of an estimator  $d(x)$  by  $R(\theta, d) = E_\theta (d(X) - \theta)^2$ . If  $G(d\theta)$  is a prior probability distribution, denote the integrated risk of  $d$  by  $r(G, d) = \int R(\theta, d)G(d\theta)$ . Let  $\delta_a$  denote point mass at  $a$  and  $\mathcal{P}([0, \infty))$  the class of probability measures supported on  $[0, \infty)$ . Consider a class of priors on  $[0, \infty)$  with an atom at 0:

$$\mathcal{G}_\epsilon = \{G = (1 - \epsilon)\delta_0 + \epsilon H : H \in \mathcal{P}([0, \infty))\}.$$

This paper is concerned with asymptotic evaluation of the (restricted Bayes) minimax risk<sup>1</sup>

$$m(\epsilon) = \sup_{\mathcal{G}_\epsilon} \inf_d r(G, d) = \inf_d \sup_{\mathcal{G}_\epsilon} r(G, d) \tag{1}$$

as  $\epsilon \rightarrow 0$ . This problem or a close relative arises in a variety of contexts, some of which we now review.

*Nearly black objects.* As part of a study of maximum entropy estimation, Donoho et. al. (1992a) consider the problem of estimating a *non-negative* vector  $\tilde{\theta} = (\theta_i)_{i=1}^n$  from noisy data

$$X_i = \theta_i + Z_i \quad i = 1, \dots, n. \tag{2}$$

---

<sup>1</sup>Here and later, footnote numbers refer to extra details found in the Appendix

where the noise terms  $Z_i$  are i.i.d.  $N(0, \sigma^2)$ . They show that a maximum entropy rule, defined as a minimiser  $\hat{\theta}_{ME, \lambda}(x)$  of

$$\sum_i (\theta_i - x_i)^2 + 2\lambda \sum_i \theta_i \log \theta_i, \quad (3)$$

achieves significant savings over linear rules in mean squared error when and only when most components  $\theta_i$  of the unknown object are nearly zero. It is then natural to define the class of “ $\epsilon$ -black images”  $\Theta_n(\epsilon)$  as the set of non-negative sequences  $(\theta_i)$  of length  $n$  satisfying  $\#\{i : \theta_i > 0\} \leq n\epsilon$ . The worst case error of an estimator  $\hat{\theta}_n(x)$  is

$$m_n(\hat{\theta}_n, \epsilon) = \sup_{\Theta_n(\epsilon)} E_\theta n^{-1} \sum_1^n (\hat{\theta}_{n,i}(X) - \theta_i)^2.$$

A benchmark against which particular estimators may be measured is the minimax risk

$$m_n(\epsilon) = \inf_{\hat{\theta}_n} m_n(\hat{\theta}_n, \epsilon). \quad (4)$$

As  $n$  becomes large,  $m_n(\epsilon)$  approaches a limit  $m(\epsilon)$ , and the evaluation of  $m(\epsilon)$  is exactly equivalent to the problem (1) (Donoho et. al., 1992a, Theorem 1).

*Robust estimation.* Denote the Fisher information of an absolutely continuous distribution  $F$  with density  $f$  by  $I(F) = \int f'^2/f$ . In Huber’s (1964, 1981) asymptotic minimax approach to robust estimation, there arises the problem of minimising Fisher information over neighborhoods of the standard Gaussian distribution. One possibility is to consider Gaussian variables that are occasionally contaminated by additive noise, i.e. distributions  $F = \Phi * G$ , where  $G$  belongs to  $\mathcal{G}'_\epsilon = \{G = (1 - \epsilon)\delta_{\{0\}} + \epsilon H, H \in \mathcal{P}((-\infty, \infty))\}$ . Here  $\Phi$  is the standard Gaussian distribution and  $*$  denotes convolution. Mallows (1978) noted that the distribution  $G_0$  minimizing  $I(\Phi * G)$  would be symmetric and discrete:  $G_0 = \sum p_j \delta_{g_j}$ , with  $G_0(\{0\}) = 1 - \epsilon$ . He further stated that “a plausible guess is that  $p_j = cp^j$ ,  $g_j = jg$  for  $j > 0$ ,” i.e., that  $G_0$  is a two-sided geometric distribution. This problem is connected with (1) through Brown’s (1971) identity

$$r(G) \triangleq \inf_d r(G, d) = 1 - I(\Phi * G).$$

Thus Mallows’ problem is identical to (1) except that  $\mathcal{G}'_\epsilon$  allows two-sided contamination distributions.

*Parametric robustness.* Bickel (1983, 1984) studied the question of optimal (minimax) estimation of  $\theta$  subject to good risk properties at  $\theta = 0$ : namely to calculate

$$\inf_d \{\sup_\theta R(\theta, d) : R(0, d) \leq s\}. \quad (5)$$

Bayes-minimax compromises (of which this is an important special case) were earlier studied by Hodges and Lehmann (1952) and Efron and Morris (1971). For related work, see also Berger (1982) and Marazzi (1980). As Bickel notes, this problem is central to estimation of a parameter  $\eta$  in the presence of a nuisance parameter  $\theta$ , when one believes that  $\theta = 0$  but desires robustness against the possibility of error.

Introducing a Lagrange multiplier  $\epsilon$  shows an equivalent form of (5) to be (1), using  $\mathcal{G}'_\epsilon$  in place of  $\mathcal{G}_\epsilon$ :

$$\begin{aligned} & \inf_d \{ (1 - \epsilon)R(0, d) + \epsilon \sup_\theta R(\theta, d) \} \\ &= \inf_d \{ (1 - \epsilon)R(0, d) + \epsilon \sup_{H \in \mathcal{P}(-\infty, \infty)} R(H, d) \} \\ &= \inf_d \sup_{\mathcal{G}'_\epsilon} R(G, d) \end{aligned} \tag{6}$$

Equivalence of the two problems means that for given  $s$  and optimal rule  $d_s^*$  in (5), there exists  $\epsilon = \epsilon(s)$  such that the optimal rule  $d_\epsilon^*$  in (6) equals  $d_s^*$ . Bickel also describes in greater detail the connections between (6), the identities of Brown and Stein, and Mallows' problem.

The remainder of this note is organised as follows. Section 2 presents a three term asymptotic expansion for the minimax risk (1) and compares it with numerical approximations. Higher order properties of the simple  $l_1$  rules studied in Donoho et. al. (1992a) are briefly described. Section 3 outlines the proof of the main result and Section 4 briefly describes connections with other constraint sets ( $l_p$  balls) and with wavelet transforms.

