

## BIAS-VARIANCE TRADEOFFS IN FUNCTIONAL ESTIMATION PROBLEMS<sup>1</sup>

BY MARK G. LOW

*University of Pennsylvania*

It is shown in infinite-dimensional Gaussian problems that affine estimators minimax the variance among all estimators of a linear functional subject to a constraint on the bias. Likewise, affine estimators also minimax the square of the bias among all estimates of a linear functional subject to a constraint on the variance.

**1. Introduction.** In nonparametric functional estimation problems, estimators which achieve good mean squared error performance usually balance bias and variance. In particular, Liu and Brown (1993) have shown that in many such singular estimation problems optimal mean squared error rates of convergence can only be attained by estimators which have the same convergence rate for both the square of the bias and the variance. For example, in these problems estimators cannot be found which attain an optimal rate of convergence for the mean squared error and have a faster rate of convergence for the square of the bias.

In the nonparametric functional estimation literature, infinite-dimensional Gaussian experiments play a central role. They capture many of the essential features of other models, such as density estimation, without as many technical difficulties. In particular, much attention has focused on problems of estimating linear functionals  $L(f)$  based on observing data  $Y$  of the form

$$(1.1) \quad Y(t) = \int_0^t Kf(s) ds + \sigma W(t), \quad 0 \leq t \leq 1,$$

where  $W(t)$  is Brownian motion,  $K$  is a linear map and  $f \in \mathcal{F}$  a convex class of functions. Donoho and Liu (1991), Sacks and Ylvisaker (1978) and Ibragimov and Hasminskii (1984) are just a sampling of the many papers which have studied essentially this same model with a variety of assumptions on both  $K$  and  $\mathcal{F}$ . In this paper we give a fairly complete analysis of the possible tradeoffs between bias and variance for problems of estimating linear functionals based on such infinite-dimensional Gaussian experiments. In particular, in these models we quantify the essentially qualitative results of Liu and Brown (1993).

---

Received July 1992; revised August 1994.

<sup>1</sup> Supported by an NSF Postdoctoral Fellowship.

1991 *subject classifications*. Primary 62C05; secondary 62E20, 62J02, 62G05, 62M99

*Key words and phrases*. Bias-variance tradeoff, Cramér-Rao inequality, minimax risk, white noise model, modulus of continuity.

The analysis relies heavily on one-dimensional subfamily theorems found in Donoho and Liu (1991) and Donoho (1994). In particular, Donoho (1994) gave a very general treatment of the Gaussian model (1.1) with  $K$  a general linear map and  $\mathcal{F}$  an arbitrary convex set under a variety of measures of performance including mean squared error (MSE), mean absolute error (MAE) and the length of confidence intervals with some minimal level of coverage probability.

It is well known that for any estimator the mean squared error can be written as the sum of the square of the bias and variance. Such a decomposition does not in general hold for the other measures unless attention is restricted to affine estimators. Then writing  $\text{Var}$  for the variance of an estimator, the mean absolute error for affine procedures can schematically be expressed as

$$(1.2) \quad \text{MAE} = \Phi(\text{Bias}, \text{Var}) = E|\text{Bias} + \sqrt{\text{Var}} Z|$$

and likewise the coverage probability of a fixed length confidence interval as

$$(1.3) \quad \alpha = \Phi(\text{Bias}, \text{Var}) = P\{|\text{Bias} + \sqrt{\text{Var}} Z| > C\},$$

where in both cases  $Z$  is a standard normal distribution. When these performance measures are increasing functions of both the absolute value of the bias and the variance, minimax theorems given in Donoho and Liu (1991) and Donoho (1994) can be exploited. These theorems break the study of estimating linear functions by affine procedures into two parts. For a linear functional  $L$ , linear map  $K$ , parameter space  $\mathcal{F}$  and a given performance measure, first find a hardest one-dimensional subfamily. Then find optimal affine estimators over these subfamilies. The minimax theorems show that these affine estimators are in fact optimal over the whole parameter space  $\mathcal{F}$  within the class of all affine estimators.

The subproblems can be identified by a modulus of continuity  $\omega(\varepsilon, L, K, \mathcal{F})$ . Write  $\|f\|_2$  for the  $L_2$  norm of a function,  $\|f\|_2^2 = \int f^2$ . Then the modulus can be written

$$(1.4) \quad \omega(\varepsilon, L, K, \mathcal{F}) = \sup\{|L(f_1) - L(f_{-1})| : \|Kf_1 - Kf_{-1}\|_2 \leq \varepsilon, f_i \in \mathcal{F}\}.$$

For each  $\varepsilon$  the affine family with endpoints  $f_1$  and  $f_{-1}$  attaining the supremum in (1.4) is a subfamily of  $\mathcal{F}$ . For each measure of loss, this subfamily is hardest for some particular noise level  $\sigma$  in the Gaussian model (1.1). By a Rao-Blackwell sufficiency argument, optimal affine procedures over these hardest subfamilies can be found by analyzing bounded normal mean problems. As just mentioned, these optimal procedures are then minimax over the whole parameter space  $\mathcal{F}$  within the class of all affine estimators. Typically, however, they are not minimax within the class of all measurable estimators. See, for example, Sacks and Strawderman (1982).

