

## ORACLE INEQUALITIES FOR INVERSE PROBLEMS

BY L. CAVALIER, G. K. GOLUBEV, D. PICARD AND A. B. TSYBAKOV

*Université Aix-Marseille 1, Université Aix-Marseille 1,  
Université Paris 7 and Université Paris 6*

We consider a sequence space model of statistical linear inverse problems where we need to estimate a function  $f$  from indirect noisy observations. Let a finite set  $\Lambda$  of linear estimators be given. Our aim is to mimic the estimator in  $\Lambda$  that has the smallest risk on the true  $f$ . Under general conditions, we show that this can be achieved by simple minimization of an unbiased risk estimator, provided the singular values of the operator of the inverse problem decrease as a power law. The main result is a nonasymptotic oracle inequality that is shown to be asymptotically exact. This inequality can also be used to obtain sharp minimax adaptive results. In particular, we apply it to show that minimax adaptation on ellipsoids in the multivariate anisotropic case is realized by minimization of unbiased risk estimator without any loss of efficiency with respect to optimal nonadaptive procedures.

**1. Introduction.** Let  $A$  be a known linear operator on a Hilbert space  $H$  with inner product  $(\cdot, \cdot)$  and norm  $\|\cdot\|$ . Let  $f \in H$  be an unknown function that we want to estimate from indirect observations

$$(1) \quad Y(g) = (Af, g) + \varepsilon\xi(g), \quad g \in H,$$

where  $0 < \varepsilon < 1$  and  $\xi(g)$  is a Gaussian random variable on a probability space  $(\Omega, \mathcal{A}, \mathbf{P})$ , with mean 0 and variance  $\|g\|^2$ , such that  $\mathbf{E}\{\xi(g)\xi(v)\} = (g, v)$ , for any  $g, v \in H$ , where  $\mathbf{E}$  is the expectation w.r.t.  $\mathbf{P}$ .

Relation (1) defines a Gaussian white noise model. Instead of using all the observations  $\{Y(g), g \in H\}$  it is usually sufficient to consider the set of values  $\{Y(\psi_k)\}$ , for some orthonormal basis  $\{\psi_k\}_{k=1}^\infty$ . Then  $\xi(\psi_k) = \xi_k$  are i.i.d. standard Gaussian random variables.

We assume that the basis  $\{\psi_k\}$  is such that  $(Af, \psi_k) = b_k\theta_k$ , where  $b_k$  are real numbers and  $\theta_k = (f, \varphi_k)$  are the Fourier coefficients of  $f$  w.r.t. some orthonormal basis  $\{\varphi_k\}$  (not necessarily  $\varphi_k = \psi_k$ ). A typical example when it occurs is that the operator  $A$  admits a singular value decomposition of the form

$$(2) \quad A\varphi_k = b_k\psi_k, \quad A^*\psi_k = b_k\varphi_k,$$

where  $A^*$  is the adjoint of  $A$ ,  $b_k$  are the singular values,  $\{\psi_k\}$  is an orthonormal basis in  $\text{Range}(A) \subset H$  and  $\{\varphi_k\}$  is the corresponding orthonormal basis in  $H$ .

---

Received July 2000; revised October 2001.

AMS 2000 subject classifications. 62G05, 62G20.

Key words and phrases. Statistical inverse problems, oracle inequalities, adaptive curve estimation, model selection, exact minimax constants.

Under these assumptions (but also in some other situations) one has the equivalent discrete sequence observation model derived from (1):

$$(3) \quad y_k = b_k \theta_k + \varepsilon \xi_k, \quad k = 1, 2, \dots,$$

where  $y_k$  stands for  $Y(\psi_k)$ . If  $b_k \neq 0$ , the model (3) can be written in the form

$$(4) \quad X_k = \theta_k + \varepsilon \sigma_k \xi_k, \quad k = 1, 2, \dots,$$

where  $X_k = y_k/b_k$  and  $\sigma_k = b_k^{-1}$ . This can also be viewed as a model with direct observations and correlated data [Johnstone (1999)]. The sequence space formulation (3) or (4) for statistical inverse problems has been studied in a number of papers [see Johnstone and Silverman (1990), Korostelev and Tsybakov (1993), Koo (1993), Donoho (1995), Mair and Ruymgaart (1996), Golubev and Khasminskii (1999, 2001), Johnstone (1999), Goldenshluger and Pereverzev (2000) and Cavalier and Tsybakov (2002) among others]. For ill-posed inverse problems we have  $|b_k| \rightarrow 0$  and  $|\sigma_k| \rightarrow \infty$ , as  $k \rightarrow \infty$ .

Let  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots)$  be an estimator of  $\theta = (\theta_1, \theta_2, \dots)$  based on the data (4). Then  $f$  is estimated by  $\hat{f} = \sum_k \hat{\theta}_k \varphi_k$ . The mean integrated squared risk (MISE) of  $\hat{f}$  is

$$\mathcal{R}(\hat{f}, f) = \mathbf{E}_f \|\hat{f} - f\|^2 = \mathbf{E}_\theta \sum_k (\hat{\theta}_k - \theta_k)^2 = \mathbf{E}_\theta \|\hat{\theta} - \theta\|^2,$$

where the notation  $\|\cdot\|$  means the  $\ell_2$ -norm when applied to  $\theta$ -vectors in the sequence space. Here and later  $\mathbf{E}_f$  and  $\mathbf{E}_\theta$  denote the expectations w.r.t.  $\{Y(g), g \in H\}$  or  $X = (X_1, X_2, \dots)$  for models (1) and (4), respectively. Analyzing the risk  $\mathcal{R}(\hat{f}, f)$  of the estimator  $\hat{f}$  is equivalent to analyzing the corresponding sequence space risk  $\mathbf{E}_\theta \|\hat{\theta} - \theta\|^2$ .

Let  $\lambda = (\lambda_1, \lambda_2, \dots)$  be a sequence of nonrandom weights, also called a *filter*. Every filter  $\lambda$  defines the estimator  $\hat{\theta}(\lambda) = (\hat{\theta}_1, \hat{\theta}_2, \dots)$ , where  $\hat{\theta}_k = \lambda_k X_k$ . Examples of commonly used weights  $\lambda_k$  are the projection weights  $\lambda_k = I(k \leq w)$ , for some integer  $w$  [where  $I(\cdot)$  denotes the indicator function], the Tikhonov–Phillips weights

$$\lambda_k = \frac{1}{1 + (k/w)^\alpha}, \quad w > 0, \alpha > 0,$$

and the Pinsker (1980) weights

$$\lambda_k = (1 - (k/w)^\alpha)_+, \quad w > 0, \alpha > 0,$$

where  $x_+ = \max(x, 0)$ . The estimator  $\hat{f}$  with Tikhonov–Phillips weights for even values  $\alpha$  is asymptotically equivalent to a smoothing spline estimator [Wahba (1977, 1990)].

Although  $\lambda$  is not finite dimensional, it is usually determined by a finite number of parameters, as in the above examples. In this paper we discuss how to choose

these parameters optimally in a data-driven way. In particular, a data-driven choice of smoothing parameters  $w$  and  $\alpha$  for the Tikhonov–Phillips method is interesting.

We suppose that there is a finite set of possible candidate filters  $\Lambda = \{\lambda^1, \dots, \lambda^N\}$ , with  $\lambda^s = (\lambda_1^s, \lambda_2^s, \dots)$ ,  $s = 1, \dots, N$ ,  $N \geq 2$ , satisfying some general conditions. These filters can be, for example, any of the three types described above, as well as pooled sets of different kinds of filters. Given the data,  $X = (X_1, X_2, \dots)$ , our aim is to select a data-dependent sequence of weights  $\lambda^* = \lambda^*(X) = (\lambda_1^*, \lambda_2^*, \dots)$ , with values in  $\Lambda$ , that has asymptotically minimal squared risk for the true  $\theta$ . We show that  $\lambda^*$  can be defined as a minimizer (with respect to  $\lambda \in \Lambda$ ) of an unbiased estimator of the risk. Optimality properties of such  $\lambda^*$  follow from the oracle inequalities that are the main result of the paper. The oracle inequalities are nonasymptotic, they are obtained under very weak conditions on  $\Lambda$  and they lead to asymptotically exact inequalities of the form

$$(5) \quad \mathcal{R}(f^*, f) \leq (1 + o(1)) \min_{\lambda \in \Lambda} \mathcal{R}(f_\lambda, f),$$

as  $\varepsilon \rightarrow 0$ , where  $f_\lambda = \sum_{k=1}^{\infty} \lambda_k X_k \varphi_k$ ,  $f^* = \sum_{k=1}^{\infty} \lambda_k^* X_k \varphi_k$  and  $o(1)$  does not depend on  $f$  but depends on the family  $\Lambda$ .

As a consequence of these inequalities, we can justify the optimal choice of smoothing parameters in the Tikhonov–Phillips and projection methods, as well as in Pinsker’s method (the last one yields as a by-product sharp minimaxity of  $\lambda^*$  on Sobolev ellipsoids). The optimality results are valid under the assumption that  $\sigma_k$  is growing as a power of  $k$ , as  $k \rightarrow \infty$ . Generality of the oracle inequalities allows us to apply them in various problems. As an example, we consider sharp adaptation on multivariate anisotropic smoothness classes. An interesting conclusion is that the adaptive estimator based on  $\lambda^*$  attains the minimax lower bound for the multivariate anisotropic case, without any loss of efficiency.

Other oracle inequalities for inverse problems have been proposed recently by Johnstone (1999) and Cavalier and Tsybakov (2002). Johnstone (1999) deals with a class of nonlinear estimators based on soft thresholding in a wavelet context. He obtains an asymptotically exact oracle inequality for this class. Cavalier and Tsybakov (2002) consider the model (4) and obtain asymptotically exact oracle inequalities of the form (5), where  $\Lambda$  is the class of all monotone weight sequences  $\lambda$  and  $f^*$  (respectively,  $\lambda^*$ ) is chosen in a different way, by application of a penalized blockwise Stein’s rule. Their method is sharp adaptive in a minimax sense on ellipsoids, but it is not suited for parameter selection in restricted classes of filters, such as the Tikhonov–Phillips one.

**2. Main results.** We deal with the model (4) and we assume the following.

ASSUMPTION 1. For any  $\lambda \in \Lambda$ ,

$$0 < \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 < \infty$$

and

$$(6) \quad \max_{\lambda \in \Lambda} \sup_i |\lambda_i| \leq 1.$$

The risk of the linear estimator  $\hat{\theta}(\lambda)$  is given by

$$(7) \quad R_\varepsilon[\lambda, \theta] = \mathbf{E}_\theta \|\hat{\theta}(\lambda) - \theta\|^2 = \sum_{i=1}^{\infty} (1 - \lambda_i)^2 \theta_i^2 + \varepsilon^2 \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2.$$

The assumption (6) is quite natural. In fact, it follows from (7) that the estimator  $\hat{\theta}(\lambda)$  with at least one  $\lambda_i \notin [0, 1]$  is inadmissible. However, we included the case of negative bounded  $\lambda_i$  since it corresponds to a number of well-known estimators such as kernel estimators with nonnegative kernels. The results below remain valid (up to a change in constants) if we replace 1 by an arbitrary constant  $C$  in (6).

Our data-driven choice of the smoothing parameter  $\lambda$  is based on the principle of unbiased risk estimation. To be more specific, recall briefly a heuristic motivation of the Mallows  $C_p$  criterion [Akaike (1973), Mallows (1973)]. The best choice of  $\lambda$  is the filter  $\lambda^0$  (called the *oracle*) that minimizes  $R_\varepsilon[\lambda, \theta]$  over  $\lambda \in \Lambda$ . The oracle  $\lambda^0$  cannot be found directly since the functional  $R_\varepsilon[\lambda, \theta]$  depends on the unknown  $\theta_i^2$ . However, an unbiased estimator of  $\theta_i^2$  is available in the form  $X_i^2 - \varepsilon^2 \sigma_i^2$ . Thus the functional

$$(8) \quad \begin{aligned} U[\lambda, X] &= \sum_{i=1}^{\infty} (\lambda_i^2 - 2\lambda_i)(X_i^2 - \varepsilon^2 \sigma_i^2) + \varepsilon^2 \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 \\ &= \sum_{i=1}^{\infty} (\lambda_i^2 - 2\lambda_i) X_i^2 + 2\varepsilon^2 \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i \end{aligned}$$

is an unbiased estimator of  $R_\varepsilon[\lambda, \theta] - \sum_{i=1}^{\infty} \theta_i^2$ :

$$(9) \quad \mathbf{E}_\theta U[\lambda, X] = R_\varepsilon[\lambda, \theta] - \sum_{i=1}^{\infty} \theta_i^2.$$

Under our assumptions the random series in (8) converges in the mean squared sense for any  $\lambda \in \Lambda$ , and the definition (8), as well as other definitions of random series appearing in the paper, are understood in this sense. Of course, in practice one does not compute infinite series. One either requires that  $\lambda_i = 0$  for all  $i$  large enough (as for the projection or Pinsker filters) or translates the definition of  $U[\lambda, X]$  into the function space to make it computable (e.g., for the Tikhonov–Phillips filters with even  $\alpha$  the computations are possible in terms of splines).

