

## STATISTICAL ESTIMATION AND OPTIMAL RECOVERY<sup>1</sup>

BY DAVID L. DONOHO

*University of California, Berkeley*

New formulas are given for the minimax linear risk in estimating a linear functional of an unknown object from indirect data contaminated with random Gaussian noise. The formulas cover a variety of loss functions and do not require the symmetry of the convex a priori class. It is shown that affine minimax rules are within a few percent of minimax even among nonlinear rules, for a variety of loss functions. It is also shown that difficulty of estimation is measured by the modulus of continuity of the functional to be estimated.

The method of proof exposes a correspondence between minimax affine estimates in the statistical estimation problem and optimal algorithms in the theory of optimal recovery.

**1. Introduction.** Suppose we observe data  $\mathbf{y}$  of the form

$$(1) \quad \mathbf{y} = K\mathbf{x} + \mathbf{z},$$

where  $\mathbf{x}$  is an element of a convex subset  $\mathbf{X}$  of  $l_2$ ,  $K$  is a linear operator and  $\mathbf{z}$  is a noise vector. We are interested in estimating the value of the linear functional  $L(\mathbf{x})$ , and we wish to do this in such a way as to minimize the error occurring at the worst  $\mathbf{x} \in \mathbf{X}$ .

When  $\mathbf{z}$  is assumed to be a zero-mean Gaussian noise with covariance  $\sigma^2\Sigma$ , this is a problem of *minimax statistical estimation*. There is a considerable literature on minimax mean-squared error estimation of linear functionals in such situations—a partial listing would include Kuks and Olman (1972), Läuter (1975), Sacks and Ylvisaker (1978), Speckman (1979), Li (1982), Ibragimov and Has'minskii (1984, 1987), Pilz (1986), Heckman (1988), Donoho and Liu (1991) and Pinelis (1991). There is also considerable literature on minimax mean-square estimation in models related to, but not identical to, (1).

When  $\mathbf{z}$  is assumed to be a vector chosen, not at random, but by an antagonistic opponent, subject to the constraint  $\langle \mathbf{z}, \Sigma^{-1}\mathbf{z} \rangle \leq \varepsilon^2$ , this is a problem of *optimal recovery* of a linear functional. The author is not able to give a complete listing of work on this topic, but is aware of, for example, Micchelli

---

Received March 1990; revised December 1992.

<sup>1</sup>Supported by NSF Grant DMS-84-51753 and by grants from Sun Microsystems, Inc., Schlumberger-Doll Research and Western Geophysical.

*AMS 1991 subject classification.* Primary 62C20, 62G07; secondary 41A25, 43A30.

*Key words and phrases.* Bounded normal mean, estimation of linear functionals, confidence statements for linear functionals, modulus of continuity, minimax risk, nonparametric regression, density estimation.

(1975), Micchelli and Rivlin (1977), Melkman and Micchelli (1979), Packel and Woźniakowski (1987), Traub, Wasilkowski and Woźniakowski (1988) and Packel (1988).

While the two problems are superficially different, there are a number of underlying similarities. Suppose that  $L, K$  and  $\Sigma$  are fixed, but we approach the problem two different ways: one time assuming the noise is random Gaussian, and the other time assuming the noise is chosen by an antagonist, subject to a quadratic constraint. In some cases both ways of stating the problem have been solved, and what happens is that while the two solutions are different in detail, they belong to the same family—that is the same family of splines, of kernel estimators or of regularized least squares estimates—only the “tuning constants” are chosen differently.

Also, a number of theoretical results in the two different fields bear resemblance. For example, Micchelli (1975) showed in the optimal recovery model that minimax linear estimates are generally minimax even among all non-linear estimates. Ibragimov and Has'minskii (1984, 1987) showed in the statistical estimation model that with  $K = I$  and squared error loss, minimax linear estimates are within some constant factor of being minimax among all estimates.

However, there are also disparities; one gets the impression that the literature on optimal recovery is more developed and intensely cultivated than the statistical estimation literature. Consequently there are a number of problems that have been treated as optimal recovery problems, and not yet as statistical estimation problems.

In previous work on statistical estimation, it has been assumed either that  $\mathbf{X}$  is ellipsoidal [cf. Kuks and Olman (1972), Läuter (1975), Speckman (1979) and Li (1982)] or hyperrectangular [cf. Sacks and Ylvisaker (1978, 1981)] or at least centrosymmetric [Ibragimov and Has'minskii (1984, 1987), Pilz (1986) and Pinelis (1991)]; an exception is Donoho and Liu (1991), where only convexity is assumed. Also, in certain instances [Ibragimov and Has'minskii (1984, 1987), Donoho and Liu (1991)], the operator  $K$  was of a very special form. Also, performance was measured via squared error loss only; an exception is Pinelis (1991). Theorems 1 and 2 of this paper give new general formulas for the minimax risk of affine estimates in the statistical estimation problem, with respect to various performance criteria. The formulas hold for general linear operators  $K$ , and without assuming more than convexity of  $\mathbf{X}$ . Our theorems may thus be viewed as the completion of a lengthy development in the statistical literature, aiming at a general characterization of minimax linear estimates of linear functionals from noisy data.

Our approach has several corollaries of immediate usefulness. Corollary 1 shows that minimax affine estimators are nearly minimax among all estimates, that is, that the minimax risk among affine estimates is within a few percent of the minimax risk among all estimates, in a variety of loss functions. We list in Section 9 a wide variety of statistical models, such as nearly linear models, semiparametric models, nonparametric regression models and signal recovery models covered by model (1). It follows that, in all these cases, min-

imax affine estimates (which are computationally tractable) are also nearly minimax among all estimates.

Corollary 2 gives relations between the modulus of continuity of the functional to be estimated and the minimax risk. It follows (see Corollary 3) that results on asymptotic behavior of minimax risk, a statistical problem, follow from asymptotic behavior of the modulus of continuity, an analytic object. The results given here form the crucial step in studying asymptotic minimax risk in a wide variety of statistical estimation problems, ranging from nonparametric and semiparametric regression to density estimation, to signal recovery. (See Theorems 3 and 4 in subsection 9.3.)

Theorems 1 and 2, and their corollaries, bring the theory of minimax linear statistical estimation to a state comparable to the theory of linear optimal recovery. This is no accident. A secondary aim of this paper is to show that, at some deeper level, the problems of statistical estimation and optimal recovery are really the same—that an estimator optimal for one problem is optimal also for the other—*provided  $\varepsilon$  and  $\sigma$  are calibrated appropriately*. This means that results obtained in one literature may be exploited in the other.

To show this, we have studied a generalization of the optimal recovery problem of Micchelli (1975). Assuming that  $\mathbf{X}$  is just convex (i.e., without assuming symmetry of  $\mathbf{X}$ ), we show in Theorem 5 the existence of affine optimal algorithms. The proof is entirely parallel to the proofs of Theorems 1 and 2; this shows that the basic results in both fields follow from the same pattern of reasoning and, in the main, from a single inequality, (50).

**2. The bounded normal mean.** The statement of our main result in Section 4 requires the introduction of some ideas and results from statistical decision theory.

Suppose we are interested in estimating the real-valued quantity  $\theta$ , from observation of the random variable  $Y = \theta + Z$ , where  $Z$  is a random variable with Gaussian distribution  $N(0, \sigma^2)$ .  $Y$  itself may be used as an estimate, of course, but suppose we know a priori that  $\theta \in [-\tau, \tau]$ , and we wish to use this a priori knowledge to do better than  $Y$ . The extent to which we can improve on  $Y$  itself depends on what measure of performance we use and on whether we use only affine (inhomogeneous linear) estimates or whether we allow the possibility of general nonlinear estimates.

Evaluate performance by worst-case mean squared error. Then the best performance among affine estimates  $cY + d$  is

$$(2) \quad \rho_A(\tau, \sigma) = \min_{c, d} \max_{\theta \in [-\tau, \tau]} E(cY + d - \theta)^2,$$

and, among nonlinear estimates  $\delta(Y)$ ,

$$(3) \quad \rho_N(\tau, \sigma) = \inf_{\delta} \max_{\theta \in [-\tau, \tau]} E(\delta(Y) - \theta)^2,$$

where the infimum is over measurable functions. These two quantities are called the *minimax affine risk* and *minimax risk*, respectively; they have

been studied by Levit (1980), Casella and Strawderman (1981), Bickel (1981) and Ibragimov and Has'minskii (1984). See also Donoho, Liu and MacGibbon (1990). They satisfy  $\rho \leq \min(\tau^2, \sigma^2)$ , the invariance  $\rho(\tau, \sigma) = \sigma^2 \rho(\tau/\sigma, 1)$  and the limiting relation  $\rho(\tau, \sigma) \rightarrow \sigma^2, \tau/\sigma \rightarrow \infty$ . Three facts are of particular interest. First, the two risks are never very different. Donoho, Liu and MacGibbon (1990) and Feldman and Brown (1989) have shown that

$$(4) \quad \rho_A(\tau, \sigma) \leq \frac{5}{4} \rho_N(\tau, \sigma).$$

Second, while there is no closed-form expression for  $\rho_N$  (various inequalities are available), for the affine risk we have

$$(5) \quad \rho_A(\tau, \sigma) = \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2}.$$

Third, the minimax affine estimator is  $c_0 Y$ , where

$$(6) \quad c_0(\tau, \sigma; l_2) = \frac{\tau^2}{\tau^2 + \sigma^2}.$$

(The  $l_2$  refers to "squared error loss" criterion.)

Suppose instead we evaluate performance by worst-case mean *absolute* error. Let  $\lambda_A(\tau, \sigma)$  denote the minimax value of  $E|cY + d - \theta|$  among affine estimates, and  $\lambda_N$  denote the minimax value among nonlinear estimates. We have not seen these discussed before in the literature, although techniques similar to those used for quadratic error may be used to study them. These measures satisfy  $\lambda \leq \min(\tau, \sqrt{2/\pi}\sigma)$ , the invariance  $\lambda(\tau, \sigma) = \sigma \lambda(\tau/\sigma, 1)$  and the limiting relation  $\lambda(\tau, \sigma) \rightarrow \sqrt{2/\pi}\sigma$ , as  $\tau/\sigma \rightarrow \infty$ . The two risks are again never very different. In unpublished work, Liu (1989) has shown (by extensive computations) that

$$(7) \quad \lambda_A(\tau, \sigma) \leq 1.23 \lambda_N(\tau, \sigma).$$

Unfortunately, there is no closed-form expression for  $\lambda_N$  or  $\lambda_A$ , although inequalities can be developed. However, the minimax affine estimator is again of the form  $c_0 Y$ , where  $c_0$  can be computed numerically, and it can be proved that

$$(8) \quad c_0(\tau, \sigma; l_1) \text{ is a monotone increasing function of } \tau/\sigma,$$

with  $0 \leq c_0 \leq 1$ .

As a third possibility, consider evaluating performance by the size of fixed-length confidence statements. That is, let  $\alpha \in [0, 1]$ , and let  $\chi_{A, \alpha}(\tau, \sigma)$  denote the smallest number  $\chi$  such that, for some  $c$  and  $d$ , we have  $P\{|cY + d - \theta| \leq \chi\} \geq 1 - \alpha$ , for all  $\theta \in [-\tau, \tau]$ . Similarly, let  $\chi_{N, \alpha}(\tau, \sigma)$  denote the smallest number  $\chi$  such that, for some  $\delta(\cdot)$ , we have  $P\{|\delta(Y) - \theta| \leq \chi\} \geq 1 - \alpha$  whatever  $\theta \in [-\tau, \tau]$  may be. These quantitative measures do not appear to have been discussed in the literature before the first draft of this manuscript. They may be analyzed by adapting techniques of Zeytinoglu and Mintz (1984, 1988). Stark

(1992) has recently given tables and other information about  $\chi_{A,\alpha}$ . Denote by  $Z_{1-\alpha}$  the 100(1 -  $\alpha$ ) percentile of the normal distribution. Both measures satisfy  $\chi \leq \min(\tau, Z_{1-\alpha/2}\sigma)$ , the invariance  $\chi(\tau, \sigma) = \sigma\chi(\tau/\sigma, 1)$  and the limiting relation  $\chi(\tau, \sigma) \rightarrow Z_{1-\alpha/2}\sigma, \tau/\sigma \rightarrow \infty$ . We also have

$$(9) \quad \chi_{N,\alpha}(\tau, 1) = \chi_{A,\alpha}(\tau, 1) = \tau, \quad \tau \leq Z_{1-\alpha}.$$

It follows that

$$(10) \quad \chi_{A,\alpha}(\tau, \sigma) \leq \frac{Z_{1-\alpha/2}}{Z_{1-\alpha}} \chi_{N,\alpha}(\tau, \sigma).$$

Hence, for  $\alpha = 0.05$ , the two risks never differ by more than 1.96/1.645 = 1.19... The minimax affine estimator is again of the form  $c_0Y$ , where  $c_0$  can be computed numerically, and it can be proved that

$$(11) \quad c_0(\tau, \sigma; \alpha) \text{ is a monotone increasing function of } \tau/\sigma.$$

with  $0 \leq c_0 \leq 1$ .

Two final, technical remarks follow: first, for each of the three criteria,

$$(12) \quad c_0(\tau, \sigma; \cdot) = o(\tau), \quad \tau \rightarrow 0.$$

This fact is apparent for the  $l_2$  measure from (6); for the other measures it may be established by analysis.

