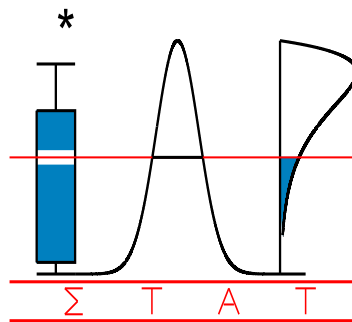


T E C H N I C A L  
R E P O R T

11036

**Adaptive Gaussian inverse regression  
with partially unknown operator**

JOHANNES, J. and M. SCHWARZ



I A P S T A T I S T I C S  
N E T W O R K

**INTERUNIVERSITY ATTRACTION POLE**

<http://www.stat.ucl.ac.be/IAP>

# Adaptive Gaussian inverse regression with partially unknown operator

JAN JOHANNES\*

MAIK SCHWARZ\*

Université catholique de Louvain

November 9, 2011

This work deals with the ill-posed inverse problem of reconstructing a function  $f$  given implicitly as the solution of  $g = Af$ , where  $A$  is a compact linear operator with unknown singular values and known eigenfunctions. We observe the coefficients of  $g$  and the singular values subject to Gaussian white noise with noise levels  $\varepsilon$  and  $\sigma$ .

We develop a minimax theory in terms of both noise levels and propose an orthogonal series estimator attaining the minimax rates. This estimator requires the optimal choice of a dimension parameter depending on certain characteristics of  $f$  and  $A$ . This work addresses the fully data-driven choice of the dimension parameter combining model selection with Lepski's method. We show that the fully data-driven estimator preserves minimax optimality over a wide range of classes for  $f$  and  $A$  and noise levels  $\varepsilon$  and  $\sigma$ . The results are illustrated considering Sobolev spaces and mildly and severely ill-posed inverse problems.

*AMS 2010 subject classifications:* 62G05, 62G08.

*Keywords:* Gaussian sequence space model, minimax theory, adaptive nonparametric estimation, model selection, Lepski's method, Sobolev spaces, mildly and severely ill-posed inverse problems

This work was supported by the IAP research network no. P6/03 the Belgian Government (Belgian Science Policy) and by the "Fonds Spéciaux de Recherche" from the Université catholique de Louvain.

## 1. Introduction

Let  $(H, \langle \cdot, \cdot \rangle_H)$  and  $(G, \langle \cdot, \cdot \rangle_G)$  be separable Hilbert spaces and  $A$  a compact linear operator from  $H$  to  $G$  with unknown singular values. This work deals with the reconstruction of a function  $f \in H$  given noisy observations of the image  $g = Af$  on the one hand and of the unknown sequence of singular values  $b = (a_j)_{j \in \mathbb{N}}$  on the other hand. In other words, we consider a statistical inverse problem with partially unknown operator. There is a vast literature on statistical inverse problems. For the case where the operator is fully known, the reader may refer to Cavalier et al. (2002), Mair and Ruymgaart (1996), Mathé and Pereverzev (2001), and Johnstone and Silverman (1990) and the references therein. A typical illustration of such a situation is a deconvolution problem (cf. Fan (1991), Ermakov (1990),

---

\*e-mail:  [{jan.johannes|maik.schwarz}@uclouvain.be](mailto:{jan.johannes|maik.schwarz}@uclouvain.be)

and Stefanski and Carroll (1990) among many others). For a more detailed discussion and motivation of the case of a partially unknown operator which we consider in this work, we refer the reader to Cavalier and Hengartner (2005). Neumann (1997) and Efromovich (1997) consider such a setting in the particular context of a deconvolution problem.

Let us describe in more detail the model we are going to consider. We suppose that  $A$  admits a singular value decomposition  $(a_j, \varphi_j, \psi_j)_{j \in \mathbb{N}}$  as follows. Denote by  $A^*$  the adjoint operator of  $A$ . Then,  $A^*A$  is a compact operator on  $H$  with eigenvalues  $(a_j^2)_{j \in \mathbb{N}}$  whose associated orthonormal basis of eigenfunctions  $\{\varphi_j\}$  we suppose to be known. Analogously, the operator  $AA^*$  has eigenvalues  $(a_j^2)_{j \in \mathbb{N}}$  and known orthonormal eigenfunctions  $\psi_j = \|A\varphi_j\|_G^{-1}A\varphi_j$  in  $G$ . Projecting the inverse problem  $g = Af$  on the eigenfunctions, we obtain the system of equations  $[g]_j := \langle g, \psi_j \rangle_G = a_j \langle f, \varphi_j \rangle_H$  for  $j \in \mathbb{N}$ . As the operator  $A$  is compact, the sequence of singular values tends to zero and the inverse problem is called ill-posed.