The principle of unbiased risk estimation suggests that we minimize over  $\lambda \in \Lambda$  the functional  $U[\lambda, X]$  in place of  $R_\varepsilon[\lambda, \theta]$ . This leads to the following data-driven choice of  $\lambda$ :

$$(10) \quad \lambda^* = \arg \min_{\lambda \in \Lambda} U[\lambda, X].$$

We show that this simple and intuitive smoothing parameter selection rule is efficient for inverse problems with power growth of  $\sigma_k$ , in the sense that it satisfies asymptotically precise oracle inequalities.

We define

$$\rho(\lambda) = \sup_k \sigma_k^2 |\lambda_k| \left\{ \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4 \right\}^{-1/2}$$

and

$$\rho = \max_{\lambda \in \Lambda} \rho(\lambda).$$

Although the main results of this paper hold for general  $\rho$ , we will typically think of  $\rho$  as being small (for small  $\varepsilon$ ). In particular, this will be explicitly assumed in the asymptotic corollaries.

Let the following assumption hold.

ASSUMPTION 2. There exists a constant  $C_1 > 0$  such that, uniformly in  $\lambda \in \Lambda$ ,

$$\sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^2 \leq C_1 \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4.$$

Assumptions 1 and 2 are very mild, and they are satisfied in most of the interesting examples. Since  $|\lambda_i| \leq 1$ , we have

$$\sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4 \leq \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^2,$$

and Assumption 2 means that both sums are of the same order. The sums  $\varepsilon^4 \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4$  and  $\varepsilon^4 \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^2$  are the main terms of the variance  $\text{Var}\{U[\lambda, X]\}$ . On the other hand  $R_\varepsilon[\lambda, \theta] \geq \varepsilon^2 \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2$  and

$$(11) \quad \frac{(\varepsilon^4 \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4)^{1/2}}{\varepsilon^2 \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2} \leq \rho.$$

Therefore, the smallness of  $\rho$  guarantees the smallness of the ratio of standard deviation to the mean error:  $\text{Var}^{1/2}\{U[\lambda, X]\}/R_\varepsilon[\lambda, \theta]$ , uniformly over  $\lambda$  and  $\theta$ .

Note also that under Assumption 2 we have

$$(12) \quad \rho(\lambda) \leq \sqrt{C_1} \quad \forall \lambda \in \Lambda.$$

Define

$$(13) \quad S = \max_{\lambda \in \Lambda} \sup_i \sigma_i^2 \lambda_i^2 / \min_{\lambda \in \Lambda} \sup_i \sigma_i^2 \lambda_i^2$$

and

$$M = \sum_{\lambda \in \Lambda} \exp\{-1/\rho(\lambda)\},$$

$$L_\Lambda = \log(NS) + \rho^2 \log^2(MS).$$

We recall that  $N$  is the number of elements of the family  $\Lambda$ . Since  $M \leq N$ , we always have  $L_\Lambda \leq \log(NS) + C_1 \log^2(NS)$ , but this bound is rough. In asymptotics, as  $\varepsilon \rightarrow 0$  and  $N \rightarrow \infty$ , the correct order is typically  $\rho = o(1)$ ,  $M \sim \text{const}$  and  $L_\Lambda \sim \log(NS)$ .

Main results of the paper are given in the next two theorems.

**THEOREM 1.** *Let Assumptions 1 and 2 hold. Then for every  $\theta \in \ell_2$ , for every  $B > B_0$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$ , where  $\lambda^*$  is defined by (10), we have*

$$(14) \quad \mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + \gamma_1 B^{-1}) \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta] + \gamma_2 B \varepsilon^2 L_\Lambda \omega(B^2 L_\Lambda),$$

where

$$\omega(x) = \max_{\lambda \in \Lambda} \sup_k \sigma_k^2 \lambda_k^2 I \left\{ \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 \leq x \sup_k \sigma_k^2 \lambda_k^2 \right\}, \quad x > 0,$$

and  $B_0 > 0$ ,  $\gamma_1 > 0$ ,  $\gamma_2 > 0$  are constants depending only on  $C_1$ .

**THEOREM 2.** *Let Assumptions 1 and 2 hold. Then there exist constants  $\gamma_3 > 0$ ,  $\gamma_4 > 0$  depending only on  $C_1$ , such that for every  $\theta \in \ell_2$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$ , where  $\lambda^*$  is defined by (10), we have*

$$\mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq [1 + \gamma_3 \rho \sqrt{L_\Lambda}] \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta],$$

provided  $\rho \sqrt{L_\Lambda} \leq \gamma_4$ .

To prove (5), in many examples it is sufficient to use Theorem 2, and even its simplified version that we are going to state now. Assume that

$$(15) \quad \lim_{\varepsilon \rightarrow 0} \rho^2 \log(NS) = 0.$$

Then  $L_\Lambda = O(\log(NS))$ , and we get the following corollary of Theorem 2.

**COROLLARY 1.** *Let Assumptions 1 and 2 hold. Then there exist constants  $C_2 > 0$ ,  $C_3 > 0$ , depending only on  $C_1$ , such that for  $\rho^2 \log(NS) < C_2$  we have*

$$\mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq \left(1 + C_3 \rho \sqrt{\log(NS)}\right) \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta],$$

for every  $\theta \in \ell_2$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$ .

Thus, condition (15) is sufficient to have asymptotically exact oracle inequalities of the form (5).

Recall that for ill-posed inverse problems  $\sigma_k \rightarrow \infty$ , as  $k \rightarrow \infty$ . A crucial restriction, implicit in Theorems 1 and 2 and Corollary 1 is that  $\sigma_k$  should grow not faster than a power of  $k$ . In fact, if  $\sigma_k$  grows exponentially, then even in the simplest case of projection weights  $\lambda_k$  we have  $\rho \neq o(1)$ , as  $\varepsilon \rightarrow 0$  (thus, Theorem 2 cannot be applied if  $N$  grows with  $\varepsilon \rightarrow 0$ ). Theorem 1 in this case can be applied but does not give correct rates.

The oracle inequality (14) is suited to obtaining minimax results or rates of convergence for classes of sequences  $\theta$ . The remainder term in the right-hand side of (14) is typically of the form  $\varepsilon^2 \log^a(1/\varepsilon)$  with  $a > 1$ . This shows limits of applications of (14): for the classes where the least favorable functions  $\theta$  satisfy  $\min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta] \ll \varepsilon^2 \log^a(1/\varepsilon)$ , the remainder term is not asymptotically negligible, and thus asymptotically sharp adaptation in a minimax sense is not possible by use of our techniques. This remark concerns the classes of analytical functions, for example. On the other hand, the remainder term is negligible compared to the minimax risk on Sobolev ellipsoids.

Theorem 2 is rather a “class of estimators” than a “class of functions” result. It allows us to get oracle inequalities of the form (5) for any fixed  $\theta \in \ell_2$ , provided the class  $\Lambda$  of estimators is small enough.

REMARK 1. Theorems 1 and 2 remain valid for non-Gaussian  $\xi_i$  such that  $\mathbf{E} \exp(C\xi_i^2) < \infty$  for some  $C > 0$ .

REMARK 2. Theorems 1 and 2 can be used not only for ill-posed, but also for well-posed inverse problems where  $\sigma_k \not\rightarrow \infty$ . For example, both theorems apply if  $\sigma_k \sim k^{-\beta}$  with  $0 \leq \beta < 1/4$  [allowing  $\rho = o(1)$ , as  $\varepsilon \rightarrow 0$ ]. For faster decreasing  $\sigma_k$  only Theorem 1 works. If  $N$  and  $S$  grow not faster than a power of  $\varepsilon^{-1}$ , the remainder term in (14) is  $O(\varepsilon^2 \log^a(1/\varepsilon))$  for some  $a > 0$ , which is only logarithmically worse than the optimal rate  $\varepsilon^2$  of the well-posed inverse problems.

REMARK 3. Consider the special case that corresponds to direct observations (i.e.,  $\sigma_k \equiv 1$ ). Here several oracle inequalities have been known previously. Theorems 1 and 2 extend these results, especially in what concerns multivariate applications.

The first oracle inequalities for the direct observations model appeared, although implicitly, in the proofs of optimality of  $C_p$ , cross-validation and related data-driven methods [Li (1986, 1987) and Polyak and Tsybakov (1990, 1992)]. They are also implicit in minimax adaptive constructions [Golubev (1987, 1992) and Golubev and Nussbaum (1992)]. Presumably, the first explicit use of oracle inequalities and its implications for minimax is due to Donoho and Johnstone

(1994) and Kneip (1994). More recent references are Donoho and Johnstone (1995, 1996), Nemirovski (2000), Birgé (2001), Birgé and Massart (2001), Cavalier and Tsybakov (2002) and Goldenshluger and Tsybakov (2001). A link of oracle inequalities to maxisets is discussed by Kerkyacharian and Picard (2002). For linear estimators, the results are often proved in the following model, which can be embedded into ours. Let  $Y = (Y_1, \dots, Y_n)^T$  be the vector of observations in the nonparametric regression model

$$(16) \quad Y_k = f(x_k) + v_k, \quad k = 1, \dots, n,$$

where  $v_k$  are i.i.d.  $\mathcal{N}(0, 1)$  random errors,  $x_k \in [0, 1]$  are nonrandom distinct points and  $f(\cdot)$  is the unknown function to be estimated. Consider the linear estimator  $\hat{f} = \mathcal{S}Y$  of the vector  $f = (f(x_1), f(x_2), \dots, f(x_n))^T$ , where  $\mathcal{S} = \sum_{i=1}^n \lambda_i u_i u_i^T$  is a symmetric  $n \times n$  matrix (smoother matrix) with eigenvalues  $\lambda_i$ , and  $\{u_i\}$  is an orthonormal basis in  $\mathbb{R}^n$ . Denoting  $X_i = n^{-1/2} u_i^T Y$ ,  $\theta_i = n^{-1/2} u_i^T f$ ,  $\xi_i = u_i^T v$  [where  $v = (v_1, \dots, v_n)^T$ ] and  $\varepsilon = n^{-1/2}$ , we rewrite the initial regression model in the equivalent form

$$X_i = \theta_i + \varepsilon \xi_i, \quad i = 1, \dots, n,$$

which is a special case of (4), modulo the fact that the  $\ell_2$ -vectors should contain zeros starting from the  $(n+1)$ th position:  $\theta = (\theta_1, \dots, \theta_n, 0, 0, \dots)$ . The linear estimator  $\hat{f}$  is translated into  $\hat{\theta}(\lambda) = (\hat{\theta}_1, \dots, \hat{\theta}_n, 0, 0, \dots)$ , where  $\hat{\theta}_i = \lambda_i X_i = n^{-1/2} \lambda_i u_i^T Y$ , and the risk is

$$\mathbf{E}_\theta \|\hat{\theta}(\lambda) - \theta\|^2 = \mathbf{E} \frac{1}{n} \sum_{i=1}^n \|\mathcal{S}Y - f\|^2.$$

Kneip (1994) studies this setup assuming that the class of filters  $\Lambda$  contains the sequences  $\lambda$  with monotone nonincreasing coefficients  $\lambda_i$  and such that for any two sequences  $\lambda, \lambda' \in \Lambda$  we have either  $\lambda_i \leq \lambda'_i, \forall i$ , or  $\lambda_i \geq \lambda'_i, \forall i$  (ordered linear smoothers). The set  $\Lambda$  in Kneip (1994) need not be finite. With the above translation into our notation, Kneip [(1994), page 844] proves the oracle inequality

$$(17) \quad \mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + O(B^{-1})) \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta] + O(B)\varepsilon^2 \quad \forall B > 0,$$

where  $\theta^*$  is the data-driven estimator (10). This is similar to (14) [but recall that Kneip's inequality (17) covers only the case  $\sigma_k \equiv 1$ ]. Another difference is that we assume finiteness of the set  $\Lambda$  (which is not restrictive in view of applications), but we drop the assumption of order. The last point is useful in multivariate models, in particular, anisotropic ones, where the ordering of the filters is not natural (cf. the example in Section 6 below).

We also note that for the regression model (16) there exist several alternatives to Mallows'  $C_p$ , for example, cross-validation (CV) and generalized cross-validation (GCV). Asymptotic optimality results [similar to (5)] for CV and GCV



in regression are given by Härdle and Marron (1985), Li (1986, 1987) and Polyak and Tsybakov (1992).

**3. Examples.** Consider some examples of application of the main results. Typical assumptions on the parameters appearing in Theorems 1 and 2 will be the following:

1. power growth of  $\sigma_k$ , as  $k \rightarrow \infty$ ;
2. power behavior of  $S$ :  $S = O(\varepsilon^{-t})$ , for some  $t > 0$ , as  $\varepsilon \rightarrow 0$ ;
3. at most power growth of  $N$ :  $N = O(\varepsilon^{-\nu})$ , for some  $\nu > 0$ , as  $\varepsilon \rightarrow 0$  [and  $N = O(\log(1/\varepsilon))$  in some examples].