Second, the minimax and minimax affine estimates for all these problems are *nonrandomized*. The reasoning is akin to "Rao-Blackwell"-ization and "Karlin-Rubin"-ization. Using terminology of Brown, Cohen and Strawderman (1976), we note that the normal location family has strict monotone likelihood ratio, and the loss functions of interest have "points of increase," so Theorem 2.1 and Remark 2.1 of Brown, Cohen and Strawderman (1976) reveal that each randomized estimator is dominated by a nonrandomized estimator. Hence if we had an opportunity to observe  $(Y, Z_2, Z_3, \dots)$ , where the  $Z_i$  are random variables whose distribution does not depend on  $\theta$  and which are stochastically independent of  $Y$ , we could always do at least as well by using a function of  $Y$  alone. Moreover, when the  $Z_i$  are, in addition, i.i.d. zero-mean normal random variables, the "Brown-Cohen-Strawderman"-ization of an affine function of  $(Y, Z_2, Z_3, \dots)$  is an affine function of  $Y$ . Therefore, given the opportunity to form an affine function of  $(Y, Z_2, Z_3, \dots)$  we could always do at least as well with an affine function of  $Y$  alone.

**3. Hardest one-dimensional subproblems.** We return now to the statistical estimation setting of the introduction. We make one specialization and one generalization. We suppose that the noise is Gaussian with covariance  $\Sigma = \sigma I$ , where  $I$  is the identity operator. We will show in Section 11 that the case of more general  $\Sigma$  is also covered by these results. We also now allow the functional  $L$  to be *affine* (inhomogeneous linear) and consider as well affine estimates of  $L$ .

We are interested in determining the minimax affine risk with squared error loss,

$$R_A^*(\sigma) = \inf_{\widehat{L} \text{ affine}} \sup_{\mathbf{x} \in \mathbf{X}} E(\widehat{L}(\mathbf{y}) - L(\mathbf{x}))^2,$$

the minimax risk with squared error loss.

$$R_N^*(\sigma) = \inf_{\widehat{L}} \sup_{\mathbf{x} \in \mathbf{X}} E(\widehat{L}(\mathbf{y}) - L(\mathbf{x}))^2,$$

and the analogous quantities for absolute error loss  $\Lambda_A^*(\sigma)$ ,  $\Lambda_N^*(\sigma)$ . We are also interested in the minimax length of fixed-length confidence statements.  $C_{\alpha, A}^*(\sigma)$  is the smallest number  $\chi$  such that, for some affine estimator  $\widehat{L}$ , the confidence interval  $[\widehat{L}(\mathbf{y}) - \chi, \widehat{L}(\mathbf{y}) + \chi]$  covers  $L(\mathbf{x})$  with probability at least  $1 - \alpha$ , for every  $\mathbf{x} \in \mathbf{X}$ . Formally,

$$C_{\alpha, A}^*(\sigma) = \inf \left\{ \chi : \exists \widehat{L} \text{ affine} \ni P(|\widehat{L}(\mathbf{y}) - L(\mathbf{x})| \leq \chi) \geq 1 - \alpha, \forall \mathbf{x} \in \mathbf{X} \right\}.$$

The definition of  $C_{\alpha, N}^*(\sigma)$  is analogous.

Suppose we knew *a priori* not just that  $\mathbf{x} \in \mathbf{X}$ , but actually that  $\mathbf{x}$  belongs to the *one-dimensional subfamily*

$$(13) \quad [\mathbf{x}_{-1}, \mathbf{x}_1] = \{t\mathbf{x}_{-1} + (1-t)\mathbf{x}_1 : t \in [0, 1]\}.$$

Set  $R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1])$  for the minimax risk in this subproblem; obviously,

$$(14) \quad R_A^*(\sigma; \mathbf{X}) \geq R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1])$$

since the additional prior information can only help. In fact,

$$(15) \quad R_A^*(\sigma; \mathbf{X}) \geq \sup \left\{ R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) : [\mathbf{x}_{-1}, \mathbf{x}_1] \subset \mathbf{X} \right\},$$

$$(16) \quad R_N^*(\sigma; \mathbf{X}) \geq \sup \left\{ R_N^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) : [\mathbf{x}_{-1}, \mathbf{x}_1] \subset \mathbf{X} \right\},$$

and similar inequalities hold for  $\Lambda_A^*(\sigma)$ ,  $\Lambda_N^*(\sigma)$  and so on. In words, the full problem is at least as hard as any one-dimensional subproblem.

We now evaluate the difficulty of a subproblem.

LEMMA 1.

$$(17) \quad R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \left( \frac{L(\mathbf{x}_1) - L(\mathbf{x}_{-1})}{\|K\mathbf{x}_1 - K\mathbf{x}_{-1}\|} \right)^2 \rho_A \left( \frac{\|K\mathbf{x}_1 - K\mathbf{x}_{-1}\|}{2}, \sigma \right)$$

and similarly for  $R_N^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1])$ ,

$$(18) \quad \Lambda_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \frac{|L(\mathbf{x}_1) - L(\mathbf{x}_{-1})|}{\|K\mathbf{x}_1 - K\mathbf{x}_{-1}\|} \lambda_A \left( \frac{\|K\mathbf{x}_1 - K\mathbf{x}_{-1}\|}{2}, \sigma \right)$$

and similarly for  $\Lambda_N^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1])$ ,  $C_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1])$  and  $C_N^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1])$ .

Let us see why. Let  $\mathbf{x}_0 = (\mathbf{x}_{-1} + \mathbf{x}_1)/2$  denote the center of the subfamily, and set  $\mathbf{w}_0 = K(\mathbf{x}_1 - \mathbf{x}_{-1})/\|K(\mathbf{x}_1 - \mathbf{x}_{-1})\|$ . Define the parameter

$$\theta = \langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_0 \rangle.$$

Consider the problem of estimating  $\theta$  from observations  $\mathbf{y}$ . We know that  $\theta \in [-\tau, \tau]$ , where  $\tau = \|K(\mathbf{x}_1 - \mathbf{x}_{-1})\|/2$ . Defining  $Y = \langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle$ , we have that  $Y$  is  $N(\theta, \sigma^2)$ . Estimation of  $\theta$  from  $Y$  was treated in Section 2, and information about minimax affine and minimax estimators was given.

Suppose  $\delta(\cdot)$  is a minimax estimator from the bounded normal mean problem for the given criterion of interest. Then  $\delta(Y)$  is obviously minimax among all functions of  $Y$ ; we claim it is minimax among all functions of  $\mathbf{y}$ . Indeed, there is an isometry mapping  $\mathbf{y}$  to  $(Y, Z_2, Z_3, \dots)$ , where  $Z_i$  are i.i.d. zero-mean Gaussian random variables, independent of  $Y - \theta$  and of each other, with probability distribution not depending on  $\theta$ . By Brown-Cohen-Strawdermanization we see that these extra, "pure noise" variables do not help us reduce the risk. Hence the minimax risk for estimating  $\theta$  from  $\mathbf{y}$  is that for estimating  $\theta$  from  $Y$ .

We now make the obvious comment that the problem of estimating  $s\theta + t$  from  $Y$  has  $s^2$  times the minimax risk of estimating  $\theta$  from  $Y$ , under quadratic loss, and  $s$  times the minimax risk of estimating  $\theta$  from  $Y$  under the absolute error or confidence-statement criterion. The restriction of  $L$  to the subfamily  $[\mathbf{x}_{-1}, \mathbf{x}_1]$  is an affine function  $L(\mathbf{x}) = L(\mathbf{x}_0) + s\theta$ . The results quoted in the lemma follow by computing  $s$ .

We now employ the lemma. Introduce the seminorm  $\|v\|_K \equiv \|Kv\|$ . The modulus of continuity of  $L$  with respect to this seminorm is defined as

$$\omega(\varepsilon; L, K, \mathbf{X}) = \sup \left\{ |L(\mathbf{x}_1) - L(\mathbf{x}_{-1})| : \|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K \leq \varepsilon \text{ and } \mathbf{x}_i \in \mathbf{X} \right\}.$$

We generally omit the secondary arguments, these being clear from context.

The modulus may be used to calculate the right-hand side of (15). Indeed,

$$\begin{aligned} \sup_{[\mathbf{x}_{-1}, \mathbf{x}_1] \in \mathbf{X}} R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) &= \sup_{\varepsilon \geq 0} \sup_{\|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K = \varepsilon} \left( \frac{L(\mathbf{x}_1) - L(\mathbf{x}_{-1})}{\varepsilon} \right)^2 \rho_A\left(\frac{\varepsilon}{2}, \sigma\right) \\ &= \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right)^2 \rho_A\left(\frac{\varepsilon}{2}, \sigma\right). \end{aligned}$$

We say that the modulus "measures the difficulty of the hardest one-dimensional subproblem." This might be an abuse of language if no such *hardest* subfamily existed (i.e., if the corresponding supremum were not attained). However, a hardest subfamily will exit in considerable generality.

We need one technical restriction on the class of problems treated.

**DEFINITION.** We say that  $L$  is *well-defined* if the modulus of continuity of  $L$  over  $\mathbf{X}$  in the usual  $l_2$  norm is continuous at 0:  $\omega(\varepsilon; L, I, \mathbf{X}) \rightarrow 0$ , as  $\varepsilon \rightarrow 0$ .

We say that the linear operator  $K$  is well-defined if the modulus of continuity of  $K$  over  $\mathbf{X}$  for the  $l_2$  norm is continuous at 0.

The condition that  $L$  and  $K$  be well-defined is much weaker than the condition that  $L$  be a bounded linear functional and  $K$  a bounded linear operator. In all the most interesting examples of Section 9, the functionals and operations involve in some way point evaluations and so are not bounded. Restricting attention to well-defined cases serves primarily to rule out consideration of nonmeasurable linear functionals and of problems where noisy data can provide essentially no information about the functional's value. This condition is satisfied by all the many examples we have looked at.

**LEMMA 2.** *If  $\mathbf{X}$  is closed, convex and bounded, if  $L$  and  $K$  are well-defined and if  $\omega(\varepsilon)$  is finite for each  $\varepsilon \geq 0$ , then the modulus of continuity is attained. That is, for each  $\varepsilon \geq 0$ , there exists a pair  $(\mathbf{x}_1, \mathbf{x}_{-1})$  such that  $\|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K \leq \varepsilon$  and*

$$|L(\mathbf{x}_1) - L(\mathbf{x}_{-1})| = \omega(\varepsilon).$$

Moreover, for each of the three performance criteria, there exists a hardest subfamily for affine estimates, that is, a family satisfying

$$R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right)^2 \rho_A \left( \frac{\varepsilon}{2}, \sigma \right),$$

a (generally different) family satisfying

$$\Lambda_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right) \lambda_A \left( \frac{\varepsilon}{2}, \sigma \right),$$

and a (still different) family satisfying

$$C_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right) \chi_{A, \alpha} \left( \frac{\varepsilon}{2}, \sigma \right).$$

**4. Main result.** The following justifies our attention to one-dimensional subproblems.

**THEOREM 1.** *Let  $\mathbf{X}$  be closed, bounded and convex, let  $L$  and  $K$  be well-defined and suppose that  $\omega(\varepsilon)$  is finite for each  $\varepsilon \geq 0$ . Then, for any of the three performance criteria, the difficulty, for affine estimates, of the full problem is **equal** to the difficulty, for affine estimates, of a hardest one-dimensional subproblem. Thus,*

$$R_A^*(\sigma) = \max_{\mathbf{x}_1, \mathbf{x}_{-1}} R_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]),$$

$$\Lambda_A^*(\sigma) = \max_{\mathbf{x}_1, \mathbf{x}_{-1}} \Lambda_A^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]),$$

$$C_{\alpha, A}^*(\sigma) = \max_{\mathbf{x}_1, \mathbf{x}_{-1}} C_{\alpha, A}^*(\sigma; [\mathbf{x}_{-1}, \mathbf{x}_1]).$$



Furthermore, the affine estimator which is minimax for a hardest subproblem is also minimax for the full problem.

The proof is given in Section 11. This theorem, together with Lemma 2, provides formulas for the minimax risk in the closed, bounded case. By an approximation argument, given in Section 14, those formulas extend to the case of general convex  $\mathbf{X}$ :

**THEOREM 2.** *Let  $L$  be affine and let  $\mathbf{X}$  be convex. If  $L$  and  $K$  are well-defined, then*

$$\begin{aligned} R_A^*(\sigma) &= \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right)^2 \rho_A \left( \frac{\varepsilon}{2}, \sigma \right), \\ \Lambda_A^*(\sigma) &= \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right) \lambda_A \left( \frac{\varepsilon}{2}, \sigma \right), \\ C_{\alpha, A}^*(\sigma) &= \sup_{\varepsilon \geq 0} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right) \chi_{A, \alpha} \left( \frac{\varepsilon}{2}, \sigma \right). \end{aligned}$$

For squared error loss, with  $L$  linear and  $K = I$ , and with  $\mathbf{X}$  centrosymmetric about  $\mathbf{0}$ , Ibragimov and Has'minskii (1984) gave the formula

$$\sup_{\mathbf{x} \in \mathbf{X}} \frac{\sigma^2 L^2(\mathbf{x})}{\sigma^2 + \|\mathbf{x}\|^2}$$

for the minimax risk of linear estimates. This may be shown to be a particular case of our formula for  $R_A^*(\sigma)$ . The formula for  $R_A^*(\sigma)$  has been proved before in special cases by Donoho and Liu (1991) and by Brown and Liu (1989). The formulas for  $\Lambda_A^*(\sigma)$  and  $C_{\alpha, A}^*(\sigma)$  are new.

