

ADAPTIVE PREDICTION AND ESTIMATION IN LINEAR REGRESSION WITH INFINITELY MANY PARAMETERS¹

BY A. GOLDENSHLUGER AND A. TSYBAKOV

University of Haifa and Université Paris VI

The problem of adaptive prediction and estimation in the stochastic linear regression model with infinitely many parameters is considered. We suggest a prediction method that is sharp asymptotically minimax adaptive over ellipsoids in ℓ_2 . The method consists in an application of blockwise Stein's rule with "weakly" geometrically increasing blocks to the penalized least squares fits of the first N coefficients. To prove the results we develop oracle inequalities for a sequence model with correlated data.

1. Introduction. Consider the regression model

$$(1) \quad y = \sum_{k=1}^{\infty} \theta_k x_k + \varepsilon,$$

where $\{x_k\}_{k=1,2,\dots}$ is a sequence of explanatory variables, y is the corresponding response, ε is the error, and $\theta = (\theta_1, \theta_2, \dots) \in \ell_2$ is an unknown regression sequence. Assume that $\{x_k\}$ and ε are random variables, and $\mathbb{E}\varepsilon = 0$ and $\mathbb{E}\varepsilon^2 = \sigma^2$. The stochastic series in (1) and later are assumed to converge in the mean squared sense. Suppose we are given n independent realizations of y and $\{x_k\}$,

$$\mathcal{U}_n = \{y(t); x_1(t), x_2(t), \dots; t = 1, \dots, n\}$$

coming from the model (1), that is,

$$y(t) = \sum_{k=1}^{\infty} \theta_k x_k(t) + \varepsilon(t), \quad t = 1, \dots, n.$$

Given $\mathcal{X}_{n+1} = \{x_1(n+1), x_2(n+1), \dots\}$, the objective is to predict the corresponding response $y(n+1)$. A *predictor* (or *prediction method*) is a random variable $\hat{y} = \hat{y}(n+1)$ measurable with respect to $(\mathcal{U}_n, \mathcal{X}_{n+1})$.

The problem of prediction in the model (1) has been considered by Shibata (1981), Breiman and Freedman (1983), Goldenshluger and Tsybakov (1999). In particular, Shibata (1981) and Breiman and Freedman (1983) study the least squares predictor of the form

$$(2) \quad \hat{y}(n+1) = \sum_{k=1}^d \hat{\theta}_k^{\text{OLS}} x_k(n+1),$$

Received June 2000; revised June 2001.

¹Supported in part by an Arc-en-ciel/Keshet grant.

AMS 2000 subject classifications. 62G05, 62G20.

Key words and phrases. Linear regression with infinitely many parameters, adaptive prediction, exact asymptotics of minimax risk, blockwise Stein's rule, oracle inequalities.

where $(\hat{\theta}_1^{\text{OLS}}, \dots, \hat{\theta}_d^{\text{OLS}})$ is the ordinary least squares (OLS) estimator of $(\theta_1, \dots, \theta_d)$ based on the reduced data $\{y(t), x_1(t), \dots, x_d(t), t = 1, \dots, n\}$. They discuss data-driven choices of d . Goldenshluger and Tsybakov (1999) suggest the predictor

$$(3) \quad \hat{y}(n+1) = \sum_{k=1}^d \lambda_k \hat{\theta}_k^{\text{P}} x_k(n+1),$$

where $(\hat{\theta}_1^{\text{P}}, \dots, \hat{\theta}_d^{\text{P}})$ is a penalized least squares estimator of $(\theta_1, \dots, \theta_d)$, and $\{\lambda_k\}$ are some weights. They show that the predictor (3) is asymptotically sharp minimax on the classes of ellipsoids in the space of coefficients θ , provided $\{\lambda_k\}$ are chosen in a proper way. In particular, the predictor (3) outperforms the OLS predictor (2) in the minimax sense. The weights $\{\lambda_k\}$ depend on the parameters of the ellipsoid, and the method (3) is not adaptive to these parameters.

In this paper we suggest an adaptive prediction method which is asymptotically as good as (3) (i.e., is asymptotically sharp minimax) on any ellipsoid within a wide scale. The method does not depend on the parameters of an ellipsoid and is not related to prior assumptions on an ellipsoidal structure. The idea is to apply the blockwise Stein rule to penalized least squares fits $\hat{\theta}_k^{\text{P}}$ of the first N coefficients θ_k . The blockwise Stein rule has been used recently to get adaptive estimators in different statistical problems [Donoho and Johnstone (1995), Johnstone (1998), Cai (1999), Cavalier and Tsybakov (2000)]. In these papers asymptotically minimax adaptivity is proved by means of oracle inequalities in sequence space. Our reasoning is similar, but the difficulty of the present setting is that the sequence space representation is non-Gaussian, correlated and biased. We get oracle inequalities that work in this situation. We also prove that, under general conditions, blockwise linear estimators are almost as good as linear monotone oracles. These results are of independent interest, and can be used in other contexts as well. The sharp minimax adaptivity is proved as a consequence of these results for a special construction of “weakly” geometrically increasing blocks. This differs from the polynomially increasing blocks as in Efromovich and Pinsker (1984, 1996), Efromovich (1999), or dyadic blocks as in the wavelet context [Donoho and Johnstone (1995), Johnstone (1998, 1999)], but is closely related to Nemirovski (2000) and Cavalier and Tsybakov (2000), where other statistical models have been studied.

The paper is organized as follows. In Section 2 we define our adaptive prediction method, analyze its properties and state our main results. Section 3 considers an equivalent sequence space model and includes basic oracle inequalities underlying our proofs. In Section 4 we prove the main results. The Appendix contains auxiliary lemmas related to properties of the Stein estimator for correlated data.

2. Adaptive prediction method. We define the predictor $\hat{y}^*(n+1)$ of $y(n+1)$ as follows. Fix a positive integer $N = N_n$, and denote $\mathcal{I} = \{1, \dots, N_n\}$.

Let $\phi_N(t) = (x_1(t), \dots, x_N(t))'$, $t = 1, \dots, n$, and consider a penalized least squares estimator

$$(4) \quad \tilde{y}_{\mathcal{J}} = \hat{\theta}_{\mathcal{J}}^P = Q_{\mathcal{J}}^{-1} \left(\frac{1}{n} \sum_{t=1}^n \phi_N(t) y(t) \right),$$

where

$$Q_{\mathcal{J}} = \frac{1}{n} \sum_{t=1}^n \phi_N(t) \phi_N'(t) + n^{-1} I$$

and I is the $N_n \times N_n$ identity matrix. Here and later v_B denotes the vector $\{v_k, k \in B\}$, where B is a set of integers.

For a monotone increasing sequence of integers $\{\kappa_j\}$ such that $\kappa_1 = 1$, we define the partition of the set $\{1, \dots, N\}$ into blocks B_j as follows

$$(5) \quad B_j = \{\kappa_j, \kappa_j + 1, \dots, \kappa_{j+1} - 1\}, \quad j = 1, \dots, J,$$

$$(6) \quad \kappa_{J+1} - 1 = N.$$

[We assume w.l.o.g. that the sequence $\{\kappa_j\}$ and N are such that (6) holds.] Denote $n_j = \kappa_j - \kappa_{j-1}$, $\theta_{(j)} = \theta_{B_j}$, and $\tilde{y}_{(j)} = \tilde{y}_{B_j}$, $j = 1, \dots, J$.

Let

$$(7) \quad \hat{\theta}^* = (\hat{\theta}_{(1)}^*, \hat{\theta}_{(2)}^*, \dots, \hat{\theta}_{(J)}^*, 0, 0, \dots),$$

where $\hat{\theta}_{(j)}^*$ is the Stein estimator for the block B_j ,

$$(8) \quad \hat{\theta}_{(j)}^* = \left(1 - \frac{\sigma^2 n_j}{n \|\tilde{y}_{(j)}\|^2} \right) \tilde{y}_{(j)}, \quad j = 1, \dots, J.$$

Here and later $\|\cdot\|$ is the Euclidean norm when applied to a finite dimensional vector.

Define the prediction method

$$(9) \quad \hat{y}^*(n+1) = \sum_{k=1}^{N_n} \hat{\theta}_k^* x_k(n+1).$$

We show that this method is asymptotically sharp adaptive in a minimax sense on the ellipsoids in the space of sequences θ . Consider the ellipsoids

$$\Theta(a, L) = \left\{ \theta \in \ell_2 : \sum_{k=1}^{\infty} a_k^2 \theta_k^2 \leq L^2 \right\},$$

where $L > 0$ and $a = \{a_k\}$ is a monotone nondecreasing positive sequence such that $a_k \rightarrow \infty$ as $k \rightarrow \infty$.