In all the cases we have  $\log(NS) = O(\log(1/\varepsilon))$ .

EXAMPLE 1 (Projection estimators). Let  $1 \leq w_1 < \dots < w_N$  be integers. Consider the projection filters  $\lambda^s = (\lambda_1^s, \lambda_2^s, \dots)$  defined by

$$(18) \quad \begin{aligned} \lambda_i^1 &= I(i \leq w_1), & \lambda_i^2 &= I(i \leq w_2), \dots, \\ \lambda_i^N &= I(i \leq w_N), & i &= 1, 2, \dots \end{aligned}$$

Throughout this section we assume a power law behavior of  $\sigma_k$ :

$$(19) \quad \sigma_{\min} k^\beta \leq |\sigma_k| \leq \sigma_{\max} k^\beta,$$

$k = 1, 2, \dots$ , for some  $\sigma_{\max} \geq \sigma_{\min} > 0, \beta \geq 0$ .

Note that Assumption 2 is satisfied with the equality and  $C_1 = 1$ . It is easy to see that

$$\rho(\lambda^s) \leq C w_s^{-1/2}, \quad \rho \leq \max_{s=1, \dots, N} C w_s^{-1/2} = C w_1^{-1/2}$$

and

$$S \leq C(w_N/w_1)^{2\beta}.$$

Here and later  $C$  stands for positive constants (possibly different on different occasions) that do not depend on  $\lambda, \theta$  and  $\varepsilon$ . Using the bound on  $\rho(\lambda)$ , we conclude that  $M$  is bounded by a constant independent of  $\varepsilon$  for any choice of the sequence  $w_i$ . Note also that

$$(20) \quad \omega(x) \leq C \sup_k k^{2\beta} I\{k^{2\beta+1} \leq Cxk^{2\beta}\} \leq Cx^{2\beta}$$

and

$$L_\Lambda \leq C(\log(Nw_N/w_1) + w_1^{-1} \log^2(w_N/w_1)).$$

Since  $\rho\sqrt{\log(NS)} \leq Cw_1^{-1/2}\sqrt{\log(Nw_N/w_1)}$ , Corollary 1 gives

$$(21) \quad \mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq \left[ 1 + C \sqrt{\frac{\log(Nw_N/w_1)}{w_1}} \right] \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta]$$

provided  $w_1^{-1} \log(Nw_N/w_1)$  is small enough. As a consequence, we get an asymptotically exact oracle inequality of the form (5):

**COROLLARY 2.** *Assume that  $\Lambda = (\lambda^1, \dots, \lambda^N)$  is the set of projection filters defined by (18) and that (19) holds. If  $N = N(\varepsilon)$  and  $w_1 = w_1(\varepsilon)$ ,  $w_N = w_N(\varepsilon)$  are such that*

$$\lim_{\varepsilon \rightarrow 0} \frac{\log(Nw_N/w_1)}{w_1} = 0,$$

*then for every  $\theta \in \ell_2$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$  we have*

$$\mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + o(1)) \inf_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta],$$

*where  $o(1) \rightarrow 0$  uniformly in  $\theta \in \ell_2$ .*

In other words, Corollary 2 states that our adaptively selected filter behaves itself asymptotically at least as well as the best projection estimator in  $\Lambda$ . For the direct case (where  $\sigma_k \equiv 1$ ) such an inequality is obtained by Birgé (2001), who uses the Lepski adaptation method rather than the Mallows  $C_p$ .

Next, consider the situation where there is no restriction on  $w_1$  except  $w_1 \geq 1$  (i.e., the class  $\Lambda$  can contain the projection filters of order less than or equal to some  $w_N$ ). Applying (21) we get an inequality with a logarithmic loss of efficiency:

**PROPOSITION 1.** *Assume that  $\Lambda = (\lambda^1, \dots, \lambda^N)$  is the set of projection filters defined by (18) and that (19) holds. For every  $\theta \in \ell_2$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$  we have*

$$\mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq C \sqrt{\log(Nw_N)} \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta],$$

*where  $C > 0$  depends only on  $\sigma_{\min}$ ,  $\sigma_{\max}$ ,  $\beta$ .*

An important special case is wavelet estimators for which we set  $w_1 = 2^{j_0}$  (where  $j_0$  is the index of the initial level),  $w_j = 2w_{j-1}$ . Typically one chooses  $2^{-j_0}$  to be decreasing as a power of  $\varepsilon$  and  $N \sim \log(1/\varepsilon)$ . It is easy to see that the result of Corollary 2 remains valid in this case. Thus, a wavelet estimator that uses our data-driven selection of  $w_j$  is asymptotically at least as good as the best linear wavelet estimator for any  $\theta$ . By taking suprema of both sides of the oracle inequality over Besov classes [for definition see Donoho and Johnstone (1994, 1995, 1996)] we

get that our adaptive estimator attains optimal rate of convergence on all the Besov classes where linear wavelet estimators attain optimal rates.

EXAMPLE 2 (Level-wise “keep-or-kill” estimators). Let  $m > 1$  and  $1 \leq w_1 < \dots < w_m$  be integers and let  $e = (e_1, \dots, e_{m-1})$  be a binary sequence of length  $m - 1$ ,  $e_k \in \{0, 1\}$ . We associate with  $e$  a filter  $\lambda(e) = (\lambda_1(e), \lambda_2(e), \dots)$  defined by

$$\lambda_i(e) = I\{i \leq w_1\} + \sum_{k=1}^{m-1} e_k I\{w_k < i \leq w_{k+1}\}, \quad i = 1, 2, \dots$$

Consider the collection of filters

$$(22) \quad \Lambda = \{\lambda(e) : e \in E\},$$

where  $E$  is the set of all binary sequences  $e$  of length  $m - 1$ . The linear estimator with weights  $\lambda_i(e)$  “keeps” the blocks of coefficients  $\{\theta_i : w_k < i \leq w_{k+1}\}$  for which  $e_k = 1$  and “kills” the blocks for which  $e_k = 0$ . Clearly,

$$N = \text{Card}(\Lambda) = 2^{m-1}.$$

As in Example 1, we get that Assumption 2 is satisfied and that  $\rho \leq Cw_1^{-1/2}$ ,  $S \leq C(w_m/w_1)^{2\beta}$ , provided (19) holds. Therefore, applying Corollary 1 we get the following result.

PROPOSITION 2. Assume that  $\Lambda$  is the set of levelwise “keep-or-kill” filters defined by (22) and that (19) holds. If  $m = m(\varepsilon)$  and  $w_1 = w_1(\varepsilon)$ ,  $w_m = w_m(\varepsilon)$  are such that

$$(23) \quad \lim_{\varepsilon \rightarrow 0} \frac{m + \log(w_m/w_1)}{w_1} = 0,$$

then for every  $\theta \in \ell_2$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$  we have

$$\mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + o(1)) \inf_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta],$$

where  $o(1) \rightarrow 0$  uniformly in  $\theta \in \ell_2$ .

Note that for wavelet bases (if we consider global thresholding level-by-level) we have  $w_1 = 2^{j_0}$ ,  $w_m = 2^{j_0+m-1}$  with an integer  $j_0$ , and condition (23) reduces to

$$\lim_{\varepsilon \rightarrow 0} 2^{-j_0} m = 0,$$

which is readily satisfied for the typical situation where  $2^{-j_0}$  decreases as a power of  $\varepsilon$  and  $m \sim \log(1/\varepsilon)$ . As a conclusion we get, in particular, that the wavelet keep-or-kill level-by-level estimator that uses our data-driven rule attains the optimal

rate of convergence on all the Besov classes where the linear wavelet level-wise keep-or-kill estimators attain optimal rates.

EXAMPLE 3 (Tikhonov–Phillips estimators). Consider the set of filters

$$(24) \quad \Lambda = \left\{ \lambda = \{\lambda_k\} : \lambda_k = \frac{1}{1 + (k/w)^\alpha}, w \in \mathcal{W}, \alpha \in \mathcal{A} \right\},$$

where  $\mathcal{A}$  is a finite set of real numbers, possibly depending on  $\varepsilon$ :

$$\mathcal{A} = \{\alpha(1), \alpha(2), \dots, \alpha(N_{\mathcal{A}})\}$$

such that  $2\beta + 1/2 < \alpha_{\min} = \alpha(1) < \alpha(2) < \dots < \alpha(N_{\mathcal{A}}) = \alpha_{\max}$ , and  $\mathcal{W}$  is a finite subset of the interval  $[w_1, w_{\max}]$ , with  $\text{Card}(\mathcal{W}) = N_{\mathcal{W}}$ , where  $N_{\mathcal{A}}, N_{\mathcal{W}}$  are integers,  $0 < w_1 < w_{\max} < \infty$ . Note that  $\text{Card}(\Lambda) = N$ , where  $N = N_{\mathcal{A}}N_{\mathcal{W}}$ .

We get by simple algebra

$$\sum_{i=1}^{\infty} \frac{i^{4\beta}}{(1 + (i/w)^\alpha)^4} \geq C_* w^{4\beta+1}, \quad \sum_{i=1}^{\infty} \frac{i^{4\beta}}{(1 + (i/w)^\alpha)^2} \leq C^* w^{4\beta+1},$$

where  $C_*$  and  $C^*$  are positive constants depending only on  $\alpha_{\min}, \alpha_{\max}, \beta$ . This and (19) guarantee that Assumption 2 is satisfied.

Next,

$$\sup_i \frac{i^{2\beta}}{[1 + (i/w)^\alpha]^2} = C(\alpha, \beta) w^{2\beta},$$

where  $C(\alpha, \beta)$  is bounded from 0 and  $\infty$  uniformly for  $\alpha_{\min} \leq \alpha \leq \alpha_{\max}$ . Therefore  $S \leq C(w_{\max}/w_1)^{2\beta}$ . Furthermore, we get similarly

$$\rho \leq C \max_{\lambda \in \Lambda} \sup_i \frac{i^{2\beta}}{1 + (i/w)^\alpha} \left( \sum_{i=1}^{\infty} \frac{i^{4\beta}}{(1 + (i/w)^\alpha)^4} \right)^{-1/2} \leq C w_1^{-1/2}.$$

Thus,

$$\rho \sqrt{\log(NS)} \leq C w_1^{-1/2} \sqrt{\log(Nw_{\max}/w_1)}$$

and, to guarantee an asymptotically exact oracle inequality, it remains to require (in view of Corollary 1) that the parameters  $w_{\max}, w_1, N_{\mathcal{W}}, N_{\mathcal{A}}$  are such that (15) holds.

PROPOSITION 3. Assume that  $\Lambda = (\lambda^1, \dots, \lambda^N)$  is the set of Tikhonov–Phillips filters defined by (24) and that (19) holds. If  $N = N_{\mathcal{A}}N_{\mathcal{W}}, w_{\max} = w_{\max}(\varepsilon), w_1 = w_1(\varepsilon)$  are such that

$$(25) \quad \lim_{\varepsilon \rightarrow 0} \frac{\log(Nw_{\max}/w_1)}{w_1} = 0,$$

then for every  $\theta \in \ell_2$  and for the estimator  $\theta^* = (\theta_1^*, \theta_2^*, \dots)$  with  $\theta_i^* = \lambda_i^* X_i$  we have

$$\mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + o(1)) \inf_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta],$$

where  $o(1) \rightarrow 0$  uniformly in  $\theta \in \ell_2$ .

Note that condition (25) is very weak and it is easily checked in typical situations. For example, one can take  $N_{\mathcal{W}} = O(\varepsilon^{-a})$ ,  $w_{\max} = O(\varepsilon^{-b})$ , for some arbitrary fixed  $a > 0$ ,  $b > 0$ , and assume that  $N_{\mathcal{A}}$ ,  $\alpha_{\max}$  do not depend on  $\varepsilon$  (it is typical to consider just a small fixed number of integers  $\alpha$ , or even one integer  $\alpha$ ). This should be completed, in order to satisfy (25), by the mild assumption  $w_1 / \log^2(1/\varepsilon) \rightarrow \infty$ . Since there is no restriction on the power  $a$ , the discrete net  $\mathcal{W}$  can be arbitrarily fine, and it is not hard to show that optimality of our discretized rule extends to the set of filters (24) where  $w$  varies continuously in the interval.

**4. Preliminary lemmas.** Let  $\xi_i$  be i.i.d.  $\mathcal{N}(0, 1)$  random variables and let  $v = \{v_i\}_1^\infty \in \ell_2$  be a random sequence measurable w.r.t.  $\{\xi_i\}_{i=1}^\infty$ . It is assumed that  $v$  takes values in a finite set  $V$  of  $\ell_2$ -sequences:  $v \in V = \{v^1, \dots, v^N\}$ .