## 2 Main Results.

Introduce, following Mallows, the geometric prior

$$G_\epsilon = (1 - \epsilon) \sum_{k=0}^{\infty} \epsilon^k \delta_{kn},$$

The lattice spacing  $n = n(\epsilon)$  is defined implicitly along with an additional parameter  $a = a(n)$  by the two equations

$$\phi(n + a) = \epsilon \phi(a) \tag{7}$$

$$(n + a)\phi(a) = 2\Phi(a) \tag{8}$$

Here  $\Phi$  and  $\phi$  are the standard Gaussian distribution and density functions respectively. The origin of these equations is discussed later, but for now we note that the orders of magnitude of  $n(\epsilon)$  and  $a_n = a(n)$  are given by

$$n \sim \sqrt{2 \log \epsilon^{-1}}, \quad a_n \sim \sqrt{2 \log c_0 n}, \quad c_0 = (2\sqrt{2\pi})^{-1}. \tag{9}$$

Define the Bayes risk of a prior  $G$ , and the maximum risk (relative to  $\mathcal{G}_\epsilon$ ) of an estimator  $d$  respectively by

$$r(G) = \inf_d r(G, d), \quad m(d, \epsilon) = \sup_{\mathcal{G}_\epsilon} r(G, d).$$

The expression (1) for minimax risk now takes the simpler form

$$\sup_{\mathcal{G}_\epsilon} r(G) = m(\epsilon) = \inf_d m(d, \epsilon).$$

**Theorem 1** *The minimax risk*

$$m(\epsilon) = \epsilon n^2 \Phi(a_n) - \frac{\pi^2}{6} \epsilon a_n \phi(a_n) + O\left(\frac{\epsilon a_n^2}{n^2}\right). \quad (10)$$

The geometric prior  $G_\epsilon$  is asymptotically least favorable and the Bayes estimator  $d_{G_\epsilon}$  is asymptotically minimax to this order:

$$m(d_{G_\epsilon}, \epsilon) - r(G_\epsilon) = O\left(\frac{\epsilon a_n^2}{n^2}\right). \quad (11)$$

Formula (10) will be checked numerically below. For theoretical interpretation, some remarks are in order. Definition (8) shows that  $n\phi(a_n) = 2 + O(a_n/n)$ , which converts the second term to  $(\pi^2/3)(\epsilon a_n/n)$  without increasing the order of error. Setting  $\tilde{\Phi} = 1 - \Phi$ , we obtain

$$m(\epsilon) = \epsilon n^2 - \epsilon n^2 \tilde{\Phi}(a_n) - \left(\frac{\pi^2}{3}\right) \frac{\epsilon a_n}{n} + O\left(\frac{\epsilon a_n^2}{n^2}\right). \quad (12)$$

Since  $n^2 \tilde{\Phi}(a_n) \sim 2n/a_n$ , this shows that (10) is actually a *third* order expansion of  $m(\epsilon)$  and that the geometric prior is *third*-order minimax. We are unaware of any previous settings in which third-order minimaxity in a parameter other than sample size has been established, though Levit(1986) gives bounds on the third term of the minimax risk for estimating a Gaussian mean restricted to an interval as the noise level decreases.

Finally, we note that, expressed in terms of  $\epsilon$ , the order of error is  $\epsilon a_n^2 n^{-2} \sim \epsilon(\log \log \epsilon^{-1} / \log \epsilon^{-1})$ .

*Simpler approximations.* The dependence of  $m(\epsilon)$  on  $\epsilon$  in (10) is implicit, since it involves the solutions of equations (7) and (8). Two cruder approximations may be derived<sup>2</sup>, the second involving only elementary functions. To this end, let  $n_0^2 = 2 \log \epsilon^{-1}$  and recall that  $c_0 = (2\sqrt{2\pi})^{-1}$ .

$$\begin{aligned} m(\epsilon) &= \epsilon n_0^2 \Phi((2 \log c_0 n_0)^{1/2}) - 2\epsilon n_0 (2 \log c_0 n_0)^{1/2} + O(\epsilon \log n_0) & (13) \\ &= \epsilon n_0^2 - 2\epsilon n_0 (2 \log c_0 n_0)^{1/2} - \epsilon n_0 (2 \log c_0 n_0)^{-1/2} + O(\epsilon n_0 (\log n_0)^{-3/2}) & (14) \end{aligned}$$

*First-order minimax rules.* In the two-sided version of (1) that uses  $\mathcal{G}'_\epsilon$  in place of  $\mathcal{G}_\epsilon$ , the first-order asymptotic minimax behavior of  $m(\epsilon)$  was described by Bickel (1983). The corresponding result for (1) is proved in Donoho et. al. (1992a) and says simply that

$$m(\epsilon) = \epsilon n_0^2 (1 + o(1)). \quad (15)$$

Denote by ‘ $\ell_1$ -rule’ an estimator of the form  $d_\lambda(x) = \max(x - \lambda, 0)$  for  $\lambda > 0$ . The name arises because  $\hat{\theta}_\lambda = (d_\lambda(x_i))$  is the co-ordinatewise solution to the problem

$$\min_{\theta \geq 0} \sum_i (\theta_i - x_i)^2 + 2\lambda \sum_i \theta_i.$$

Both Bickel (1983) and Donoho et. al. (1992a) show for their respective settings that  $\ell_1$ -rules with  $\lambda(\epsilon) \sim n_0(\epsilon) = (2 \log \epsilon^{-1})^{1/2}$  are first-order asymptotically minimax:

$$m(d_{\lambda_\epsilon}, \epsilon) \sim m(\epsilon).$$

However, as Bickel notes, the first-order approximation is rather crude and not practically useful. The next result gives higher-order behavior of the best  $\ell_1$ -rule.

**Theorem 2** *Let  $\lambda_\epsilon$  minimize  $m(d_\lambda, \epsilon)$ . Then*

$$\lambda_\epsilon = (n_0^2 - 6 \log n_0 - \log 2\pi)^{1/2} + O(n_0^{-3} \log n_0) \quad \text{and} \quad (16)$$

$$m(d_{\lambda_\epsilon}, \epsilon) = \epsilon \left[ n_0^2 - 6 \log n_0 - \log 2\pi + 3 + 18n_0^{-2} \log n_0 + O(n_0^{-2}) \right]. \quad (17)$$

*Discussion of (17).* A simple relation between  $n_0$  and  $n$  results from putting the definition  $\epsilon = e^{-n_0^2/2}$  into (7):

$$n_0^2 = n^2 + 2na_n. \quad (18)$$

Substituting (18) into (17) and comparing with (10) shows immediately that the best  $\ell_1$ -rule is not even second-order minimax. Although hardly needed here for theoretical purposes, the extra terms in expansion (17) occur naturally in the proof and are retained to provide possibly greater numerical accuracy.

*Numerical evaluation.* Table 1 compares the various approximations for a range of values of  $\epsilon$ . Note, for example, that  $\epsilon = 10^{-6}$  corresponds to a single non-zero pixel in a  $1000 \times 1000$  screen image. To calculate approximation (10), values of  $n$  and  $a_n$  were obtained by solving equations (7) and (8) numerically.

As a check on these values, one can take geometric priors  $G_{\nu, \epsilon} = (1 - \epsilon) \sum_0^\infty \epsilon^k \delta_{k\nu}$  and compute the Bayes risk  $r(G_{\nu, \epsilon})$  by numerical integration and summation. By numerical minimization over  $\nu$ , one obtains an optimal spacing  $\nu = n_{lo}(\epsilon)$  and thus a lower bound  $m_{lo}(\epsilon)$  to  $m(\epsilon)$  that Mallows' conjecture suggests ought to be quite sharp (and would in fact equal  $m(\epsilon)$  were the conjecture to be exactly true). The results for selected values of  $\epsilon$  are displayed in Table 2.

A corresponding numerical upper bound for  $m(\epsilon)$  was obtained by locating the maximum  $\theta_{max}$  of the risk function of the Bayes rule corresponding to  $G_{n_{lo}(\epsilon), \epsilon}$  and evaluating the left side of (6) using this Bayes rule and  $\theta_{max}$  to obtain the upper bound  $m_{up}(\epsilon)$ . Table 2 shows that even at  $\epsilon = .2$  the upper and lower bounds differ by only 2.5%, and the bounds become tighter as  $\epsilon$  decreases.