The prediction error of any predictor \hat{y} is defined as $E[\hat{y}(n+1) - y(n+1)]^2$. Note that this error cannot be arbitrarily small; it is at least σ^2 for large n , because of the nonvanishing innovation component $\varepsilon(n+1)$ independent of

$(\mathcal{X}_n, \mathcal{X}_{n+1})$. We therefore consider the difference $\mathbb{E}[\hat{y}(n+1) - y(n+1)]^2 - \sigma^2$, and define the maximal risk over $\Theta(a, L)$ in the form

$$\mathcal{R}[\hat{y}; \Theta(a, L)] = \sup_{\theta \in \Theta(a, L)} \mathbb{E}[\hat{y}(n+1) - y(n+1)]^2 - \sigma^2.$$

We show that $\hat{y}^* = \hat{y}^*(n+1)$ is an asymptotically minimax predictor; that is, it minimizes the maximal prediction error,

$$\mathcal{R}[\hat{y}^*; \Theta(a, L)] = \inf_{\hat{y}} \mathcal{R}[\hat{y}; \Theta(a, L)](1 + o(1)), \quad n \rightarrow \infty,$$

where inf is taken over all possible prediction methods based on the observations $(\mathcal{X}_n, \mathcal{X}_{n+1})$.

The following assumptions will be used.

ASSUMPTION 1. The random variables $\{x_k\}_{k=1,2,\dots}$ are independent, $\mathbb{E}x_k = 0$, $\mathbb{E}x_k^2 = 1$, and there exist constants $H > 0$ and $c_* > 0$ such that

$$(10) \quad \mathbb{E} \exp\{\lambda x_k^2\} \leq c_* < \infty, \quad |\lambda| < H, \quad k = 1, 2, \dots$$

The assumption that $\{x_k\}$ are uncorrelated zero mean random variables with variance 1 is quite natural in the prediction context, since typically $\{x_k\}$ are considered as “principal components” of some original random covariates [see the discussion in Breiman and Freedman (1983) and Goldenshluger and Tsybakov (1999)]. In particular, Breiman and Freedman (1983) work with i.i.d. standard Gaussian $\{x_k\}$.

Note also that under Assumption 1 there is a simple and natural relationship between prediction risk and the ℓ_2 -risk in estimation of the regression coefficients. In this case the prediction risk $\mathbb{E}[\hat{y}(n+1) - y(n+1)]^2 - \sigma^2$ of any predictor of the type $\hat{y}(n+1) = \sum_{k=1}^\infty \hat{\theta}_k x_k(n+1)$ with $\hat{\theta}_k$'s based on \mathcal{X}_n only, coincides with $\mathbb{E}\|\hat{\theta} - \theta\|^2$. Thus, the prediction problem is equivalent to estimating coefficients of the corresponding regression model with explanatory variables satisfying Assumption 1.

ASSUMPTION 2. The random variable ε is Gaussian $\mathcal{N}(0, \sigma^2)$, and ε is independent of $\{x_k\}_{k=1,2,\dots}$.

ASSUMPTION 3. The partition (5) and (6) satisfies

$$\kappa_1 = 1, \kappa_2 = \nu_n + 1, \kappa_j = \kappa_{j-1} + \lfloor (1 + \rho_n)^{j-1+\nu_n} \rfloor, \quad j = 3, \dots, J,$$

where $\rho_n = (\ln \ln n)^{-1}$, $\nu_n = \lfloor \rho_n^{-1} \ln \rho_n^{-1} \rfloor$, and J is such that $N_n \asymp \sqrt{n}(\ln n)^{-1}$.

Let c_n denote the solution to the equation [cf. Pinsker (1980)]

$$\sigma^2 n^{-1} \sum_{k=1}^\infty a_k (1 - c_n a_k)_+ = c_n L^2$$

(note that the solution is unique since $a_k \rightarrow \infty$), and let

$$r_n = r_n(\Theta(a, L)) = \sigma^2 n^{-1} \sum_{k=1}^{\infty} (1 - c_n a_k)_+.$$

For the i.i.d. sequence space model, the value r_n is exactly the minimax linear risk on the ellipsoid $\Theta(a, L)$, and asymptotically r_n equals to the minimax risk among all estimators on $\Theta(a, L)$ [Pinsker (1980); see also Belitser and Levit (1995)]. As shown in Goldenshluger and Tsybakov (1999), the value r_n gives also a lower bound for the minimax risk in our problem, and this bound cannot be improved among all prediction methods. The next theorem is a corollary of the lower bound in Goldenshluger and Tsybakov (1999).

THEOREM 1. *Let Assumptions 1 and 2 hold. Assume that the ellipsoid $\Theta(a, L)$ is such that the sequence $\{a_k\}$ is monotone nondecreasing and there exist $\beta > 1/2$, and the positive constants a_{\min}, a_{\max} such that*

$$a_{\min} k^\beta \leq a_k \leq a_{\max} k^\beta, \quad k = 1, 2, \dots$$

Then for every prediction method $\hat{y} = \hat{y}(n + 1)$ one has

$$(11) \quad \mathcal{R}[\hat{y}; \Theta(a, L)] \geq r_n(1 + o(1)), \quad n \rightarrow \infty.$$

Now we state the main results of the paper. First, we claim that the lower bound (11) is attained by the predictor \hat{y}^* ; that is, this predictor is asymptotically minimax sharp adaptive on the scale of ellipsoids satisfying the assumptions of Theorem 1.

THEOREM 2. *Let Assumptions 1–3 hold. Let $\{a_k\}$ satisfy the assumptions of Theorem 1. Then the prediction method defined in (9) satisfies*

$$(12) \quad \mathcal{R}[\hat{y}^*; \Theta(a, L)] \leq r_n(1 + o(1)), \quad n \rightarrow \infty.$$

Next, the result of Theorem 2 can be extended to a larger class of ellipsoids, and the $o(1)$ term in (12) is uniformly small over this class as stated in the following theorem.

THEOREM 3. *Let Assumptions 1–3 hold. Let $A(\beta_0, \beta_1)$ be the set of monotone nondecreasing sequences $a = \{a_k\}$ such that*

$$a_{\min} k^{\beta_0} \leq a_k \leq a_{\max} k^{\beta_1}, \quad k = 1, 2, \dots,$$

where $2\beta_1/(2\beta_1 + 1) < \beta_0 \leq \beta_1 < \infty$. For given numbers $0 < L_{\min} \leq L_{\max} < \infty$, let \mathcal{A} denote the collection of pairs (a, L) such that $a \in A(\beta_0, \beta_1)$ and $L \in [L_{\min}, L_{\max}]$. Then the prediction method defined in (9) satisfies

$$\sup_{(a, L) \in \mathcal{A}} \{\mathcal{R}[\hat{y}^*; \Theta(a, L)]/r_n\} = 1 + o(1), \quad n \rightarrow \infty.$$

Proofs of Theorems 2 and 3 are given in Section 4.

REMARK. Theorems 2 and 3 can be extended in a standard way to the case where the variance σ^2 is not known. It suffices to replace σ^2 in the definition of $\hat{\theta}^*$ by a consistent estimator $\hat{\sigma}^2$. This can be a standard estimator of variance based on the sum of squares of the OLS residuals.

3. The main tools.

3.1. *Equivalent sequence model.* It follows from (4) that

$$(13) \quad \tilde{y}_{\mathcal{J}} = \theta_{\mathcal{J}} + \delta_{\mathcal{J}} + \frac{\sigma}{\sqrt{n}} \xi_{\mathcal{J}}, \quad \mathcal{J} = \{1, \dots, N_n\},$$

where

$$(14) \quad \begin{aligned} \delta_{\mathcal{J}} &= \mathbf{Q}_{\mathcal{J}}^{-1} \left(-\frac{\theta_{\mathcal{J}}}{n} + \frac{1}{n} \sum_{t=1}^n \phi_N(t) \sum_{k=N_n+1}^{\infty} \theta_k x_k(t) \right), \\ \xi_{\mathcal{J}} &= \mathbf{Q}_{\mathcal{J}}^{-1} \left(\frac{1}{\sqrt{n}} \sum_{t=1}^n \phi_N(t) \tilde{\varepsilon}(t) \right), \end{aligned}$$

and $\tilde{\varepsilon} = \sigma^{-1} \varepsilon \sim \mathcal{N}(0, 1)$.