LEMMA 1. For any  $K \geq 1$ ,

$$\mathbf{E} \left| \sum_{i=1}^\infty v_i \xi_i \right| \leq \sqrt{2 \log(NK)} (\mathbf{E} \|v\| + \sqrt{2 \mathbf{E} \|v\|^2 / K}).$$

PROOF. It suffices to consider the case where  $\|v^j\| \neq 0, \forall v^j \in V$ . Denote  $\zeta_v = \|v\|^{-1} \sum_{i=1}^\infty v_i \xi_i$ , where the sums  $\|v\|^2$  and  $\sum_{i=1}^\infty v_i \xi_i$  are understood in the sense of mean squared convergence. By the Cauchy–Schwarz inequality,

$$\begin{aligned} \mathbf{E} \left| \sum_{i=1}^\infty v_i \xi_i \right| &\leq \mathbf{E} \|v\| \max_{v \in V} |\zeta_v| \leq \mathbf{E} \|v\| \max_{v \in V} |\zeta_v| I \left\{ \max_{v \in V} |\zeta_v| \leq \sqrt{2 \log(NK)} \right\} \\ &\quad + \mathbf{E} \|v\| \max_{v \in V} |\zeta_v| I \left\{ \max_{v \in V} |\zeta_v| > \sqrt{2 \log(NK)} \right\} \\ &\leq \sqrt{2 \log(NK)} \mathbf{E} \|v\| \\ &\quad + (\mathbf{E} \|v\|^2)^{1/2} \left( \mathbf{E} \max_{v \in V} |\zeta_v|^2 I \left\{ \max_{v \in V} |\zeta_v| > \sqrt{2 \log(NK)} \right\} \right)^{1/2}. \end{aligned}$$

Now, for the function  $F(t) = t^2 I\{t > \sqrt{2 \log(NK)}\}$  we use that

$$F\left(\max_{v \in V} |\zeta_v|\right) \leq \sum_{v \in V} F(|\zeta_v|),$$

and since  $\zeta_{v^1}, \dots, \zeta_{v^N}$  are identically distributed standard Gaussian variables, we get

$$\begin{aligned} \mathbf{E} \left| \sum_{i=1}^{\infty} v_i \xi_i \right| &\leq \sqrt{2 \log(NK)} \mathbf{E} \|v\| \\ &\quad + (\mathbf{E} \|v\|^2)^{1/2} \left( N \mathbf{E} |\zeta_{v^1}|^2 I \left\{ |\zeta_{v^1}| > \sqrt{2 \log(NK)} \right\} \right)^{1/2} \\ &\leq \sqrt{2 \log(NK)} (\mathbf{E} \|v\| + \sqrt{2 \mathbf{E} \|v\|^2 / K}), \end{aligned}$$

where for the last inequality we used

$$\mathbf{E} |\zeta_{v^1}|^2 I \{ |\zeta_{v^1}| > x \} \leq 2(2\pi)^{-1/2} (x + x^{-1}) \exp(-x^2/2) \quad \forall x > 0,$$

and  $NK \geq 2$ .  $\square$

LEMMA 2. Let  $\|v\| \neq 0$ ,  $\forall v \in V$ , and denote  $m(v) = \sup_i |v_i| / \|v\|$ ,  $m_V = \max_{v \in V} m(v)$ ,

$$(26) \quad M(q) = \sum_{v \in V} \exp\{-q/m(v)\},$$

where  $q > 0$ . Then for any  $K \geq 1$  we have

$$\mathbf{E} \left| \sum_{i=1}^{\infty} v_i (\xi_i^2 - 1) \right| \leq D \left( \sqrt{\log(NK)} + m_V \log(M(q)K) \right) (\mathbf{E} \|v\| + \sqrt{\mathbf{E} \|v\|^2 / K}),$$

where  $D$  is a constant depending only on  $q$ .

PROOF. Let  $\eta_v = (\sqrt{2}\|v\|)^{-1} \sum_{i=1}^{\infty} v_i (\xi_i^2 - 1)$ , where the sums  $\|v\|^2$  and  $\sum_{i=1}^{\infty} v_i (\xi_i^2 - 1)$  are understood in the sense of mean squared convergence. Using the Markov inequality and the formula

$$-\log(1-x) = \sum_{k=1}^{\infty} \frac{x^k}{k}$$

one obtains, for any  $0 < t < [\sqrt{2}m(v)]^{-1}$ ,

$$\begin{aligned} \mathbf{P}\{\eta_v > x\} &\leq \exp(-tx) \mathbf{E} \exp(t\eta_v) \\ &= \exp(-tx) \prod_{i=1}^{\infty} \exp\left\{ -\frac{tv_i}{\sqrt{2}\|v\|} - \frac{1}{2} \log\left(1 - \frac{\sqrt{2}tv_i}{\|v\|}\right) \right\} \\ &= \exp(-tx) \exp\left\{ \sum_{k=2}^{\infty} \sum_{i=1}^{\infty} \frac{1}{2k} \left( \frac{\sqrt{2}tv_i}{\|v\|} \right)^k \right\} \end{aligned}$$

$$\begin{aligned} &\leq \exp(-tx) \exp\left\{\frac{1}{m^2(v)} \sum_{k=2}^{\infty} \frac{1}{2k} [\sqrt{2tm(v)}]^k\right\} \\ &\leq \exp(-tx) \exp\left\{-\frac{1}{2m^2(v)} \log[1 - \sqrt{2tm(v)}] - \frac{t}{\sqrt{2m(v)}}\right\}. \end{aligned}$$

Minimization of the last expression with respect to  $t$  yields

$$\mathbf{P}\{\eta_v > x\} \leq \exp[\varphi_v(x)], \quad \varphi_v(x) = \frac{1}{2m^2(v)} \log[1 + \sqrt{2xm(v)}] - \frac{x}{\sqrt{2m(v)}}.$$

Note that for  $u \geq 0$  we have

$$\log(1+u) - u = u \int_0^1 \left(-\frac{\tau u}{1+\tau u}\right) d\tau \leq -\int_0^1 \frac{\tau u^2}{1+u} d\tau = -\frac{u^2}{2(1+u)}.$$

Thus

$$\varphi_v(x) \leq -\frac{x^2}{2(1 + \sqrt{2xm(v)})},$$

and we conclude that

$$(27) \quad \mathbf{P}\{|\eta_v| > x\} \leq 2 \exp\left\{-\frac{x^2}{2(1 + \sqrt{2xm(v)})}\right\} \quad \forall x > 0.$$

Using (27), we find, for any  $Q > 0$ ,

$$\begin{aligned} \mathbf{E}\eta_v^2 I\{|\eta_v| > Q\} &= 2 \int_Q^\infty x \mathbf{P}\{|\eta_v| > x\} dx \\ &\leq 4 \int_Q^\infty x \exp\left\{-\frac{x^2}{2(1 + \sqrt{2xm(v)})}\right\} dx. \end{aligned}$$

It is easy to see that

$$-\frac{x^2}{2(1 + \sqrt{2xm(v)})} \leq \begin{cases} -\frac{x^2}{4}, & \sqrt{2m(v)}x \leq 1, \\ -\frac{x}{\sqrt{32m(v)}}, & \sqrt{2m(v)}x > 1, \end{cases}$$

and we get by simple algebra, for  $Q > q\sqrt{32}$ ,

$$(28) \quad \mathbf{E}\eta_v^2 I\{|\eta_v| > Q\} \leq C \exp\left(-\frac{Q^2}{4}\right) + CQ \exp\left\{-\frac{Q}{\sqrt{32}m(v)}\right\}.$$

In view of (26) we have, for any  $Q > q\sqrt{32}$ ,

$$(29) \quad \sum_{v \in V} \exp\left(-\frac{Q}{\sqrt{32}m(v)}\right) \leq M(q) \exp\left(-\frac{Q/\sqrt{32}-q}{m_V}\right).$$

Now, acting as in the proof of Lemma 1 and using (28), (29) we obtain

$$\begin{aligned} & \mathbf{E} \left| \sum_{i=1}^{\infty} v_i (\xi_i^2 - 1) \right| \\ & \leq \sqrt{2} \mathbf{E} \|v\| \max_{v \in V} |\eta_v| \\ & \leq \sqrt{2} \mathbf{E} \|v\| \max_{v \in V} |\eta_v| I\left\{ \max_{v \in V} |\eta_v| \leq Q \right\} \\ & \quad + \sqrt{2} \mathbf{E} \|v\| \max_{v \in V} |\eta_v| I\left\{ \max_{v \in V} |\eta_v| > Q \right\} \\ & \leq \sqrt{2} Q \mathbf{E} \|v\| + \sqrt{2 \mathbf{E} \|v\|^2} \left( \mathbf{E} \max_{v \in V} \eta_v^2 I\left\{ \max_{v \in V} |\eta_v| > Q \right\} \right)^{1/2} \\ & \leq \sqrt{2} Q \mathbf{E} \|v\| + \sqrt{2 \mathbf{E} \|v\|^2} \left( \sum_{v \in V} \mathbf{E} \eta_v^2 I\{|\eta_v| > Q\} \right)^{1/2} \\ & \leq \sqrt{2} Q \mathbf{E} \|v\| + C \sqrt{\mathbf{E} \|v\|^2} \left[ N \exp\left(-\frac{Q^2}{4}\right) + M(q) Q \exp\left(-\frac{Q}{\sqrt{32}m_V}\right) \right]^{1/2}, \end{aligned}$$

for any  $Q > q\sqrt{32}$ . Choosing  $Q = 2\sqrt{\log(NK)} + \sqrt{32}m_V \log(M(q)K) + q\sqrt{32}$  one gets the lemma.  $\square$

Consider the estimator

$$\tilde{\theta}_i = \lambda_i(X) X_i,$$

where  $\lambda_i = \lambda_i(X) \in [-1, 1]$  depends on the data  $X$  [not necessarily  $\lambda_i(X) = \lambda_i^*(X)$ ]. We assume that the filter  $\lambda(X) = (\lambda_1(X), \lambda_2(X), \dots)$  takes values in the set of candidate filters  $\Lambda$ . In the next lemma we give a bound for the risk of this estimator. We need the following notation:

$$\Delta^\varepsilon[\lambda] = \varepsilon^2 L_\Lambda \sup_i \sigma_i^2 \lambda_i^2 \quad \forall \lambda \in \Lambda.$$



LEMMA 3. For any  $B > 0$  we have

$$\mathbf{E}_\theta \|\tilde{\theta} - \theta\|^2 \leq (1 + 2B^{-1})\mathbf{E}_\theta R_\varepsilon[\lambda(X), \theta] + CB\mathbf{E}_\theta \Delta^\varepsilon[\lambda(X)],$$

where  $C > 0$  is an absolute constant.

PROOF. Write

$$\begin{aligned} \mathbf{E}_\theta \|\tilde{\theta} - \theta\|^2 &= \mathbf{E}_\theta \left[ \sum_{i=1}^{\infty} (1 - \lambda_i(X))^2 \theta_i^2 + \varepsilon^2 \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) \right] \\ (30) \quad &\quad - 2\varepsilon \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i(X)) \theta_i \lambda_i(X) \sigma_i \xi_i \\ &\quad + \varepsilon^2 \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) (\xi_i^2 - 1). \end{aligned}$$

Using Lemma 1 with  $K = S$  we get

$$\begin{aligned} &\left| \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i(X)) \theta_i \lambda_i(X) \sigma_i \xi_i \right| \\ &\leq \sqrt{2 \log(NS)} \mathbf{E}_\theta \left[ \sum_{i=1}^{\infty} (1 - \lambda_i(X))^2 \theta_i^2 \lambda_i^2(X) \sigma_i^2 \right]^{1/2} \\ &\quad + 2\sqrt{\log(NS)} S^{-1/2} \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i(X))^2 \theta_i^2 \lambda_i^2(X) \sigma_i^2 \right]^{1/2} \\ &\leq \sqrt{2 \log(NS)} \mathbf{E}_\theta \left[ \sup_i |\sigma_i| |\lambda_i(X)| \left( \sum_{i=1}^{\infty} (1 - \lambda_i(X))^2 \theta_i^2 \right)^{1/2} \right] \\ &\quad + 2\sqrt{\log(NS)} S^{-1/2} \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i(X))^2 \theta_i^2 \right]^{1/2} \max_{\lambda \in \Lambda} \sup_i |\sigma_i| |\lambda_i|. \end{aligned}$$

Now (13) and the elementary inequality  $2ab \leq B^{-1}a^2 + Bb^2$ ,  $\forall B > 0$ , yield, for any  $B > 0$ ,

$$\begin{aligned} &\varepsilon \mathbf{E}_\theta \left| \sum_{i=1}^{\infty} [1 - \lambda_i(X)] \theta_i \lambda_i(X) \sigma_i \xi_i \right| \\ (31) \quad &\leq B^{-1} \mathbf{E}_\theta \sum_{i=1}^{\infty} [1 - \lambda_i(X)]^2 \theta_i^2 \\ &\quad + \varepsilon^2 B \log(NS) \mathbf{E}_\theta \sup_i \sigma_i^2 \lambda_i^2(X) + \varepsilon^2 B \log(NS) \min_{\lambda \in \Lambda} \sup_i \sigma_i^2 \lambda_i^2 \\ &\leq B^{-1} \mathbf{E}_\theta \sum_{i=1}^{\infty} [1 - \lambda_i(X)]^2 \theta_i^2 + 2B\mathbf{E}_\theta \Delta^\varepsilon[\lambda(X)]. \end{aligned}$$