In this paper we consider a new optimization problem where affine procedures are in fact minimax within the class of all measurable procedures. Find procedures which minimax the variance subject to a constraint on the bias. Likewise find estimators which minimax the square of the bias among all

estimators with a constraint on the variance. These optimization problems can be written in terms of bias and variance as

$$(1.5) \quad V_B(\text{Bias}, \text{Var}) = \begin{cases} \infty, & \text{when } |\text{Bias}| > B, \\ \text{Var}, & \text{when } |\text{Bias}| \leq B, \end{cases}$$

and

$$(1.6) \quad B_V^2(\text{Bias}, \text{Var}) = \begin{cases} \infty, & \text{when } |\text{Var}| > V, \\ \text{Bias}^2, & \text{when } |\text{Var}| \leq V. \end{cases}$$

Just as for the other performance measures, the analysis can be broken into two parts. In Section 2 we give a complete analysis of the possible bias–variance tradeoffs possible in the one-dimensional normal problem. It relies heavily on the use of the Cramér–Rao inequality to bound these tradeoffs. See Hall (1989) and Brown and Farrell (1990) for similar arguments used to bound the mean squared error. This analysis is carried over to the infinite-dimensional setting in Section 3. It allows for a new optimality interpretation of some well known linear estimators such as those of Sacks and Ylvisaker (1981) and Epanechnikov (1969).

**2. Bounded normal mean problem.** As mentioned in the Introduction, the problem of estimating linear functionals based on the infinite-dimensional Gaussian experiments (1.1) with a performance given by an increasing function of the absolute bias and variance can be reduced to the study of hardest one-dimensional subfamilies. By a Rao–Blackwell sufficiency argument these one-dimensional subfamilies are equivalent to estimating the mean  $\theta$  of a normal distribution when  $\theta$  is known to lie in some closed interval. In this section we concentrate on the bias–variance tradeoff problem for this one-dimensional Gaussian experiment. The connection to the infinite-dimensional problem is made in the next section.

Suppose that  $X \sim N(\theta, \sigma^2)$ . Write  $\text{Bias}_\sigma(\delta(X), \theta) = E\delta(X) - \theta$  for the bias of an estimator  $\delta(X)$  when the noise level is  $\sigma$  and the mean is  $\theta$ . Likewise write  $\text{Var}_\sigma(\delta(X), \theta) = E(\delta(X) - E\delta(X))^2$  for the variance of the estimator  $\delta(X)$  at the same noise level  $\sigma$  and mean  $\theta$ . For notational convenience we shall also write  $\text{Var}_\sigma \delta(X)$  for  $\sup_{|\theta - \theta_0| \leq \tau} \text{Var}_\sigma(\delta(X), \theta)$ . Likewise write  $|\text{Bias}_\sigma(\delta(X))|$  for  $\sup_{|\theta - \theta_0| \leq \tau} |\text{Bias}_\sigma(\delta(X), \theta)|$ . If it is known that the mean  $\theta$  lies in an interval, say  $|\theta - \theta_0| \leq \tau$ , then the minimax mean squared error procedure for estimating  $\theta$  is not an affine procedure. For small values of  $\tau/\sigma$ , Casella and Strawderman (1981) in fact found the minimax procedure. In general the minimax mean squared error over the class of all affine procedures is only a small multiple of the minimax mean squared error over the class of all measurable procedures. See, for example, Donoho and Liu (1991) and Brown and Feldman (1989).

In this context the study of the possible bias–variance tradeoffs becomes that of finding procedures which minimax the variance subject to a constraint on the bias and likewise to find procedures which minimax the square of the

bias subject to a constraint on the variance. These optimization problems can be written as

$$(2.1) \quad \beta^2(v, \sigma, \tau) = \inf_{\text{Var}_\sigma(\delta(X)) \leq v} \sup_{|\theta - \theta_0| \leq \tau} \text{Bias}_\sigma^2(\delta(X), \theta)$$

and

$$(2.2) \quad v(\beta, \sigma, \tau) = \inf_{|\text{Bias}_\sigma \delta(X)| \leq \beta} \sup_{|\theta - \theta_0| \leq \tau} \text{Var}_\sigma(\delta(X), \theta).$$

The following theorem then shows that from this point of view there are affine estimators which are minimax over the class of all measurable procedures. Moreover, such minimax estimators are essentially unique.