**5. Near-minimaxity of affine estimates.** These formulas imply that affine estimators cannot be improved on much by nonlinear estimators. Indeed, using Theorem 2 and (4), we have

$$\begin{aligned} R_A^*(\sigma) &= \sup_{\varepsilon \geq 0} \frac{\omega(\varepsilon)^2}{\varepsilon^2} \rho_A \left( \frac{\varepsilon}{2}, \sigma \right) \\ &\leq \frac{5}{4} \sup_{\varepsilon \geq 0} \frac{\omega(\varepsilon)^2}{\varepsilon^2} \rho_N \left( \frac{\varepsilon}{2}, \sigma \right) \\ &\leq \frac{5}{4} R_N^*(\sigma). \end{aligned}$$

Arguing similarly for the other measures of performance and using the facts (7) and (10) gives the following corollary.

COROLLARY 1. *Under the assumptions of Theorem 2,*

$$\begin{aligned} R_A^*(\sigma) &\leq 1.25R_N^*(\sigma), \\ \Lambda_A^*(\sigma) &\leq 1.23\Lambda_N^*(\sigma), \\ C_{\alpha,A}^*(\sigma) &\leq \frac{Z_{1-\alpha/2}}{Z_{1-\alpha}}C_{\alpha,N}^*(\sigma). \end{aligned}$$

Hence, quite generally and with respect to several worst-case performance measures, affine estimators cannot be dramatically improved upon by nonlinear estimators.

Previous work has assumed the squared error loss criterion. Sacks and Strawderman (1982) had shown that in some case the minimax linear risk was strictly larger than the minimax risk; Ibragimov and Has'minskii (1984) had shown, under the assumptions  $K = I$  and  $\mathbf{X}$  centrosymmetric, that the ratio of the minimax linear risk and minimax nonlinear risk was less than some unknown, finite positive constant. This "Ibragimov-Has'minskii constant" has been shown by Donoho, Liu and MacGibbon (1990) and by Feldman and Brown (1989) to be less than  $\frac{5}{4}$ .

Here we see that, for general  $K$ , without any assumption of symmetry, and in several different performance measures, the minimax affine estimator must be quantitatively quite close to minimax.

**6. The minimax affine estimator.** For this section, fix one of the three performance criteria. Suppose that a hardest subfamily for affine estimates  $[\mathbf{x}_{-1}, \mathbf{x}_1]$  exists under that criterion (e.g., if  $\mathbf{X}$  is closed and norm-bounded). Define the parameters  $\mathbf{x}_0$  and  $\mathbf{w}_0$  as in the proof of Lemma 1. For estimating the parameter  $\theta = \langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_0 \rangle$  in the subfamily  $[\mathbf{x}_{-1}, \mathbf{x}_1]$ , the minimax affine estimator is unique: it is just  $\hat{\theta} = c_0 \langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle$  (here  $c_0$  depends on the performance criterion we have chosen). The restriction of  $L$  to the family is affine,  $L(\mathbf{x}) = L(\mathbf{x}_0) + s\theta$ , with slope  $s = (L(\mathbf{x}_1) - L(\mathbf{x}_{-1})) / \|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K$ . Hence the unique minimax affine estimator in the subfamily is

$$L_0(\mathbf{y}) = L(\mathbf{x}_0) + s\hat{\theta}.$$

Theorem 1 says that the minimax affine risk of this subproblem is the minimax affine risk of the full problem, so there is an affine estimator for the full problem which is also minimax affine for the subproblem. However,  $L_0$  is uniquely the minimax affine estimator for the subproblem. This forces  $L_0$  to be minimax affine for the *full* problem.

The formula for  $L_0$  can be rewritten as

$$(19) \quad L_0(\mathbf{y}) = L(\mathbf{x}_0 + c_0 \langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle \cdot \mathbf{u}_0),$$

where  $\mathbf{u}_0 = (\mathbf{x}_1 - \mathbf{x}_{-1}) / \|\mathbf{x}_1 - \mathbf{x}_{-1}\|$ . This says that the minimax affine estimator has the form of projecting the data orthogonally onto the hardest subfamily,

shrinking toward the center of the subfamily by a factor  $c_0$  and evaluating  $L$  on the projected, shrunken result.

The shrinkage coefficient  $c_0$  has an interesting form. Assume that the hardest subproblem has length  $\|K(\mathbf{x}_1 - \mathbf{x}_{-1})\| = \varepsilon_0$ . One can calculate formally that

$$(20) \quad c_0 = \frac{\varepsilon_0 \omega'(\varepsilon_0)}{\omega(\varepsilon_0)};$$

we will prove this later. Thus, if  $\omega(\varepsilon) = A\varepsilon^r$ , for  $\varepsilon > 0$  and  $r \in (0, 1]$ , then  $c_0 = r$ . So in this case the estimator reduces to shrinkage by the rate exponent in the modulus of continuity.

**7. Risk and modulus.** The modulus of continuity of a linear functional over a convex set is subadditive. Hence  $\omega(\varepsilon)/\varepsilon$  is a decreasing function of  $\varepsilon$ . It follows that

$$\sup_{\varepsilon \geq \sigma} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right)^2 \rho_A \left( \frac{\varepsilon}{2}, \sigma \right) \leq \left( \frac{\omega(\sigma)}{\sigma} \right)^2 \sup_{\varepsilon \geq \sigma} \rho_A \left( \frac{\varepsilon}{2}, \sigma \right) = \omega^2(\sigma).$$

On the other hand,  $\omega(\varepsilon)$  is monotone increasing, so

$$\sup_{\varepsilon \leq \sigma} \left( \frac{\omega(\varepsilon)}{\varepsilon} \right)^2 \rho_A \left( \frac{\varepsilon}{2}, \sigma \right) \leq \omega^2(\sigma) \sup_{\varepsilon \leq \sigma} \varepsilon^{-2} \rho_A \left( \frac{\varepsilon}{2}, \sigma \right) \leq \omega^2(\sigma).$$

Combining these displays,  $R_A^*(\sigma) \leq \omega^2(\sigma)$ . Continuing in this fashion and using Theorem 2 for lower bounds, one proves the following corollary.

**COROLLARY 2.** *Under the assumptions of Theorem 2,*

$$\begin{aligned} \rho_N \left( \frac{1}{2}, 1 \right) \omega^2(\sigma) &\leq R_N^*(\sigma) \leq R_A^*(\sigma) \leq \omega^2(\sigma), \\ \lambda_N \left( \frac{1}{2}, 1 \right) \omega(\sigma) &\leq \Lambda_N^*(\sigma) \leq \Lambda_A^*(\sigma) \leq \omega(\sigma), \\ \omega(2Z_{1-\alpha}\sigma) &\leq C_{\alpha,N}^*(\sigma) \leq C_{\alpha,A}^*(\sigma) \leq \omega(2Z_{1-\alpha/2}\sigma). \end{aligned}$$

So the modulus of continuity determines, to within reasonable constant factors, the behavior of the minimax risks.

**8. Asymptotics as  $\sigma \rightarrow 0$ .** If we use the notation  $\omega(\varepsilon) \asymp \varepsilon^r$  to mean that  $c_0\varepsilon^r \leq \omega(\varepsilon) \leq c_1\varepsilon^r$  as  $\varepsilon \rightarrow 0$ , then the preceding inequalities show that  $\omega(\varepsilon) \asymp \varepsilon^r$  implies  $R_N^*(\sigma) \asymp \sigma^{2r}$ , as well as  $\Lambda_N^*(\sigma) \asymp \sigma^r$  and  $C_{\alpha,N}^*(\sigma) \asymp \sigma^r$ . Hence Corollary 2 shows that asymptotics of the modulus control asymptotics of the risk.

\* We say that  $\omega(\varepsilon)$  has exponent  $r$  if  $\omega(\varepsilon) = A\varepsilon^r + o(\varepsilon^r)$ . When this condition holds, it is possible to make precise asymptotic statements.

**COROLLARY 3.** *Suppose that Theorem 2 applies and that the modulus of*

continuity has exponent  $r$ . Then

$$\begin{aligned} R_A^*(\sigma) &= \xi_{2,A}(r)\omega^2(\sigma)(1+o(1)), \\ \Lambda_A^*(\sigma) &= \xi_{1,A}(r)\omega(\sigma)(1+o(1)), \\ C_{\alpha,A}^*(\sigma) &= \xi_{\alpha,A}(r)\omega(\sigma)(1+o(1)), \end{aligned}$$

as  $\sigma \rightarrow 0$ , where

$$\begin{aligned} \xi_{2,A}(r) &= \sup_{v \geq 0} v^{2r-2} \rho_A(v/2, 1), \\ \xi_{1,A}(r) &= \sup_{v \geq 0} v^{r-1} \lambda_A(v/2, 1), \\ \xi_{\alpha,A}(r) &= \sup_{v \geq 0} v^{r-1} \chi_{A,\alpha}(v/2, 1). \end{aligned}$$

Also

$$\begin{aligned} R_N^*(\sigma) &\geq \xi_{2,N}(r)\omega^2(\sigma)(1+o(1)), \\ \Lambda_N^*(\sigma) &\geq \xi_{1,N}(r)\omega(\sigma)(1+o(1)), \\ C_{\alpha,N}^*(\sigma) &\geq \xi_{\alpha,N}(r)\omega(\sigma)(1+o(1)), \end{aligned}$$

as  $\sigma \rightarrow 0$ , where

$$\begin{aligned} \xi_{2,N}(r) &= \sup_{v \geq 0} v^{2r-2} \rho_N(v/2, 1), \\ \xi_{1,N}(r) &= \sup_{v \geq 0} v^{r-1} \lambda_N(v/2, 1), \\ \xi_{\alpha,N}(r) &= \sup_{v \geq 0} v^{r-1} \chi_{N,\alpha}(v/2, 1). \end{aligned}$$

Calculus gives the closed-form expression

$$\xi_{2,A} = 2^{2r-2} r^r (1-r)^{1-r}.$$

For all the other quantities, it is necessary to get bounds via computational means.

It follows from these formulas that

$$(21) \quad \lim_{\sigma \rightarrow 0} \frac{R_A^*(\sigma)}{R_N^*(\sigma)} \leq \frac{\xi_{2,A}(r)}{\xi_{2,N}(r)},$$

$$(22) \quad \lim_{\sigma \rightarrow 0} \frac{\Lambda_A^*(\sigma)}{\Lambda_N^*(\sigma)} \leq \frac{\xi_{1,A}(r)}{\xi_{1,N}(r)},$$

$$(23) \quad \lim_{\sigma \rightarrow 0} \frac{C_{\alpha,A}^*(\sigma)}{C_{\alpha,N}^*(\sigma)} \leq \frac{\xi_{\alpha,A}(r)}{\xi_{\alpha,N}(r)}.$$

These may be used to give somewhat tighter bounds than those proved in Corollary 1. For example, under squared error loss, in problems with  $r = \frac{1}{2}$ , (21) shows that minimax affine estimates can be improved upon by at most 7% as  $\sigma \rightarrow 0$ . See, for example, Donoho and Liu [(1991), Table 1].

Another form of asymptotic relationship can be deduced.

**COROLLARY 4.** *Suppose the modulus of continuity has exponent  $r$  and Theorem 2 applies. Then, for each of the three performance criteria, if  $c_0$  and  $\varepsilon_0$  refer to the shrinkage coefficient in a minimax affine estimator for that criterion and the length of a hardest subfamily for that criterion,*

$$(24) \quad c_0(\varepsilon_0/2, \sigma; \cdot) \rightarrow r \quad \text{as } \sigma \rightarrow 0.$$

Moreover, if  $v_r$  denotes the solution of  $c_0(v, 1; \cdot) = r$  for the criterion of interest, then

$$(25) \quad \varepsilon_0 = 2v_r\sigma(1 + o(1)) \quad \text{as } \sigma \rightarrow 0.$$

In other words, the shrinkage coefficient in the minimax affine estimator tends to  $r$ , and the length of the hardest subproblem behaves like a fixed constant times the noise level. For the  $l_2$  criterion we have, by calculus,

$$(26) \quad v_{2,r} = \sqrt{\frac{r}{1-r}}.$$

The other quantities  $v_{1,r}$  and  $v_{\alpha,r}$  must be found numerically.

**9. Applications.** We now briefly point out some of the different areas in which results given above can be applied.

**9.1. Some familiar statistical models.** The model with observations (1) subsumes many situations familiar to statisticians. In view of this, minimax affine estimators for such models are nearly minimax among all estimates.

*Approximately linear models* [Sacks and Ylvisaker (1978)]. Let  $t_i$  be  $n$  fixed numbers, and suppose we observe

$$(27) \quad Y_i = a + \beta t_i + \delta_i + z_i, \quad i = 1, \dots, n,$$

where  $a$  and  $\beta$  are unknown real numbers and the  $\delta_i$  are unknown, but they are known to satisfy

$$(28) \quad |\delta_i| \leq c_i, \quad i = 1, \dots, n,$$

with the  $c_i$  known constants. The  $z_i$  are, as usual, a  $N(0, \sigma^2)$  Gaussian white noise. We are interested in the value of  $\beta$ . Except for the perturbations  $\delta_i$ , this model posits a linear relation between  $Y_i$  and  $t_i$ —hence the term “approximately linear model.” Sacks and Ylvisaker (1978) have developed a complete treatment of minimax mean square estimation in this model.