Considering $\tilde{y}_{\mathcal{J}}$ as “new observations,” we note that the model (13) is non-gaussian with correlated nonzero mean errors. We note, however, that conditionally on the σ -algebra $\mathcal{F}_x^n = \sigma(\{x_k(t)\}_{k=1,2,\dots,t=1,\dots,n})$ the random vector $\tilde{y}_{\mathcal{J}}$ is Gaussian. Below we state a lemma showing that on a set of “large” probability (13) can be regarded as a Gaussian model with small correlations.

Note that

$$(15) \quad \mathbb{E}(\xi_{\mathcal{J}} | \mathcal{F}_x^n) = 0, \quad S_{\mathcal{J}} = \mathbb{E}(\xi_{\mathcal{J}} \xi'_{\mathcal{J}} | \mathcal{F}_x^n) = \mathbf{Q}_{\mathcal{J}}^{-1} - n^{-1} \mathbf{Q}_{\mathcal{J}}^{-2}.$$

Let $v_{ij} = n^{-1} \sum_{t=1}^n x_i(t)x_j(t) - \delta_{ij}$, where $i, j = 1, \dots, N_n$, and δ_{ij} is the Kronecker delta.

LEMMA 1. *Let Assumptions 1 and 2 hold. Let $\alpha \in (0, 1)$. Then there exists a constant $q = q(c_*, H) > 0$ such that*

$$(16) \quad \mathbb{P}(\Omega_{\alpha}) \geq 1 - \alpha,$$

where

$$(17) \quad \Omega_{\alpha} = \left\{ \omega \in \Omega : \max_{i, j=1, \dots, N_n} |v_{ij}| \leq \sqrt{\frac{q}{n} \ln \frac{2N_n^2}{\alpha}} \right\}.$$

Furthermore, let

$$(18) \quad N_n \mu_n(\alpha) < \frac{1}{2},$$

where

$$\mu_n(\alpha) = \sqrt{\frac{q}{n} \ln \frac{2N_n^2}{\alpha}} + \frac{1}{n}.$$

Then on the event Ω_α we have $S_{\mathcal{J}} = I - A_{\mathcal{J}}$, where $A_{\mathcal{J}}$ is an $N_n \times N_n$ matrix satisfying

$$(19) \quad \max_{i, j=1, \dots, N_n} |[A_{\mathcal{J}}]_{ij}| \leq 3\mu_n(\alpha),$$

and $[\cdot]_{ij}$ stands for the (i, j) th entry of a matrix.

Proof of the lemma is given in Section 4.

In view of Lemma 1 conditionally on \mathcal{F}_x^n on the event Ω_α the model (13) can be regarded as a Gaussian model with small correlations. Our proof of Theorems 2 and 3 is based on combining several oracle inequalities. The first inequality shows that the risk of the blockwise Stein estimator is almost as small as the risk of the blockwise linear oracle in the Gaussian model with small correlations. The second inequality uses the first one and guarantees that the effect of the bias $\delta_{\mathcal{J}}$ is asymptotically negligible. Finally, the third inequality allows linking the risk of the blockwise linear oracle to that of the monotone linear one.

3.2. *An oracle inequality for a sequence model with correlated errors.* Denote $\varepsilon = \sigma n^{-1/2}$. Consider the model

$$(20) \quad y_{\mathcal{J}} = \theta_{\mathcal{J}} + \varepsilon \xi_{\mathcal{J}},$$

where $\xi_{\mathcal{J}}$ is the Gaussian vector with zero mean, and the covariance matrix $Q_{\mathcal{J}} = I - A_{\mathcal{J}}$ such that $\max_{i, j} |[A_{\mathcal{J}}]_{ij}| \leq \mu = \mu_\varepsilon < 1$.

Denote $y_{(j)} = y_{B_j}$, $j = 1, \dots, J$, and introduce the Stein estimators,

$$\hat{\theta}_{(j)} = \left(1 - \frac{\varepsilon^2 n_j}{\|y_{(j)}\|^2}\right) y_{(j)}, \quad j = 1, \dots, J.$$

The estimate $\hat{\theta}_{\mathcal{J}}$ of the sequence $\theta_{\mathcal{J}}$ is given by

$$(21) \quad \hat{\theta}_{\mathcal{J}} = (\hat{\theta}_{(1)}, \hat{\theta}_{(2)}, \dots, \hat{\theta}_{(J)}).$$

Consider the ideal blockwise linear risk

$$r_n^{\text{BL}} = \sum_{j=1}^J \frac{\|\theta_{(j)}\|^2 \varepsilon^2 n_j}{\|\theta_{(j)}\|^2 + \varepsilon^2 n_j}.$$

It is well known [see, e.g., Efroimovich and Pinsker (1984), Johnstone (1998)] that in the sequence space model with independent errors, r_n^{BL} represents the minimal risk of linear estimators whose weights are constant on the blocks $\{B_j\}$. This minimum is attained on a “pseudo-estimator” that depends on θ (called a *blockwise linear oracle*).

LEMMA 2. *Let $\{B_j\}$ be the partition (5) and (6). Let $n_j > 4$ and $\mu < 1/6 - 2/(3n_j)$ for all $j = 1, \dots, J$. Then for the estimate (21) in the model (20) we*

have

$$(22) \quad \mathbb{E} \|\hat{\theta}_{\mathcal{J}} - \theta_{\mathcal{J}}\|^2 \leq r_n^{\text{BL}} + 7\varepsilon^2 \mu N + 4\varepsilon^2 J.$$

Proof of the lemma is given in the Appendix.

3.3. *Passage to infinite sequences and the i.i.d. approximation.* Denote by $\bar{\theta}(h)$ a linear estimator of $\theta \in \ell_2$ with nonrandom weights $h = \{h_k\}$, that is,

$$\bar{\theta}(h) = (h_1 y_1, h_2 y_2, \dots).$$

Let \mathbb{E}_* be the expectation with respect to $\{y_k\}$ satisfying

$$(23) \quad y_k = \theta_k + \varepsilon \xi_k, \quad k = 1, 2, \dots, \quad \xi_k \text{ i.i.d. } \mathcal{N}(0, 1).$$

The risk of the estimator $\bar{\theta}(h)$ in this model is

$$\mathcal{R}_\varepsilon(h, \theta) = \mathbb{E}_* \|\bar{\theta}(h) - \theta\|^2 = \sum_k (1 - h_k)^2 \theta_k^2 + \varepsilon^2 \sum_k h_k^2,$$

where $\|\cdot\|$ is the ℓ_2 -norm.

Let \mathcal{H} be a set of sequences $\{h_k\}$ piecewise constant on the blocks B_j :

$$\mathcal{H} = \{\{h_k\}: h_k = h_{\kappa_j}, \forall k \in B_j, j = 1, \dots, J, \text{ and } h_k = 0, \forall k > N\}.$$

It is easy to see that

$$\inf_{h \in \mathcal{H}} \mathcal{R}_\varepsilon(h, \theta) = \sum_{j=1}^J \frac{\|\theta_{(j)}\|^2 \varepsilon^2 n_j}{\|\theta_{(j)}\|^2 + \varepsilon^2 n_j} + \sum_{k=N+1}^\infty \theta_k^2 = r_n^{\text{BL}} + \sum_{k=N+1}^\infty \theta_k^2.$$

LEMMA 3. *Let Assumptions 1 and 2 hold. Assume that the partition (5) and (6) satisfies the following conditions:*

$$(24) \quad n_j > 4, \quad j = 1, \dots, J_n, \quad N_n \asymp \sqrt{n}(\ln n)^{-1}.$$

Then for $\hat{\theta}^*$ defined in (7) and (8) and for every $\theta \in \ell_2$ such that $\|\theta\| \leq L$ we have

$$(25) \quad \mathbb{E} \|\hat{\theta}^* - \theta\|^2 \leq \left(\inf_{h \in \mathcal{H}} \mathcal{R}_{\sigma/\sqrt{n}}(h, \theta) + \frac{4\sigma^2 J_n}{n} \right) (1 + o(1)), \quad n \rightarrow \infty,$$

here $o(1)$ is uniform over $\|\theta\| \leq L$.

The proof is given in Section 4.