**THEOREM 1.** *If  $X \sim N(\theta, \sigma^2)$ , then*

$$(2.3) \quad \beta^2(v, \sigma, \tau) = \left( \left( \frac{\sqrt{v}}{\sigma} \wedge 1 \right) - 1 \right)^2 \tau^2$$

and the affine procedure

$$(2.4) \quad \delta_v(X) = \left( \frac{\sqrt{v}}{\sigma} \wedge 1 \right) (x - \theta_0) + \theta_0$$

is essentially the unique procedure satisfying

$$(2.5) \quad \sup_{|\theta - \theta_0| \leq \tau} \text{Var}(\delta_v(X), \theta) \leq v \wedge \sigma^2$$

and

$$(2.6) \quad \beta^2(v, \sigma, \tau) = \sup_{|\theta - \theta_0| \leq \tau} \text{Bias}_\sigma^2(\delta_v(X), \theta).$$

Likewise

$$(2.7) \quad v(\beta, \sigma, \tau) = \left( \frac{\sigma}{\tau} \right)^2 ([\tau - \beta]_+)^2$$

and the affine procedure

$$(2.8) \quad \delta_v(X) = \frac{\sqrt{v(\beta, \sigma, \tau)}}{\sigma} (X - \theta_0) + \theta_0$$

is essentially the unique procedure satisfying

$$(2.9) \quad \sup_{|\theta - \theta_0| \leq \tau} \text{Bias}_\sigma^2(\delta_v(X), \theta) \leq \beta^2 \wedge \tau^2$$

and

$$(2.10) \quad v(\beta, \sigma, \tau) = \sup_{|\theta - \theta_0| \leq \tau} \text{Var}(\delta_v(X), \theta).$$

**PROOF.** We shall only prove (2.3), (2.5) and (2.6) as the proof of (2.7), (2.9) and (2.10) is essentially the same. Note that if  $v \geq \sigma^2$ , then the unbiased

estimator  $\delta(X) = X$  is essentially the unique unbiased estimator with variance  $\text{Var}_\sigma \delta(X) \leq \sigma^2$  and the theorem trivially holds. So assume that  $v \leq \sigma^2$  and note that  $\sqrt{v}/\sigma \wedge 1 = \sqrt{v}/\sigma$ .

The Fisher information for  $\theta$  based on  $X$  is given by  $I(\theta)$ , where

$$(2.11) \quad I(\theta) = \frac{1}{\sigma^2}, \quad |\theta - \theta_0| \leq \tau.$$

For simplicity rewrite the bias of an estimator as

$$(2.12) \quad b(\theta) = E\delta(X) - \theta.$$

The Cramér–Rao inequality then yields

$$(2.13) \quad v \geq \text{Var } \delta \geq \frac{(1 + b'(\theta))^2}{1/\sigma^2}, \quad |\theta - \theta_0| \leq \tau,$$

which can be rewritten as

$$(2.14) \quad b'(\theta) \leq \frac{\sqrt{v}}{\sigma} - 1, \quad |\theta - \theta_0| \leq \tau.$$

Hence

$$(2.15) \quad b(\theta_0 + \tau) - b(\theta_0 - \tau) \leq \left( \frac{\sqrt{v}}{\sigma} - 1 \right) 2\tau$$

and it follows that

$$(2.16) \quad \max(b^2(\theta_0 + \tau), b^2(\theta_0 - \tau)) \geq \left( \frac{\sqrt{v}}{\sigma} - 1 \right)^2 \tau^2.$$

Equation (2.3) is then an immediate consequence of (2.16).

The proof of (2.5) and (2.6) is a straight calculation which we leave to the reader. We now prove the essential uniqueness of  $\delta_v(X)$  when  $v \leq \sigma^2$ . Suppose in this case that  $\delta(X)$  and  $\tilde{\delta}(X)$  satisfy (2.5) and (2.6). It then follows that (2.14)–(2.16) must hold with equality for both  $\delta(X)$  and  $\tilde{\delta}(X)$ . In particular, the bias functions of these two estimators must be equal and thus  $E_\theta \delta(X) = E_\theta \tilde{\delta}(X)$  for all  $\theta$ . It is then easy to check that  $\varphi(X) = \frac{1}{2}(\delta(X) + \tilde{\delta}(X))$  also satisfies (2.5) and (2.6). By (2.7) it follows that

$$(2.17) \quad \sup_{|\theta - \theta_0| \leq \tau} \text{Var}_\sigma(\varphi(X), \theta) \geq \left( \frac{\sigma}{\tau} \right)^2 \left( \tau - \left( 1 - \frac{\sqrt{v}}{\sigma} \right) \tau \right)^2 = v.$$

Moreover, simple computations show that

$$(2.18) \quad \begin{aligned} \text{Var } \varphi(X) &= \frac{1}{4}(\text{Var } \delta(X) + \text{Var } \tilde{\delta}(X) + 2 \text{Cov}(\delta(X), \tilde{\delta}(X))) \\ &\leq \frac{1}{4}(\text{Var } \delta(X) + \text{Var } \tilde{\delta}(X) + 2[\text{Var } \delta(X) \text{Var } \tilde{\delta}(X)]^{1/2}) \\ &\leq \frac{1}{4}(v + v + 2(vv)^{1/2}) \\ &= v, \end{aligned}$$

where we have suppressed the dependence of  $\text{Var } \varphi(X)$  on  $\theta$  and  $\sigma$ . Hence it follows from (2.17) and (2.18) that

$$(2.19) \quad \sup_{|\theta - \theta_0| \leq \tau} \text{Var } \varphi(X) = v.$$