This model is a particular instance of ours. Define  $\mathbf{x} = (a, \beta, \delta_1, \dots, \delta_n)$  and  $(K\mathbf{x})_i = a + \beta t_i + \delta_i$ ,  $i = 1, \dots, n$ . Then with  $\Delta$  the hyperrectangular set defined by (28), and with  $\mathbf{X} = \mathbf{R}^2 \times \Delta$  we get precisely a problem of the form mentioned in the Introduction, with  $L(\mathbf{x}) \equiv \beta$ . Of course, our framework handles generalizations of the original Sacks–Ylvisaker model, by defining  $\Delta$  differently—as an ellipsoid, for example, or some other convex set. For example, one might impose monotonicity constraints or moment conditions on the  $(\delta_i)$ .

*Semiparametric models* [Heckman (1988)]. Suppose we observe

$$(29) \quad y_i = \beta t_i + f(u_i) + z_i, \quad i = 1, \dots, n,$$

where  $t_i$  and  $u_i$  are fixed constants,  $u_i \in [0, 1]$ , say, and  $f$  is unknown but is known to lie in a convex function class  $\mathcal{F}$ . Again  $(z_i)$  is Gaussian white noise. We are again interested in estimating  $\beta$ , but  $f$  represents a nuisance which affects our measurements in an unknown but smooth fashion. Set  $\delta_i = f(u_i)$  and

$$\Delta = \{(\delta_i): \delta_i = f(u_i), f \in \mathcal{F}\}.$$

Because convexity of  $\mathcal{F}$  implies that of  $\Delta$ , we have an instance of the (generalized) approximately linear model previously mentioned.

*Nonparametric regression* [Speckman (1979), Li(1982)]. Here we have

$$(30) \quad y_i = f(t_i) + z_i, \quad i = 1, \dots, n,$$

where now  $f \in \mathcal{F}$ , a convex function class on domain  $\mathcal{D} \subset \mathbf{R}^d$  and the samples  $t_i$  are taken at points of  $D$ . We pedantically spell out the representation as a problem of the form (1). We are interested in estimating functionals  $T(f)$  such as  $T_0(f) = f(t_0)$  or  $T_1(f) = f'(t_0)$ , and so on. Let  $(\phi_j(\cdot))$  be an orthonormal basis for  $L_2(\mathcal{D})$ , let  $x_j = \int f \phi_j$ ,  $\mathbf{x} = (x_j)$  and set  $\mathbf{X} = \{(x_j(f)): f \in \mathcal{F}\}$ . Finally, set  $(K\mathbf{x})_i = \sum_j x_j \phi_j(t_i) = f(t_i)$  and  $L(\mathbf{x}) \equiv T(f)$ . This is a problem of our type.

REMARK. The functionals  $T$  one might be interested in estimating here, such as  $T_0$  and  $T_1$ , are not bounded linear functionals, nor is the operator  $K$  a bounded linear operator. However, if the class  $\mathcal{F}$  consists of sufficiently smooth functions, both  $T$  and  $K$  will be well-defined in the sense of our definition.

*Inverse problems* [O’Sullivan (1986)]. Here we have

$$(31) \quad y_i = (Pf)(t_i) + z_i, \quad i = 1, \dots, n,$$

where  $P$  is a linear operator, such as Radon transform, Abel transform, convolution transform and so on. This is again a problem of our type; the setup is as in nonparametric regression, only  $K$  has changed:  $(K\mathbf{x})_i = \sum_j x_j (P\phi_j)(t_i) = (Pf)(t_i)$ .

*Signal recovery* [Hall (1990)]. Here we have noisy, filtered observations of a signal  $\mathbf{x} = (x_i)$ , where now  $i$  ranges over the lattice  $\mathbf{Z}^2$ :

$$(32) \quad y_i = \sum_j k_{i-j} x_j + z_i, \quad i, j \in \mathbf{Z}^2.$$

The noise is i.i.d. Gaussian with variance  $\sigma^2$ , and we wish to recover  $L(\mathbf{x}) = x_0$ . We know a priori that the signal  $x_i$  is slowly changing in  $i$ ; this is expressed by the constraint  $\mathbf{x} \in \mathbf{X}$ , with  $\mathbf{X}$  a certain convex class.

*White noise model* [Ibragimov and Has'minskii (1984), Donoho and Liu (1991)]. We observe

$$(33) \quad Y(t) = \int_{-a}^t f(u) du + \sigma W(t), \quad t \in [-a, a],$$

where  $W(t)$  is a (two-sided) Wiener process [ $W(-a) = 0$ ]. [This is a rigorous way of writing  $dY(t) = f(t) + \sigma dW(t)$ , hence the term "observations in white noise".] We wish to estimate the linear functional  $T(f)$ , and we know a priori that  $f \in \mathcal{F}$ , a convex subset of  $L_2[-a, a]$ .

This reduces to a model of our type with  $K = I$ . With  $\{\phi_i\}_{i=1}^\infty$  a complete orthonormal basis for  $L_2[-a, a]$ , let  $x_i = x_i(f)$  denote the  $i$ th Fourier-Bessel coefficient of  $f$  with respect to this basis, so that  $f \sim \sum_{i=1}^\infty x_i \phi_i$ . Then set  $\mathbf{X} =$  the set of coefficient sequences  $\mathbf{x} = (x_i)$  of members of  $\mathcal{F}$ , and set  $L(\mathbf{x}) = T(f)$  whenever  $\mathbf{x} = \mathbf{x}(f)$ . Observing  $Y$  is equivalent to observing the Fourier-Bessel coefficient sequence  $\mathbf{y} = (y_i)$  where  $y_i = \int \phi_i Y(dt)$ . However, for this we have the observation equation  $y_i = x_i + z_i$ ,  $i = 1, 2, \dots$ , with  $(z_i)$  i.i.d.  $N(0, \sigma^2)$ . Thus the mapping from functions to their coefficient sequences maps the white noise model (33) onto the present one.

It follows from Corollary 1 that in all the models just mentioned, *minimax affine estimates are nearly minimax among all estimates*.

**9.2. Deriving minimax affine estimates.** Our theory may be used to derive new approaches to the models just mentioned. For example, in Heckman's treatment of the semiparametric model, only two particular function classes  $\mathcal{F}$  are considered, and minimax linear estimators are derived for those two cases. Our approach would allow derivation of parametric quadratic programming algorithms to design estimators useful for convex function classes other than the two considered by Heckman; for example, for classes of smooth monotone functions. However, for reasons of space we turn to other matters.

**9.3. Asymptotic statistical theory.** The results of Sections 7 and 8 allow us to derive, by simple, general techniques, relatively precise results on the behavior of asymptotic minimax risk in statistical estimation problems with increasing sample size. In essence, the risk in problems such as semiparametric and nonparametric estimation with  $n \rightarrow \infty$  is equivalent to the risk in white noise problems with  $\sigma \rightarrow 0$ . This principle has been formulated for local asymptotic minimax risk in Low (1988), for minimax affine risk in Donoho and Liu (1991) and in Donoho and Low (1990), and for minimax risk in Brown and Low (1990).

Theorems 1 and 2 and their corollaries provide asymptotics in the white noise problem as  $\sigma \rightarrow 0$  and thereby give asymptotics in the statistical problems as  $n \rightarrow \infty$ . Thus the results of this paper, together with approximation

arguments developed elsewhere, give a variety of results in asymptotic decision theory. We mention three examples.

*Optimal rates of convergence in nonparametric regression* [Donoho and Low (1990)]. In the nonparametric regression model mentioned earlier, suppose that the evaluation points  $t_i$  are a random sample from the uniform distribution on  $\mathcal{D}$ .

We are interested in estimating the affine functional  $T(f)$ . An affine rule for this problem is any rule of the form  $\widehat{T}(\mathbf{y}) = e + \sum_j l_j y_j$ , where the  $l_j$  are allowed to depend on the  $(t_i)$  but not the  $(y_i)$ . Denote by  $R_A(n)$  the minimax risk of an affine procedure based on  $n$  observations, with respect to squared error loss. Define  $\Lambda_A(n)$  and  $C_{A,\alpha}(n)$  similarly. Combining results in Donoho and Low (1990) with those of Section 7, we get the following theorem.

**THEOREM 3.** *Let  $\omega(\varepsilon)$  be the  $L_2(\mathcal{D})$  modulus of continuity of the functional  $T$  over the class  $\mathcal{F}$ . Suppose that the function class  $\mathcal{F}$  consists of elements all bounded by  $M$  in supremum norm. Let  $\tau = \sqrt{\sigma^2 + M^2}$ . Then*

$$\begin{aligned} \omega^2\left(\frac{\sigma/\sqrt{n}}{5}\right) &\leq R_A(n) \leq \omega^2\left(\frac{\tau}{\sqrt{n}}\right), \\ \omega\left(\frac{\sigma/\sqrt{n}}{2}\right) &\leq \Lambda_A(n) \leq \omega\left(\frac{\tau}{\sqrt{n}}\right), \\ \omega\left(2Z_{1-\alpha}\frac{\sigma}{\sqrt{n}}\right) &\leq C_{A,\alpha}(n) \leq \omega\left(2Z_{1-\alpha/2}\frac{\tau}{\sqrt{n}}\right), \end{aligned}$$

for all  $n$ . Hence, the modulus of continuity  $\omega(\varepsilon) \asymp A\varepsilon^r$  as  $\varepsilon \rightarrow 0$  iff

$$\begin{aligned} R_A(n) &\asymp n^{-r}, \\ \Lambda_A(n) &\asymp n^{-r/2}, \\ C_{A,\alpha}(n) &\asymp n^{-r/2}. \end{aligned}$$

In other words, determining the rate at which the minimax risk converges to zero as  $n \rightarrow \infty$  is *completely equivalent* to determining the exponent in the modulus of continuity of  $T$  over  $\mathcal{F}$ .

*Minimax risk in density estimation* [Donoho and Liu (1991)]. Suppose we observe  $X_i, i = 1, \dots, n$ , independent and identically distributed  $F$ , where the distribution  $F$  is unknown but assumed to have a density  $f = F'$  in a class  $\mathcal{F}$ , and we wish to estimate the linear functional  $T(f) = f(0)$ . Suppose  $\mathcal{F}$  is the class of decreasing, Lipschitz densities defined by

$$\begin{aligned} \mathcal{F} = \left\{ f: 1 \geq f(-1) \geq f(t) \geq f(1) \geq 0, \quad \text{for } t \in [-1, 1], \right. \\ \left. \text{and } 0 \leq f(t) - f(t+h) \leq Ch, \quad \text{for } h > 0, \text{ and } \int_{-1}^1 f = 1 \right\}. \end{aligned}$$

This class is convex asymmetric.



Donoho and Liu (1991) studied this problem from the minimax mean squared error viewpoint. Their calculations, combined with Section 8 of this paper, give results for other performance measures. Some terminology follows: an affine procedure is any rule of the form  $e + (nh_n)^{-1} \sum_i k(X_i/h_n)$ , a "kernel estimate". Let  $\Lambda_A(n)$  denote the minimax expected absolute error for estimating  $T$  by an affine procedure using  $n$  observations, and define the confidence statement measure  $C_{A,\alpha}(n)$  similarly.

**THEOREM 4.** *The triangular kernel  $k(t) = (1 - |t|)_+$  asymptotically minimax among kernel estimates for estimating  $T(f) = f(0)$  over  $\mathcal{F}$  for each of our loss functions, when the bandwidth is chosen appropriately. For absolute error loss, the optimal choice of bandwidth is*

$$h_n = v_{1,2/3}^{2/3} 6^{1/3} C^{-2/3} n^{-1/3},$$

and to get asymptotic minimaxity for  $(1 - \alpha)$  confidence statement length, we should use bandwidth

$$h_n = v_{\alpha,2/3}^{2/3} 6^{1/3} C^{-2/3} n^{-1/3}.$$

Moreover, the optimally tuned triangular kernel is within 23% of minimax (absolute error loss) and 19% of minimax (95% confidence statement loss). Finally,

$$\Lambda_A(n) = \xi_{1,A}(2/3) (6C)^{1/3} n^{-1/3} (1 + o(1))$$

and

$$C_{A,\alpha}(n) = \xi_{\alpha,A}(2/3) (6C)^{1/3} n^{-1/3} (1 + o(1)).$$

The results of Section 8 play an integral role in this result, which explains the appearance of the constant  $v$  and  $\xi$ , and the figures 19% and 23%. For this application it is important that our theorems hold for convex, asymmetric  $\mathbf{X}$ .

The approach is, of course, not limited to this one example; it can easily give minimax risk for  $l_1$  and confidence statement loss in many other problems of density estimation, in particular, all those discussed in Donoho and Liu (1991).

*Minimax quadratic estimation of a quadratic functional* [Donoho and Nussbaum (1990)]. Suppose we have nonparametric regression data  $y_i = f(t_i) + z_i$ , with the  $t_i$  equispaced on  $[-\pi, \pi]$ . We are interested in estimating the quadratic functional  $\int_{-\pi}^{\pi} (f^{(k)}(t))^2 dt$  using a quadratic rule  $e + \langle \mathbf{y}, \mathbf{M}\mathbf{y} \rangle$ . We know a priori that  $f^{(l)}$  is periodic and absolutely continuous for  $0 \leq l < m$  and that  $\int_{-\pi}^{\pi} (f^{(m)}(t))^2 dt \leq 1$ .

While this is a quadratic, rather than linear, problem, Donoho and Nussbaum exhibit a transformation which allows a solution using the methods developed here. They derive a formula for the asymptotic minimax risk among quadratic estimates and a formula for a computationally effective quadratic estimator attaining this asymptotic minimax risk. Theorems 1 and 2 and their corollaries, play a key role in this solution. For this application it is crucial that these theorems hold for *asymmetric* convex sets  $\mathbf{X}$ .