To bound the last term in (30) we use Lemma 2 with  $v_i = \sigma_i^2 \lambda_i^2(X)$ ,  $K = S$ ,  $q = 1$ . The inequalities  $|\lambda_i(X)| \leq 1$  and (12) entail  $m(v) \leq \rho(\lambda(X)) \leq \sqrt{C_1}$ . Thus,  $M(1) = \sum_{v \in V} \exp\{-1/m(v)\} \leq M$ . Note also that  $m_V \leq \rho$  and  $\sqrt{\log(NS)} + m_V \log(M(q)S) \leq \sqrt{2L_\Lambda}$ . Therefore, application of Lemma 2 yields

$$\begin{aligned} & \varepsilon^2 \mathbf{E}_\theta \left| \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) (\xi_i^2 - 1) \right| \\ & \leq D\sqrt{2L_\Lambda} \varepsilon^2 \left\{ \mathbf{E}_\theta \left[ \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4(X) \right]^{1/2} + S^{-1/2} \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4(X) \right]^{1/2} \right\}. \end{aligned}$$

Next, by (13),

$$\begin{aligned} S^{-1/2} \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4(X) \right]^{1/2} & \leq \min_{\lambda \in \Lambda} \sup_i \sigma_i \lambda_i \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) \right]^{1/2} \\ & \leq \left[ \mathbf{E}_\theta \sup_i \sigma_i^2 \lambda_i^2(X) \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) \right]^{1/2}. \end{aligned}$$

Hence, for any  $B > 0$ ,

$$\begin{aligned} & \varepsilon^2 \mathbf{E}_\theta \left| \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) (\xi_i^2 - 1) \right| \\ & \leq D\sqrt{2L_\Lambda} \varepsilon^2 \mathbf{E}_\theta \left[ \sup_i |\sigma_i| |\lambda_i(X)| \left( \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) \right)^{1/2} \right] \\ (32) \quad & + D\sqrt{2L_\Lambda} \varepsilon^2 \left[ \mathbf{E}_\theta \sup_i \sigma_i^2 \lambda_i^2(X) \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) \right]^{1/2} \\ & \leq 2\varepsilon^2 B^{-1} \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) + D^2 L_\Lambda \varepsilon^2 B \mathbf{E}_\theta \sup_i \sigma_i^2 \lambda_i^2(X) \\ & = 2\varepsilon^2 B^{-1} \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2(X) + D^2 B \mathbf{E}_\theta \Delta^\varepsilon[\lambda(X)]. \end{aligned}$$

Combining (30)–(32) we complete the proof.  $\square$

**5. Proof of the theorems.** Denote by  $\lambda^0 = (\lambda_1^0, \lambda_2^0, \dots)$  the oracle  $\lambda^0 = \arg \min_{\lambda \in \Lambda} R_\varepsilon[\lambda, \theta]$ . We have

$$\begin{aligned}
\mathbf{E}_\theta U[\lambda^*, X] &= \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] - \sum_{i=1}^{\infty} \theta_i^2 \\
&\quad + 2\varepsilon \mathbf{E}_\theta \sum_{i=1}^{\infty} (\lambda_i^{*2} - 2\lambda_i^*) \theta_i \sigma_i \xi_i + \varepsilon^2 \mathbf{E}_\theta \sum_{i=1}^{\infty} (\lambda_i^{*2} - 2\lambda_i^*) \sigma_i^2 (\xi_i^2 - 1) \\
(33) \quad &= \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] - \sum_{i=1}^{\infty} \theta_i^2 \\
&\quad + 2\varepsilon \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i^*)^2 \theta_i \sigma_i \xi_i + \varepsilon^2 \mathbf{E}_\theta \sum_{i=1}^{\infty} (\lambda_i^{*2} - 2\lambda_i^*) \sigma_i^2 (\xi_i^2 - 1).
\end{aligned}$$

We now bound the last two terms in (33). First, note that

$$\begin{aligned}
(34) \quad [(1 - \lambda_i^*)^2 - (1 - \lambda_i^0)^2]^2 &= [(1 - \lambda_i^*) + (1 - \lambda_i^0)]^2 [\lambda_i^* - \lambda_i^0]^2 \\
&\leq 2[(1 - \lambda_i^*)^2 + (1 - \lambda_i^0)^2] (\lambda_i^{*2} + \lambda_i^{02}).
\end{aligned}$$

Then we have, by Lemma 1 with  $K = S$  and (34),

$$\begin{aligned}
&\varepsilon \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i^*)^2 \theta_i \sigma_i \xi_i \\
&= \varepsilon \mathbf{E}_\theta \sum_{i=1}^{\infty} [(1 - \lambda_i^*)^2 - (1 - \lambda_i^0)^2] \theta_i \sigma_i \xi_i \\
&\geq -\varepsilon \mathbf{E}_\theta \left| \sum_{i=1}^{\infty} [(1 - \lambda_i^*)^2 - (1 - \lambda_i^0)^2] \theta_i \sigma_i \xi_i \right| \\
&\geq -\varepsilon \sqrt{2 \log(NS)} \mathbf{E}_\theta \left\{ \sum_{i=1}^{\infty} [(1 - \lambda_i^*)^2 - (1 - \lambda_i^0)^2]^2 \theta_i^2 \sigma_i^2 \right\}^{1/2} \\
&\quad - 2\varepsilon \sqrt{\log(NS)/S} \left\{ \mathbf{E}_\theta \sum_{i=1}^{\infty} [(1 - \lambda_i^*)^2 - (1 - \lambda_i^0)^2]^2 \theta_i^2 \sigma_i^2 \right\}^{1/2} \\
&\geq -2\varepsilon \sqrt{\log(NS)} \mathbf{E}_\theta \left\{ \sum_{i=1}^{\infty} [(1 - \lambda_i^*)^2 + (1 - \lambda_i^0)^2] (\lambda_i^{*2} + \lambda_i^{02}) \theta_i^2 \sigma_i^2 \right\}^{1/2} \\
&\quad - 2\varepsilon \sqrt{2 \log(NS)/S} \left\{ \mathbf{E}_\theta \sum_{i=1}^{\infty} [(1 - \lambda_i^*)^2 + (1 - \lambda_i^0)^2] (\lambda_i^{*2} + \lambda_i^{02}) \theta_i^2 \sigma_i^2 \right\}^{1/2}.
\end{aligned}$$

The argument as in the proof of Lemma 3 gives, for any  $B > 0$ ,

$$\begin{aligned} & 2\varepsilon\sqrt{\log(NS)}\mathbf{E}_\theta\left\{\sum_{i=1}^{\infty}[(1-\lambda_i^*)^2+(1-\lambda_i^0)^2](\lambda_i^{*2}+\lambda_i^{02})\theta_i^2\sigma_i^2\right\}^{1/2} \\ & \leq \frac{1}{2B}\mathbf{E}_\theta\sum_{i=1}^{\infty}[(1-\lambda_i^*)^2+(1-\lambda_i^0)^2]\theta_i^2 \\ & \quad + 2B\varepsilon^2\log(NS)\mathbf{E}_\theta\left\{\sup_i\sigma_i^2\lambda_i^{*2}+\sup_i\sigma_i^2\lambda_i^{02}\right\} \end{aligned}$$

and

$$\begin{aligned} & 2\varepsilon\sqrt{2\log(NS)/S}\left\{\mathbf{E}_\theta\sum_{i=1}^{\infty}[(1-\lambda_i^*)^2+(1-\lambda_i^0)^2](\lambda_i^{*2}+\lambda_i^{02})\theta_i^2\sigma_i^2\right\}^{1/2} \\ & \leq 2\varepsilon\sqrt{2\log(NS)/S}\left\{\mathbf{E}_\theta\sum_{i=1}^{\infty}[(1-\lambda_i^*)^2+(1-\lambda_i^0)^2]\theta_i^2\right\}^{1/2} \\ & \quad \times \left(\max_{\lambda\in\Lambda}\sup_i\{\sigma_i^2\lambda_i^2+\sigma_i^2\lambda_i^{02}\}\right)^{1/2} \\ & \leq \frac{1}{2B}\mathbf{E}_\theta\sum_{i=1}^{\infty}[(1-\lambda_i^*)^2+(1-\lambda_i^0)^2]\theta_i^2 \\ & \quad + 4B\varepsilon^2\log(NS)\mathbf{E}_\theta\left\{\sup_i\sigma_i^2\lambda_i^{*2}+\sup_i\sigma_i^2\lambda_i^{02}\right\}. \end{aligned}$$

Putting together these inequalities, we find

$$\begin{aligned} & \varepsilon\mathbf{E}_\theta\sum_{i=1}^{\infty}(1-\lambda_i^*)^2\theta_i\sigma_i\xi_i \\ & \geq -B^{-1}\mathbf{E}_\theta\sum_{i=1}^{\infty}(1-\lambda_i^*)^2\theta_i^2 \\ (35) \quad & - B^{-1}\sum_{i=1}^{\infty}(1-\lambda_i^0)^2\theta_i^2 - 6B\varepsilon^2\log(NS)\mathbf{E}_\theta\left\{\sup_i\sigma_i^2\lambda_i^{*2}+\sup_i\sigma_i^2\lambda_i^{02}\right\} \\ & \geq -B^{-1}\mathbf{E}_\theta\sum_{i=1}^{\infty}(1-\lambda_i^*)^2\theta_i^2 - B^{-1}R_\varepsilon[\lambda^0,\theta] - 6B\mathbf{E}_\theta\Delta^\varepsilon[\lambda^*] - 6B\Delta^\varepsilon[\lambda^0]. \end{aligned}$$

Now, we bound the last term in (33) using Lemma 2 with  $K = S$ ,  $q = 3$  and  $v_i = (\lambda_i^{*2} - 2\lambda_i^*)\sigma_i^2$ . Note that

$$\lambda_i^2 \leq (\lambda_i^2 - 2\lambda_i)^2 \leq 9\lambda_i^2 \quad \forall |\lambda_i| \leq 1.$$

This and (12) entail

$$m(v) = \frac{\sup_i |\lambda_i^2 - 2\lambda_i| \sigma_i^2}{[\sum_{i=1}^{\infty} \sigma_i^4 (\lambda_i^2 - 2\lambda_i)^2]^{1/2}} \leq 3\rho(\lambda) \leq 3\sqrt{C_1}.$$

Hence,  $M(3) = \sum_{v \in V} \exp\{-3/m(v)\} \leq M$ . Furthermore,  $m_V \leq 3\rho$ , and

$$\sqrt{\log(NS)} + m_V \log(M(3)S) \leq C\sqrt{L_\Lambda}.$$

Therefore, by Lemma 2, we obtain

$$\begin{aligned} & \varepsilon^2 \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i^*)^2 \sigma_i^2 (\xi_i^2 - 1) \\ & \geq -C\sqrt{L_\Lambda} \varepsilon^2 \mathbf{E}_\theta \left[ \sum_{i=1}^{\infty} (\lambda_i^{*2} - 2\lambda_i^*)^2 \sigma_i^4 \right]^{1/2} \\ & \quad - C\sqrt{L_\Lambda} \varepsilon^2 S^{-1/2} \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} (\lambda_i^{*2} - 2\lambda_i^*)^2 \sigma_i^4 \right]^{1/2} \\ & \geq -C\sqrt{L_\Lambda} \varepsilon^2 \left\{ \mathbf{E}_\theta \left[ \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^{*4} \right]^{1/2} + S^{-1/2} \left[ \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^{*4} \right]^{1/2} \right\}. \end{aligned}$$

Here and below in this section  $C$  is a generic notation for positive constants that depend only on  $C_1$ . Repeating the argument of (32) to bound the last expression we finally get

$$\begin{aligned} & \varepsilon^2 \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i^*)^2 \sigma_i^2 (\xi_i^2 - 1) \\ (36) \quad & \geq -2\varepsilon^2 B^{-1} \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^{*2} - C\varepsilon^2 B L_\Lambda \mathbf{E}_\theta \sup_k \sigma_k^2 \lambda_k^{*2} \\ & = -2\varepsilon^2 B^{-1} \mathbf{E}_\theta \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^{*2} - C B \mathbf{E}_\theta \Delta_\varepsilon[\lambda^*]. \end{aligned}$$

Now we are ready to complete the proof of Theorems 1 and 2. From (35), (36) we have

$$\begin{aligned} & 2\varepsilon \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i^*)^2 \theta_i \sigma_i \xi_i + \varepsilon^2 \mathbf{E}_\theta \sum_{i=1}^{\infty} (1 - \lambda_i^*)^2 \sigma_i^2 (\xi_i^2 - 1) \\ & \geq -2B^{-1} R_\varepsilon[\lambda^*, \theta] - 2B^{-1} R_\varepsilon[\lambda^0, \theta] - C B \mathbf{E}_\theta \Delta_\varepsilon[\lambda^*] - C B \Delta_\varepsilon[\lambda^0]. \end{aligned}$$

This and (33) yield

$$(37) \quad \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] \leq \mathbf{E}_\theta U[\lambda^*, X] + \sum_{i=1}^{\infty} \theta_i^2 + 2B^{-1} \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] \\ + C B \mathbf{E}_\theta \Delta_\varepsilon[\lambda^*] + 2B^{-1} R_\varepsilon[\lambda^0, \theta] + C B \Delta^\varepsilon[\lambda^0].$$