Lemma 3 shows that, up to the factor $1 + o(1)$, the behavior of our estimator $\hat{\theta}^*$ is asymptotically at least as good as that of the blockwise Stein estimator for the i.i.d. model (23). In fact, in view of (22), we have for the model (23) that the risk of the blockwise Stein estimator $(\hat{\theta}_{(1)}, \dots, \hat{\theta}_{(J)}, 0, 0, \dots)$ is bounded by $\inf_{h \in \mathcal{H}} \mathcal{R}_\varepsilon(h, \theta) + 4\varepsilon^2 J$.

3.4. *The blockwise linear oracle is almost as good as the monotone linear one.* The oracle inequality (25) shows that our estimator is asymptotically almost as good as the best blockwise linear oracle for the i.i.d. Gaussian sequence space model. Here we will see that a similar result is valid for the linear monotone oracle in place of the blockwise linear oracle. This follows from Lemma 3 and the next lemma.

Denote \mathcal{H}_{mon} the set of monotone nonincreasing sequences $h = \{h_k\} \in \ell_2$,

$$\mathcal{H}_{\text{mon}} = \{h = \{h_k\} \in \ell_2: 1 \geq h_1 \geq h_2 \geq \dots \geq 0\}.$$

LEMMA 4. *For any $\theta \in \ell_2$ and any partition (5) and (6) such that*

$$(26) \quad \frac{\text{card}(B_{j+1})}{\text{card}(B_j)} \leq 1 + \eta,$$

where $\eta = \eta_\varepsilon > 0$ we have

$$(27) \quad \inf_{h \in \mathcal{H}} \mathcal{R}_\varepsilon(h, \theta) \leq (1 + \eta) \inf_{h \in \mathcal{H}_{\text{mon}}} \mathcal{R}_\varepsilon(h, \theta) + \varepsilon^2 \text{card}(B_1) + \sum_{k=N+1}^\infty \theta_k^2.$$

PROOF. It suffices to prove that for any $h \in \mathcal{H}_{\text{mon}}$ there exists $\bar{h} \in \mathcal{H}$ such that

$$(28) \quad \mathcal{R}_\varepsilon(\bar{h}, \theta) \leq (1 + \eta) \mathcal{R}_\varepsilon(h, \theta) + \varepsilon^2 \text{card}(B_1) + \sum_{k=N+1}^\infty \theta_k^2.$$

Given $h \in \mathcal{H}_{\text{mon}}$, define $\bar{h} \in \mathcal{H}$ by

$$\bar{h}_k = \begin{cases} 1, & k \in B_1, \\ h_{\kappa_j}, & k \in B_j, j = 2, \dots, J, \\ 0, & k > N. \end{cases}$$

Since $0 \leq h_k \leq \bar{h}_k \leq 1$ for $k \leq N$ we have

$$\begin{aligned} \mathcal{R}_\varepsilon(\bar{h}, \theta) &= \sum_{k=1}^\infty (1 - \bar{h}_k)^2 \theta_k^2 + \varepsilon^2 \sum_{k=1}^\infty \bar{h}_k^2 \\ &\leq \sum_{k=1}^N (1 - h_k)^2 \theta_k^2 + \varepsilon^2 \sum_{k=1}^N \bar{h}_k^2 + \sum_{k=N+1}^\infty \theta_k^2. \end{aligned}$$

Hence to show (28) it suffices to prove that

$$\varepsilon^2 \sum_{k=1}^N \bar{h}_k^2 \leq (1 + \eta) \varepsilon^2 \sum_{k=1}^N h_k^2 + \varepsilon^2 \text{card}(B_1).$$

Note that $\sum_{k=1}^N \bar{h}_k^2 = \text{card}(B_1) + \sum_{j=2}^J h_{\kappa_j}^2 \text{card}(B_j)$. By (26) and monotonicity of $\{h_k\}$,

$$\begin{aligned} \sum_{j=2}^J h_{\kappa_j}^2 \text{card}(B_j) &\leq (1 + \eta) \sum_{j=1}^{J-1} h_{\kappa_{j+1}}^2 \text{card}(B_j) \\ &\leq (1 + \eta) \sum_{j=1}^{J-1} \sum_{k \in B_j} h_k^2 \leq (1 + \eta) \sum_{k=1}^N h_k^2. \quad \square \end{aligned}$$

REMARK. Inequality (27) can be used to prove minimax adaptivity of block rules for different statistical models after reduction to the i.i.d. Gaussian sequence model. In fact, minimaxity of block rules for a class of sequences Θ follows from Lemma 4 under the condition that the minimax estimator on Θ is linear and has the form $\bar{\theta}(h^*)$ with $h^* \in \mathcal{H}_{\text{mon}}$.

4. Proofs.

PROOF OF LEMMA 1. First we will prove (16). For any $s > 0$ we have

$$\begin{aligned} \mathbb{P} \left\{ \max_{i, j=1, \dots, N_n} |v_{ij}| \geq s \right\} &\leq \mathbb{P} \left\{ \max_{i=1, \dots, N_n} \left| \frac{1}{n} \sum_{t=1}^n x_i^2(t) - 1 \right| \geq s \right\} \\ &\quad + \mathbb{P} \left\{ \max_{i, j=1, \dots, N_n, i \neq j} \left| \frac{1}{n} \sum_{t=1}^n x_i(t)x_j(t) \right| \geq s \right\} \\ &\equiv P_1 + P_2. \end{aligned}$$

It follows from (10) that there exist positive constants H_1, q_1 and q_2 such that for $k, j = 1, \dots, N_n, k \neq j$,

$$\begin{aligned} \mathbb{E} \exp\{\lambda(x_k^2 - 1)\} &\leq \exp\{q_1 \lambda^2\}, \\ \mathbb{E} \exp\{\lambda x_k x_j\} &\leq \exp\{q_2 \lambda^2\} \quad \text{for } |\lambda| < H_1. \end{aligned}$$

Therefore by Theorem 2.6 from Petrov (1995) we have for $s < 2H_1 q_1$ and any $k = 1, \dots, N_n$,

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n (x_k^2(t) - 1) \right| \geq s \right\} \leq 2 \exp \left\{ -\frac{ns^2}{4q_1} \right\} \quad \text{and} \quad P_1 \leq 2N_n \exp \left\{ -\frac{ns^2}{4q_1} \right\}.$$

The term P_2 is bounded similarly. Finally, we conclude that for $s \leq 2H_1 \times \min(q_1, q_2)$

$$\mathbb{P} \left\{ \max_{i, j=1, \dots, N_n} |v_{ij}| \geq s \right\} \leq 2N_n^2 \exp \left\{ -\frac{ns^2}{4q_3} \right\},$$

where $q_3 = q_1 \wedge q_2$. This completes the proof of (16).

Let $C_{\mathcal{J}} = I - Q_{\mathcal{J}}$; then on the set Ω_α ,

$$|[C_{\mathcal{J}}]_{ij}| \leq |v_{ij}| + \delta_{ij} n^{-1} \leq \mu_n(\alpha).$$

By (18), $\|C_{\mathcal{J}}\| \leq N_n \mu_n(\alpha) < 1/2$, where $\|\cdot\|$ denotes the standard spectral matrix norm. Hence on the set Ω_α ,

$$(29) \quad Q_{\mathcal{J}}^{-1} = (I - C_{\mathcal{J}})^{-1} = I + D_{\mathcal{J}}, \quad D_{\mathcal{J}} = \sum_{k=1}^{\infty} C_{\mathcal{J}}^k.$$

We note also that

$$\begin{aligned} \max_{i, j=1, \dots, N} |[C_{\mathcal{J}}^k]_{ij}| &\leq N_n \mu_n(\alpha) \max_{i, j=1, \dots, N_n} |[C_{\mathcal{J}}^{k-1}]_{ij}| \\ &\leq N_n^{k-1} \mu_n^k(\alpha), \quad k = 1, 2, \dots, \end{aligned}$$

and therefore by (29) and (18),

$$(30) \quad \max_{i, j=1, \dots, N_n} |[D_{\mathcal{J}}]_{ij}| \leq \sum_{k=1}^{\infty} N_n^{k-1} \mu_n^k(\alpha) \leq 2\mu_n(\alpha).$$

Now define $A_{\mathcal{J}} = -D_{\mathcal{J}}(1 - 2n^{-1}) + D_{\mathcal{J}}^2 n^{-1} + n^{-1}I$; then by definition $S_{\mathcal{J}} = I - A_{\mathcal{J}}$. It follows from (30) and (18) that

$$\begin{aligned} \max_{i, j=1, \dots, N_n} |[A_{\mathcal{J}}]_{ij}| &\leq 2\mu_n(\alpha) \left(1 - \frac{2}{n}\right) + \frac{4}{n} N_n \mu_n^2(\alpha) + \frac{1}{n} \\ &\leq 2\mu_n(\alpha) - \frac{2\mu_n(\alpha)}{n} + \frac{1}{n} \leq 3\mu_n(\alpha). \end{aligned}$$

This completes the proof. \square

The following auxiliary result will be used in the proof of Lemma 3.