The agreement between Table 2 and (10) is remarkable at  $\epsilon = .02$  and  $\epsilon = .01$ , suggesting that (10) is likely to be quite accurate for smaller  $\epsilon$ . Over the range of numerical calculations, approximation (10) is typically smaller than  $m_{lo}(\epsilon)$ , but by an amount less than the difference between the upper and lower bounds  $m_{up} - m_{lo}$ . This implies a relative error in  $m(\epsilon)$  of about 6% at  $\epsilon = .05$ , dropping to about 2.5% at  $\epsilon = .001$ . By contrast, the first order expression (15) is too large by a factor of 2 for plausible values of  $\epsilon$ , and the approximations (13) and (14), while considerably better for small  $\epsilon$ , are useless above  $\epsilon \sim 10^{-6}$ , due to the logarithms being undefined.

Table 3 shows the approximations to  $\lambda_\epsilon$  and  $m(d_{\lambda_\epsilon}, \epsilon)$  as given by Theorem 2. Formulas (34) and (44) in the proof of Theorem 2 below show that it is in principle not difficult to calculate  $m(d_{\lambda_\epsilon}, \epsilon)$  directly by evaluating the critical  $\lambda_\epsilon$  numerically and then substituting into (34). This was done to provide a check on (17). For  $\epsilon \geq .02$ , the approximation is undefined, and indeed the asymptotic formula is barely satisfactory over the range shown. It would presumably improve for  $\epsilon \leq 10^{-6}$ . In any case, the table shows clearly the higher-order suboptimality of the  $\ell_1$ -rule over a wide range of  $\epsilon$ .

< Tables 1, 2, and 3 about here >

### 3 Approximating minimax rules

This section outlines the proofs of Theorems 1 and 2, with some details deferred to the Appendix. The usual device for proving that the geometric prior  $G_\epsilon$  is asymptotically minimax is to show that its Bayes risk is close to the maximum risk  $m(d_{G_\epsilon}, \epsilon)$  of the Bayes rule corresponding to  $G_\epsilon$ . For  $G_\epsilon$ , the Bayes risk

$$r(G_\epsilon) = (1 - \epsilon)R(0, d_{G_\epsilon}) + \epsilon(1 - \epsilon) \sum_{k=1}^{\infty} \epsilon^{k-1} R(kn, d_{G_\epsilon}) \quad (19)$$

Since the “maximum risk”

$$m(d, \epsilon) = \sup\{(1 - \epsilon)R(0, d) + \epsilon \int R(\theta, d)H(d\theta) : H \in \mathcal{P}([0, \infty))\},$$

it may be rewritten for the Bayes estimator as

$$m(d_{G_\epsilon}, \epsilon) = (1 - \epsilon)R(0, d_{G_\epsilon}) + \epsilon \sup_{\theta \geq 0} R(\theta, d_{G_\epsilon}). \quad (20)$$

To establish (11), the method is to construct a lattice spacing  $n$  so that up to terms of order  $O(a_n^2/n^2)$ , the maximum of  $R(\theta, d_{G_\epsilon})$  is attained at each of the support points  $\{kn; k = 1, 2, \dots\}$  of the positive component of  $G_\epsilon$ .

To this end, approximations to the risk of  $d_{G_\epsilon}$  are useful. The posterior distribution  $G_\epsilon(\{kn\} | x)$  is proportional to  $\epsilon^k \phi(x - kn)$ , and for  $n$  in the range of interest ( $\geq 2.5$ ), this is almost entirely concentrated on at most two points. So, introduce the change of variables  $x = k_0 n + z$  and note that

$$G_\epsilon(\{kn\} | x) \propto \begin{cases} \epsilon^{-1} \phi(z + n) & k = k_0 - 1 \\ \phi(z) & \text{as } k = k_0 \\ \epsilon \phi(z - n) & k = k_0 + 1. \end{cases}$$

In fact, for most  $z$ , a single value of  $k$  dominates, so that the Bayes estimator  $d_{G_\epsilon}$ , being the mean of the posterior distribution, approximately equals  $kn$ . The contributions from the support points  $(k_0 - 1)n$  and  $k_0 n$  balance when  $\epsilon^{-1} \phi(z + n) = \phi(z)$ , that is, when  $\log \epsilon^{-1} - nz - n^2/2 = 0$ . The defining equation (7) shows that this occurs when  $z = a_n$ , that is, when  $x = k_0 n + a_n$ . One finds similarly that  $k_0 n$  and  $(k_0 + 1)n$  balance when  $z = a_n + n$ . Thus the posterior distribution is

approximately periodic in  $x$  with period  $n$ , at least for  $x$  situated away from the left endpoint of the support of  $G_\epsilon$ .

For  $x$  within  $n/2$  standard deviations of  $k_0n + a_n$ , we therefore approximate  $d_{G_\epsilon}(x)$  by the Bayes rule for a two point prior putting mass proportional to 1 at  $k_0n$  and to  $\epsilon^{-1}$  at  $(k_0 - 1)n$ . The approximation is slightly modified for  $x \leq n/2 + a$  to the Bayes rule for the two-point prior  $\delta_0 + \epsilon\delta_n$ .

Explicitly, let  $z \sim N(\zeta, 1)$  and  $\zeta \sim \delta_0 + \epsilon^{-1}\delta_{-n}$ . Then

$$d_0(z) = E(\zeta | z) = \frac{-ne^{-nz}}{e^{-na} + e^{-nz}}. \quad (21)$$

Let  $k_0(x)$  denote the *positive* integer  $k$  for which  $kn + a$  is closest to  $x$ , and set

$$d_\epsilon(x) = nk_0(x) + d_0(x - nk_0(x)). \quad (22)$$

Estimator  $d_\epsilon$  uniformly approximates  $d_{G_\epsilon}$  both pointwise and in risk<sup>3</sup> :

$$|d_{G_\epsilon}(x) - d_\epsilon(x)| \leq Mne^{-n^2/2} \quad (23)$$

$$|R(\theta, d_{G_\epsilon}) - R(\theta, d_\epsilon)| \leq Mn^2e^{-n^2/2} \quad (24)$$

[Here and below  $M$  denotes a generic constant, not necessarily the same at each appearance.]

Figure 1 displays the approximate form of  $d_\epsilon$  and its risk function. Thus<sup>4</sup> any maxima of  $\theta \rightarrow R(\theta, d_\epsilon)$  with values greater than  $(1 - \delta)n^2$  occur only within intervals of the form  $[kn - \beta n, kn + \beta n]$  for  $k \geq 1$  and a small positive constant  $\beta$ . Set  $\gamma = \frac{1}{2} - \beta$ . Now replace consideration of estimator  $d_\epsilon$  on each of these intervals by a single two point prior Bayes rule: the inequality<sup>5</sup>

$$\sup \{|R(\theta, d_\epsilon) - R(\theta - k_0n, d_0)| : |\theta - k_0n| \leq \beta n\} \leq Mn\phi(\gamma n - a_n) \quad (25)$$

means we need only look at the maximum of the risk function

$$R_0(\zeta) = E_\zeta (d_0(z) - \zeta)^2$$

on the single interval  $|\zeta| \leq \beta n$ .

The logistic form (21) of  $d_0$  still precludes explicit evaluation of  $R_0(\zeta)$ , so consider as a further approximation the step function

$$d_{00}(z) = -nI(z < a_n), \quad (26)$$

with easily evaluated risk function and derivative

$$R_{00}(\zeta) = E_\zeta (d_{00}(z) - \zeta)^2 = (n^2 + 2n\zeta)\Phi(a_n - \zeta) + \zeta^2 \quad (27)$$

$$R'_{00}(\zeta) = 2n\Phi(a_n - \zeta) - (n^2 + 2n\zeta)\phi(a_n - \zeta) + 2\zeta. \quad (28)$$

For large  $n$ , Taylor expansions give the error in these approximations<sup>6</sup> :

$$R_0(\zeta) - R_{00}(\zeta) = -n\phi(a_n - \zeta)\left[1 + \frac{\pi^2}{6}\frac{a_n - \zeta}{n} + O(n^{-2}(a_n - \zeta)^2)\right] \quad (29)$$

$$R'_0(\zeta) - R'_{00}(\zeta) = -n(a_n - \zeta)\phi(a_n - \zeta)\left[1 + O(n^{-1}|a_n - \zeta|)\right], \quad (30)$$



valid uniformly in  $|\zeta| \leq \text{const}$ .