By definition of  $\lambda^*$  and (9),

$$\mathbf{E}_\theta U[\lambda^*, X] \leq \mathbf{E}_\theta U[\lambda^0, X] = R_\varepsilon[\lambda^0, \theta] - \sum_{i=1}^{\infty} \theta_i^2.$$

This and (37) imply

$$(38) \quad (1 - 2B^{-1}) \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] \leq (1 + 2B^{-1}) R_\varepsilon[\lambda^0, \theta] \\ + C B \mathbf{E}_\theta \Delta_\varepsilon[\lambda^*] + C B \Delta^\varepsilon[\lambda^0].$$

Next, by Lemma 3 for any  $B > 0$ ,

$$(39) \quad \mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + 2B^{-1}) \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] + C B \mathbf{E}_\theta \Delta^\varepsilon[\lambda^*].$$

Inequalities (38) and (39) entail Theorems 1 and 2. In fact, to get Theorem 1 we use the following argument. Note that, for any  $x > 0$ ,

$$\sup_i \sigma_i^2 \lambda_i^2 = \sup_i \sigma_i^2 \lambda_i^2 I \left\{ x \sup_i \sigma_i^2 \lambda_i^2 < \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 \right\} \\ + \sup_i \sigma_i^2 \lambda_i^2 I \left\{ x \sup_i \sigma_i^2 \lambda_i^2 \geq \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 \right\} \\ \leq \frac{1}{x} \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 + \max_{\lambda \in \Lambda} \sup_i \sigma_i^2 \lambda_i^2 I \left\{ x \sup_i \sigma_i^2 \lambda_i^2 \geq \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 \right\} \\ = \frac{1}{x} \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 + \omega(x) \leq \frac{1}{x \varepsilon^2} R_\varepsilon[\lambda, \theta] + \omega(x) \quad \forall \theta \in \ell_2.$$

This inequality with  $x = B^2 L_\Lambda$ , the definition of  $\Delta^\varepsilon[\lambda]$  and (38), (39) yield Theorem 1.

PROOF OF THEOREM 2. Note first that, in view of (11),

$$\frac{\sup_i \sigma_i^2 \lambda_i^2}{\sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2} \leq \rho \frac{\sup_i \sigma_i^2 \lambda_i^2}{(\sum_{i=1}^{\infty} \sigma_i^4 \lambda_i^4)^{1/2}} \leq \rho^2.$$

Thus,

$$(40) \quad \Delta^\varepsilon[\lambda] \leq \rho^2 L_\Lambda R_\varepsilon[\lambda, \theta] \quad \forall \lambda \in \Lambda, \theta \in \ell_2.$$

Substitution of (40) into (38) and (39) gives, respectively,

$$(41) \quad \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta] \leq \frac{1 + 2B^{-1} + \gamma_5 B \rho^2 L_\Delta}{1 - 2B^{-1} - \gamma_6 B \rho^2 L_\Delta} R_\varepsilon[\lambda^0, \theta]$$

(provided the denominator here is positive) and

$$(42) \quad \mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + 2B^{-1} + \gamma_7 B \rho^2 L_\Delta) \mathbf{E}_\theta R_\varepsilon[\lambda^*, \theta].$$

The constants  $\gamma_i$  here are positive and depend only on  $C_1$ . There exists  $\gamma_4 > 0$  small enough such that if  $\rho^2 L_\Delta < \gamma_4$ , the choice of  $B = (\rho^2 L_\Delta)^{-1/2}$  satisfies the inequality  $2B^{-1} + \gamma_6 B \rho^2 L_\Delta < 1/2$ . With this choice of  $B$ , (41) and (42) entail Theorem 2.  $\square$

**6. Application to sharp adaptive estimation.** In this section we apply the oracle inequalities of Section 2 to show that sharp minimax adaptive estimators for inverse problems can be obtained by the principle of unbiased risk estimation. We study the problem where sharp adaptive estimators were not known previously, namely a recovery of *anisotropic* smooth functions from indirect noisy data. For brevity, we restrict the discussion to a specific example (measuring the temperature of the earth). However, the key elements of the proofs are given under general assumptions and the result can be easily extended.

Let  $\phi = (\phi_1, \phi_2)$  be the polar coordinates of a point on the surface of the earth (we suppose for simplicity that the earth is a sphere). Consider the problem of measuring the temperature  $t(\phi)$  at the point  $\phi$ . The function  $t(\phi)$  is sufficiently smooth (since it is a solution of a thermal conductivity equation) and, of course, periodic. More specifically, we assume that  $t(\phi)$  belongs to the following anisotropic Sobolev ball:

$$(43) \quad W_2^m(p) = \left\{ t : \int_T \left[ p_1^2 \left( \frac{\partial^{m_1} t}{\partial \phi_1^{m_1}} \right)^2 + p_2^2 \left( \frac{\partial^{m_2} t}{\partial \phi_2^{m_2}} \right)^2 \right] d\phi \leq 1 \right\},$$

where  $T = [0, 2\pi] \times [0, 2\pi]$  and the parameters  $m = (m_1, m_2)$ ,  $p = (p_1, p_2)$  characterize the smoothness of the function  $t(\phi)$  with respect to  $\phi_1$  and  $\phi_2$  ( $m_i \in \mathbb{N}$ ,  $p_i > 0$ ). The anisotropic character of smoothness here is important and reflects the fact that the temperature changes more rapidly along the meridians than along the parallels.

Next, it is assumed that temperature is measured by means of remote infrared detectors located for instance on satellites or on planes. It means that one cannot measure temperature at a given point directly, but one rather measures an average temperature in a vicinity of this point, that is,

$$Et(\phi) = \frac{1}{C(\alpha_0)} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\{-2\pi\alpha_0\|\phi - \phi'\|\} t(\phi'_1, \phi'_2) d\phi'_1 d\phi'_2.$$

Here  $C(\alpha_0)$  is the normalizing constant,

$$C(\alpha_0) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\{-2\pi\alpha_0\|\phi\|\} d\phi_1 d\phi_2$$

and  $\|\phi\|^2 = \phi_1^2 + \phi_2^2$ . The parameter  $\alpha_0^{-1}$  characterizes the size of the vicinity of the point where we measure the temperature. (We consider here and below integration over the whole line instead of integration over  $T$ , in order to simplify the technical details and the resulting sequence model.)

The measurements are contaminated by a noise inherent to infrared detection. We may write the resulting model as

$$(44) \quad y(\phi) = Et(\phi) + \varepsilon n(\phi), \quad \phi \in T,$$

where  $n(\phi)$  is a standard white Gaussian noise in the Hilbert space of periodic functions. Writing (44) means that for every  $g \in L_2(T)$  we observe the random variable  $Y(g) = (g, Et) + \varepsilon \xi(g)$ , where  $\xi(g) \sim \mathcal{N}(0, \|g\|_2)$  and  $\mathbf{E}(\xi(g)\xi(g')) = 0$  if  $(g, g') = 0$  [here  $(\cdot, \cdot)$  and  $\|\cdot\|$  are the scalar product and the norm in  $L_2(T)$ ].

Since the function  $t(\phi_1, \phi_2)$  is periodic it is convenient to consider the statistical model (44) in the Fourier domain. Denote by

$$\theta_{kl} = \int_T \exp(ik\phi_1 + il\phi_2)t(\phi) d\phi, \quad k, l = 0, \pm 1, \pm 2, \dots,$$

the Fourier coefficients of  $t(\phi)$ . Then (44) is equivalent to the model

$$(45) \quad z_{kl} = \theta_{kl} + \varepsilon \sigma_{kl} \xi_{kl},$$

where  $z_{kl} = Y(\exp(ik\phi_1 + il\phi_2))$ ,  $\xi_{kl}$  are i.i.d.  $\mathcal{N}(0, 1)$  random variables and

$$\sigma_{kl}^{-2} = \frac{1}{C(\alpha_0)} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\{ik\phi_1 + il\phi_2 - 2\pi\alpha_0\|\phi\|\} d\phi = \frac{\alpha_0^3}{(\alpha_0^2 + k^2 + l^2)^{3/2}}.$$

[For the last equality see, for instance, Stein and Weiss (1971), Theorem 1.14.] Note also that the image of Sobolev's ball  $W_2^m(p)$  in the Fourier domain has the form

$$(46) \quad \Theta_2^m(p) = \left\{ \{\theta_{kl}\} : \sum_{k,l=-\infty}^{\infty} \theta_{kl}^2 a_{kl}^2 \leq 1 \right\},$$

where  $a_{kl} = [p_1^2(2\pi k)^{2m_1} + p_2^2(2\pi l)^{2m_2}]^{1/2}$ . Thus our problem is equivalent to recovering the Fourier coefficients  $\theta = \{\theta_{kl}\}$  based on the data (45) under the prior information provided by (46).

Recall that an estimator  $\widehat{\theta}^M$  is called asymptotically minimax on  $\Theta_2^m(p)$  if, as  $\varepsilon \rightarrow 0$ ,

$$r^\varepsilon[\Theta_2^m(p)] = \inf_{\widehat{\theta}} \sup_{\theta \in \Theta_2^m(p)} \mathbf{E}\|\widehat{\theta} - \theta\|^2 = (1 + o(1)) \sup_{\theta \in \Theta_2^m(p)} \mathbf{E}\|\widehat{\theta}^M - \theta\|^2,$$

where the infimum is taken over all estimators. The next proposition derived from Pinsker's (1980) theorem gives an asymptotically minimax estimator and its risk.

Denote by  $\mu$  a root of the equation

$$(47) \quad \frac{\varepsilon^2}{\mu} \sum_{k,l} \sigma_{kl}^2 a_{kl} [1 - \mu a_{kl}]_+ = 1.$$



PROPOSITION 4. *The linear estimator*

$$(48) \quad \widehat{\theta}_{kl}^M = [1 - \mu a_{kl}]_+ z_{kl}, \quad k, l = 0, \pm 1, \pm 2, \dots,$$

is asymptotically minimax on  $\Theta_2^m(p)$  with risk

$$\sup_{\theta \in \Theta_2^m(p)} \mathbf{E} \|\widehat{\theta}^M - \theta\|^2 = \varepsilon^2 \sum_{k,l} \sigma_{kl}^2 [1 - \mu a_{kl}]_+.$$

Using this proposition one can compute the asymptotic minimax risk with the exact constant. The calculations are cumbersome and we do not reproduce them here, since the knowledge of the exact asymptotic minimax risk is not needed to prove sharp adaptivity of our estimator. We need only to evaluate the order of  $r^\varepsilon[\Theta_2^m(p)]$ :

LEMMA 4. *Let  $m_2 \geq m_1$ . Then, as  $\varepsilon \rightarrow 0$ ,*

$$r^\varepsilon[\Theta_2^m(p)] \leq C (\alpha_0^3 p_1^{4/m_1} p_2^{1/m_2})^{-\gamma} \varepsilon^{2\gamma}$$

and

$$\mu = O(\varepsilon^\gamma),$$

where the constant  $C > 0$  does not depend on  $\alpha_0, p_1, p_2$  and

$$\gamma = \frac{2m_1 m_2}{2m_1 m_2 + m_1 + 4m_2}.$$

REMARK 4. We can write  $\gamma = 2s/(2s + 1)$ , where

$$\frac{1}{s} = \frac{1}{m_2} + \frac{4}{m_1} = \frac{1}{\max(m_1, m_2)} + \frac{4}{\min(m_1, m_2)}.$$

It is remarkable that the respective roles of  $m_1$  and  $m_2$  are not equivalent, unlike the case where  $\sigma_{kl} \equiv 1$ . The coefficient  $\sigma_{kl}$  is symmetric in  $k, l$ , and nevertheless the rate is asymmetric. This effect is due to the anisotropy of the class  $\Theta_2^m(p)$ .