**10. Optimal recovery.** Our inequalities between minimax risk and the modulus of continuity have a deeper explanation—they express a close connection between the problem of optimal recovery and that of statistical estimation.

Suppose that we have data of the form (1), where  $\mathbf{z}$  is assumed to satisfy only  $\|\mathbf{z}\| \leq \varepsilon$ . Our measure of performance is the worst-case error:

$$E(\widehat{L}, \mathbf{x}) = \sup_{\|\mathbf{z}\| \leq \varepsilon} |\widehat{L}(\mathbf{y}) - L(\mathbf{x})|.$$

This problem setting has been treated by many authors, for example, Micchelli (1975), Micchelli and Rivlin (1977), and Traub, Wasilkowski and Woźniakowski (1983, 1988). See these sources for further references, going back to the 1965 Moscow dissertation of Smolyak and the seminal paper of Golomb and Weinberger (1959).

For the sake of later sections, we pedantically spell out our approach to the problem. We are interested in the minimax error, either over affine estimators or over general nonlinear estimators. Hence, set

$$E_A^*(\varepsilon) = \inf_{\widehat{L} \text{ affine}} \sup_{\mathbf{x} \in \mathbf{X}} E(\widehat{L}, \mathbf{x}),$$

$$E_N^*(\varepsilon) = \inf_{\widehat{L}} \sup_{\mathbf{x} \in \mathbf{X}} E(\widehat{L}, \mathbf{x}).$$

We consider lower bounds based on hardest subproblem arguments. Begin with the analog of the bounded normal mean. Suppose that we are interested in estimation of the scalar  $\theta$  from data  $y = \theta + z$ ; we know that  $|\theta| \leq \tau$  and that  $|z| \leq \varepsilon$ . If  $\tau < \varepsilon$ , a minimax procedure is  $\widehat{\theta} = 0$ . If  $\tau > \varepsilon$ , a minimax procedure is to estimate  $\widehat{\theta} = y$ . If  $\tau = \varepsilon$ , any procedure  $c y$  with  $c \in [0, 1]$  is minimax. Thus, the minimax errors satisfy

$$(34) \quad e_N(\tau, \varepsilon) = e_A(\tau, \varepsilon) = \min(\tau, \varepsilon).$$

Now suppose we wish to estimate  $L(\mathbf{x})$ , for  $\mathbf{x}$  known to lie in  $[\mathbf{x}_{-1}, \mathbf{x}_1]$ . The minimax errors satisfy

$$(35) \quad E_N^*(\varepsilon; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \frac{|L(\mathbf{x}_1) - L(\mathbf{x}_{-1})|}{\|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K} e_N\left(\frac{\|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K}{2}, \varepsilon\right),$$

and so on. The difficulty of a hardest subproblem is

$$(36) \quad \sup_{\mathbf{x}_1, \mathbf{x}_{-1} \in \mathbf{X}} E_N^*(\varepsilon; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \sup_{\delta \geq 0} \frac{\omega(\delta)}{\delta} e_N\left(\frac{\delta}{2}, \varepsilon\right).$$

Now  $\omega$  is monotone, so

$$\sup_{\delta \leq 2\varepsilon} \frac{\omega(\delta)}{\delta} e_N\left(\frac{\delta}{2}, \varepsilon\right) = \sup_{\delta \leq 2\varepsilon} \frac{\omega(\delta)}{\delta} \frac{\delta}{2} = \frac{\omega(2\varepsilon)}{2},$$

and it is subadditive, so

$$\sup_{\delta \geq 2\varepsilon} \frac{\omega(\delta)}{\delta} e_N\left(\frac{\delta}{2}, \varepsilon\right) = \varepsilon \sup_{\delta \geq 2\varepsilon} \frac{\omega(\delta)}{\delta} = \frac{\omega(2\varepsilon)}{2}.$$

Hence

$$(37) \quad \sup_{\mathbf{x}_1, \mathbf{x}_{-1} \in \mathbf{X}} E_N^*(\varepsilon; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \frac{\omega(2\varepsilon)}{2}.$$

On the other hand, the nonlinear procedure

$$\begin{aligned} \widehat{L}^*(\mathbf{y}) &= \frac{1}{2} \sup \{L(\mathbf{x}): \|\mathbf{y} - K\mathbf{x}\| \leq \varepsilon, \mathbf{x} \in \mathbf{X}\} \\ &\quad + \frac{1}{2} \inf \{L(\mathbf{x}): \|\mathbf{y} - K\mathbf{x}\| \leq \varepsilon, \mathbf{x} \in \mathbf{X}\} \end{aligned}$$

[called the *central algorithm* in Traub, Wasilkowski and Woźniakowski (1983)] attains, as one can check,

$$(38) \quad \sup_{\mathbf{x} \in \mathbf{X}} E(\widehat{L}^*, \mathbf{x}) = \omega(2\varepsilon)/2.$$

So, in our terminology, the difficulty of a hardest subproblem for nonlinear estimates is equal to the difficulty of the full problem, and

$$(39) \quad E_N^*(\varepsilon) = \omega(2\varepsilon)/2.$$

Micchelli (1975) and Micchelli and Rivlin (1977) showed that if  $\mathbf{X}$  is centrosymmetric about  $\mathbf{0}$ , there exists a linear optimal algorithm. Since we have

$$(40) \quad \sup_{\mathbf{x}_1, \mathbf{x}_{-1} \in \mathbf{X}} E_A^*(\varepsilon; [\mathbf{x}_{-1}, \mathbf{x}_1]) = \frac{\omega(2\varepsilon)}{2},$$

existence of linear optimal algorithms is equivalent to the statement that the difficulty, for linear estimates, of the full problem is the same as the difficulty, for linear estimates, of a hardest one-dimensional subproblem.

It is possible to generalize the optimal recovery theorem. Assuming just convexity of  $\mathbf{X}$ , but not centrosymmetry, we can say that *affine* optimal algorithms exist.

**THEOREM 5.** *Let  $\mathbf{X}$  be convex, closed and bounded, and let  $L$  and  $K$  be well-defined. Then the difficulty of a hardest one-dimensional subproblem is equal to the difficulty of the full problem*

$$E_A^*(\varepsilon, \mathbf{X}) = \max_{\mathbf{x}_1, \mathbf{x}_{-1}} E_A^*(\varepsilon, [\mathbf{x}_{-1}, \mathbf{x}_1]).$$

Even if we assume only that  $\mathbf{X}$  is convex and that  $L$  and  $K$  are well-defined, we may still conclude that there exists an affine estimator which attains the minimax error and that

$$(41) \quad E_A^*(\varepsilon) = \frac{\omega(2\varepsilon)}{2}.$$

Combining (41) with Corollary 2 yields

$$(42) \quad (E_A^*(\sigma))^2/4 \leq R_A^*(\sigma) \leq (E_A^*(\sigma))^2.$$

In other words, if we equate noise levels  $\varepsilon = \sigma$ , then  $(E^*(\varepsilon))^2$  is approximately  $R_A(\sigma)$ . The connection between the optimal recovery and statistical estimation model will be further spelled out in Section 12.

**11. Proofs of Theorems 1 and 5.** We may assume that  $\omega(\varepsilon)$  is finite for all  $\varepsilon$  else the subadditivity of  $\omega(\varepsilon)$  implies that there are one-dimensional subproblems with arbitrarily high difficulty, in which case the theorems are trivially true.

**PROOF OF THEOREM 1.** We claim that, for a specific subproblem  $[\mathbf{x}_{-1}, \mathbf{x}_1]$  and a specific choice of  $d$ , the affine estimator

$$(43) \quad L_0(\mathbf{y}) = L(\mathbf{x}_0) + d\langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle$$

has two properties:

**PROPERTY  $P_0$ .**  $L_0$  is affine minimax for the subproblem  $[\mathbf{x}_{-1}, \mathbf{x}_1]$ .

**PROPERTY  $P_1$ .**  $L_0$  attains its worst-case risk, over all of  $\mathbf{X}$ , in the subproblem  $[\mathbf{x}_{-1}, \mathbf{x}_1]$ .

Thus, the difficulty of the full problem for this particular estimator is no more than the difficulty of the subproblem, but as this estimator is affine minimax for the subproblem, it must also be affine minimax for the full problem. We conclude that the subproblem is a hardest one-dimensional subproblem and that the difficulty of the subproblem is equal to that of the full problem.

The proof shows that the set  $\mathcal{E}_0$  of estimators of the form (43) which have Property  $P_0$  and the set  $\mathcal{E}_1$  of estimators which have Property  $P_1$  have nonempty intersection. This shows there is an estimator (43) with both properties and implies Theorem 1.

To begin, we state without proof the following result, which follows by easy calculations and standard results in Rockefellar (1970).

**LEMMA 3.** *The modulus of continuity of an affine functional over a convex set is a concave function of  $\varepsilon$ . It is nonnegative and, if it is bounded on an*

interval  $[0, \varepsilon]$ , it is locally Lipschitz continuous at  $\delta$  to that interval. It has a bounded superdifferential  $\partial\omega(\delta)$  at each  $\delta$  interior to that interval. That is, let  $\partial\omega(\delta)$  denote the set of slopes of lines passing through  $(\delta, \omega(\delta))$  which lie above the graph of  $\omega(\delta)$ ,

$$\partial\omega(\delta) = \{d: \omega(\varepsilon) \leq \omega(\delta) + d(\varepsilon - \delta), \varepsilon > 0\}.$$

Then  $\partial\omega(\delta)$  is a nonempty, closed, bounded, convex subset of  $\mathbf{R}$ . Viewed as a curve in the plane  $(\varepsilon, \partial\omega(\varepsilon))_{\varepsilon > 0}$  is a connected, monotone nonincreasing curve.

Define the set  $\Gamma_1(\varepsilon) = \varepsilon\partial\omega(\varepsilon)/\omega(\varepsilon)$ . Then  $\Gamma_1 = \bigcup_{\varepsilon > 0} (\{\varepsilon\} \times \Gamma_1(\varepsilon))$  is a subset of  $[0, \infty] \times [0, 1]$ . By Lemma 3,  $(\varepsilon, \partial\omega(\varepsilon))_{\varepsilon > 0}$  and hence also  $\Gamma_1$  are connected subsets of  $\mathbf{R}^2$ .

Under our assumptions,

$$\varepsilon^* = \sup \|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K \leq \infty.$$

As  $\omega(\varepsilon) = \omega(\varepsilon^*)$ , for  $\varepsilon \geq \varepsilon^*$ , we have  $0 \in \partial\omega(\varepsilon), \varepsilon \geq \varepsilon^*$ . Thus

$$(44) \quad 0 \in \Gamma_1(\varepsilon^*).$$

Also, as  $L$  is nonconstant (otherwise the theorem is trivially true),

$$(45) \quad \liminf_{\varepsilon \rightarrow 0} \frac{\inf \partial\omega(\varepsilon)}{\omega(\varepsilon)} > \text{const.} > 0.$$

For the criterion of interest, define

$$\Gamma_0 = \bigcup_{\varepsilon > 0} (\{\varepsilon\} \times \{c_0(\varepsilon/2, \sigma; \cdot)\}).$$

As  $c_0$  is monotone increasing and continuous for whichever of the three criteria we have chosen [see Section 2, equations (6), (8) and (11)],  $\Gamma_0$  is a connected monotone increasing curve of  $\mathbf{R}^2$ .

It follows from (45) and (12) that, for all sufficiently small  $\varepsilon$ ,

$$\inf \Gamma_1(\varepsilon) > c_0(\varepsilon/2, \sigma; \cdot).$$

However, by (44), (6), (8) and (11),

$$0 = \inf \Gamma_1(\varepsilon^*) < c_0(\varepsilon^*/2, \sigma; \cdot).$$

Hence, by connectedness of  $\Gamma_1$  and  $\Gamma_0$ , these two curves “cross”:

$$\Gamma_1 \cap \Gamma_0 \neq \emptyset.$$

The crossing of the curves implies that, for some  $\varepsilon_0 \in [0, \varepsilon^*]$ ,

$$(46) \quad c_0(\varepsilon_0/2, \sigma; \cdot) \in \Gamma_1(\varepsilon_0).$$

Let  $(\mathbf{x}_{-1}, \mathbf{x}_1)$  attain the modulus at  $\varepsilon_0$ . Define

$$(47) \quad L_0(\mathbf{y}) = L(\mathbf{x}_0) + c_0 \left( \frac{\varepsilon_0}{2}, \sigma; \cdot \right) \frac{\omega(\varepsilon_0)}{\varepsilon_0} \langle \mathbf{w}_0, \mathbf{y} - \mathbf{x}_0 \rangle.$$

We claim that  $L_0$  has the two properties desired.

Property  $P_0$  follows from the use of  $c_0$  in (47) and from the discussion in Section 6.

As for Property  $P_1$ , let  $\mathcal{L}_{\mathbf{x}}V$  denote the probability law of the random variable  $V$  when  $\mathbf{x}$  is the true object. Set

$$\mathcal{P}(L_0, [\mathbf{x}_{-1}, \mathbf{x}_1]) = \{ \mathcal{L}_{\mathbf{x}}(L_0(\mathbf{y}) - L(\mathbf{x})): \mathbf{x} \in [\mathbf{x}_{-1}, \mathbf{x}_1] \}$$

and

$$\mathcal{P}(L_0, \mathbf{X}) = \{ \mathcal{L}_{\mathbf{x}}(L_0(\mathbf{y}) - L(\mathbf{x})): \mathbf{x} \in \mathbf{X} \}.$$

We wish to show that the full problem is no harder than the subproblem, so that

$$(48) \quad \mathcal{P}(L_0, \mathbf{X}) = \mathcal{P}(L_0, [\mathbf{x}_{-1}, \mathbf{x}_1]).$$

Note that

$$\mathcal{L}_{\mathbf{x}}(L_0(\mathbf{y}) - L(\mathbf{x})) = N(\text{Bias}(L_0, \mathbf{x}), d^2 \sigma^2).$$

Thus it is enough to show that

$$(49) \quad |\text{Bias}(L_0, \mathbf{x}_1)| \geq |\text{Bias}(L_0, \mathbf{x})|, \quad \mathbf{x} \in \mathbf{X}.$$

We now exploit the intersection condition (46). By definition of  $\Gamma_1$ , this implies that  $L_0$  is of the form (43), where  $d \in \partial\omega(\varepsilon_0)$ . The following result is fundamental to the article and is proved in Section 14.