The solution  $f$  is characterized by its coefficients  $[f]_j := \langle f, \varphi_j \rangle_H$ . Our objective is their estimation based on the following observations:

$$Y_j = [g]_j + \sqrt{\varepsilon} \xi_j = a_j [f]_j + \sqrt{\varepsilon} \xi_j \quad \text{and} \quad X_j = a_j + \sqrt{\sigma} \eta_j \quad (j \in \mathbb{N}), \quad (1.1)$$

where the  $\xi_j, \eta_j$  are iid. standard normally distributed random variables and  $\varepsilon, \sigma \in (0, 1)$  are noise levels. Thus we represent the problem at hand as a hierarchical Gaussian sequence space model. Of course  $f$  can only be reconstructed from such observations if all the  $a_j$  are non-zero which is the case if and only if the operator  $A$  is injective. We assume this from now on, which allows us to write  $f = \sum_{j=1}^{\infty} [g]_j a_j^{-1} \varphi_j$ . Hence, an orthogonal series estimator of  $f$  seems to be a natural approach:

$$\hat{f}_k := \sum_{j=1}^k \frac{Y_j}{X_j} \mathbf{1}_{[X_j^2 \geq \sigma]} \varphi_j.$$

The stabilizing threshold on the random denominator  $X_j$  corresponds to its noise level as an estimator of  $a_j$ . Note that  $\hat{f}_k$  depends on a dimension parameter  $k$  whose choice essentially determines the estimation accuracy. Its optimal choice generally depends on both unknown sequences  $([f]_j)$  and  $(a_j)$ . Our purpose is to establish an adaptive estimation procedure for the function  $f$  which does not depend on these sequences. More precisely, assuming that the solution and the operator belong to given classes  $f \in \mathcal{F}$  and  $A \in \mathcal{A}$ , respectively, we shall measure the accuracy of an estimator  $\tilde{f}$  of  $f$  by the maximal weighted risk  $\mathcal{R}_\omega(\tilde{f}, \mathcal{F}, \mathcal{A}) := \sup_{f \in \mathcal{F}} \sup_{A \in \mathcal{A}} \mathbb{E} \|\tilde{f} - f\|_\omega^2$  defined with respect to some weighted norm  $\|\cdot\|_\omega := \sum_{j \in \mathbb{N}} \omega_j |[\cdot]_j|^2$ , where  $\omega := (\omega_j)_{j \in \mathbb{N}}$  is a strictly positive weight sequences. This allows us to quantify the estimation accuracy in terms of the mean integrated square error (MISE) not only of  $f$  itself, but as well of its derivatives, for example. Given observations  $Y = (Y_j)_{j \in \mathbb{N}}$  and  $X = (X_j)_{j \in \mathbb{N}}$  with respective noise levels  $\varepsilon$  and  $\sigma$  according to (1.1), the minimax risk with respect to the classes  $\mathcal{F}$  and  $\mathcal{A}$  is then defined as  $\mathcal{R}_\omega^*(\varepsilon, \sigma, \mathcal{F}, \mathcal{A}) := \inf_{\tilde{f}} \mathcal{R}_\omega(\tilde{f}, \mathcal{F}, \mathcal{A})$ , where the infimum is taken over all possible estimators  $\tilde{f}$  of  $f$ . An estimator  $\hat{f}$  is said to attain the minimax rate or to be minimax optimal with respect to  $\mathcal{F}$  and  $\mathcal{A}$  if there is a constant  $C > 0$  depending on the classes only such that  $\mathcal{R}_\omega(\hat{f}, \mathcal{F}, \mathcal{A}) \leq C \mathcal{R}_\omega^*(\varepsilon, \sigma, \mathcal{F}, \mathcal{A})$  for all  $\varepsilon, \sigma \in (0, 1)$ . An estimation procedure which is fully data-driven and minimax optimal for a wide range of classes  $\mathcal{F}$  and  $\mathcal{A}$  is called *adaptive*.

In the next section, we show that for a wide range of classes  $\mathcal{F}$  and  $\mathcal{A}$  the orthogonal series estimator  $\hat{f}_{k_\varepsilon^*}$  attains the minimax rate for an optimal choice  $k_\varepsilon^*$  of the dimension parameter. We illustrate this result considering subsets of Sobolev spaces for  $\mathcal{F}$  and distinguishing two types of operator classes  $\mathcal{A}$  specifying the decay of the singular values: If  $(a_j)$  decays polynomially, the inverse problem is called mildly ill-posed and severely ill-posed if they decay exponentially. However,  $k_\varepsilon^*$  is chosen subject to a classical variance-squared-bias tradeoff and depends on properties of both classes  $\mathcal{F}$  and  $\mathcal{A}$  which are unknown in general.