PROOF OF LEMMA 4. Denote  $V_i = (2\pi)^{-1}(p_i \mu)^{-1/m_i}, i = 1, 2$ . Since  $\sigma_{kl}^2$  is increasing in  $k, l$  we have

$$(49) \quad \varepsilon^2 \sum_{k,l} \sigma_{kl}^2 [1 - \mu a_{kl}]_+ \leq C \varepsilon^2 V_1 V_2 \sigma_{V_1 V_2}^2.$$

Here and later we denote by  $C$  positive constants, possibly different on different occasions. It is clear from (47) that  $\mu \rightarrow 0$  (and thus  $V_1 \rightarrow \infty, V_2 \rightarrow \infty$ ), as

$\varepsilon \rightarrow 0$ . Consequently, for  $\varepsilon$  small enough we get

$$\begin{aligned}
 & \sum_{k \geq V_1/2, l \geq V_2/2} a_{kl} [1 - \mu a_{kl}]_+ \\
 & \geq 2^{-m_2} \mu^{-1} \sum_{k \geq V_1/2, l \geq V_2/2} \left[ 1 - \sqrt{(k/V_1)^{2m_1} + (l/V_2)^{2m_2}} \right]_+ \\
 (50) \quad & \geq 2^{-m_2} \mu^{-1} \int_{x \geq V_1/2+1} \int_{y \geq V_2/2+1} \left[ 1 - \sqrt{(x/V_1)^{2m_1} + (y/V_2)^{2m_2}} \right]_+ dx dy \\
 & \geq 2^{-m_2} V_1 V_2 \mu^{-1} \int_{x \geq 1/2+1/V_1} \int_{y \geq 1/2+1/V_2} \left[ 1 - \sqrt{x^2 + y^2} \right]_+ dx dy \\
 & \geq C V_1 V_2 \mu^{-1},
 \end{aligned}$$

where we used that  $m_1 \geq 1, m_2 \geq 1$ . Noting that  $\sigma_{\lfloor k/2 \rfloor \lfloor l/2 \rfloor}^2 \geq c \sigma_{kl}^2$ , where  $c > 0$  does not depend on  $k, l$ , and using (50) one obtains

$$\begin{aligned}
 1 &= \frac{\varepsilon^2}{\mu} \sum_{k,l} \sigma_{kl}^2 a_{kl} [1 - \mu a_{kl}]_+ \geq \frac{c \varepsilon^2}{\mu} \sum_{|k| \geq V_1/2, |l| \geq V_2/2} \sigma_{V_1 V_2}^2 a_{kl} [1 - \mu a_{kl}]_+ \\
 &\geq C \varepsilon^2 \mu^{-2} V_1 V_2 \sigma_{V_1 V_2}^2 = C \alpha_0^{-3} \varepsilon^2 \mu^{-2} V_1 V_2 (\alpha_0^2 + V_1^2 + V_2^2)^{3/2}.
 \end{aligned}$$

The proof follows now from this inequality, (49), Proposition 4 and simple algebra.  $\square$

The minimax estimator  $\hat{\theta}^M$  defined by (48) depends on the parameters  $(m, p)$  of the functional class  $\Theta_2^m(p)$ , which are usually not known in practice. To overcome this drawback of  $\hat{\theta}^M$  we construct another estimator which is asymptotically minimax but does not depend on the parameters of the Sobolev ball. This estimator is called sharp adaptive. We show that a sharp adaptive estimator can be obtained by using the filter (10) with appropriately chosen set  $\Lambda$ .

Let  $\Lambda_0$  be the set of all filters with weights  $\lambda_{kl}$  having the form

$$\lambda_{kl}(W, \beta) = \lambda_{kl}(W_1, W_2, \beta_1, \beta_2) = \left[ 1 - \sqrt{(k/W_1)^{2\beta_1} + (l/W_2)^{2\beta_2}} \right]_+,$$

where  $W \in [\log(1/\varepsilon), \infty) \times [\log(1/\varepsilon), \infty)$ ,  $\beta \in [1, \infty) \times [1, \infty)$ . The minimax filter belongs to this set for  $\varepsilon$  small enough [cf. (48) and note that  $\mu = O(\varepsilon^\gamma)$  by Lemma 4]. Unfortunately  $\Lambda_0$  contains infinitely many elements and we cannot apply Theorem 1 for  $\Lambda = \Lambda_0$ . Note also that the direct numerical minimization of the unbiased risk estimator over  $\Lambda_0$  is very time consuming. Therefore we look for a finite subfamily which approximates well the filters in  $\Lambda_0$ . To make the estimator numerically feasible we will pick an approximating set containing the ‘‘minimal’’ number of elements.

Let  $\delta = \log^{-1}(1/\varepsilon)$ . Denote  $w_\delta(k) = (1 + \delta)^k$ ,  $k = 1, 2, \dots$ , and  $\beta_\delta(k) = (1 - \delta k)^{-1}$ ,  $k = 1, 2, \dots, \lfloor 1/\delta \rfloor$ . Define the approximating finite subfamily  $\Lambda_\delta$  as

a set of filters of the following form:  $\lambda = \{\lambda_{kl}(w_\delta(i), w_\delta(j), \beta_\delta(s), \beta_\delta(p))\}$ , where  $i, j, s, p \in \mathbb{N}$  are such that  $(w_\delta(i), w_\delta(j)) \in [\log(1/\varepsilon), \varepsilon^{-1}] \times [\log(1/\varepsilon), \varepsilon^{-1}]$  and  $(\beta_\delta(s), \beta_\delta(p)) \in (1, \sqrt{\log(1/\varepsilon)}] \times (1, \sqrt{\log(1/\varepsilon)}]$ .

Although, for the sake of simplicity, we assumed here a specific form of  $\sigma_{kl}$ , our results hold for the general situation where  $\sigma_{kl}$  satisfy the following assumption.

**ASSUMPTION 3.** The sequence  $\sigma_{kl}^2$  is positive, monotone nondecreasing in  $k$  and  $l$ , and there exists  $c_* > 0$  such that  $\sigma_{\lfloor k/2 \rfloor \lfloor l/2 \rfloor}^2 \geq c_* \sigma_{kl}^2$  for all  $k, l$ .

**LEMMA 5.** *Let Assumption 3 hold. For any  $\lambda \in \Lambda_0$  there exists a filter  $\lambda' \in \Lambda_\delta$  such that, for all  $\theta \in l_2$ ,*

$$(51) \quad \mathbf{E}_\theta \|\widehat{\theta}(\lambda') - \theta\|^2 \leq (1 + C\delta) \mathbf{E}_\theta \|\widehat{\theta}(\lambda) - \theta\|^2,$$

where  $C > 0$  does not depend on  $\theta, \lambda, \lambda'$ .

**PROOF.** Since

$$\mathbf{E}_\theta \|\widehat{\theta}(\lambda) - \theta\|^2 = \sum_{k,l} (1 - \lambda_{kl})^2 \theta_{kl}^2 + \varepsilon^2 \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2$$

it suffices to show that for any  $\lambda \in \Lambda_0$  there exists a filter  $\lambda' \in \Lambda_\delta$  such that  $\lambda_{kl} \leq \lambda'_{kl}$  for all  $k, l$  and

$$\sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2 \leq (1 + C\delta) \sum_{k,l} \sigma_{kl}^2 \lambda'_{kl}^2.$$

To prove this, first note that the componentwise inequality  $(W, \beta) \leq (W', \beta')$  implies

$$(52) \quad \lambda_{kl}(W, \beta) \leq \lambda_{kl}(W', \beta').$$

Fix  $(W, \beta) = (W_1, W_2, \beta_1, \beta_2)$  and define

$$\begin{aligned} i_0 &= \max\{i : w_\delta(i) \leq W_1\}, & j_0 &= \max\{j : w_\delta(j) \leq W_2\}, \\ s_0 &= \max\{s : \beta_\delta(s) \leq \beta_1\}, & p_0 &= \max\{p : \beta_\delta(p) \leq \beta_2\}. \end{aligned}$$

Set

$$\begin{aligned} (W^0, \beta^0) &= (w_\delta(i_0), w_\delta(j_0), \beta_\delta(s_0), \beta_\delta(p_0)), \\ (W^1, \beta^1) &= (w_\delta(i_0 + 1), w_\delta(j_0 + 1), \beta_\delta(s_0 + 1), \beta_\delta(p_0 + 1)). \end{aligned}$$

We have  $(W^0, \beta^0) \leq (W, \beta) \leq (W^1, \beta^1)$ , and thus

$$(53) \quad \lambda_{kl}(W^0, \beta^0) \leq \lambda_{kl}(W, \beta) \leq \lambda_{kl}(W^1, \beta^1).$$

Set  $\lambda'_{kl} = \lambda_{kl}(W^1, \beta^1)$ . In view of (53) it suffices to show that

$$(54) \quad \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2(W^1, \beta^1) \leq (1 + C\delta) \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2(W^0, \beta^0).$$

Let  $V(W_1, W_2, \beta_1, \beta_2) = \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2(W_1, W_2, \beta_1, \beta_2)$ . From Lemma 6 of the Appendix we have

$$\begin{aligned} & V(W^1, \beta^1) - V(W^0, \beta^0) \\ & \leq CV(W^0, \beta^0) \max_k \left( \frac{w_\delta(k+1) - w_\delta(k)}{w_\delta(k)} + \frac{\beta_\delta(k+1) - \beta_\delta(k)}{\beta_\delta^2(k+1)} \right). \end{aligned}$$

This and the obvious inequalities

$$\beta_\delta(k+1) - \beta_\delta(k) \leq \delta \beta_\delta^2(k+1), \quad w_\delta(k+1) - w_\delta(k) = \delta w_\delta(k)$$

imply (54). Lemma 5 is proved.  $\square$

Now, choose the data-driven filter  $\lambda^* = (\lambda_{kl}^*, k, l = 0, \pm 1, \pm 2, \dots)$  in the finite family of filters  $\Lambda = \Lambda_\delta$  using the unbiased risk estimator as in (10):

$$\lambda^* = \arg \min_{\lambda \in \Lambda_\delta} \left\{ \sum_{k,l} (\lambda_{kl}^2 - 2\lambda_{kl}) z_{kl}^2 + 2\varepsilon^2 \sum_{k,l} \sigma_{kl}^2 \lambda_{kl} \right\}.$$

**THEOREM 3.** *The estimator  $\theta^* = (\theta_{kl}^*, k, l = 0, \pm 1, \pm 2, \dots)$ , where  $\theta_{kl}^* = \lambda_{kl}^* z_{kl}$ , is minimax sharp adaptive on the scale of Sobolev classes; that is, for every  $(m, p)$  such that  $m_i \in \mathbb{N}$ ,  $p_i > 0$  we have*

$$(55) \quad \sup_{\theta \in \Theta_2^m(p)} \mathbf{E}_\theta \|\theta^* - \theta\|^2 \leq (1 + o(1)) \inf_{\tilde{\theta}} \sup_{\theta \in \Theta_2^m(p)} \mathbf{E}_\theta \|\tilde{\theta} - \theta\|^2, \quad \varepsilon \rightarrow 0,$$

where  $\inf_{\tilde{\theta}}$  is the infimum over all estimators.

**PROOF.** We apply Theorem 1. First, we check Assumptions 1 and 2. Assumption 1 is obvious; Assumption 2 follows from the inequalities [cf. (59) of the Appendix]

$$(56) \quad \sum_{k,l} \sigma_{kl}^4 \lambda_{kl}^4 \geq c_* \sigma_{W_1 W_2}^4 \sum_{|k| \geq W_1/2, |l| \geq W_2/2} \lambda_{kl}^4 \geq C W_1 W_2 \sigma_{W_1 W_2}^4$$

and

$$\sum_{k,l} \sigma_{kl}^4 \lambda_{kl}^2 \leq \sum_{|k| \leq W_1, |l| \leq W_2} \sigma_{kl}^4 \leq C W_1 W_2 \sigma_{W_1 W_2}^4,$$

where the constants  $C > 0$  do not depend on  $(W, \beta)$ . Again using (56), we find

$$\rho(\lambda) \leq \sup_{|k| \leq W_1, |l| \leq W_2} \sigma_{kl}^2 \left( \sum_{k,l} \sigma_{kl}^4 \lambda_{kl}^4 \right)^{-1/2} \leq C (W_1 W_2)^{-1/2} \leq C \log^{-1}(1/\varepsilon),$$

where the constants  $C > 0$  do not depend on  $(W, \beta)$ . Furthermore, using (59) of the Appendix, we obtain

$$\begin{aligned} \omega(x) &\leq \max_{\lambda \in \Lambda_0} \max_{k,l} \sigma_{kl}^2 \lambda_{kl}^2 I \left\{ \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2 \leq x \max_{k,l} \sigma_{kl}^2 \lambda_{kl}^2 \right\} \\ &\leq \max_{W_1, W_2} \max_{|k| \leq W_1, |l| \leq W_2} \sigma_{kl}^2 I \left\{ C W_1 W_2 \sigma_{W_1 W_2}^2 \leq x \max_{|k| \leq W_1, |l| \leq W_2} \sigma_{kl}^2 \right\} \\ &\leq C \max_{W_1, W_2} [W_1^2 + W_2^2]^{3/2} I \{W_1 W_2 \leq Cx\} \leq Cx^3 \log^{-3}(1/\varepsilon), \end{aligned}$$

where we used that  $W_i \geq \log(1/\varepsilon)$ . Next,  $N = \text{Card}(\Lambda_\delta)$  is of the order  $\log^6(1/\varepsilon)$  and the parameter  $S$  has the order  $\varepsilon^{-3}$ . Thus by Theorem 1 and Lemma 5 one concludes that for  $\varepsilon$  small enough the estimator  $\theta^*$  satisfies the oracle inequality

$$\begin{aligned} \mathbf{E}_\theta \|\theta^* - \theta\|^2 &\leq (1 + CB^{-1}) \min_{\lambda \in \Lambda_\delta} \mathbf{E}_\theta \|\widehat{\theta}(\lambda) - \theta\|^2 + C\varepsilon^2 B^7 \log(1/\varepsilon) \\ (57) \quad &\leq (1 + CB^{-1})(1 + C\delta) \min_{\lambda \in \Lambda_0} \mathbf{E}_\theta \|\widehat{\theta}(\lambda) - \theta\|^2 + C\varepsilon^2 B^7 \log(1/\varepsilon) \\ &= (1 + CB^{-1})(1 + C\delta) \mathbf{E}_\theta \|\widehat{\theta}^M - \theta\|^2 + C\varepsilon^2 B^7 \log(1/\varepsilon), \end{aligned}$$

where  $\widehat{\theta}^M$  is the asymptotically minimax estimator defined in (48). For the last equality we used that this estimator belongs to  $\Lambda_0$  for  $\varepsilon$  small enough [in view of (48) and the fact that  $\mu = O(\varepsilon^\gamma)$  by Lemma 4]. Choose  $B = \varepsilon^{-r/4}$ , where  $r = 1 - \gamma > 0$  and  $\gamma$  is defined in Lemma 4. Taking the supremum of (57) w.r.t.  $\theta \in \Theta_2^m(p)$ , using Proposition 4 and observing that by Lemma 4 the minimax risk over  $\Theta_2^m(p)$  has the order  $\varepsilon^{2-2r}$ , we get (55).  $\square$

Note that Theorem 3 remains valid if  $m$  in the definition of the ellipsoid  $\Theta_2^m(p)$  is real-valued. Furthermore, it is easy to see from the proof that the  $o(1)$  in (55) converges to 0 uniformly over a large set of values  $m, p$ .