Since the set  $\{\theta: |\theta - \theta_0| \leq \tau\}$  is compact, there is a point  $\theta_c$  such that  $|\theta_c - \theta_0| \leq \tau$  and

$$(2.20) \quad \text{Var}(\varphi(X), \theta_c) = v.$$

Hence at this point from (2.18) and (2.19),  $\text{Cov}(\delta(X), \tilde{\delta}(X)) = (\text{Var } \delta(X) \text{Var } \tilde{\delta}(X))^{1/2}$  and since as mentioned earlier  $E\delta(X) = E\tilde{\delta}(X)$ , it follows that  $\delta = \tilde{\delta}$  almost surely.  $\square$

**3. Infinite-dimensional Gaussian models.** We now return to the infinite-dimensional Gaussian experiment given in (1.1). This experiment contains a large variety of statistical models including nonparametric regression, semiparametric models, inverse problems and white noise models. In these models consistent estimators of the linear functional  $L$  can only exist as the noise level  $\sigma \rightarrow 0$  if the modulus of continuity  $\omega(\varepsilon, L, I, \mathcal{F}) \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , where the modulus  $\omega$  is defined by (1.4) and  $I$  is the identity operator. Following Donoho (1994) it is natural to call such functionals well defined since it is possible to extend these functionals to the  $L_2$  completion of  $\mathcal{F}$  in the following way. Write  $\bar{\mathcal{F}}$  for the closure of  $\mathcal{F}$  in the  $L_2$  norm. Then for  $f \in \bar{\mathcal{F}}$ , let  $f_n \in \mathcal{F}$ ,  $n = 1, 2, \dots$ , be a sequence such that  $\lim_{n \rightarrow \infty} \|f_n - f\|_2 = 0$ . Define  $Lf$  by  $Lf = \lim_{n \rightarrow \infty} Lf_n$ . It is easy to check that this limit exists and hence  $Lf$  is well defined as long as  $\omega(\varepsilon, L, I, \mathcal{F}) \rightarrow 0$  when  $\varepsilon \rightarrow 0$ .

Similarly, it is natural to call the linear map  $K$  well defined whenever  $\lim_{\varepsilon \rightarrow 0} \sup\{\|Kf - Kg\|_2: \|f - g\|_2 \leq \varepsilon, f \in \mathcal{F}, g \in \mathcal{F}\} = 0$ . Then the bias and variance of any estimator are continuous in the  $L_2$  topology of the parameter space. If  $L$  and  $K$  are well defined, then for any estimator  $\hat{L}$  the supremum of the bias and variance over  $\mathcal{F}$  is equal to the supremum over the closure  $\bar{\mathcal{F}}$  of  $\mathcal{F}$ . The results of Theorem 1 together with hardest one-dimensional subfamily arguments can then yield the possible tradeoffs of bias and variance available in these problems.

For any linear functional  $L$ , estimator  $\hat{L}$  and parameter  $f \in \mathcal{F}$ , write  $\text{Bias}_\sigma(\hat{L}, L, f)$  for the bias  $E\hat{L} - Lf$  when the noise level is  $\sigma$ . Likewise write  $\text{Var}_\sigma(\hat{L}, f)$  for the variance  $E(\hat{L} - E\hat{L})^2$  when  $f$  is the true parameter and  $\sigma$  is the noise level. As in Section 2 write  $\text{Var}_\sigma(\hat{L})$  for  $\sup_{\mathcal{F}} \text{Var}_\sigma(\hat{L}, f)$  and  $|\text{Bias}_\sigma(\hat{L})|$  for  $\sup_{\mathcal{F}} |\text{Bias}_\sigma(\hat{L}, L, f)|$ .

Then the bias–variance optimization problem can be written

$$(3.1) \quad B^2(V, \sigma, L, \mathcal{F}) = \inf_{\text{Var}_\sigma(\hat{L}) \leq V} \sup_{\mathcal{F}} \text{Bias}_\sigma^2(\hat{L}, L, f)$$

and

$$(3.2) \quad V(B, \sigma, L, \mathcal{F}) = \inf_{|\text{Bias}_\sigma(\hat{L})| \leq B} \sup_{\mathcal{F}} \text{Var}_\sigma(\hat{L}, f).$$

When both  $K$  and  $L$  are well defined, we may by the remarks in the previous paragraph restrict attention to parameter spaces which are closed. The following theorem then gives affine estimators which achieve the optimal trades of bias and variance.