The aim now is to choose  $a_n$  so that the risk function  $R_0(\zeta)$  attains its maximum at a support point of the prior, namely 0. A sign change argument<sup>7</sup> shows that in fact  $R_0(\zeta)$  has a single local maximum. Together, (28) and (30) show that  $R'_0(0) \approx 0$  if  $a$  is required to satisfy the second defining equation (8). Further approximation<sup>8</sup> in (30) shows that for this choice of  $a_n$ ,

$$R'_0(-n^{-1}) \sim na_n\phi(a_n) > 0, \quad R'_0(0) \sim -\frac{1}{6}\pi^2 a_n^2 \phi(a_n) < 0. \quad (31)$$

Thus for large  $n$ , the maximum of  $R_0(\zeta)$  occurs at some point  $\xi^* \in (-n^{-1}, 0)$ . From (8) and (9) follows  $n\phi(a_n) \sim 2$ , and since  $R_0(\zeta)$  turns out to be concave on  $(-n^{-1}, 0)$ ,

$$R_0(\xi^*) - R_0(0) \leq n^{-1}|R'_0(0)| = O(a_n^2/n^2). \quad (32)$$

which is the order of error claimed in Theorem 1.

Together, (24), (25) and (32), imply that

$$\begin{aligned} \sup R(\theta, d_{G_\epsilon}) &= R_0(0) + O(n^2 e^{-n^2/2} + n\phi(\gamma n - a_n) - a_n^2/n^2), \\ R(kn, d_{G_\epsilon}) &= R_0(0) + O(n^2 e^{-n^2/2} + n\phi(\gamma n - a_n)), \quad k \geq 1. \end{aligned}$$

Thus (19) and (20) show that, up to terms of order  $O(\epsilon a_n^2/n^2)$ ,  $d_{G_\epsilon}$  is asymptotically minimax. This establishes (11), the first part of Theorem 1.

The approximation (10) to the minimax risk requires also the risk of  $d_{G_\epsilon}$  at  $\theta = 0$ . Replacing  $d_{G_\epsilon}$  by the two point prior  $\delta_0 + \epsilon\delta_n$ , yields an approximation<sup>9</sup>

$$R(0, d_{G_\epsilon}) = \epsilon n\phi(a_n) \left[ 1 + O(a_n^2/n^2) \right]. \quad (33)$$

Since  $r(G_\epsilon) \leq m(\epsilon) \leq m(d_{G_\epsilon}, \epsilon)$ , we now evaluate the minimax risk from (27), (29) and (33) as

$$m(\epsilon) = (1 - \epsilon)\epsilon n\phi(a_n) + \epsilon \left\{ n^2 \Phi(a_n) - n\phi(a_n) \left( 1 + \frac{\pi^2 a_n}{6n} \right) + O(\epsilon a_n^2/n^2) \right\},$$

which reduces to (10).

*Remark.* It does not seem likely that the geometric prior is exactly least favorable for  $\epsilon > 0$  in this setting – it would be necessary to choose the lattice spacing  $n$  so that the risk of  $d_{G_\epsilon}$  was constant at  $\theta = kn$  for  $k = 1, 2, \dots$ . However, the geometric prior is probably asymptotically minimax to higher orders: this might be shown using a three piece linear approximation to  $d_0$ , tangent to  $d_0$  at 0, in place of (26).

*Theorem 2: Asymptotics for  $\ell_1$  rule.* These are fairly straightforward, since the maximum risk of  $\theta \rightarrow R(\theta, d_\lambda)$  on  $[0, \infty)$  occurs<sup>10</sup> at  $\infty$  for any  $\lambda > 0$ , and equals  $1 + \lambda^2$ . Evaluating also  $R(0, d_\lambda)$ , we obtain

$$m(d_{\lambda_\epsilon}, \epsilon) = \inf_\lambda H(\lambda) = \inf_\lambda (1 - \epsilon) \left[ (\lambda^2 + 1)\tilde{\Phi}(\lambda) - \lambda\phi(\lambda) \right] + \epsilon(1 + \lambda^2). \quad (34)$$

The derivatives of  $H$  are easily found, and in particular show that  $H$  is convex. The form (16) of  $\lambda_\epsilon$  is obtained<sup>11</sup> by a one-step approximation to the initial value  $n_0^2 = 2 \log \epsilon^{-1}$  and substitution into  $H$  yields (17).

## 4 Discussion

The problem (1) studied in this paper is essentially identical to Mallows' problem involving two sided contamination. Problem (1) is technically simpler, since one doesn't have to deal with the effect of prior probability mass at negative atoms, but it seems likely that the techniques developed here would readily yield corresponding expansions and third order results for Mallows' model.

Most of the recent work on minimax properties of various models for sparsity (references below) has concentrated on first-order risk behavior, in part for reasons set out in Donoho, Johnstone, Kerkyacharian and Picard (1993). Two and three point asymptotically least favorable prior distributions arise commonly in this work. Thus, an additional contribution of this paper is to provide tools that might be adapted to study of higher order risk properties in these related settings.

*Sparsity and wavelet bases.* As preparation for studying the beneficial properties of wavelet bases in function estimation, Donoho and Johnstone (1992b,c) have investigated minimax estimation in Gaussian white noise of a mean vector known to lie in a finite dimensional  $l_p$  ball  $0 < p < \infty$ . Since  $\|\theta\|_{n,p}^p = \sum_1^n |\theta_i|^p \rightarrow \#\{i : |\theta_i| > 0\}$  as  $p \rightarrow 0$ , the 'nearly black' conditions studied here may be regarded in some sense as  $l_0$ -ball constraints.

Indeed, a characteristic property of the wavelet transform is that the wavelet coefficients of smooth or piecewise smooth functions are *typically* sparse – at higher resolution levels, only those coefficients in the vicinity of a discontinuity of the function or its derivatives are significantly non-zero.

Not surprisingly, therefore, the ideas and methods of each paper are related. Thus, an asymptotically least favorable distribution over  $\{\theta \geq 0 : \|\theta\|_{n,p} \leq r\}$  (as  $n \rightarrow \infty$  in model (2) ) is  $(1 - \epsilon)\delta_\mu + \epsilon\delta_0$ , where  $\epsilon$  and  $\mu$  are determined by  $\epsilon\mu^p = n^{-1}(r/\sigma)^p$  and equations (7) and (8) (with  $n$  replaced by  $\mu$ ). Indeed a two point prior of this form is all that is needed to establish first-order asymptotic minimaxity in the nearly black setting in Donoho et. al. (1992a). The  $l_p$ -balls ( and their weak analogs) with  $0 < p < 1$  arise naturally in the study of optimal spatially adaptive function estimates (such as, for example, variable kernel estimators – Donoho (1992), Johnstone (1993)).