LEMMA 5. *Let Assumption 1 hold. Denote $U_n(\alpha) = 1 + 2N_n \mu_n(\alpha)$. Then for the subvectors $\delta_{(j)} = \delta_{B_j}$, $j = 1, \dots, J$ we have*

(i)

$$(31) \quad \mathbb{E}[\|\delta_{(j)}\|^2 \mathbf{1}(\Omega_\alpha)] \leq 2U_n^2(\alpha) \left[\frac{\|\theta_{(j)}\|^2}{n^2} + \frac{\sqrt{2qn_j}}{n} \sum_{k=N_n+1}^{\infty} \theta_k^2 \right];$$

(ii)

$$(32) \quad \|\mathbb{E}[\delta_{(j)} \mathbf{1}(\Omega_\alpha)]\| \leq U_n(\alpha) \|\theta_{(j)}\| n^{-1}.$$

PROOF. (i) By (29) and (30) on the set Ω_α we have $\|Q_{\mathcal{J}}^{-1}\| \leq U_n(\alpha)$. Then (31) is easily obtained as in the proof of Lemma 2 in Goldenshluger and Tsybakov (1999).

(ii) By (29) and (30) on the set Ω_α we have $Q_{\mathcal{J}}^{-1} = I + D_{\mathcal{J}}$. Therefore, conditioning on $\{x_k(t); t = 1, \dots, n; k = 1, \dots, N\}$ and using independence of

$x_k(t)$, $k = 1, 2, \dots$ and the fact that, by Assumption 1, $\mathbb{E}(x_k(t)) = 0$, we obtain

$$\begin{aligned} \mathbb{E}[\delta_{\mathcal{J}} \mathbf{1}(\Omega_\alpha)] &= -\frac{\theta_{\mathcal{J}}}{n} + \mathbb{E}\left[\frac{1}{n} \sum_{t=1}^n \phi_N(t) \sum_{k=N_n+1}^\infty \theta_k x_k(t) \mathbf{1}(\Omega_\alpha)\right] \\ &+ \mathbb{E}\left[D_{\mathcal{J}} \left(-\frac{\theta_{\mathcal{J}}}{n} + \frac{1}{n} \sum_{t=1}^n \phi_N(t) \sum_{k=N_n+1}^\infty \theta_k x_k(t)\right) \mathbf{1}(\Omega_\alpha)\right] \\ &= -\frac{\theta_{\mathcal{J}}}{n} (I + \mathbb{E}[D_{\mathcal{J}} \mathbf{1}(\Omega_\alpha)]). \end{aligned}$$

Then (32) follows from (30). \square

PROOF OF LEMMA 3. We proceed in steps.

Step 1'. Let Ω_α be defined by (17), and choose $\alpha = \alpha_* = 2N_n^2 n^{-8}$. Denote $\Omega_* = \Omega_{\alpha_*}$ and write

$$(33) \quad \mathbb{E}[\|\hat{\theta}^* - \theta\|^2] = \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] + \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\bar{\Omega}_*)] + \sum_{k=N_n+1}^\infty \theta_k^2.$$

Now we establish upper bounds on the first two terms in the r.h.s. of (33).

Step 2'. First, note that Ω_* is \mathcal{F}_x^n measurable and

$$\mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] = \mathbb{E}\{\mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mid \mathcal{F}_x^n] \mathbf{1}(\Omega_*)\}.$$

Now consider the sequence model (13), and observe that conditionally on \mathcal{F}_x^n the vector $\tilde{y}_{\mathcal{J}}$ is Gaussian. Note that with our choice $\alpha = \alpha_*$ we have $\mu_n(\alpha_*) = \mu_n^* = 2\sqrt{2qn^{-1} \ln n} + n^{-1}$. Thus, (18) is fulfilled for n large enough. If the event Ω_* holds, then by Lemma 1 the conditional covariance matrix of $\xi_{\mathcal{J}}$ equals $I - A_{\mathcal{J}}$, and inequality (19) is valid. Therefore, due to (24), we can apply Lemma 2 with $\mu_\varepsilon = 3\mu_n^*$ conditionally on \mathcal{F}_x^n to estimate the vector $\theta_{\mathcal{J}} + \delta_{\mathcal{J}}$.

Thus, on the event Ω_* one has

$$\begin{aligned} \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}} - \delta_{\mathcal{J}}\|^2 \mid \mathcal{F}_x^n] &\leq \sum_{j=1}^{J_n} \frac{\|\theta_{(j)} + \delta_{(j)}\|^2 (\sigma^2/n) n_j}{\|\theta_{(j)} + \delta_{(j)}\|^2 + (\sigma^2/n) n_j} \\ &+ \frac{\sigma^2}{n} (21\mu_n^* N_n + 4J_n). \end{aligned}$$

Taking expectation and using the Jensen inequality, we obtain

$$\begin{aligned} \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}} - \delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] &\leq \sum_{j=1}^{J_n} \frac{\mathbb{E}[\|\theta_{(j)} + \delta_{(j)}\|^2 \mathbf{1}(\Omega_*)] (\sigma^2/n) n_j}{\mathbb{E}[\|\theta_{(j)} + \delta_{(j)}\|^2 \mathbf{1}(\Omega_*)] + (\sigma^2/n) n_j} \\ &+ \frac{\sigma^2}{n} (21\mu_n^* N_n + 4J_n). \end{aligned}$$

Let $V_{j,n}(\alpha)$ denote the expression on the r.h.s. of (31). By Lemma 5,

$$\begin{aligned} \mathbb{E}[\|\theta_{(j)} + \delta_{(j)}\|^2 \mathbf{1}(\Omega_*)] &\leq \|\theta_{(j)}\|^2 \mathbb{P}(\Omega_*) + \mathbb{E}[\|\delta_{(j)}\|^2 \mathbf{1}(\Omega_*)] + 2U_n^* n^{-1} \|\theta_{(j)}\|^2 \\ &\leq \|\theta_{(j)}\|^2 (1 + 2U_n^* n^{-1}) + V_{j,n}^* \end{aligned}$$

[here we use the notation $U_n^* = U_n(\alpha_*)$ and $V_{j,n}^* = V_{j,n}(\alpha_*)$]. On the other hand,

$$\mathbb{E}[\|\theta_{(j)} + \delta_{(j)}\|^2 \mathbf{1}(\Omega_*)] \geq \|\theta_{(j)}\|^2 (1 - \alpha_* - 2U_n^* n^{-1}).$$

Thus,

$$\begin{aligned} &\mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}} - \delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] \\ &\leq \sum_{j=1}^{J_n} \frac{\{\|\theta_{(j)}\|^2 (1 + 2U_n^* n^{-1}) + V_{j,n}^*\} (\sigma^2/n) n_j}{\|\theta_{(j)}\|^2 (1 - \alpha_* - 2U_n^* n^{-1}) + (\sigma^2/n) n_j} + \frac{\sigma^2}{n} (21\mu_n^* N_n + 4J_n) \\ &\leq \sum_{j=1}^{J_n} \left(\frac{\|\theta_{(j)}\|^2 (\sigma^2/n) n_j}{\|\theta_{(j)}\|^2 (1 - \alpha_* - 2U_n^* n^{-1}) + (\sigma^2/n) n_j} + 4U_n^* \frac{\sigma^2 n_j}{n^2} + V_{j,n}^* \right) \\ (34) \quad &+ \frac{\sigma^2}{n} (21\mu_n^* N_n + 4J_n) \\ &\leq \sum_{j=1}^{J_n} \frac{\|\theta_{(j)}\|^2 (\sigma^2/n) n_j}{\|\theta_{(j)}\|^2 + (\sigma^2/n) n_j} \left(1 + \frac{\alpha_* + 2U_n^* n^{-1}}{1 - \alpha_* - 2U_n^* n^{-1}} \right) \\ &+ 2U_n^* \frac{\sigma^2 N_n}{n^2} + \sum_{j=1}^{J_n} V_{j,n}^* + \frac{\sigma^2}{n} (21\mu_n^* N_n + 4J_n) \\ &= \left(r_n^{\text{BL}} + \frac{4\sigma^2 J_n}{n} \right) (1 + o(1)) + o\left(\sum_{k=N_n+1}^{\infty} \theta_k^2 \right), \quad n \rightarrow \infty. \end{aligned}$$