**THEOREM 2.** *Suppose that  $L$  and  $K$  are well defined, that  $\omega(\varepsilon, L, K, \mathcal{F}) \rightarrow 0$  as  $\varepsilon \rightarrow 0$  and that  $\mathcal{F}$  is closed and convex. Then the minimum square of the bias over the class of estimators with variance bounded by  $V \geq 0$  is given by*

$$(3.3) \quad B^2(V, \sigma, L, \mathcal{F}) = 4^{-1} \sup_{\varepsilon > 0} \left( \left[ \omega(\varepsilon, L, K, \mathcal{F}) - \varepsilon\sqrt{V}/\sigma \right]_+ \right)^2.$$

*Likewise the minimum variance over the class of estimators with the absolute value of the bias bounded by  $B \geq 0$  is given by*

$$(3.4) \quad V(B, \sigma, L, \mathcal{F}) = \sup_{\varepsilon > 0} (\sigma/\varepsilon)^2 \left( \left[ \omega(\varepsilon, L, K, \mathcal{F}) - 2B \right]_+ \right)^2.$$

*Affine estimators which yield this optimal tradeoff of bias and variance are given as follows. Suppose that  $\varepsilon_V = \arg \max_{\varepsilon \geq 0} (\left[ \omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{V}\varepsilon/\sigma \right]_+)$  exists such that  $0 < \varepsilon_V < \infty$ . Then there are  $f_{-1} \in \mathcal{F}$ ,  $f_1 \in \mathcal{F}$ , such that  $\|Kf_1 - Kf_{-1}\|_2 = \varepsilon_V$  and  $Lf_1 - Lf_{-1} = \omega(\varepsilon_V, L, K, \mathcal{F})$ . Write  $Kf_0 \equiv \frac{1}{2}(Kf_1 + Kf_{-1})$  for the center and  $Ku \equiv (Kf_1 - Kf_{-1})/\varepsilon_V$  for the direction of the affine family joining  $Kf_{-1}$  and  $Kf_1$ . Then the estimator*

$$(3.5) \quad \hat{L}_V = Lf_0 + \frac{\sqrt{V}}{\sigma} \int Ku(t)(Y(dt) - Kf_0(t) dt)$$

*has constant variance*

$$(3.6) \quad \text{Var}_\sigma(\hat{L}_V, f) = V$$

*and maximum bias*

$$(3.7) \quad \sup_{\mathcal{F}} \text{Bias}_\sigma^2(\hat{L}_V, L, f) = B^2(V, \sigma, L, \mathcal{F}).$$

*Writing  $V$  for  $V(B, \sigma, L, \mathcal{F})$  given in (3.4) and once again assuming that  $\varepsilon_V = \arg \max_{\varepsilon \geq 0} (\left[ \omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{V}\varepsilon/\sigma \right]_+)$  exists such that  $0 < \varepsilon_V < \infty$ , then the estimator  $\hat{L}_V$  defined by (3.5) satisfies*

$$(3.8) \quad \sup_{\mathcal{F}} \text{Bias}_\sigma^2(\hat{L}_V, L, f) = B^2(V, \sigma, L, \mathcal{F})$$

*and*

$$(3.9) \quad \text{Var}_\sigma(\hat{L}_V, f) = V.$$

**PROOF.** The proof of (3.4), (3.8) and (3.9) is essentially the same as that of (3.3), (3.6) and (3.7); hence, we shall only prove the latter. Moreover we shall only prove (3.3) for the case when  $\varepsilon_V$  exists, although we do not assume that there is such an  $\varepsilon_V = \arg \max_{\varepsilon \geq 0} (\left[ \omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{V}\varepsilon/\sigma \right]_+)$  which is strictly positive. The extension to the general case follows by the same approximation arguments used to prove Theorem 2 of Donoho (1994).

First assume that (3.3), (3.6) and (3.7) hold whenever there exists an  $\varepsilon_V = \arg \max_{\varepsilon \geq 0} ([\omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{V}\varepsilon/\sigma]_+)$  such that  $0 < \varepsilon_V < \infty$ . We shall then show that (3.3) holds when  $\sup_{\varepsilon} (\omega(\varepsilon) - \sqrt{V}\varepsilon/\sigma) = 0$  and there does not exist an  $\varepsilon_V = \arg \max_{\varepsilon \geq 0} ([\omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{V}\varepsilon/\sigma]_+)$  such that  $\varepsilon_V > 0$ . In that case, let

$$(3.10) \quad \limsup_{\varepsilon \rightarrow 0} \frac{\omega(\varepsilon, L, I, \mathcal{F})}{\varepsilon} = a.$$

It then follows that  $V \geq a^2$  and also that  $\omega(\varepsilon) < \sqrt{V}\varepsilon/\sigma$  for all  $\varepsilon > 0$ . Now let  $\nu_n \geq 0$  be an increasing sequence where  $\nu_n \rightarrow a^2$ . It is then easy to check that for each  $n$ ,  $\varepsilon_{\nu_n} = \arg \max_{\varepsilon \geq 0} ([\omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{\nu_n}\varepsilon/\sigma]_+)$  exists with  $0 < \varepsilon_{\nu_n} < \infty$ . Moreover

$$(3.11) \quad \lim_{n \rightarrow \infty} B^2(\nu_n, \sigma, L, \mathcal{F}) = 0$$

and

$$(3.12) \quad \lim_{n \rightarrow \infty} \text{Var } \hat{L}_{\nu_n} = a^2 \leq V.$$

Hence (3.3) once again holds.