*Related work.* In problem (2) the large sample ( $n \rightarrow \infty$ ) limit plus constraints led to a "restricted Bayes" minimax problem (1) for estimating a Gaussian mean subject to constraints on the class of prior distributions. (For more on the large sample limiting process, see Johnstone (1993)). A number of other such restricted minimax problems have been studied: as noted above, Donoho and Johnstone (1992b) study moment constraints, see also Feldman (1991). The limit as  $p \rightarrow \infty$  yields the case when the mean is known to be bounded in absolute value, by  $\eta$  say. See, e.g. Casella and Strawderman, (1981), Donoho, Liu and MacGibbon (1990). Casella and Strawderman gave the exact minimax value for  $\eta < 1.05$ , when a symmetric two point prior  $\frac{1}{2}\delta_\eta + \frac{1}{2}\delta_{-\eta}$  is exactly least favorable. Although this paper also studies situations in which "most of the mass is small", the situation here is of large, infrequent non-zero components, as compared to uniformly non-zero, but small components.

The relation between the “sparse” and “dense” settings becomes clearer in Donoho and Johnstone (1992b), where the  $L_p$  moment constraint is combined with an  $l_q$  loss function (here of course  $q = 2$ ). As  $\eta \rightarrow 0$ , the “sparse” situation arises when  $p < q$  (e.g.  $p = 0$ ), and the “dense” case when  $p \geq q$  (e.g.  $p = \infty$ ).

**Acknowledgements.** This work, begun in 1986, has been supported in part by NSF DMS 84-51750 and 92-09130, NIH PHS GM 21215-12 and the Sloan Foundation. My thanks go to David Donoho and Boris Levit for some helpful conversations, and to the referees and Associate Editor for their suggestions which improved an earlier version of the manuscript.

## References

- [1] Berger, J. (1982) Estimation in continuous exponential families: Bayesian estimation subject to risk restrictions and inadmissibility results. *Statistical Decision Theory and Related Topics III*. S.S. Gupta and J. Berger, eds. Academic, New York.
- [2] Bickel, P.J. (1983). Minimax estimation of the mean of a normal distribution subject to doing well at a point. *Recent Advances in Statistics* (Chernoff volume) 511–528. Academic Press.
- [3] Bickel, P.J. (1984) Parametric robustness : small biases can be worthwhile. *Ann. Statist.*, **12**. 864–879.
- [4] Brown, L.D. (1971) Admissible estimators, recurrent diffusions and insoluble boundary value problems. *Ann. Math. Statist.*, **42**, 855–903. Correction, *Ann. Statist.*, **1**, 594–596.
- [5] Casella, G. and Strawderman, W.E. Estimating a bounded normal mean. *Ann. Stat.* **9**, 870-878. (1981).
- [6] Donoho, D.L. (1992) De-Noising via Soft-Thresholding. Technical Report, Department of Statistics, Stanford University.
- [7] Donoho, D.L., Johnstone, I.M., Hoch, J.C. and Stern, A.S. (1992a). Maximum Entropy and the Nearly Black Object. *J. Roy. Stat. Soc. Ser. B* (with discussion) , **54**, 41 – 81.
- [8] Donoho, D. L. and Johnstone, I.M. (1992b) Minimax risk over  $l_p$  balls for  $l_q$  error. Technical Report No. 401., Department of Statistics, Stanford University.
- [9] Donoho, D. L. and Johnstone, I. M (1992c) Minimax Estimation via Wavelet shrinkage. Technical Report, Department of Statistics, Stanford University.
- [10] Donoho, D. L. and Johnstone, I. M (1992d) Ideal Spatial Adaptation via Wavelet Shrinkage. Technical Report, Department of Statistics, Stanford University. Tentatively accepted, *Biometrika*.

- [11] Donoho, D. L., Johnstone, I. M., Kerkyacharian, G. and Picard, D (1993) Wavelet Shrinkage: Asymptopia? Technical Report No. 419, Department of Statistics, Stanford University.
- [12] Donoho, D. L., Liu, R. C. and MacGibbon, K. B. Minimax risk over hyperrectangles, and implications. *Ann. Statist.*, **18**, 1416–1437. (1990).
- [13] Efron, B., and Morris C. (1971) Limiting the risk of Bayes and empirical Bayes estimates: Part I. The Bayes case. *J. Amer. Statist. Assoc.* **66** 807–815.
- [14] Feldman, I. (1991) Constrained minimax estimation of the mean of the normal distribution. *Ann. Statist.* **19**. 2259–2265.
- [15] B.R. Frieden. (1972) Restoring with Maximum Entropy II. Superresolution of photographs with diffraction-blurred impulses. *Journal of the Optical Society of America* **62**, 1202-1210,
- [16] S.F. Gull and G.J. Daniell. (1978) Image Reconstruction from incomplete and noisy data. *Nature* **272** 686-690.
- [17] Hodges, J.L. Jr. and Lehmann, E.L. (1952) The use of previous experience in reaching statistical decisions. *Ann. Math. Statist.* **23** 396–407.
- [18] Huber, P.J. (1964). Robust estimation of a location parameter. *Ann. Math. Statist.*, **35**, 73–101.
- [19] Huber, P.J. (1981). *Robust Statistics*. Wiley. New York.
- [20] Johnstone, I.M. (1993) Minimax Bayes, asymptotic minimax and sparse wavelet priors. To appear. *Proc. Fifth Purdue International Symposium on Decision Theory and Related Topics*.
- [21] Levit, B. Ya. (1986) On bounds for the minimax risk. *Prob. Theory and Math. Stat.* **2** 203–216. Prohorov et al., (eds) VNU Science Press.
- [22] Mallows, C.L. (1978) Problem 78–4. *SIAM Review*, 183.
- [23] Marazzi, A. (1980) Robust Bayesian estimation for the linear model. Technical Report No. 27 E.T.H. Zürich.
- [24] Stein, C. M. (1981) Estimation of the mean of a multivariate normal distribution. *Ann. Statist.* **9**. 1135-1151.

## 5 Appendix

1. *Note on terminology.* The use of the term “minimax risk” in the title of the paper derives from the motivating problem (4) of estimating a high dimensional vector with few non-zero components. The reduced form (1) poses a minimax problem in which Nature chooses from a restricted class of prior distributions on a single Gaussian mean. After noting connections with the restricted Bayes ideas of Hodges and

Lehmann (1952), Bickel (1983) refers to (1) as a restricted minimax problem. The fuller term restricted Bayes minimax emphasises that the payoff function  $r(G, d)$  involves expectation with respect to the prior distribution as well as the data distribution.

2. The derivation of approximations (13) and (14) involves three steps. First, define  $a^0 = (2 \log c_0 n)^{1/2}$  (compare (9)) and a one-step Newton approximation  $a^n$  to the solution of (8) starting from  $a^0$ . Second, express the right side of (10) in terms of  $n$  and  $a^0$ . Finally, obtain an expression in terms of  $n_0$  rather than  $n$  by exploiting (18) derived from (7). A sequence of approximations results, of which (13) and (14) are the crudest. The rather tedious details are omitted.