Here we used that $2U_n^* N_n n^{-2} = o(n^{-1})$, $\mu_n^* = O(\sqrt{\ln n/n})$, $N_n = O(\sqrt{n}/\ln n)$, and

$$\begin{aligned} \sum_{j=1}^{J_n} V_{j,n}^* &= 2(U_n^*)^2 \left[\frac{\|\theta_{\mathcal{J}}\|^2}{n^2} + \frac{N_n}{n} \sum_{k \geq N_n+1} \theta_k^2 \right] \\ (35) \quad &= O\left(\frac{1}{\sqrt{n} \ln n} \right) \sum_{k=N_n+1}^{\infty} \theta_k^2 + O\left(\frac{1}{n^2} \right). \end{aligned}$$

On the other hand, we have

$$\begin{aligned} &\mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] \\ (36) \quad &\leq \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}} - \delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] + \mathbb{E}[\|\delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] \\ &+ 2(\mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}} - \delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)])^{1/2} (\mathbb{E}[\|\delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)])^{1/2}. \end{aligned}$$

It follows from (31) and (35) that

$$\mathbb{E}[\|\delta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] \leq \sum_{j=1}^{J_n} V_{j,n}^* = O\left(\frac{1}{\sqrt{n}}\right) \left[\left(r_n^{\text{BL}} + \frac{4\sigma^2 J_n}{n} \right) + o\left(\sum_{k=N_n+1}^{\infty} \theta_k^2 \right) \right].$$

Combining this with (34) and (36) we conclude that

$$(37) \quad \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\Omega_*)] \leq \left(r_n^{\text{BL}} + \frac{4\sigma^2 J_n}{n} \right) (1 + o(1)) \\ + o\left(\sum_{k=N_n+1}^{\infty} \theta_k^2 \right), \quad n \rightarrow \infty.$$

Step 3'. Now we establish an upper bound on the second term in the r.h.s. of (33). We have

$$(38) \quad \mathbb{E}[\|\hat{\theta}_{\mathcal{J}}^* - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\bar{\Omega}_*)] = \sum_{j=1}^{J_n} \mathbb{E}[\|\hat{\theta}_{(j)} - \theta_{(j)}\|^2 \mathbf{1}(\bar{\Omega}_*)] \leq K_1 + K_2,$$

where

$$K_1 = 2 \sum_{j=1}^{J_n} \mathbb{E}[\|\tilde{y}_{(j)} - \theta_{(j)}\|^2 \mathbf{1}(\bar{\Omega}_*)], \\ K_2 = \frac{2\sigma^4}{n^2} \sum_{j=1}^{J_n} \left(n_j^2 \mathbb{E}[\|\tilde{y}_{(j)}\|^{-2} \mathbf{1}(\bar{\Omega}_*)] \right).$$

It follows immediately from (14) and definition of the set Ω_* that

$$(39) \quad \mathbb{E}[\|\delta_{\mathcal{J}}\|^2 \mathbf{1}(\bar{\Omega}_*)] \\ \leq 2n^2 \left(n^{-1} \|\theta_{\mathcal{J}}\|^2 \mathbb{P}(\bar{\Omega}_*) + \mathbb{E} \left[\left\| \frac{1}{n} \sum_{t=1}^n \phi_N(t) \sum_{k=N_n+1}^{\infty} \theta_k x_k(t) \right\|^2 \mathbf{1}(\bar{\Omega}_*) \right] \right) \\ \leq 2n \|\theta_{\mathcal{J}}\| \alpha_* + 2n^2 \left[\mathbb{E} \left\| \frac{1}{n} \sum_{t=1}^n \phi_N(t) \sum_{k=N_n+1}^{\infty} \theta_k x_k(t) \right\|^4 \right]^{1/2} \sqrt{\alpha_*}.$$

For every $k = 1, \dots, N$, due to Assumption 1,

$$\mathbb{E} \left(\frac{1}{n} \sum_{t=1}^n x_k(t) \sum_{j=N+1}^{\infty} \theta_j x_j(t) \right)^4 \\ = \frac{1}{n^4} \mathbb{E} \left[\sum_{t,s=1}^n x_k^2(t) x_k^2(s) \left(\sum_{j=N+1}^{\infty} \theta_j x_j(t) \right)^2 \left(\sum_{j=N+1}^{\infty} \theta_j x_j(s) \right)^2 \right] \\ \leq \frac{1}{n^2} \left(\sum_{j=N+1}^{\infty} \theta_j^2 \right)^2 + \frac{C_1}{n^4} \sum_{t=1}^n \mathbb{E} \left(\sum_{j=N+1}^{\infty} \theta_j x_j(t) \right)^4 \\ \leq \frac{L^4}{n^2} + \frac{C_1}{n^3} \sum_{j_1, \dots, j_4=N+1}^{\infty} \theta_{j_1} \theta_{j_2} \theta_{j_3} \theta_{j_4} \mathbb{E}[x_{j_1}(t) x_{j_2}(t) x_{j_3}(t) x_{j_4}(t)] \leq C_2 n^{-2}$$

(here $C_1 > 0$, $C_2 > 0$ are constants depending on L , c_* and H only). Using this bound and (39) we conclude that

$$\mathbb{E}[\|\delta_{\mathcal{J}}\|^2 \mathbf{1}(\bar{\Omega}_*)] \leq 2n\|\theta_{\mathcal{J}}\|\alpha_* + 2C_2N^2\sqrt{\alpha_*} = o(n^{-2}), \quad n \rightarrow \infty.$$

Further, acting as in Lemma 3 in Goldenshluger and Tsybakov (1999) we get $\mathbb{E}(\|\xi_{\mathcal{J}}\|^4)^{1/2} \leq C_3N_n$, where $C_3 > 0$ depends on c_* and H only. Thus,

$$\begin{aligned} K_1 &= \mathbb{E}[\|\tilde{y}_{\mathcal{J}} - \theta_{\mathcal{J}}\|^2 \mathbf{1}(\bar{\Omega}_*)] \\ (40) \quad &\leq 2\mathbb{E}[\|\delta_{\mathcal{J}}\|^2 \mathbf{1}(\bar{\Omega}_*)] + \frac{2\sigma^2}{n} \mathbb{E}[\|\xi_{\mathcal{J}}\|^4]^{1/2} \sqrt{\alpha_*} \\ &= o(n^{-2}), \quad n \rightarrow \infty. \end{aligned}$$

It remains to bound K_2 . We note again that conditionally on \mathcal{F}_x^n , $\tilde{y}_{(j)}$ is a Gaussian vector with mean $\theta_{(j)} + \delta_{(j)}$ and covariance matrix $\sigma^2 n^{-1} S_{(j)}$, where $S_{(j)}$ denotes the corresponding principle submatrix of the matrix $S_{\mathcal{J}}$ [see (15)]. Further, write $\tilde{y}_{(j)} = \sigma n^{-1/2} S_{(j)}^{1/2} \tilde{z}_{(j)}$, where $\tilde{z}_{(j)} | \mathcal{F}_x^n \sim \mathcal{N}(\sigma^{-1} n^{1/2} S_{(j)}^{-1/2} [\theta_{(j)} + \delta_{(j)}], I)$. Thus, we have

$$\begin{aligned} (41) \quad \mathbb{E}[\|\tilde{y}_{(j)}\|^{-2} | \mathcal{F}_x^n] &\leq \frac{n}{\sigma^2 \lambda_{\min}[S_{(j)}]} \mathbb{E}[\|\tilde{z}_{(j)}\|^{-2} | \mathcal{F}_x^n] \\ &\leq \frac{n}{\sigma^2 \lambda_{\min}[S_{(j)}] (n_j - 2)}. \end{aligned}$$