The proof will be complete if we show that (3.3), (3.6) and (3.7) hold whenever there exists an  $\varepsilon_V = \arg \max_{\varepsilon \geq 0} (\omega(\varepsilon, L, K, \mathcal{F}) - \sqrt{V}\varepsilon/\sigma)_+$  such that  $0 < \varepsilon_V < \infty$ . For such an  $\varepsilon_V$  the existence of functions  $f_1$  and  $f_{-1}$  with  $\|Kf_1 - Kf_{-1}\|_2 = \varepsilon_V$  follows from the fact that the parameter space is closed. It then follows that the estimator  $\hat{L}_V$  given by (3.5) is well defined.

We now show that this  $\hat{L}_V$  satisfies (3.6) and (3.7) and also that equation (3.3) holds. This follows essentially from Theorem 1 plus hardest one-dimensional subfamily theorems of Donoho (1994) and Donoho and Liu (1991). The hardest one-dimensional subfamily theorems show that the maximum absolute bias over  $\mathcal{F}$  of the affine estimator  $\hat{L}_V$  is attained at the endpoints  $f_{-1}$  and  $f_1$  of the affine family joining  $f_{-1}$  and  $f_1$ . This can also be easily checked as follows. Let  $g$  be any other element of  $\mathcal{F}$ . The affine family joining  $f_1$  and  $g$  is given by  $(1 - \theta)f_1 + \theta g$ ,  $0 \leq \theta \leq 1$ . Let

$$(3.13) \quad J(\theta) = \left( L((1 - \theta)f_1 + \theta g) - Lf_{-1} - \frac{\sqrt{V}}{\sigma} \|K((1 - \theta)f_1 + \theta g) - Kf_{-1}\|_2 \right).$$

Since  $f_1$  and  $f_{-1}$  are the extremal functions attaining the supremum on the right-hand side of (3.3) and  $Lf_1 - Lf_{-1} > 0$ , it follows that  $J'(0) \leq 0$  and hence by a simple computation that

$$(3.14) \quad Lg - Lf_1 - \frac{\sqrt{V}}{\sigma} \int Ku(t)(Kg(t) - Kf_1(t)) dt \leq 0.$$

Now

$$(3.15) \quad \text{Bias}_{\sigma}(\hat{L}_V, L, f_1) = Lf_0 + \frac{\sqrt{V}}{\sigma} \int Ku(t)(Kf_1(t) - Kf_0(t)) dt - Lf_1$$



and

$$(3.16) \quad \text{Bias}_\sigma(\hat{L}_V, L, g) = Lf_0 + \frac{\sqrt{V}}{\sigma} \int Ku(t)(Kg(t) - Kf_0(t)) dt - Lg.$$

It then follows from (3.14)–(3.16) that

$$(3.17) \quad \text{Bias}_\sigma(\hat{L}_V, L, f_1) - \text{Bias}_\sigma(\hat{L}_V, L, g) \leq 0.$$

Likewise it is easy to show that

$$(3.18) \quad \text{Bias}_\sigma(\hat{L}_V, L, f_{-1}) - \text{Bias}_\sigma(\hat{L}_V, L, g) \geq 0.$$

It then follows that the maximum absolute bias over  $\mathcal{F}$  of the estimator  $\hat{L}_V$  is attained at the parameter points  $f_1$  and  $f_{-1}$ .

Now a simple computation shows that

$$(3.19) \quad \text{Var}(\hat{L}_V, f) = V$$

and that

$$(3.20) \quad \text{Bias}_\sigma^2(\hat{L}_V, L, f_1) = 4^{-1} \sup_{\varepsilon > 0} \left( \left[ \omega(\varepsilon, L, K, \mathcal{F}) - \frac{\varepsilon\sqrt{V}}{\sigma} \right]_+ \right)^2.$$

Hence to finish the proof we need only show that

$$(3.21) \quad B^2(V, \sigma, L, \mathcal{F}) \geq 4^{-1} \sup_{\varepsilon > 0} \left( \left[ \omega(\varepsilon, L, K, \mathcal{F}) - \frac{\varepsilon\sqrt{V}}{\sigma} \right]_+ \right)^2.$$

Now let  $\phi(\theta) = Lf_0 + \theta(Lf_1 - Lf_0)$  and  $f_\theta = f_0 + \theta(f_1 - f_0)$ . Over the family  $\{f_\theta: -1 \leq \theta \leq 1\}$ , the estimator  $\hat{L}_V$  defined by (3.5) is sufficient for  $\phi(\theta)$ . Also since  $\{f_\theta: -1 \leq \theta \leq 1\} \in \mathcal{F}$ ,

$$(3.22) \quad B^2(V, \sigma, L, \mathcal{F}) \geq B^2(V, \sigma, L, \{f_\theta: -1 \leq \theta \leq 1\}).$$