3. *Proof of (23) and (24)* Denote the unnormalised posterior  $\epsilon^k \phi(x - kn)$  by  $\phi_k = \phi_k(x)$ . Fix  $k_0 \geq 1$  and consider  $x \in [(k_0 - 1)n + a + n/2, k_0 n + a + n/2]$ . In this interval, setting  $z = x - k_0 n$ , we have

$$\begin{aligned} n^{-1}[d_{G_\epsilon}(x) - nk_0] &= \sum_{-k_0}^{\infty} j\phi_j / \sum_{-k_0}^{\infty} \phi_j, \\ n^{-1}[d_\epsilon(x) - nk_0] &= -\phi_{-1}/(\phi_0 + \phi_{-1}). \end{aligned}$$

We apply the equality

$$\frac{a + \delta_1}{b + \delta_2} = \frac{a}{b} + \left( \frac{\delta_1}{b} - \frac{a\delta_2}{b^2} \right) \left( 1 + \frac{\delta_2}{b} \right)^{-1} \quad (35)$$

for

$$\begin{aligned} a &= -\phi_{-1}, & b &= \phi_0 + \phi_1, \\ \delta_1 &= \sum_{j \neq 0, -1} j\phi_j, & \delta_2 &= \sum_{j \neq 0, -1} \phi_j. \end{aligned}$$

For  $|z - a| \leq n/2$ , the rapid decay of Gaussian tails ensures that  $\delta_1$  and  $\delta_2$  are at most a constant multiple of the leading terms  $\delta_3 = \phi_{-2} + \phi_1$ . Add the fact that  $|a/b| \leq 1$ , and it will be enough to show that  $\delta_3/b$  is  $O(e^{-n^2/2})$ . But simple algebra using (7), namely  $\log \epsilon = -n^2/2 - an$ , yields

$$\begin{aligned} \log \phi_1/\phi_0 &= (z - a)n - n^2 \leq -n^2/2 \\ \log \phi_{-2}/\phi_{-1} &= -(z - a)n - n^2 \leq -n^2/2. \end{aligned}$$

For  $k = 0$  (so that  $x = z$ ) and  $x \leq n/2 + a$ , apply a similar argument using  $a = \phi_1$  and  $b = \phi_0 + \phi_1$ . In this case,

$$\log \phi_2/\phi_0 = 2n(z - a_n) - 3n^2 \leq -2n^2$$

and so one obtains the stronger estimate

$$d_{G_\epsilon}(x) = d_\epsilon(x) + O(ne^{-2n^2}). \quad (36)$$

To establish (23), note the global bounds

$$|d_{G_\epsilon}(x) - x_+| \leq n, \quad |d_\epsilon(x) - x_+| \leq n, \quad (37)$$

where  $x_+ = x \vee 0$ . From (23),

$$\begin{aligned} R(\theta, d_{G_\epsilon}) - R(\theta, d_\epsilon) &\leq E_\theta |d_{G_\epsilon} - d_\epsilon| |d_{G_\epsilon} + d_\epsilon - 2\theta| \\ &\leq O(ne^{-n^2/2}) E_\theta [2n + 2|x_+ - \theta|] \\ &= O(n^2 e^{-n^2/2}). \end{aligned}$$

4. *Maxima of  $R(\theta, d_\epsilon)$*  The key to showing that  $R(\theta, d_\epsilon)$  has maxima only within intervals  $[kn - \beta n, kn + \beta n]$  is to show, for  $k_0 \geq 1$  and  $\theta - k_0 n \in [\beta n, (1 - \beta)n]$ , that  $d_\epsilon(x) - k_0 n$  is essentially zero over a somewhat larger range of  $x$ . More precisely, one can pick  $\gamma < \beta/2$  and note that  $d_\epsilon(x) - k_0 n$  is an odd function about  $x = k_0 n + a + n/2$  to verify that

$$\sup \{|d_\epsilon(x) - k_0 n| : x - k_0 n - a \in [\gamma n, (1 - \gamma)n]\} \leq e^{-\gamma n^2} \quad (38)$$

Using (38) on the set  $\{z = x - k_0 n \in [a + \gamma n, a + (1 - \gamma)n]\}$  and (37) plus bounds on Gaussian tails on the complement of the set leads to

$$E_\theta [d_\epsilon(x) - \theta]^2 \leq n^2 [(1 - \beta)^2 + M\tilde{\Phi}(\beta n/4)]. \quad (39)$$

To establish the claim, it remains to establish (39) also for  $\theta \leq n - \beta n$ . This is done analogously, but now using the inequality  $d_\epsilon(x) \leq n[1 + e^{n^2\gamma}]^{-1}$  valid for  $x \leq a + (1 - \gamma)n$ .

5. *Proof of (25)* The bound (25) is easily established after noting that  $d_\epsilon(k_0 n + z) = d_0(z)$  for  $|z - a| \leq n/2$ , so that (if  $\zeta = \theta - k_0 n \in [-\beta n, \beta n]$ ,  $z = x - k_0 n$ )

$$\begin{aligned} |R(\theta, d_\epsilon) - R(\zeta, d_0)| &\leq E_\theta \left\{ (d_\epsilon(x) - \theta)^2 + (d_0(x) - \theta)^2, |z - a| > n/2 \right\} \\ &\leq 2E_\zeta [n^2 + (z - \zeta)^2, |z - a| > n/2] \\ &\leq Mn\phi(\gamma n - a). \end{aligned}$$

6. *Proof of (29) and (30)* Suppose first that  $\gamma$  is an even function and set  $\gamma_k = \int_{-\infty}^{\infty} v^k \gamma(v) dv$ . Then if  $z \sim N(\zeta, 1)$ ,

$$\begin{aligned} E_\zeta \gamma [n(z - a_n)] &= n^{-1} \int_{-\infty}^{\infty} \gamma(v) \phi(a_n - \zeta + n^{-1}v) dv \\ &= \gamma_0 n^{-1} \phi(a_n - \zeta) + \frac{1}{2} \gamma_2 n^{-3} \phi''(a_n - \zeta) + \dots \end{aligned}$$

Using the derivatives  $\phi''(x) = (x^2 - 1)\phi(x)$  and  $\phi'''(x) = x(x^2 - 3)\phi(x)$ , and arguing analogously for odd functions, we obtain

**Lemma 3** *Suppose  $\int |v|^k \gamma(v) dv < \infty$  for  $k \leq 4$ , and that  $\gamma_n(w) = c_n \gamma(nw)$ . According as  $\gamma$  is even or odd*

$$\begin{aligned} E_\zeta \gamma_n(z - a_n) &= c_n n^{-1} \phi(a_n - \zeta) [\gamma_0 + O(n^{-2}[(a_n - \zeta)^2 - 1])], \\ E_\zeta \gamma_n(z - a_n) &= -c_n n^{-2} (a_n - \zeta) \phi(a_n - \zeta) [\gamma_1 + O(n^{-2}[(a_n - \zeta)^2 - 3])]. \end{aligned}$$

Before applying the lemma, we note the following relations between estimators  $d_0$ ,  $d_{00}$  and their risk functions:

$$d_0(z) = d_{00}(z) + \psi(z); \quad \psi(z) = n\tilde{\psi}(nw); \quad w = z - a_n$$

where  $\tilde{\psi}(w)$  is an *odd* function defined for  $w > 0$  by

$$\tilde{\psi}(w) = -e^{-w}/(1 + e^{-w}).$$

Using the risk (27) and its derivative (28) of  $d_{00}$ , and the identity  $\partial/\partial\zeta E_\zeta g(z) = E_\zeta(z - \zeta)g(z)$ :

$$R_0(\zeta) = R_{00}(\zeta) + E_\zeta [\psi^2 + 2\psi(d_{00} - \zeta)] \quad (40)$$

$$R'_0(\zeta) = R'_{00}(\zeta) + E_\zeta(z - \zeta) [\psi^2 + 2\psi(d_{00} - \zeta)] - 2E_\zeta\psi. \quad (41)$$

Applying Lemma 3 to the various components of (40) and (41) yields (showing leading terms only, and setting  $r = \phi(a_n - \zeta)$ )

$$\begin{aligned} r^{-1}E\psi^2 &= n\gamma_{a,0} + \dots \\ r^{-1}2E\psi d_{00} &= -n\gamma_{b,0} + (\zeta - a_n)\gamma_{b,1} + \dots \\ r^{-1}E\psi &= n^{-1}(\zeta - a_n)\gamma_{e,1} + \dots \\ r^{-1}E(z - a_n)\psi^2 &= n^{-1}(\zeta - a_n)\gamma_{d,1} + \dots \\ r^{-1}2E(z - a_n)\psi d_{00} &= -\gamma_{e,0} - n^{-1}(\zeta - a_n)\gamma_{ee,1} + \dots \\ r^{-1}E(z - a_n)\psi &= n^{-1}\gamma_{f,0} + \dots \end{aligned} \quad (42)$$

where

$$\begin{aligned} \gamma_{a,0} &= 2 \int_0^\infty [e^{-x}/(1 + e^{-x})]^2 dx = 2 \log 2 - 1 \\ \gamma_{b,0} &= 2 \int_0^\infty e^{-x}/(1 + e^{-x}) dx = 2 \log 2 \\ \gamma_{b,1} &= 2 \int_0^\infty xe^{-x}/(1 + e^{-x}) dx = \pi^2/6, \end{aligned}$$

and it is unnecessary to compute the remaining constants. Combining the terms in (42) in accordance with (40) and (41), and tracking the error terms provided by Lemma 3 leads to (29) and (30).