The last inequality is a consequence of the following. Note that $\mathbb{E}[\|\tilde{z}_{(j)}\|^{-2} | \mathcal{F}_x^n] \leq \mathbb{E}[\|\tilde{z}_{(j)}\|^{-2} | \mathcal{F}_x^n]$ [cf. Lehmann and Casella (1998), page 355], where $\tilde{z}_{(j)} | \mathcal{F}_x^n \sim \mathcal{N}(0, I)$. Then (41) follows from the fact that $\|\tilde{z}_{(j)}\|^2 | \mathcal{F}_x^n \sim \chi_{(n_j)}^2$ and $n_j > 4$. Thus,

$$\begin{aligned} (42) \quad \mathbb{E}[\|\tilde{y}_{(j)}\|^{-2} \mathbf{1}(\bar{\Omega}_*)] &= \mathbb{E}\{\mathbb{E}[\|\tilde{y}_{(j)}\|^{-2} | \mathcal{F}_x^n] \mathbf{1}(\bar{\Omega}_*)\} \\ &\leq \frac{n}{\sigma^2 (n_j - 2)} \mathbb{E}\{\lambda_{\min}^{-1}[S_{(j)}] \mathbf{1}(\bar{\Omega}_*)\} \\ &\leq \frac{n}{\sigma^2 (n_j - 2)} \{E\lambda_{\min}^{-2}[S_{(j)}]\}^{1/2} \sqrt{\alpha_*}. \end{aligned}$$

Since $S_{(j)}$ is a principle submatrix of $S_{\mathcal{J}}$, we have $\lambda_{\min}[S_{(j)}] \geq \lambda_{\min}[S_{\mathcal{J}}]$ [see, e.g., Marcus and Minc (1992), Chapter III, Section 3.6.5]. By definition of $Q_{\mathcal{J}}$ and independently of the event Ω_* , we have $\lambda_{\min}[Q_{\mathcal{J}}^{-1}] = \lambda_{\max}^{-1}[Q_{\mathcal{J}}] \geq (\sqrt{2q}N_n + n^{-1})^{-1}$. Therefore, if n is large enough, then $\lambda_{\min}[S_{\mathcal{J}}] \geq C_4 N_n^{-1}$, where $C_4 > 0$ depends only on q . Combining these inequalities with (42), we obtain

$$\mathbb{E}[\|\tilde{y}_{(j)}\|^{-2} \mathbf{1}(\bar{\Omega}_*)] \leq \frac{C_5 n N_n^2}{\sigma^2 (n_j - 2) n^4},$$

where $C_5 > 0$. This entails $K_2 = o(n^{-2})$, $n \rightarrow \infty$. Combining this with (33), (37), (38) and (40) we complete the proof. \square

PROOF OF THEOREMS 2 AND 3. We prove only Theorem 2 because the proof of Theorem 3 is quite similar. We indicate how the statement of Theorem 3 follows from the proof of Theorem 2.

First note that Assumption 3 on the blocks guarantees that (24) is satisfied for n large enough. Note also that under this assumption,

$$J_n = O(\rho_n^{-1} \ln N_n) \quad \text{and} \quad \frac{\text{card}(B_{j+1})}{\text{card}(B_j)} \leq 1 + O(\rho_n).$$

Therefore, using Lemmas 3 and 4 we get for $\theta \in \Theta(a, L)$ and for $\hat{\theta}^*$ defined in (7) and (8),

$$\begin{aligned} & \mathbb{E} \|\hat{\theta}^* - \theta\|^2 \\ & \leq \left(\inf_{h \in \mathcal{H}} \mathcal{R}_{\sigma/\sqrt{n}}(h, \theta) + \frac{4\sigma^2 J_n}{n} \right) (1 + o(1)) \\ & \leq \left([1 + O(\rho_n)] \inf_{h \in \mathcal{H}_{\text{mon}}} \mathcal{R}_{\sigma/\sqrt{n}}(h, \theta) + \frac{\sigma^2}{n} (4J_n + \nu_n) + \sum_{k=N+1}^{\infty} \theta_k^2 \right) (1 + o(1)), \end{aligned}$$

where the $o(1)$ term is uniform in $\|\theta\| \leq L$. Observe that

$$\mathbb{E} [\hat{y}^*(n+1) - y(n+1)]^2 - \sigma^2 = \mathbb{E} \|\hat{\theta}^* - \theta\|^2.$$

Hence

$$\begin{aligned} & \mathbb{E} [\hat{y}^*(n+1) - y(n+1)]^2 - \sigma^2 \\ (43) \quad & \leq \left(\inf_{h \in \mathcal{H}_{\text{mon}}} \mathcal{R}_{\sigma/\sqrt{n}}(h, \theta) + \frac{\sigma^2}{n} (4\rho_n^{-1} \ln N_n + \nu_n) + \frac{L^2}{a_{\min}^2} N_n^{-2\beta} \right) (1 + o(1)) \\ & \leq (1 + o(1)) \inf_{h \in \mathcal{H}_{\text{mon}}} \mathcal{R}_{\sigma/\sqrt{n}}(h, \theta) + O\left(\frac{\ln n}{n\rho_n} + (\ln n)^{2\beta} n^{-\beta}\right), \end{aligned}$$

where we used that $\sum_{k=N_n+1}^{\infty} \theta_k^2 \leq L^2/a_{N_n}^2 \leq L_{\max}^2 N_n^{-2\beta}/a_{\min}^2$. Note that both $o(1)$ and $O(\cdot)$ in (43) are uniform over $\theta \in \cup_{(a, L) \in \mathcal{A}} \Theta(a, L)$. Taking the supremum of both sides of (43) over $\theta \in \Theta(a, L)$ and using the monotonicity of $\{a_k\}$ we get

$$\begin{aligned} & \mathcal{R}[\hat{y}^*; \Theta(a, L)] \leq \inf_{\bar{\theta}} \sup_{\theta \in \Theta(a, L)} \mathbb{E}_* \|\bar{\theta} - \theta\|^2 [1 + o(1)] \\ (44) \quad & \quad \quad \quad + O\left(\frac{\ln n}{n\rho_n} + (\ln n)^{2\beta} n^{-\beta}\right). \end{aligned}$$

By Pinsker's (1980) theorem [see also Belitser and Levit (1995)], we get

$$(45) \quad \inf_{\bar{\theta}} \sup_{\theta \in \Theta(a, L)} \mathbb{E}_* \|\bar{\theta} - \theta\|^2 = r_n (1 + o(1)) \geq cn^{-2\beta/(2\beta+1)}$$

and $c > 0$ depends only on β, a and L . Recalling that $\beta > 1/2$ and combining the last two inequalities we get the result of Theorem 2.

The proof of Theorem 3 follows the same lines as the above proof. We need to modify only the last step. Since $\Theta(a, L) \supseteq \Theta(a^1, L_{\min}), \forall (a, L) \in \mathcal{A}$, where $a^1 = \{a_{\max} k^{\beta_1}\}$, instead of (45) we may write

$$(46) \quad \inf_{\tilde{\theta}} \sup_{\theta \in \Theta(a, L)} \mathbb{E}_* \|\tilde{\theta} - \theta\|^2 \geq \inf_{\tilde{\theta}} \sup_{\theta \in \Theta(a^1, L_{\min})} \mathbb{E}_* \|\tilde{\theta} - \theta\|^2 \geq c' n^{-2\beta_1/(2\beta_1+1)},$$

where $c' > 0$ depends only on β_1, a_{\max} and L_{\min} . On the other hand, $O(\frac{\ln n}{n\rho_n} + (\ln n)^{2\beta} n^{-\beta}) = O((\ln n)^{2\beta_1} n^{-\beta_0})$, since $\beta_0 \leq \beta \leq \beta_1$. This and (46) show that under the condition $2\beta_1/(2\beta_1 + 1) < \beta_0 \leq \beta_1$ the first term in the r.h.s. of (44) asymptotically dominates the second one uniformly over all pairs $(a, L) \in \mathcal{A}$. \square

APPENDIX

PROOF OF LEMMA 2. The proof is obtained as an immediate consequence of Lemma 7 given below. First we present an auxiliary lemma concerning the estimation of a d -dimensional Gaussian mean.

Let $y_k = \theta_k + \varepsilon \xi_k, k = 1, \dots, d$, where $y = (y_1, \dots, y_d)'$ is the observation vector, and $\theta = (\theta_1, \dots, \theta_d)'$ is the parameter to be estimated. We assume that $\xi = (\xi_1, \dots, \xi_d)' \sim \mathcal{N}_d(0, Q)$, where Q is a positive definite $d \times d$ matrix.