Now note that on the family  $\{f_\theta: -1 \leq \theta \leq 1\}$ ,

$$(3.23) \quad \hat{L}_V \sim N(\phi(\theta) - \theta\sqrt{\gamma}, V),$$

where we have set

$$(3.24) \quad \gamma = 4^{-1} \sup_{\varepsilon > 0} \left( \left[ \omega(\varepsilon, L, K, \mathcal{F}) - \frac{\varepsilon\sqrt{V}}{\sigma} \right]_+ \right)^2.$$

If we put

$$(3.25) \quad c = \frac{Lf_1 - Lf_0}{Lf_1 - Lf_0 - \sqrt{\gamma}},$$

it follows that

$$(3.26) \quad (\hat{L}_V - Lf_0)c + Lf_0 \sim N(\phi(\theta), Vc^2).$$

A simple computation using (2.3) then yields

$$(3.27) \quad B^2(V, \sigma, L, \{f_\theta: -1 \leq \theta \leq 1\}) \geq \gamma,$$

which by (3.20) and (3.24) yields (3.3) and hence also (3.7). This completes the proof of the theorem.  $\square$

Theorem 2 can be used to provide a new interpretation of the optimality of some well known linear estimators. For example, the following theorem of Speckman (1979) is quoted in Donoho (1994).

**THEOREM 3.** *Let  $y_i = f(t_i) + z_i$ ,  $i = 1, \dots, n$ , where  $t_i \in [0, 1]$ ,  $z_i$  are i.i.d.  $N(0, \sigma^2)$  and where the function  $f$  is known to satisfy  $\int (f''(t))^2 dt \leq C^2$ . Let  $g_\mu$  be the solution to*

$$\min_{g \in \mathcal{F}} \sum_i (g(t_i) - y_i)^2 + \mu \int_0^1 (g''(t))^2 dt.$$

*Then  $g_\mu$  is a cubic spline. Let  $L$  be a linear functional with finite minimax mean squared error. Then with  $\mu = \sigma^2/C^2$ , the estimate*

$$L_0(y_1, y_2, \dots, y_n) = L(g_\mu)$$

*is the minimax linear estimator of  $L$  under squared error loss.*

This theorem shows that over the ellipsoid  $\int f''^2 \leq C^2$ , minimax linear estimators of a linear functional under squared error loss are given by applying that linear functional to a cubic smoothing spline. Donoho (1994) has shown that these minimax linear estimators are in fact also minimax linear estimators under absolute error loss at another noise level and are also the center of shortest fixed length affine confidence intervals at yet another noise level. Theorem 2 of this paper shows that Speckman's estimates are not only minimax affine estimates, but also in fact minimax the variance among all measurable estimators with a particular bound on the variance. This new optimality property also holds for other well known estimators such as those of Sacks and Ylvisaker (1981), Epanechnikov (1969) and Li (1982).

Theorem 2 also gives a precise quantification for the possible tradeoffs of bias and variance. In the following discussion write  $\omega(\varepsilon)$  for  $\omega(\varepsilon, L, K, \mathcal{F})$ . First note that if an estimator has maximum variance much smaller than  $\omega^2(\sigma)$ , then the maximum squared bias must be much larger than  $\omega^2(\sigma)$ . Likewise an estimator with a small maximum squared bias compared to  $\omega^2(\sigma)$  must have a large maximum variance compared to  $\omega^2(\sigma)$ . This is summarized in the following corollary which follows almost immediately from Theorem 2. It also follows from results of Liu and Brown (1993).

**COROLLARY 1.** *Suppose that  $L$  is well defined and that*

$$(3.28) \quad \limsup_{C \rightarrow \infty} \liminf_{\sigma \rightarrow 0} \frac{\omega(C\sigma)}{\omega(\sigma)} = \infty.$$

It follows that if  $\hat{L}_\sigma$  is a sequence of estimators such that

$$(3.29) \quad \liminf_{\sigma \rightarrow 0} \sup_{\mathcal{F}} \frac{\text{Var}_\sigma(\hat{L}_\sigma, f)}{\omega^2(\sigma)} = 0,$$

then

$$(3.30) \quad \limsup_{\sigma \rightarrow 0} \sup_{\mathcal{F}} \frac{\text{Bias}_\sigma^2(\hat{L}_\sigma, L, f)}{\omega^2(\sigma)} = \infty.$$

Likewise if

$$(3.31) \quad \limsup_{\varepsilon \rightarrow 0} \frac{\omega(\varepsilon)}{\varepsilon} = \infty$$

and

$$(3.32) \quad \liminf_{\sigma \rightarrow 0} \sup_{\mathcal{F}} \frac{\text{Bias}_\sigma^2(\hat{L}_\sigma, L, f)}{\omega^2(\sigma)} = 0,$$

then

$$(3.33) \quad \limsup_{\sigma \rightarrow 0} \sup_{\mathcal{F}} \frac{\text{Var}_\sigma(\hat{L}_\sigma, f)}{\omega^2(\sigma)} = \infty.$$