## APPENDIX

LEMMA 6. *Let Assumption 3 hold. Let  $(W_1, W_2, \beta_1, \beta_2)$  belong to the set  $\mathcal{W} = [\log(1/\varepsilon), \infty) \times [\log(1/\varepsilon), \infty) \times [1, \sqrt{\log(1/\varepsilon)}] \times [1, \sqrt{\log(1/\varepsilon)}]$ . The partial derivatives of  $\log V$  satisfy uniformly over  $\mathcal{W}$  the inequalities*

$$(58) \quad \left| \frac{\partial}{\partial W_i} \log V \right| \leq C W_i^{-1}, \quad \left| \frac{\partial}{\partial \beta_i} \log V \right| \leq C \beta_i^{-2}, \quad i = 1, 2.$$

PROOF. Using Assumption 3 and acting as in (50) we get, for  $W_i \geq \log(1/\varepsilon)$ ,

$$\begin{aligned}
& V(W_1, W_2, \beta_1, \beta_2) \\
&= \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2(W_1, W_2, \beta_1, \beta_2) \\
&\geq c_* \sigma_{W_1 W_2}^2 \sum_{|k| \geq W_1/2, |l| \geq W_2/2} \lambda_{kl}^2(W_1, W_2, \beta_1, \beta_2) \\
(59) \quad &\geq c_* \sigma_{W_1 W_2}^2 \int_{|x| \geq W_1/2+1} \int_{|y| \geq W_2/2+1} \left[1 - \sqrt{(x/W_1)^{2\beta_1} + (y/W_2)^{2\beta_2}}\right]_+^2 dx dy \\
&\geq c_* W_1 W_2 \sigma_{W_1 W_2}^2 \int_{|x| \geq 1/2+1/W_1} \int_{|y| \geq 1/2+1/W_2} \left[1 - \sqrt{x^2 + y^2}\right]_+^2 dx dy \\
&\geq C W_1 W_2 \sigma_{W_1 W_2}^2,
\end{aligned}$$

where the constant  $C > 0$  does not depend on  $(W, \beta)$ . For brevity write  $A(x, y) = (x^{2\beta_1} + y^{2\beta_2})^{-1/2}$ . Since  $A(x, y) \leq x^{-\beta_1}$ ,  $\forall x > 0, y > 0$ , one obtains

$$\begin{aligned}
(60) \quad & \left| \frac{\partial}{\partial W_1} V \right| = \left| \frac{\partial}{\partial W_1} \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2(W_1, W_2, \beta_1, \beta_2) \right| \\
&= 4\beta_1 W_1^{-1} \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}(W_1, W_2, \beta_1, \beta_2) A\left(\frac{k}{W_1}, \frac{l}{W_2}\right) \left(\frac{k}{W_1}\right)^{2\beta_1-1} \\
&\leq 4\beta_1 W_1^{-1} \sigma_{W_1 W_2}^2 \sum_{|k| \leq W_1, |l| \leq W_2} \left(\frac{k}{W_1}\right)^{\beta_1-1} \\
&\leq C\beta_1 W_1^{-1} \sigma_{W_1 W_2}^2 W_2 \left[ W_1 \int_0^1 x^{\beta_1-1} dx + 1 \right] \leq C W_2 \sigma_{W_1 W_2}^2,
\end{aligned}$$

where the constants  $C > 0$  do not depend on  $(W, \beta)$ , and we used that  $W_1 \geq \log(1/\varepsilon)$ ,  $\beta_1 \leq \sqrt{\log(1/\varepsilon)}$ . Combining (59) and (60) we get the first inequality in (58). The second inequality in (58) is checked similarly. In fact, it suffices to use (59) and the relation

$$\begin{aligned}
& \left| \frac{\partial}{\partial \beta_1} V \right| = \left| \frac{\partial}{\partial \beta_1} \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}^2(W_1, W_2, \beta_1, \beta_2) \right| \\
&= 4 \sum_{k,l} \sigma_{kl}^2 \lambda_{kl}(W_1, W_2, \beta_1, \beta_2) A\left(\frac{k}{W_1}, \frac{l}{W_2}\right) \left(\frac{k}{W_1}\right)^{2\beta_1} \log \frac{k}{W_1} \\
&\leq 4 \sigma_{W_1 W_2}^2 \sum_{k,l} \left[1 - \sqrt{\left(\frac{k}{W_1}\right)^{2\beta_1} + \left(\frac{l}{W_2}\right)^{2\beta_2}}\right]_+ \left(\frac{k}{W_1}\right)^{\beta_1} \left| \log \frac{k}{W_1} \right|
\end{aligned}$$

$$\begin{aligned} &\leq CW_2\sigma_{W_1W_2}^2 \sum_{|k|\leq W_1} \left(\frac{k}{W_1}\right)^{\beta_1} \left|\log \frac{k}{W_1}\right| \\ &\leq CW_2\sigma_{W_1W_2}^2 \left[ W_1 \int_0^1 x^{\beta_1} |\log(x)| dx + 1 \right] \leq CW_1W_2\sigma_{W_1W_2}^2\beta_1^{-2}, \end{aligned}$$

where the constants  $C > 0$  do not depend on  $(W, \beta)$ , and we used that  $W_1 \geq \log(1/\varepsilon)$ ,  $\beta_1 \leq \sqrt{\log(1/\varepsilon)}$ .  $\square$

## REFERENCES

- AKAIKE, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Proc. 2nd Internat. Symp. Inform. Theory, Budapest* (B. N. Petrov and F. Csaki, eds.) 267–281. Akademiai Kiado, Budapest.
- BIRGÉ, L. (2001). An alternative point of view on Lepski's method. In *State of the Art in Probability and Statistics. Festschrift for W. R. van Zwet* (M. de Gunst, C. Klaassen and A. van der Vaart, eds.) 113–133. IMS, Hayward, CA.
- BIRGÉ, L. and MASSART, P. (2001). Gaussian model selection. *J. Eur. Math. Soc. (JEMS)* **3** 203–268.
- CAVALIER, L. and TSYBAKOV, A. B. (2002). Sharp adaptation for inverse problems with random noise. *Probab. Theory Related Fields*. To appear. Available at [www.proba.jussieu.fr](http://www.proba.jussieu.fr).
- DONOHO, D. L. (1995). Nonlinear solution of linear inverse problems by wavelet–vaguelette decomposition. *Appl. Comput. Harmon. Anal.* **2** 101–126.
- DONOHO, D. L. and JOHNSTONE, I. M. (1994). Ideal spatial adaptation via wavelet shrinkage. *Biometrika* **81** 425–455.
- DONOHO, D. L. and JOHNSTONE, I. M. (1995). Adapting to unknown smoothness via wavelet shrinkage. *J. Amer. Statist. Assoc.* **90** 1200–1224.
- DONOHO, D. L. and JOHNSTONE, I. M. (1996). Neoclassical minimax problems, thresholding and adaptive function estimation. *Bernoulli* **2** 39–62.
- GOLDENSHLUGER, A. and PEREVERZEV, S. V. (2000). Adaptive estimation of linear functionals in Hilbert scales from indirect white noise observations. *Probab. Theory Related Fields* **118** 169–186.
- GOLDENSHLUGER, A. and TSYBAKOV, A. B. (2001). Adaptive prediction and estimation in linear regression with infinitely many parameters. *Ann. Statist.* **29** 1601–1619.
- GOLUBEV, G. K. (1987). Adaptive asymptotically minimax estimates of smooth signals. *Problems Inform. Transmission* **23** 57–67.
- GOLUBEV, G. K. (1992). Nonparametric estimation of smooth probability densities in  $L_2$ . *Problems Inform. Transmission* **28** 44–54.
- GOLUBEV, G. K. and KHASHMINSKII, R. Z. (1999). A statistical approach to some inverse boundary problems for partial differential equations. *Problems Inform. Transmission* **35** 136–149.
- GOLUBEV, G. K. and KHASHMINSKII, R. Z. (2001). A statistical approach to the Cauchy problem for the Laplace equation. In *State of the Art in Probability and Statistics. Festschrift for W. R. van Zwet* (M. de Gunst, C. Klaassen and A. van der Vaart, eds.) 419–433. IMS, Hayward, CA.
- GOLUBEV, G. K. and NUSSBAUM, M. (1992). Adaptive spline estimates in a nonparametric regression model. *Theory Probab. Appl.* **37** 521–529.
- HÄRDLE, W. and MARRON, J. S. (1985). Optimal bandwidth selection in nonparametric regression function estimation. *Ann. Statist.* **13** 1465–1481.

- JOHNSTONE, I. M. (1999). Wavelet shrinkage for correlated data and inverse problems: adaptivity results. *Statist. Sinica* **9** 51–83.
- JOHNSTONE, I. M. and SILVERMAN, B. W. (1990). Speed of estimation in positron emission tomography and related inverse problems. *Ann. Statist.* **18** 251–280.
- KERKYACHARIAN, G. and PICARD, D. (2002). Minimax or maxisets? *Bernoulli* **8** 219–253.
- KNEIP, A. (1994). Ordered linear smoothers. *Ann. Statist.* **22** 835–866.
- KOO, J.-Y. (1993). Optimal rates of convergence for nonparametric statistical inverse problems. *Ann. Statist.* **21** 590–599.
- KOROSTELEV, A. P. and TSYBAKOV, A. B. (1993). *Minimax Theory of Image Reconstruction. Lecture Notes in Statist.* **82**. Springer, New York.
- LI, K.-C. (1986). Asymptotic optimality of  $C_L$  and generalized cross-validation in ridge regression with application to spline smoothing. *Ann. Statist.* **14** 1101–1112.
- LI, K.-C. (1987). Asymptotic optimality of  $C_p$ ,  $C_L$ , cross-validation and generalized cross-validation: discrete index set. *Ann. Statist.* **15** 958–976.
- MAIR, B. and RUYMGAART, F. H. (1996). Statistical inverse estimation in Hilbert scales. *SIAM J. Appl. Math.* **56** 1424–1444.
- MALLOWS, C. L. (1973). Some comments on  $C_p$ . *Technometrics* **15** 661–675.
- NEMIROVSKI, A. (2000). Topics in nonparametric statistics. *Ecole d'Été de Probabilités de St. Flour XXVIII. Lecture Notes in Math.* **1738** 85–277. Springer, New York.
- PINSKER, M. S. (1980). Optimal filtering of square integrable signals in Gaussian white noise. *Problems Inform. Transmission* **16** 120–133.
- POLYAK, B. T. and TSYBAKOV, A. B. (1990). Asymptotic optimality of the  $C_p$ -test for the orthogonal series estimation of regression. *Theory Probab. Appl.* **35** 293–306.
- POLYAK, B. T. and TSYBAKOV, A. B. (1992). A family of asymptotic optimal methods for choosing the estimate order in orthogonal series regression. *Theory Probab. Appl.* **37** 471–481.
- STEIN, E. and WEISS, G. (1971). *Introduction to Fourier Analysis on Euclidean Spaces*. Princeton Univ. Press.
- WAHBA, G. (1977). Practical approximate solutions to linear operator equations when the data are noisy. *SIAM J. Numer. Anal.* **14** 651–667.
- WAHBA, G. (1990). *Spline Models for Observational Data*. SIAM, Philadelphia.

L. CAVALIER  
G. K. GOLUBEV  
CMI  
UNIVERSITÉ AIX-MARSEILLE I  
39 RUE F. JOLIOT-CURIE  
F-13453 MARSEILLE CEDEX 13  
FRANCE

D. PICARD  
A. B. TSYBAKOV  
LABORATOIRE DE PROBABILITÉS  
ET MODÈLES ALÉATOIRES  
UNIVERSITÉS PARIS 7–PARIS 6  
4 PL. JUSSIEU, BP 188  
F-75252 PARIS CEDEX 05  
FRANCE  
E-MAIL: tsybakov@ccr.jussieu.fr