**LEMMA 4.** *Let the modulus be attained at  $\varepsilon_0$  by  $(\mathbf{x}_{-1}, \mathbf{x}_1)$ , and let  $d \in \partial\omega(\varepsilon_0)$ . Suppose that labels are chosen so that  $L(\mathbf{x}_1) > L(\mathbf{x}_{-1})$ . Then, for every  $\mathbf{x} \in \mathbf{X}$ ,*

$$(50) \quad \begin{aligned} L(\mathbf{x}) - L(\mathbf{x}_1) &\leq d \langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle, \\ L(\mathbf{x}) - L(\mathbf{x}_{-1}) &\geq d \langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_{-1} \rangle. \end{aligned}$$

To apply this lemma, note that  $\text{Bias}(L_0, \mathbf{x})$  is an affine functional with  $\text{Bias}(L_0, \mathbf{x}_0) = 0$ . Thus  $\text{Bias}$  takes opposite signs at  $\mathbf{x}_1$  and  $\mathbf{x}_{-1}$ . Our choice of labels  $L(\mathbf{x}_1) > L(\mathbf{x}_{-1})$  forces  $\text{Bias}(L_0, \mathbf{x}_1) \leq 0$ . Then, using (50),

$$\begin{aligned} \text{Bias}(L_0, \mathbf{x}_1) - \text{Bias}(L_0, \mathbf{x}) &= L_0(K\mathbf{x}_1) - L_0(K\mathbf{x}) - L(\mathbf{x}_1) + L(\mathbf{x}) \\ &= d \langle \mathbf{w}_0, K(\mathbf{x}_1 - \mathbf{x}) \rangle - L(\mathbf{x}_1) + L(\mathbf{x}) \\ &\leq 0 \quad (\text{by 50}). \end{aligned}$$

On the other hand, our assumption forces  $\text{Bias}(L_0, \mathbf{x}_{-1}) \geq 0$ , and again using (50),

$$\begin{aligned} \text{Bias}(L_0, \mathbf{x}_{-1}) - \text{Bias}(L_0, \mathbf{x}) &= d\langle \mathbf{w}_0, K(\mathbf{x}_{-1} - \mathbf{x}) \rangle - L(\mathbf{x}_{-1}) + L(\mathbf{x}) \\ &\geq 0. \end{aligned}$$

Finally, as  $|\text{Bias}(L_0, \mathbf{x}_1)| = |\text{Bias}(L_0, \mathbf{x}_{-1})|$ , we have (49) and the proof is complete.  $\square$

**PROOF OF THEOREM 5.** Theorem 5 makes two statements. The first statement is analogous to Theorem 1 and will be proven in a moment. The second statement, which is analogous to Theorem 2, follows by applying the proof of Theorem 2 word-for-word, only using the auxiliary function  $m(a, b) = a + \varepsilon|b|$ .

The proof of Theorem 1 is also a proof of the optimal recovery theorem. To see this, define  $\Gamma_0 = (0, 2\varepsilon) \times \{0\} \cup \{2\varepsilon\} \times [0, 1] \cup (2\varepsilon, \infty) \times \{1\}$ , and define  $\Gamma_1$  exactly as before. Then, just as before,  $\Gamma_0 \cap \Gamma_1 \neq \emptyset$ ; in fact the two sets intersect at  $\varepsilon_0 = 2\varepsilon$ . Let  $(\mathbf{x}_{-1}, \mathbf{x}_1)$  attain the modulus at  $\varepsilon_0$  and let  $d \in \partial\omega(\varepsilon_0)$ . Define the rule

$$L_0(\mathbf{y}) = L(\mathbf{x}_0) + d\langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle.$$

We claim that  $L_0$  has two key properties:

**PROPERTY  $P_0$ .** It is minimax for the subproblem  $[\mathbf{x}_{-1}, \mathbf{x}_1]$ .

**PROPERTY  $P_1$ .** It attains its worst error over all  $\mathbf{X}$  in the subproblem  $[\mathbf{x}_{-1}, \mathbf{x}_1]$ .

As a result, in the optimal recovery setting, the estimator is minimax for the full problem, the family  $[\mathbf{x}_{-1}, \mathbf{x}_1]$  is a hardest subproblem, the difficulty of the hardest subproblem is equal to that of the full problem. This establishes the first part of Theorem 5, as desired.

We verify that  $L_0$  has Property  $P_0$ . For estimating  $\theta = \langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_0 \rangle$ , from  $\mathbf{y} = \langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle$  we have that  $|\theta| \leq \varepsilon$  and that  $z = \mathbf{y} - \theta$  has  $|z| \leq \varepsilon$  also. By an earlier comment, any estimator  $c\mathbf{y}$  with  $c \in [0, 1]$  is minimax for estimating  $\theta$ . It follows that any estimator

$$\widehat{L}(\mathbf{y}) = L(\mathbf{x}_0) + c \frac{\omega(2\varepsilon)}{2\varepsilon} \langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle,$$

with  $c \in [0, 1]$  is minimax for estimating  $L$  in the subproblem. By monotonicity and subadditivity of  $\omega$ , any element  $d \in \partial\omega(2\varepsilon)$  satisfies  $0 \leq d \leq \omega(2\varepsilon)/(2\varepsilon)$ , that is, we can write  $d = c\omega(2\varepsilon)/2\varepsilon$  with  $c \in [0, 1]$ . So our choice of  $d$  makes  $L_0$  a minimax estimator in the subproblem.

Finally, we verify that  $L_0$  has Property  $P_1$ . Write

$$\begin{aligned} L_0(\mathbf{y}) - L(\mathbf{x}) &= L_0(K\mathbf{x}) - L(\mathbf{x}) + L_0(\mathbf{y}) - L_0(K\mathbf{x}) \\ &= \text{Bias}(L_0, \mathbf{x}) + d\langle \mathbf{w}_0, \mathbf{z} \rangle. \end{aligned}$$

Picking the noise  $\mathbf{z}$  aligned with  $\mathbf{w}_0$  [i.e.,  $\langle \mathbf{w}_0, \mathbf{z} \rangle = \text{sgn}(\text{Bias}(L_0, \mathbf{x})) \cdot \varepsilon$ ], we see that

$$\sup_{\|\mathbf{z}\| \leq \varepsilon} |L_0(\mathbf{y}) - L(\mathbf{x})| = |\text{Bias}(L_0, \mathbf{x})| + d\varepsilon.$$

In this expression only Bias depends on  $\mathbf{x}$ . In this case (50) again implies (49), and so

$$|\text{Bias}(L_0, \mathbf{x}_1)| = \sup_{\mathbf{x} \in \mathbf{X}} |\text{Bias}(L_0, \mathbf{x})|,$$

which implies

$$\sup_{\mathbf{x} \in \mathbf{X}} |E(L_0, \mathbf{x})| = \sup_{\mathbf{x} \in [\mathbf{x}_{-1}, \mathbf{x}_1]} E(L_0, \mathbf{x}),$$

that is, property  $P_1$ .  $\square$

**11.1. Other situations.** The proof of Theorem 1 may obviously be adapted to other performance measures in the statistical estimation problem. The fundamental issue is that the minimax affine estimator for  $|\theta| \leq \tau$  be linear and that the coefficient  $c_0$  which furnishes the minimax affine estimator be continuous in  $\tau$ . This will hold for many other loss functions.

In other words, a single proof idea handles various performance criteria in the statistical estimation problem and also the optimal recovery problem.

**12. Correspondence theorem.** The preceding proof also establishes the following.

**COROLLARY 5.** *Let the assumptions of Theorem 1 hold. Choose any one of the three performance criteria in the statistical estimation problem. Let a hardest subfamily for affine estimates under that criterion have length  $\varepsilon_0$ . Then the estimator*

$$(51) \quad L_0(\mathbf{y}) = L(\mathbf{x}_0) + d\langle \mathbf{w}_0, \mathbf{y} - K\mathbf{x}_0 \rangle,$$

where  $d \in \partial\omega(\varepsilon_0)$  is an affine minimax estimator for the statistical problem and is also an optimal algorithm for the optimal recovery problem at noise level  $\varepsilon = \varepsilon_0/2$ .

In words, if we calibrate noise levels so that the hardest one-dimensional subproblems for optimal recovery and for statistical estimation have the same length, then they have optimal estimators in common.

Here is a simple illustration. Speckman (1979) proved the following result, which expresses the minimaxity of cubic smoothing splines. [For extensions of this result, see Li (1982).]



**THEOREM 6** [Speckman (1979)]. *Let  $y_i = f(t_i) + z_i$ ,  $i = 1, \dots, n$ , where  $t_i \in [0, 1]$  and the  $z_i$  are i.i.d.  $N(0, \sigma^2)$  and where the function  $f$  is known to satisfy  $\int_0^1 (f''(t))^2 dt \leq C^2$ . Let  $g_\mu$  be the solution to*

$$\min_g \sum_i (g(t_i) - y_i)^2 + \mu \int_0^1 (g''(t))^2 dt.$$

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

$$L_0(\mathbf{y}) = L(g_\mu)$$

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

Now consider the associated optimal recovery problem, with observations  $y_i = f(t_i) + z_i$ ,  $i = 1, \dots, n$ ,  $\int_0^1 (f''(t))^2 dt \leq C^2$ , where now the  $z_i$  are nonstochastic and are known only to satisfy  $\sum_i z_i^2 \leq \varepsilon^2$ . Speckman's theorem and our Corollary 5 imply that, for some  $\mu_{\text{or}} = \mu_{\text{or}}(\varepsilon, C)$ , the cubic-spline-based estimator  $L_0(\mathbf{y}) = L(g_{\mu_{\text{or}}})$  is an optimal recovery algorithm—a fact due, essentially, to Schoenberg (1964a, b). In other words, Speckman's theorem implies Schoenberg's. And, of course, vice versa.

In the other direction, consider the prototypical problem of optimal recovery: estimating the integral  $L(f) = \int_0^1 f(t) dt$  from data  $y_i = f(t_i) + z_i$ ,  $i = 1, \dots, n$ . Here we take  $t_i = (i - 0.5)/n$ . We know a priori only that  $f$  belongs to  $\mathcal{F} = \{f: |f(s) - f(t)| \leq C|s - t|\}$  and the nonstochastic noise satisfies  $\sum_i z_i^2 \leq \varepsilon^2$ . Then the modulus is attained with  $f_{-1} = -f_1$ , where  $f_1$  is the sawtooth function

$$f_1(t) = \min_i \left( \frac{\varepsilon}{\sqrt{n}} + C|t - t_i| \right).$$

We get  $\omega(\varepsilon) = \varepsilon/\sqrt{n} + C/(n - 1)$  and that  $L_0(\mathbf{y}) = (1/n)\sum_{i=1}^n y_i$  is an optimal algorithm, for each  $\varepsilon > 0$ . Turning to the associated statistical estimation problem, where the noise is i.i.d.  $N(0, \sigma^2)$ , we note that the formula  $\sup_\varepsilon (\omega(\varepsilon)/\varepsilon)^2 \rho_A(\varepsilon/2, \sigma)$  has its maximum at some  $\varepsilon_0 \in (0, \infty)$ , and it follows that the associated  $L_0$  is minimax affine for the statistical estimation problem. A side calculation gives  $R_A^*(\sigma) = C^2/(16n^2) + \sigma^2/n$ .

In short, if a problem has been solved in one of the two literatures, that solution may be considered as a solution of the problem in the other literature.

We also have correspondence between the solutions to the statistical estimation problem with different loss criteria.

**COROLLARY 6.** *Under the assumptions of Theorem 1 and 2, there exist monotone, continuous functions  $\sigma_1(\sigma)$  and  $\sigma_2(\sigma)$  (which depend on  $L, K$  and  $\mathbf{X}$ ) so that an affine estimator can be found which is affine minimax for squared error loss at noise level  $\sigma$ , for absolute error loss at  $\sigma_1(\sigma)$  and for the confidence statement criterion at  $\sigma_2(\sigma)$ .*

In situations where asymptotics as  $\sigma \rightarrow 0$  make sense, of course, Corollary 4 shows that we must have the relationships

$$\begin{aligned}\sigma_1 &= \frac{v_{2,r}}{v_{1,r}}\sigma(1 + o(1)), \\ \sigma_\alpha &= \frac{v_{2,r}}{v_{\alpha,r}}\sigma(1 + o(1)).\end{aligned}$$

Speckman's theorem, quoted previously, shows that cubic-spline-based estimates of a linear functional are, under certain assumptions, minimax among linear estimates under squared error loss. Corollary 6 says that the same estimates will also be minimax for absolute error and confidence statements measures, at certain noise levels. For example, with absolute error loss, let  $\sigma_1^{-1}(\sigma)$  denote the solution to  $\sigma_1(s) = \sigma$ . If the true noise level is  $\sigma$ , we set  $\mu_1 = (\sigma_1^{-1}(\sigma)/C)^2$  and set  $L_0(\mathbf{y}) = L(g_{\mu_1})$ ; this is affine minimax for absolute error loss.