7. *Single maximum for  $R_0(\zeta)$*  Using Stein's unbiased estimate of risk (Stein, 1981), we may write (setting  $z = x - a$ ,  $\zeta = \eta - a$  and  $\gamma(x) = -n(1 + e^{nx})^{-1}$ ),

$$\begin{aligned} R_0(\zeta) - c &= E_\eta(\gamma(x) + \eta - a)^2 - c \\ &= E_\eta[(\gamma(x) + x - a)^2 - 2\gamma'(x) - (1 + c)]. \end{aligned}$$

For large values of  $n$  and  $a = a(n)$ , the integrand turns out to have at most four sign changes. Since the Gaussian kernel is totally positive of all orders,  $R_0(\zeta) - c$  can have at most four sign changes. Since  $R_0(\zeta) \nearrow \infty$  as  $|\zeta| \nearrow \infty$  and  $R_0$  dips down to values that are at most  $O(a_n^2)$  near  $\zeta = a + c \log n/n$  and  $\zeta = -n + a$ , it follows that  $R_0$  can have at most one maximum of order  $n^2$ .

8. *Proof of (31).* Here it is necessary to retain one extra term in the approximation to  $R'_0(\zeta)$  derived from (41) and (42):

$$R'_0(\zeta) = - \left[ n^2 + 2n\zeta + n(a_n - \zeta) + (a_n - \zeta)^2 \gamma_{b,1} + \gamma_{e,0} \right] \phi(a_n - \zeta) + 2n\Phi(a_n - \zeta) + 2\zeta + o\left((a_n - \zeta)^2 \phi(a_n - \zeta)\right).$$

This expansion is valid at least for  $|\zeta| \leq \text{constant}$ , and so (31) follows by substitution. Differentiation shows that

$$R''(\zeta) \sim -n^2 a_n \phi(a_n - \zeta) < 0$$

on  $[-n^{-1}, 0]$ , so that  $R_0(\zeta)$  lies below its tangent at  $\zeta = 0$ . This establishes the bound (32).

9. *Proof of (33).* Let  $d_1(x) = n\epsilon\phi(x - n)/[\phi(x) + \epsilon\phi(x - n)]$  be the Bayes rule for the two point prior  $\delta_0 + \epsilon\delta_n$ . Set  $\Delta = d_{G_\epsilon}^2 - d_1^2 = (d_{G_\epsilon} - d_1)(d_{G_\epsilon} + d_1)$ . Let  $I_0, I_1$  and  $I_2$  be the events that  $x$  lies in  $(-\infty, n/2 + a)$ ,  $[n/2 + a, 3n/2 + a]$  and  $(3n/2 + a, \infty)$  respectively. Using (23), (36), (37), bounds on the tail of the Gaussian distribution, and finally that  $|d_{G_\epsilon}(x)| \leq |x| + n$ ,  $|d_1(x)| \leq n$ ,

$$\begin{aligned} E_0[\Delta, I_0] &\leq Mne^{-2n^2} E_0[|x| + n, I_0] \leq Mn^2 e^{-2n^2} \\ E_0[\Delta, I_1] &\leq Mn^2 e^{-n^2/2} E_0[|x| + n, I_0^c] \leq Mn^2 e^{-n^2/2} \phi(n/2) \\ E_0[\Delta, I_2] &\leq ME_0[|x|^2 + n^2, I_2] \leq Mn\phi(3n/2). \end{aligned}$$

Using equation (7),  $\epsilon = e^{-n^2/2 - na_n}$ , we have

$$\begin{aligned} E_0 d_1^2(x) &= \int_{-\infty}^{\infty} \frac{n^2}{[1 + e^{n(n+a-x)}]^2} \phi(x) dx \\ &= n\epsilon\phi(a_n) \left[ \int \frac{e^w dw}{(1 + e^w)^2} + \frac{a_n}{n} \int \frac{we^w}{(1 + e^w)^2} dw + O\left(\frac{a_n}{n}\right)^2 \right], \end{aligned}$$

after making the substitution  $w = n(n + a - x)$ . This establishes (33) since the first integral equals one and the transformation  $x = e^w[1 + e^w]^{-1}$  reduces the second integral to  $\int_0^1 w(x) dx$ , which vanishes since  $w(x) = \log x/(1 - x)$  is odd about  $x = 1/2$ .

10. *Maximum risk of  $\ell_1$ -rule.* The  $\ell_1$ -rule  $d_\lambda(x) = (x - \lambda)_+$  may be written in the form  $x + \psi(x)$  with  $\psi(x) = -xI\{x \leq \lambda\} - \lambda I\{x > \lambda\}$ . Applying Stein's (1981) unbiased estimate of risk,

$$\begin{aligned} R(\theta, d_\lambda) &= 1 + E_\theta \left[ 2\psi'(x) + \psi^2(x) \right] \\ &= 1 + E_\theta \left[ -2I(x < \lambda) + x^2 I(x \leq \lambda) + \lambda^2 I(x > \lambda) \right] \end{aligned} \quad (43)$$

The integrand in (43) crosses any horizontal line at most twice, so it follows from the variation diminishing property of the Gaussian kernel that  $R(\theta, d_\lambda)$  attains its maximum on  $[0, \infty)$  at either 0 or  $+\infty$ .



11. *Minimum of  $H(\lambda)$ .* We record the following derivatives and approximations:

$$H'(\lambda) = 2(1 - \epsilon) [\lambda \tilde{\Phi}(\lambda) - \phi(\lambda)] + 2\lambda\epsilon \quad (44)$$

$$H''(\lambda) = 2(1 - \epsilon) \tilde{\Phi}(\lambda) + 2\epsilon \quad (45)$$

$$\lambda \tilde{\Phi}(\lambda) - \phi(\lambda) = -\lambda \int_{\lambda}^{\infty} x^{-2} \phi(x) dx = \lambda^{-2} \phi(\lambda) [1 + O(\lambda^{-2})] \quad (46)$$

$$(\lambda^2 + 1) \tilde{\Phi}(\lambda) - \lambda \phi(\lambda) = \int_{\lambda}^{\infty} x^{-2} (x^2 - \lambda^2) \phi(x) dx = 2\lambda^{-3} \phi(\lambda) [1 + O(\lambda^{-2})] \quad (47)$$

To approximately locate a zero of (44), set  $\lambda^2(a) = n_0^2 - 3 \log a n_0^2$ , where  $n_0^2 = 2 \log \epsilon^{-1}$ . The choices  $a_0 = (2\pi)^{1/3}$  and  $a_1 = a_0(1 - c n_0^{-2} \log n_0^2)$  (for  $c > 0$  large) lead to

$$H'(\lambda(a_0)) \sim -9\epsilon \lambda(a_0) n_0^{-2} \log a_0 n_0^2 < 0, \quad H'(\lambda(a_1)) \sim (c - 9)\epsilon n_0^{-2} \log a_0 n_0^2 > 0.$$