It is also possible to give more precise results. In many problems the modulus of continuity has an exact power law  $\omega(\varepsilon, L, K, \mathcal{F}) = A\varepsilon^r$  with  $0 < r < 1$  at least for  $\varepsilon < \varepsilon_0$ . Many such examples are given in Donoho and Liu (1991), Donoho and Low (1992) and Low (1992). Then if  $0 < \delta \ll 1$ , where  $\delta$  can depend on  $\sigma$ ,

$$(3.34) \quad V(\delta\omega(\sigma), \sigma, L, \mathcal{F}) = \omega^2(\sigma)(2\delta^2)^{1-1/r}(1-r)^{2/r}\left(\frac{r}{1-r}\right)^2.$$

Since  $1 - 1/r < 0$  it follows that the term  $(\delta^2)^{1-1/r}$  blows up as  $\delta \rightarrow 0$ . Likewise

$$(3.35) \quad \begin{aligned} B^2(\delta^2\omega^2(\sigma), \sigma, L, \mathcal{F}) \\ = 4^{-1}\omega^2(\sigma)(\delta^2)^{(1-1/(1-r))}\{r^{r/(1-r)} - r^{1/(1-r)}\}^2. \end{aligned}$$

In most other problems, even when an exact power law does not hold, the modulus usually has an asymptotic relation of the form  $\omega(\varepsilon, L, K, \mathcal{F}) \sim A\varepsilon^r$ , where  $0 < r < 1$  and once again it follows that

$$(3.36) \quad \begin{aligned} V(\delta\omega(\sigma), \sigma, L, \mathcal{F}) \\ \sim \omega^2(\sigma)((2\delta^2))^{1-1/r}(1-r)^{2/r}\left(\frac{r}{1-r}\right)^2, \quad \sigma \rightarrow 0, \end{aligned}$$

and

$$(3.37) \quad \begin{aligned} B^2(\delta^2\omega^2(\sigma), \sigma, L, \mathcal{F}) \\ \sim 4^{-1}\omega^2(\sigma)(\delta^2)^{(1-1/(1-r))}\{r^{r/(1-r)} - r^{1/(1-r)}\}^2, \quad \sigma \rightarrow 0. \end{aligned}$$

**Acknowledgments.** The author would like to thank David Donoho, Iain Johnstone and Larry Brown for extremely useful communications. The final manuscript also greatly benefitted from very helpful referee reports. A particularly detailed report formed the basis for the present Introduction and also the notation used in this final manuscript.

## REFERENCES

- BROWN, L. D. and FARRELL, R. H. (1990). A lower bound for the risk in estimating the value of a probability density. *J. Amer. Statist. Assoc.* **85** 1147–1153.
- BROWN, L. D. and FELDMAN, I. (1989). The minimax risk for estimating a bounded normal mean. Technical Report, Cornell Statistics Center.
- CASELLA, G. and STRAWDERMAN, W. E. (1981). Estimating a bounded normal mean. *Ann. Statist.* **9** 870–878.
- DONOHO, D. L. (1994). Statistical estimation and optimal recovery. *Ann. Statist.* **22** 238–270.
- DONOHO, D. L. and LIU, R. C. (1991). Geometrizing rates of convergence, III. *Ann. Statist.* **19** 668–701.
- DONOHO, D. L. and LOW, M. G. (1992). Renormalization exponents and optimal pointwise rates of convergence. *Ann. Statist.* **20** 944–970.
- EPANECHNIKOV, V. (1969). Nonparametric estimates of a multivariate probability density. *Theory Probab. Appl.* **14** 153–158.
- HALL, P. (1989). On convergence rates in nonparametric problems *Internat. Statist. Rev.* **57** 45–58.
- IBRAGIMOV, I. A. and HASMINSKII, R. Z. (1984). On nonparametric estimation of the value of a linear functional in Gaussian white noise. *Theory Probab. Appl.* **29** 18–32.
- LI, K. C. (1982). Minimality of the method of regularization on stochastic processes *Ann. Statist.* **10** 937–942.
- LIU, R. C. and BROWN, L. D. (1993). Non-existence of informative unbiased estimators in singular problems. *Ann. Statist.* **21** 1–14.
- LOW, M. G. (1992). Renormalization and white noise approximation for nonparametric functional estimation problems *Ann. Statist.* **20** 545–554.
- SACKS, J. and STRAWDERMAN, W. (1982). Improvements on linear minimax estimates. In *Statistical Decision Theory and Related Topics 3* (S. S. Gupta and J. O. Berger, eds.) **2** 287–304. Academic Press, New York.
- SACKS, J. and YLVIKAKER, D. (1978). Linear estimation for approximately linear models. *Ann. Statist.* **6** 1122–1137.
- SACKS, J. and YLVIKAKER, D. (1981). Asymptotically optimum kernels for density estimation at a point. *Ann. Statist.* **9** 334–346.
- SPECKMAN, P. (1979). Minimax estimates of linear functionals in a Hilbert space. Unpublished manuscript.

DEPARTMENT OF STATISTICS  
 THE WHARTON SCHOOL  
 UNIVERSITY OF PENNSYLVANIA  
 PHILADELPHIA, PENNSYLVANIA 19104-6302