Even without recalibration, the solution to one problem furnishes a fairly good solution to any one of the others. For example, suppose we know how to design an affine optimal algorithm  $L_0$ , for the optimal recovery model at noise level  $\varepsilon$ . We pick  $\varepsilon = \sigma$  and we apply the resulting  $L_0$  in a statistical estimation problem with noise level  $\sigma$ . With respect to the squared error loss criterion, a simple analysis will show that

$$\sup_{\mathbf{x} \in \mathbf{X}} E(L_0(\mathbf{y}) - L(\mathbf{x}))^2 \leq \omega(2\sigma)^2/4,$$

whereas, by Theorem 2,  $R_N^*(\sigma) \geq \omega(\sigma)^2/5$ . Hence the optimal algorithm, although designed for deterministic noise, is within a factor of about 4 of minimax in MSE for the statistical estimation problem.

Much the same story holds for other performance measures. Consider confidence statement length. Set  $\varepsilon = Z_{1-\alpha/2}\sigma$ , and obtain an  $L_0$  which is an affine optimal algorithm for deterministic noise of norm  $\varepsilon$ . Apply this estimator in the statistical estimation problem with noise level  $\sigma$ . One calculates that the interval

$$L_0(\mathbf{y}) \pm \omega(2Z_{1-\alpha/2}\sigma)/2$$

covers the true  $L(\mathbf{x})$  with at least  $1-\alpha$  coverage probability for any  $\mathbf{x} \in \mathbf{X}$ . Thus this optimal algorithm for dealing with deterministic noise may be used to design a valid fixed-width  $1-\alpha$  confidence interval. Moreover, by our preceding results, any fixed-width interval which is a measurable function of the data and which has at least  $1-\alpha$  coverage probability must be at least a factor  $Z_{1-\alpha}/Z_{1-\alpha/2}$  as long. So the interval is within a few percent of efficient.

### 13. Discussion.

13.1. *Nonwhite noise.* A certain class of problems with nonwhite noise can be mapped onto the present one. If our observations (1) have  $\mathbf{z}$  with nonwhite

covariance, and if the covariance is an operator with a bounded inverse, then we can transform the observations via  $\mathbf{y}' = \Sigma^{-1/2}\mathbf{y}$ , giving data

$$\mathbf{y}' = K'\mathbf{x} + \mathbf{z},$$

where now  $\mathbf{z}$  is white and  $K' = \Sigma^{-1/2}K$ . Proceeding as before, we define the modulus with respect to the seminorm defined by  $K'$ , and the formulas from before all continue to apply. In this way we could recapture results of not only Ibragimov and Has'minskii (1987), but others as well, since our results allow indirect observations ( $K \neq I$ ), asymmetry of  $\mathbf{X}$ , various loss functions and so on. Also, we could demonstrate a close mathematical connection between estimation in nonwhite noise and in the optimal recovery model with constraint  $\langle \mathbf{z}, \Sigma^{-1}\mathbf{z} \rangle \leq \varepsilon^2$ .

**13.2. Nonlinear functionals.** We have shown here a close connection between the modulus of continuity and the difficulty of estimation of linear functionals from incomplete data with Gaussian noise. The connection between the modulus and difficulty of estimation need not persist when we consider estimation of nonlinear functionals. Some basic information about estimation of nonlinear functionals in white noise is given in, for example Ibragimov, Nemirovskii and Has'minskii (1986) and Fan (1991). Donoho and Nussbaum (1990) show that in such problems the minimax risk may go to zero much more slowly than the rate at which the modulus goes to zero.

In contrast, in the optimal recovery model, under very mild conditions, the modulus of continuity measures the difficulty of estimation quite precisely for general nonlinear functionals, that is, the "central algorithm" described in Section 10 can be used for general nonlinear functionals; it gives the worst-case error  $\omega(2\varepsilon)/2$  for quite a wide variety of situations, and this can be shown to be the minimax error. Compare Traub, Wasilkowski and Woźniakowski (1983, 1988).

Thus the connection we are describing between optimal recovery and statistical estimation need not persist when we consider estimating nonlinear functionals.

However, the results of this paper are still useful in nonlinear cases, as we have suggested in Section 9.3.

**13.3. Estimating the whole object.** If, rather than estimating just a single linear functional of the object, we were estimating the whole object  $\mathbf{x}$  with, say,  $l_2$  norm loss, statistical estimation and optimal recovery would no longer, in general, have a close connection. In general, minimax linear statistical estimation is connected with minimizing the Hilbert–Schmidt norm of the estimator, subject to a side constraint on the norm of the bias, while linear optimal recovery is connected with minimizing the operator norm of the estimator, subject to a constraint on the norm of the bias. Of course for estimators with one-dimensional range, that is, *functionals*, Hilbert–Schmidt and operator norms are the same, which explains why the connection holds for one-dimensional functionals and not for more general objects.

13.4. *Other norms.* The basic theorem of linear optimal recovery is not restricted to use of the  $l_2$  norm in specifying the constraint  $\|z\| \leq \varepsilon$ . For example, Micchelli and Rivlin (1977) showed that one can use any Banach space norm for the error norm, and there will still exist an optimal linear algorithm under quite general conditions. However, optimal recovery under these other error norms does not necessarily relate to statistical estimation.

One exception is when one has in the optimal recovery model an  $l_p$  error norm  $\|z\|_{l_p} \leq \varepsilon$ , for  $p \in [2, \infty]$ . This corresponds to linear statistical estimation with a white symmetric stable noise of index  $\alpha$  conjugate to  $p$  ( $1/p + 1/\alpha = 1$ ). Of course,  $p = \alpha = 2$  is the case we have covered in this paper: the case  $p = \infty, \alpha = 1$  might be an interesting one to consider. It connects deterministic noise small in supremum norm with stochastic noise following a Cauchy distribution.

14. **Proofs.** Note that we omit detailed proofs of Corollaries 1, 2, 5 and 6; these follow from Theorems 1 and 2 and from other information, such as the discussion of Section 2 or the proof of Theorem 2.

PROOF OF LEMMA 2. The whole result follows once we know that the modulus of continuity is attained. For, by Lemma 4, when the modulus is finite, it is concave and continuous; the suprema over  $\varepsilon$  in the formulas are really therefore suprema of continuous functions of  $\varepsilon$ . Moreover, under the assumptions, only a finite range  $[0, \varepsilon^*]$  need be considered, where  $\varepsilon^* = \sup_{\mathbf{x}_1, \mathbf{x}_{-1}} \|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K < \infty$ . A continuous function on a compact set takes on its maximum, and so in each of the formulas the supremum is attained at some  $\varepsilon_0$ . The family that attains the modulus at that  $\varepsilon_0$  is the hardest one-dimensional subfamily for that criterion.

Suppose now that  $\{\mathbf{x}_{-1,n}, \mathbf{x}_{1,n}\}$  is a sequence of subfamilies of  $\mathbf{X}$ , with  $\|\mathbf{x}_{1,n} - \mathbf{x}_{-1,n}\|_K \leq \varepsilon$  but  $L(\mathbf{x}_{1,n}) - L(\mathbf{x}_{-1,n}) \rightarrow \omega(\varepsilon)$ . By hypothesis,  $\mathbf{X}$  is weakly compact, and along a subsequence  $\mathbf{x}_{1,n}$  and  $\mathbf{x}_{-1,n}$  have weak limits  $\mathbf{x}_1$  and  $\mathbf{x}_{-1}$  in  $\mathbf{X}$ . As  $K$  is well-defined, the restriction to  $\mathbf{V} = \mathbf{X} - \mathbf{X}$  of the seminorm  $J(\mathbf{v}) = \|\mathbf{v}\|_K$  is continuous for  $l_2$ -convergence;  $J$  is also convex. By (52) in Lemma 5, the seminorm is lower semicontinuous for weak convergence. It follows that  $\|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K \leq \varepsilon$ . As  $L$  is well-defined, the restriction of the functional  $J(\mathbf{v}) = L(\mathbf{v})$  to  $\mathbf{V}$  is both convex and concave and is continuous for  $l_2$ -convergence. Hence, by (52), both  $L(\mathbf{x}_{1,n}) \rightarrow L(\mathbf{x}_1)$  and  $L(\mathbf{x}_{-1,n}) \rightarrow L(\mathbf{x}_{-1})$ , as  $n \rightarrow \infty$ . Thus

$$L(\mathbf{x}_1) - L(\mathbf{x}_{-1}) = \omega(\varepsilon);$$

the modulus is attained by  $(\mathbf{x}_{-1}, \mathbf{x}_1)$ .  $\square$

LEMMA 5. Let  $J(\mathbf{v})$  be a convex functional on a norm-closed, norm-bounded convex set  $\mathbf{V}$  which is continuous for strong  $l_2$ -convergence. Let  $\mathbf{v}_n$  be a sequence of elements in  $\mathbf{V}$  converging weakly to  $\mathbf{v}$ . Then

$$(52) \quad J(\mathbf{v}) \leq \liminf_{n \rightarrow \infty} J(\mathbf{v}_n).$$

PROOF. Define  $V_{\leq j} = \{\mathbf{v} \in \mathbf{V}: J(\mathbf{v}) \leq j\}$ . This set is convex and bounded. Because  $J$  is continuous, the set is strongly closed, hence (by convexity) weakly closed.

Suppose that along a subsequence, for all sufficiently large  $n$ ,  $J(\mathbf{v}_n) \leq j$ , that is,  $\mathbf{v}_n \in V_{\leq j}$ . As  $V_{\leq j}$  is weakly closed, the weak limit  $\mathbf{v} \in V_{\leq j}$ . Hence  $J(\mathbf{v}) \leq j$ . Inequality (52) follows.  $\square$

PROOF OF THEOREM 2. Let  $\text{MaxRisk}(\widehat{L}, \mathbf{X})$  denote the *supremum* risk of  $\widehat{L}$  over  $\mathbf{X}$ , according to whichever loss criterion we are considering. We note that

$$(53) \quad \text{MaxRisk}(\widehat{L}, \mathbf{X}) = m(\text{MaxBias}(\widehat{L}, \mathbf{X}), \|\widehat{L}'\|),$$

where  $\text{MaxBias}$  denotes the *supremum* of the absolute value of the bias of  $\widehat{L}$  over  $\mathbf{X}$ , and  $\widehat{L}'$  is the homogeneous linear part of  $\widehat{L}$ . Here the function  $m$  depends on the loss criterion. For example, if loss is squared error,  $m(a, b) = a^2 + \sigma^2 b^2$ . In any event,

$$(54) \quad m(a, b) \text{ is a continuous function, monotone increasing in each argument separately.}$$

We now explain why closedness of  $\mathbf{X}$  is not necessary for the formulas to work. Indeed, the minimax affine risk is unaffected by taking the  $l_2$ -closure. For an affine estimator  $\widehat{L}$  with finite minimax risk,

$$(55) \quad \text{MaxRisk}(\widehat{L}, \mathbf{X}) = \text{MaxRisk}(\widehat{L}, \text{cl}(\mathbf{X})).$$

Indeed, by (53) the  $l_2$  norm of the homogeneous linear part  $\widehat{L}'$  of  $\widehat{L}$  is finite. As  $K$  is well-defined,  $\widehat{L}(K\mathbf{x})$  is a uniformly continuous function of  $\mathbf{x} \in \mathbf{X}$ , with a unique continuous extension to  $\text{cl}(\mathbf{X})$ . As  $L$  is well-defined, it too is a uniformly continuous function of  $\mathbf{x} \in \mathbf{X}$  with unique extension. We conclude that  $\text{Bias}(\widehat{L}, \mathbf{x}) = \widehat{L}(K\mathbf{x}) - L(\mathbf{x})$  is a uniformly continuous function of  $\mathbf{x}$  with unique extension, and so

$$\text{MaxBias}(\widehat{L}, \mathbf{X}) = \text{MaxBias}(\widehat{L}, \text{cl}(\mathbf{X})).$$

From this and from (53) and (54), (55) follows.

We now explain why norm-boundedness is unnecessary for the formulas to work. We assume that the supremum,  $M$ , say, of minimax risks of all one-dimensional subfamilies is finite; otherwise there is nothing to prove. Let  $\mathbf{X}_k$  denote the set  $\text{cl}(\mathbf{X} \cap B(0, k))$  (restricting attention to only those  $k \geq k_0$  for which the set is nonempty).  $\mathbf{X}_k$  is a closed, convex, norm-bounded set. By Theorem 1, there exists an affine estimator  $L_k$ , say, which is affine minimax for estimation of  $L$  over  $\mathbf{X}_k$ . Let  $M_k$  denote the affine minimax risk. Fix  $x_0$  in every  $\mathbf{X}_k$ ,  $k \geq k_0$ , and set  $l_k = L_k(Kx_0)$ . Let  $L'_k$  be the homogeneous linear part of  $L_k$ .

The sequence of norms ( $\|L'_k\|$ ) is bounded because

$$\begin{aligned} M &\geq M_k = \text{MaxRisk}(L_k, \mathbf{X}_k) \\ &= m(\text{MaxBias}(L_k, \mathbf{X}_k), \|L'_k\|) \\ &\geq m(0, \|L'_k\|). \end{aligned}$$

We can extract a weak limit  $L'_0$  from the norm-bounded sequence  $(L'_k)$ . By weak semicontinuity of the norm,

$$(56) \quad \|L'_0\| \leq \liminf_k \|L'_k\|,$$

where  $k$  is along the subsequence which gives rise to  $L'_0$ .