From (45),  $H$  is convex, and so  $\lambda_{\epsilon}$  is bracketed between  $\lambda(a_0)$  and  $\lambda(a_1)$ . Since  $\lambda^2(a_1) - \lambda^2(a_0) = 3 \log(a_0/a_1) = O(n_0^{-2} \log n_0)$ , (16) follows. From (45),  $H$  is convex, and since  $H'(\lambda(a_0)) < 0$ ,

$$\begin{aligned} 0 \leq H(\lambda(a_0)) - H(\lambda_{\epsilon}) &\leq (\lambda(a_0) - \lambda_{\epsilon}) H'(\lambda(a_0)) \\ &\leq |\lambda(a_0) - \lambda(a_1)| |H'(\lambda(a_0))| \\ &= O(\epsilon n_0^{-4} \log^2 n_0). \end{aligned} \quad (48)$$

Finally, using (34) and (47), and setting  $\lambda_1 = \lambda(a_0)$ ,

$$\begin{aligned} H(\lambda_1) &= (1 - \epsilon) \left\{ 2\lambda_1^{-3} \phi(\lambda_1) [1 + O(\lambda_1^{-2})] \right\} + \epsilon(1 + \lambda_1^2) \\ &= \epsilon \left[ \lambda_1^2 + 3 + 9n_0^{-2} \log(2\pi)^{1/3} n_0^2 + O(\lambda_1^{-2}) \right]. \end{aligned}$$

In view of error bound (48), this establishes (17) and completes the proof of Theorem 2. In fact, the bound (48) suggests that (17) could easily be improved by adding the next term in expansion (47).

Table 1: Various approximations to the minimax risk  $m(\epsilon)$ .  $n$  and  $a_n$  are the solutions to equations (7) and (8).  $n_0 = (2 \log \epsilon^{-1})^{1/2}$ . Equation (10) is the third-order approximation provided by Theorem 1, (15) is the first-order approximation, and (13), (14) are simpler, but less accurate, versions of (10).

$\epsilon$	$n$	$n_0$	$a_n$	$m(\epsilon)(10)$	$m(\epsilon)(15)$	$m(\epsilon)(13)$	$m(\epsilon)(14)$
7.5e-2	2.40	2.27	-0.12	2.00e-1	3.88e-1	NA	NA
5.0e-2	2.47	2.44	-0.02	1.50e-1	2.99e-1	NA	NA
4.0e-2	2.52	2.53	0.01	1.28e-1	2.57e-1	NA	NA
3.0e-2	2.58	2.64	0.06	1.04e-1	2.10e-1	NA	NA
2.0e-2	2.66	2.79	0.13	7.69e-2	1.56e-1	NA	NA
1.0e-2	2.81	3.03	0.23	4.52e-2	9.21e-2	NA	NA
5.0e-3	2.96	3.25	0.30	2.62e-2	5.29e-2	NA	NA
2.0e-3	3.15	3.52	0.39	1.25e-2	2.48e-2	NA	NA
1.0e-3	3.29	3.71	0.45	7.04e-3	1.38e-2	NA	NA
1.0e-4	3.73	4.29	0.59	9.78e-4	1.84e-3	NA	NA
1.0e-5	4.15	4.79	0.69	1.26e-4	2.30e-4	NA	NA
1.0e-6	4.53	5.25	0.77	1.57e-5	2.76e-5	1.39e-5	7.31e-6
1.0e-7	4.89	5.67	0.84	1.88e-6	3.22e-6	1.66e-6	1.51e-6
1.0e-8	5.23	6.06	0.90	2.19e-7	3.68e-7	1.94e-7	1.95e-7
1.0e-9	5.55	6.43	0.94	2.52e-8	4.14e-8	2.24e-8	2.32e-8
1.0e-10	5.86	6.78	0.98	2.84e-9	4.60e-9	2.54e-9	2.67e-9
1.0e-11	6.16	7.11	1.02	3.18e-10	5.06e-10	2.85e-10	3.02e-10
1.0e-12	6.45	7.43	1.05	3.51e-11	5.52e-11	3.17e-11	3.36e-11
1.0e-15	7.24	8.31	1.14	4.54e-14	6.90e-14	4.14e-14	4.40e-14
1.0e-18	7.98	9.10	1.20	5.60e-17	8.28e-17	5.16e-17	5.46e-17

Table 2: Numerically computed approximations to the minimax risk  $m(\epsilon)$  and risk of the best  $\ell_1$ -rule. Compare with asymptotic approximations in Table 1 (Col (10)) and Table 3 (Col (17)) respectively.  $n_{lo}(\epsilon)$  is the numerical approximation to  $n(\epsilon)$  corresponding to  $m_{lo}(\epsilon)$ .

$\epsilon$	$m_{lo}(\epsilon)$	$m_{up}(\epsilon)$	$m(d_{\lambda_\epsilon}, \epsilon)$	$n_{lo}(\epsilon)$	$\theta_{max}$
0.2000	0.3907	0.4100	0.4100	2.4781	4.2655
0.1000	0.2481	0.2604	0.2600	2.5890	4.6248
0.0500	0.1534	0.1602	0.1670	2.7018	4.9684
0.0200	0.0784	0.0811	0.0873	2.8564	5.4028
0.0100	0.0461	0.0474	0.0522	2.9776	5.7182
0.0050	0.0267	0.0272	0.0306	3.1012	6.0227
0.0020	0.0127	0.0128	0.0148	3.2662	6.4096
0.0010	0.0071	0.0072	0.0084	3.3911	6.6913

Table 3: Maximum risk over  $\epsilon$ -black objects for the best  $\ell_1$ -rule.  $\lambda_\epsilon$  and  $m(d_{\lambda_\epsilon}, \epsilon)$  are the approximations given in Theorem 2 by (16) and (17) respectively.  $\lambda_t$  and  $m(d_{\lambda_t}, \epsilon)$  denote the numerically determined optimal values ( $\lambda_t$  is the root of (44) and  $m(d_{\lambda_t}, \epsilon) = H(\lambda_t)$ .) For comparison,  $m(\epsilon)$  given by (10), is the third order approximation to minimax risk.

$\epsilon$	$\lambda_\epsilon$	$\lambda_t$	$m(d_{\lambda_\epsilon}, \epsilon)$	$m(d_{\lambda_t}, \epsilon)$	$m(\epsilon)$
7.5e-2	NA	1.00	NA	0.219	0.195
5.0e-2	NA	1.15	NA	0.167	0.148
4.0e-2	NA	1.23	NA	0.143	0.126
3.0e-2	NA	1.34	NA	0.116	0.103
2.0e-2	NA	1.48	NA	8.73e-02	7.65e-2
1.0e-2	0.84	1.72	5.87e-2	5.22e-02	4.51e-2
5.0e-3	1.29	1.94	3.33e-2	3.06e-02	2.62e-2
2.0e-3	1.74	2.23	1.56e-2	1.48e-02	1.25e-2
1.0e-3	2.02	2.43	8.81e-3	8.41e-03	7.04e-3
5.0e-4	2.28	2.63	4.89e-3	4.72e-03	3.92e-3
2.0e-4	2.58	2.88	2.22e-3	2.17e-03	1.79e-3
1.0e-4	2.80	3.06	1.22e-3	1.19e-03	9.78e-4
5.0e-5	3.00	3.23	6.67e-4	6.53e-04	5.32e-4