LEMMA 6. Let $g: \mathbb{R}^d \rightarrow \mathbb{R}^d$ and

$$\mathbb{E} \left(|y_k| |g_k(y)| + \left| \frac{\partial g_k(y)}{\partial y_k} \right| \right) < \infty, \quad k = 1, \dots, d.$$

Then

$$\mathbb{E} \|y + g(y) - \theta\|^2 = \varepsilon^2 \text{tr}(Q) + 2\varepsilon^2 \text{Etr}[QD_g(y)] + \mathbb{E} \|g(y)\|^2,$$

where $D_g = D_g(y) = \{\partial g_i / \partial y_j\}_{i, j=1, \dots, d}$.

The proof is an easy modification of Stein's (1981) theorem [e.g., Johnstone (1998)].

Assume that $Q = I - A$, where A is a symmetric $d \times d$ matrix such that $\max_{i, j=1, \dots, d} |[A]_{ij}| \leq \mu < 1$. Consider the Stein estimator

$$\tilde{\theta} = \left(1 - \frac{\varepsilon^2 d}{\|y\|^2} \right) y.$$

LEMMA 7. Assume that $d > 4$, and $\mu < 1/6 - 2/(3d)$. Then

$$(47) \quad \mathbb{E} \|\tilde{\theta} - \theta\|^2 \leq \frac{\|\theta\|^2 \varepsilon^2 d}{\|\theta\|^2 + \varepsilon^2 d} + 7\varepsilon^2 \mu d + 4\varepsilon^2.$$

For $\mu = 0$ this inequality is given in Donoho and Johnstone (1995).

PROOF. We apply Lemma 6 with $g(y) = -\varepsilon^2 d \|y\|^{-2} y$. First we note that

$$(48) \quad \frac{\partial g_i(y)}{\partial y_j} = \begin{cases} -\varepsilon^2 d \left(\frac{1}{\|y\|^2} - \frac{2y_i^2}{\|y\|^4} \right), & i = j, \\ 2\varepsilon^2 d \frac{y_i y_j}{\|y\|^4}, & i \neq j. \end{cases}$$

Thus, $\mathbb{E} \operatorname{tr}[D_g(y)] = -\varepsilon^2 d(d-2)\mathbb{E}\|y\|^{-2}$, and simple algebra leads to

$$(49) \quad \begin{aligned} \mathbb{E}\|\tilde{\theta} - \theta\|^2 &= \varepsilon^2 \operatorname{tr}(I - A) + 2\varepsilon^2 \mathbb{E} \operatorname{tr}[(I - A)D_g(y)] + \mathbb{E}\|g(y)\|^2 \\ &= \varepsilon^2 d - \varepsilon^4 d(d-4)\mathbb{E}\|y\|^{-2} - \varepsilon^2 \operatorname{tr}(A) - 2\varepsilon^2 \mathbb{E} \operatorname{tr}(AD_g). \end{aligned}$$

Further,

$$(50) \quad \mathbb{E} \operatorname{tr}(AD_g) = \operatorname{tr}(A\mathbb{E}(D_g)) \leq \mu \mathbb{E} \left(\sum_{i,j=1}^d |[D_g]_{ij}| \right).$$

In view of (48), $\mathbb{E}[D_g]_{ii} \leq \varepsilon^2 d \mathbb{E}(\|y\|^{-2} + 2y_i^2 \|y\|^{-4})$, and $\mathbb{E}[D_g]_{ij} \leq 2\varepsilon^2 d \mathbb{E}(|y_i y_j| \|y\|^{-4})$ for $i \neq j$. Therefore,

$$\begin{aligned} \mathbb{E} \left(\sum_{i,j=1}^d |[D_g]_{ij}| \right) &\leq \varepsilon^2 d \sum_{i=1}^d [\mathbb{E}\|y\|^{-2} + 2\mathbb{E}(y_i^2 \|y\|^{-4})] + 2\varepsilon^2 d \sum_{i,j=1, i \neq j}^d \mathbb{E}(|y_i y_j| \|y\|^{-4}) \\ &= \varepsilon^2 (d^2 + 2d) \mathbb{E}\|y\|^{-2} + 2\varepsilon^2 d \mathbb{E} \left[\|y\|^{-4} \left(\left(\sum_{i=1}^d |y_i| \right)^2 - \sum_{i=1}^d |y_i|^2 \right) \right] \\ &\leq \varepsilon^2 (d^2 + 2d) \mathbb{E}\|y\|^{-2} + 2\varepsilon^2 d(d-1) \mathbb{E} \left(\|y\|^{-4} \sum_{i=1}^d |y_i|^2 \right) \\ &= 3\varepsilon^2 d^2 \mathbb{E}\|y\|^{-2}. \end{aligned}$$

Taking into account (50) and (49) we come to

$$(51) \quad \begin{aligned} \mathbb{E}\|\tilde{\theta} - \theta\|^2 &\leq \varepsilon^2 d - \varepsilon^4 d(d-4)\mathbb{E}\|y\|^{-2} + \varepsilon^2 \mu d + 6\varepsilon^4 \mu d^2 \mathbb{E}\|y\|^{-2} \\ &= \varepsilon^2 d - \varepsilon^4 d(d-4)\mathbb{E}\|y\|^{-2} \left(1 - \frac{6\mu d}{d-4} \right) + \varepsilon^2 \mu d. \end{aligned}$$

By Jensen's inequality $\mathbb{E}\|y\|^{-2} \geq (\mathbb{E}\|y\|^2)^{-1} = (\|\theta\|^2 + \varepsilon^2 d)^{-1}$. This along with (51) and $6\mu d < d-4$, leads to (47). \square

REFERENCES

- BELITSER, E. N. and LEVIT, B. YA. (1995). On minimax filtering on ellipsoids. *Math. Method Statist.* **4** 259–273.
 BREIMAN, L. and FREEDMAN, D. (1983). How many variables should be entered in a regression equation? *J. Amer. Statist. Assoc.* **78** 131–136.

- CAI, T. (1999). Adaptive wavelet estimation: a block thresholding and oracle inequality approach. *Ann. Statist.* **27** 898–924.
- CAVALIER, L. and TSYBAKOV, A. (2000). Sharp adaptation for inverse problems with random noise. *Probab. Theory Related Fields*. To appear. Available at www.proba.jussieu.fr.
- DONOHO, D. and JOHNSTONE, I. M. (1995). Adapting to unknown smoothness via wavelet shrinkage. *J. Amer. Statist. Assoc.* **90** 1200–1224.
- EFROMOVICH, S. (1999). *Nonparametric Curve Estimation*. Springer, New York.
- EFROMOVICH, S. YU. and PINSKER, M. S. (1984). Learning algorithm for nonparametric filtering. *Automat. Remote Control*, **11** 1434–1440.
- EFROMOVICH, S. and PINSKER, M. S. (1996). Sharp-optimal and adaptive estimation for heteroscedastic nonparametric regression. *Statist. Sinica* **6** 925–942.
- GOLDENSHLUGER, A. and TSYBAKOV, A. (1999). Optimal prediction for linear regression with infinitely many parameters. Available at www.proba.jussieu.fr.
- JOHNSTONE, I. M. (1998). Function estimation in Gaussian noise: sequence models. Available at www-stat.stanford.edu/.
- JOHNSTONE, I. M. (1999). Wavelet shrinkage for correlated data and inverse problems: adaptivity results. *Statist. Sinica* **9** 51–83.
- LEHMANN, E. and CASELLA, G. (1998). *Theory of Point Estimation*. Springer, New York.
- MARCUS, M. and MINC, H. (1992). *A Survey of Matrix Theory and Matrix Inequalities*. Dover, New York.
- NEMIROVSKI, A. (2000). *Topics in Non-Parametric Statistics. Ecole d'été de Probabilités Saint Flour XXVIII. Lecture Notes in Math.* **1738** 89–277. Springer, Berlin.
- PETROV, V. V. (1995). *Limit Theorems of Probability Theory*. Clarendon Press, Oxford.
- PINSKER, M. S. (1980). Optimal filtering of square integrable signals in Gaussian white noise. *Problems Inform. Transmission* **16** 120–133.
- SHIBATA, R. (1981). An optimal selection of regression variables. *Biometrika* **68** 45–54.
- STEIN, CH. (1981). Estimation of the mean of a multivariate normal distribution. *Ann. Statist.* **9** 1135–1151.

DEPARTMENT OF STATISTICS
UNIVERSITY OF HAIFA
HAIFA 31905
ISRAEL
E-MAIL: goldensh@stat.haifa.ac.il

LABORATOIRE DE PROBABILITÉS
ET MODÈLES ALÉATOIRES
UNIVERSITÉ PARIS VI
4PL. JUSSIEU
PARIS 75252
FRANCE