The sequence  $(l_k)$  is bounded because

$$\begin{aligned} M &\geq M_k = \text{MaxRisk}(L_k, \mathbf{X}_k) = m(\text{MaxBias}(L_k, \mathbf{X}_k), \|L'_k\|) \\ &\geq m(|\text{Bias}(L_k, \mathbf{x}_0)|, 0) = m(|l_k - L(\mathbf{x}_0)|, 0). \end{aligned}$$

We can extract, along a further subsequence of the initial subsequence, a limit  $l_0$ .

Define  $L_0(\mathbf{y}) = l_0 + L'_0(\mathbf{y} - \mathbf{x}_0)$ . This is affine and has  $\text{Bias}(L_0, \mathbf{x}) = L_0(K\mathbf{x}) - L(\mathbf{x})$ . Pick any  $\mathbf{x} \in \mathbf{X}$ . Now

$$\text{Bias}(L_0, \mathbf{x}) - \text{Bias}(L_k, \mathbf{x}) = L'_0(\mathbf{x}) - L'_k(\mathbf{x}) + l_0 - l_k.$$

Along the second subsequence, the right-hand side tends to the limit 0, and so

$$(57) \quad \text{Bias}(L_0, \mathbf{x}) = \lim_k \text{Bias}(L_k, \mathbf{x}) \leq \limsup_k \text{MaxBias}(L_k, \mathbf{X}_k),$$

where the last step follows from the fact that  $\mathbf{x} \in \mathbf{X}_k$  as soon as  $k \geq \|\mathbf{x}\|$ .

It follows from (56), (57), (53), (54) and Theorem 1 that

$$\begin{aligned} m(\text{MaxBias}(L_0, \mathbf{X}), \|L'_0\|) &\leq \limsup_k m(\text{MaxBias}(L_k, \mathbf{X}_k), \|L'_k\|) \\ &= \limsup_k M_k = M. \end{aligned}$$

In other words,

$$\text{MaxRisk}(L_0, \mathbf{X}) \leq M.$$

Recalling that  $M$  is the supremum of the difficulties of all one-dimensional subproblems, we show (by exhibiting the estimator  $L_0$ ) that the difficulty of the full problem is not harder. The formulas follow.  $\square$

**PROOF OF LEMMA 4.** We present only the argument for the first inequality; the second is similar. Suppose that, for a given  $d$ , we have

$$(58) \quad L(\mathbf{x}) - L(\mathbf{x}_1) > (d + \delta)\langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle$$

for some  $\mathbf{x} \in \mathbf{X}$ , which remains fixed throughout the proof. We will show that  $d \notin \partial\omega(\varepsilon_0)$ .

Set  $\mathbf{x}_h = (1 - h)\mathbf{x}_1 + h\mathbf{x}$ . Using the definition of  $\mathbf{x}_1$  and  $\mathbf{x}_{-1}$ , we have  $L(\mathbf{x}_1) - L(\mathbf{x}_{-1}) = \omega(\varepsilon_0)$ , and so

$$L(\mathbf{x}_h) - L(\mathbf{x}_{-1}) > \omega(\varepsilon_0) + h(d + \delta)\langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle.$$

Now

$$\begin{aligned} \|\mathbf{x}_h - \mathbf{x}_{-1}\|_K^2 &= \|\mathbf{x}_1 - \mathbf{x}_{-1} + \mathbf{x}_h - \mathbf{x}_1\|_K^2 \\ &= \varepsilon_0^2 + 2\langle K\mathbf{x}_1 - K\mathbf{x}_{-1}, K\mathbf{x}_h - K\mathbf{x}_1 \rangle + \|\mathbf{x}_h - \mathbf{x}_1\|_K^2 \\ &= \varepsilon_0^2 + 2h\varepsilon_0\langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle + h^2\|\mathbf{x} - \mathbf{x}_1\|_K^2. \end{aligned}$$

Note that  $|\langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle| > 0$ . Otherwise, we would have  $\|\mathbf{x}_h - \mathbf{x}_{-1}\|_K = \|\mathbf{x}_1 - \mathbf{x}_{-1}\|_K + o(h)$ . But (58) shows  $|L(\mathbf{x}_h) - L(\mathbf{x}_{-1})| > |L(\mathbf{x}_1) - L(\mathbf{x}_{-1})| + \text{const. } h$ , which contradicts the assumption that  $(\mathbf{x}_1, \mathbf{x}_{-1})$  attain the modulus.

It follows that  $h^2\|\mathbf{x} - \mathbf{x}_1\|_K^2 = O(h^2|\langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle|^2)$ . Setting

$$\eta = h\langle \mathbf{w}_0, K\mathbf{x} - K\mathbf{x}_1 \rangle,$$

we have

$$(59) \quad \omega(\varepsilon_0 + \eta + O(\eta^2)) > \omega(\varepsilon_0) + (d + \delta)\eta.$$

On the other hand, by definition, for any  $d \in \partial\omega(\varepsilon_0)$  we must have

$$(60) \quad \omega(\varepsilon_0 + \eta) \leq \omega(\varepsilon_0) + d\eta,$$

for all admissible  $\eta$ . However, (59) makes (60) impossible. As  $d$  does not satisfy (60), it cannot belong to  $\partial\omega(\varepsilon_0)$ .  $\square$

**PROOF OF COROLLARY 3.** The result follows by plugging in  $A\varepsilon^r + o(\varepsilon^r)$  in place of  $\omega(\varepsilon)$  in earlier results, and bounding remainder terms.  $\square$

**PROOF OF COROLLARY 4.** Under the hypothesis that the modulus has exponent  $r$ , it follows from concavity of the modulus that we have the set convergence

$$\frac{\varepsilon \partial\omega(\varepsilon)}{\omega(\varepsilon)} \rightarrow r, \quad \text{as } \varepsilon \rightarrow 0.$$

In the context of the proof of Theorem 1, this means that asymptotically, for small  $\varepsilon$ , we have  $\Gamma_1(\varepsilon) \sim r$  for small  $\varepsilon$ . It follows that asymptotically, as  $\sigma \rightarrow 0$ ,  $\Gamma_0$  intersects  $\Gamma_1$ , where both take approximately the  $y$ -value  $r$ . Hence  $c_0 \approx r$ , and the other formulas all follow.  $\square$

**Acknowledgments.** The author would like to thank L. D. Brown, Jianqing Fan, Roger Farrell, Iain Johnstone, Lucien Le Cam, Richard Liu, Mark Low, Max Mintz, Alex Samarov, Paul Speckman, Philip Stark, Joseph Traub, Hans Weinberger, Henryk Woźniakowski and the referees for helpful correspondence.

## REFERENCES

- BICKEL, P. J. (1981). Minimax estimation of the mean of a normal distribution when the parameter space is restricted. *Ann. Statist.* **9** 1301–1309.
- BROWN, L. D., COHEN, A. and STRAWDERMAN, W. E. (1976). A complete class theorem for strict monotone likelihood ratio with applications. *Ann. Statist.* **4** 712–722.
- BROWN, L. D. and LIU, R. (1989). A sharpened inequality concerning the hardest affine subproblem. Unpublished manuscript.
- BROWN, L. D. and LOW, M. G. (1990). Asymptotic equivalence of nonparametric regression and white noise. Unpublished manuscript.
- CASELLA, G. and STRAWDERMAN, W. E. (1981). Estimating a bounded normal mean. *Ann. Statist.* **9** 870–878.
- DONOHO, D. L. and LIU, R. C. (1991). Geometrizing rates of convergence, III. *Ann. Statist.* **19** 668–701.
- DONOHO, D. L., LIU, R. C. and MACGIBBON, K. B. (1990). Minimax risk over hyperrectangles, and implications. *Ann. Statist.* **18** 1416–1437.
- DONOHO, D. L. and LOW, M. G. (1990). White noise approximation and minimax risk. Technical report, Dept. Statistics, Univ. California, Berkeley.
- DONOHO, D. L. and NUSSBAUM, M. (1990). Minimax quadratic estimation of quadratic functionals. *J. Complexity* **6** 290–323.
- FAN, J. (1991). Nonparametric estimation of quadratic functionals in Gaussian white noise. *Ann. Statist.* **19** 1273–1295.
- FELDMAN, I. and BROWN, L. D. (1989). Unpublished manuscript.
- GOLOMB, M. and WEINBERGER, H. F. (1959). Optimal approximation and error bounds. In *On Numerical Approximation* (R. E. Langer, ed.) 117–190. Univ. Wisconsin Press.
- HALL, P. (1990). Optimal convergence rates in signal recovery. *Ann. Probab.* **18** 887–900.
- HECKMAN, N. (1988). Minimax estimation in a semiparametric model. *J. Amer. Statist. Assoc.* **83** 1090–1096.
- IBRAGIMOV, I. A. and HAS'MINSKII, R. Z. (1984). On nonparametric estimation of values of a linear functional in a Gaussian white noise. *Teor. Veroyatnost. i Primenen.* **29** 19–32. (In Russian.)
- IBRAGIMOV, I. A. and HAS'MINSKII, R. Z. (1987). On estimating linear functionals in a Gaussian noise. *Teor. Veroyatnost. i Primenen.* **32** 35–44. (In Russian.)
- IBRAGIMOV, I. A., NEMIROVSKII, A. S. and HAS'MINSKII, R. Z. (1986). Some problems on nonparametric estimation in Gaussian white noise. *Teor. Veroyatnost. i Primenen.* **31** 391–406. (In Russian.)
- KARLIN, S. and RUBIN, H. (1956). The theory of decision procedures for distributions with monotone likelihood ratio. *Ann. Math. Statist.* **27** 272–299.
- KUJS, J. A. and OLMAN, V. (1972). A minimax estimator of regression coefficients. *Izv. Akad. Nauk Eston. SSR* **21** 66–72. (In Russian.)
- LÄUTER, H. (1975). A minimax linear estimator for linear parameters under restrictions in form of inequalities. *Math. Operationsforsch. Statist.* **6** 689–695.
- LEVIT, B. Y. (1980). On asymptotic minimax estimates of the second order. *Theory Probab. Appl.* **25** 552–568.
- LI, K. C. (1982). Minimality of the method of regularization on stochastic processes. *Ann. Statist.* **10** 937–942.



- LIU, R. (1989). Unpublished manuscript.
- LOW, M. D. (1988). Towards a unified theory of asymptotic minimax estimation. Ph.D. dissertation, Cornell Univ.
- MELKMAN, A. A. and MICCHELLI, C. A. (1979). Optimal estimation of linear operators in Hilbert spaces from inaccurate data. *SIAM J. Numer. Anal.* **16** 87–105.
- MICCHELLI, C. A. (1975). Optimal estimation of linear functionals. IBM Research Report 5729. IBM, Armonk, NY.
- MICCHELLI, C. A. and RIVLIN, T. J. (1977). A survey of optimal recovery. In *Optimal Estimation in Approximation Theory* (C. A. Micchelli and T. J. Rivlin, eds.) 1–54. Plenum, New York.
- O'SULLIVAN, F. (1986). A statistical perspective on ill-posed inverse problems (with discussion). *Statist. Sci.* **1** 502–527.
- PACKEL, E. W. (1988). Do linear problems have linear optimal algorithms? *SIAM Rev.* **30** 388–403.
- PACKEL, E. W. and WOŹNIAKOWSKI, H. (1987). Recent developments in information-based complexity. *Bull. Amer. Math. Soc.* **17** 9–36.
- PILZ, J. (1986). Minimax linear regression estimation with symmetric parameter restrictions. *J. Statist. Plann. Inf.* **13** 297–318.
- PINELIS, I. F. (1991). On minimax risk. *Theory Probab. Appl.* **35** 104–109.
- ROCKEFELLAR, R. T. (1970). *Convex Analysis*. Princeton Univ. Press.
- 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, New York.
- SACKS, J. and YLVISAKER, D. (1978). Linear estimation for approximately linear models. *Ann. Statist.* **6** 1122–1137.
- SACKS, J. and YLVISAKER, D. (1981). Asymptotically optimum kernels for density estimation at a point. *Ann. Statist.* **9** 334–346.
- SCHOENBERG, I. J. (1964a). Spline interpolation and best quadrature for mullae. *Bull. Amer. Math. Soc.* **70** 143–148.
- SCHOENBERG, I. J. (1964b). On best approximation of linear operators. *Indag. Math.* **26** 155–163.
- SPECKMAN, P. (1979). Minimax estimates of linear functionals in a Hilbert space. Unpublished manuscript.
- STARK, P. B. (1992). Affine minimax confidence intervals for a bounded normal mean. *Statist. Probab. Lett.* **13** 39–44.
- TRAUB, J. F., WASILKOWSKI, G. W. and WOŹNIAKOWSKI, H. (1983). *Information, Uncertainty, Complexity*. Addison-Wesley, Reading, MA.
- TRAUB, J. F., WASILKOWSKI, G. W. and WOŹNIAKOWSKI, H. (1988). *Information-Based Complexity*. Addison-Wesley, Reading, MA.
- ZEYTINGLU, M. and MINTZ, M. (1984). Optimal fixed-size confidence procedures for a restricted parameter space. *Ann. Statist.* **12** 945–957.
- ZEYTINGLU, M. and MINTZ, M. (1988). Robust fixed-size confidence procedures for a restricted parameter space. *Ann. Statist.* **16** 1241–1253.

DEPARTMENT OF STATISTICS  
 STANFORD UNIVERSITY  
 STANFORD, CALIFORNIA 94305