The last section is devoted to the development of a data-driven choice  $\widehat{k}$  of  $k$ , using a combination of the model selection scheme (Barron et al., 1999, cf.) with Lepski's procedure which is inspired by the work of Goldenshluger and Lepski (2011) who consider bandwidth selection for kernel estimators. Given a random sequence  $(\widehat{\text{pen}}_k)_{k \geq 1}$  of penalties, a random set  $\{1, \dots, \widehat{K}\}$  of admissible dimension parameters and the random sequence of contrasts

$$\widehat{\Psi}_k := \max_{k \leq j \leq \widehat{K}} \left\{ \|\widehat{f}_j - \widehat{f}_k\|_\omega^2 - \widehat{\text{pen}}_j \right\} \quad (k \in \mathbb{N}), \quad (1.2)$$

the dimension parameter is selected as the minimizer<sup>1</sup> of a penalized contrast

$$\widehat{k} := \operatorname{argmin}_{1 \leq k \leq \widehat{K}} \left\{ \widehat{\Psi}_k + \widehat{\text{pen}}_k \right\}. \quad (1.3)$$

We assess the accuracy of the fully data-driven estimator  $\widehat{f}_{\widehat{k}}$  deriving an upper bound for  $\mathcal{R}_\omega(\widehat{f}_{\widehat{k}}, \mathcal{F}, \mathcal{A})$ . Obviously this upper bound heavily depends the random sequence  $(\widehat{\text{pen}}_k)$  and the random upper bound  $\widehat{K}$ . However, we construct these objects in such a way that the resulting fully data-driven estimator  $\widehat{f}_{\widehat{k}}$  is minimax optimal in many cases and thus adaptive.

Adaptive estimation in a hierarchical Gaussian sequence space model has previously been considered by Cavalier and Hengartner (2005). Though, the authors restrict their investigation to the mildly ill-posed case and to noise levels satisfying  $\sigma \leq \varepsilon$ . The new approach presented in this paper has the advantage of not requiring such restrictions. On the contrary, the influence of the two noise levels on the estimation accuracy is characterized. Moreover, the estimator presented in this paper is can attain optimal convergence rates independently of whether the underlying inverse problem is mildly or severely ill-posed, for example, even when  $\varepsilon \ll \sigma$ .

The more technical proofs and some auxiliary results are deferred to the appendix.

## 2. Minimax

In this section we develop a minimax theory for Gaussian inverse regression with respect to the classes

$$\mathcal{F}_\gamma^r := \left\{ h \in H \mid \sum_{j \in \mathbb{N}} \gamma_j |[h]_j|^2 =: \|h\|_\gamma^2 \leq r \right\} \text{ and}$$

$$\mathcal{A}_\lambda^d := \left\{ T \in C(H, G) \mid \text{The eigenvalues } \{u_j\} \text{ of } T^*T \text{ satisfy } 1/d \leq \frac{u_j^2}{\lambda_j} \leq d \quad \forall j \in \mathbb{N} \right\},$$

where  $C(H, G)$  denotes the set of all compact linear operators from  $H$  to  $G$  having  $\{\varphi_j\}$  and  $\{\psi_j\}$  as eigenfunctions, respectively. The minimal regularity conditions on the solution, the operator and the weighted norm  $\|\cdot\|_\omega$  which we need in this section are summarized in the following assumption.

**Assumption 2.1** Let  $\gamma := (\gamma_j)_{j \in \mathbb{N}}$ ,  $\omega := (\omega_j)_{j \in \mathbb{N}}$  and  $\lambda := (\lambda_j)_{j \in \mathbb{N}}$  be strictly positive sequences of weights with  $\gamma_1 = \omega_1 = \lambda_1 = 1$  such that  $\omega/\gamma$  and  $\lambda$  are non-increasing, respectively.

**Illustration 2.2** As an illustration of the results below, we will consider weight sequences  $\gamma_j = j^{2p}$ , for which  $\mathcal{F}_\gamma^r$  is a Sobolev space of  $p$ -times differentiable functions if we consider the trigonometric basis in  $H = L^2[0, 1]$ . As for the operator, we will distinguish the cases  $\lambda_j = j^{-2b}$ , referred to as *mildly ill-posed* (**[m]**) and  $\lambda_j = \exp(1 - j^{2b})$ , the *severely ill-posed* case (**[s]**). Concerning the weighted norm, we will consider sequences<sup>2</sup>  $\omega_j \sim j^{2s}$ , such that  $\|f\|_\omega = \|f^{(s)}\|_{L^2}$  for all  $f \in \mathcal{F}_\gamma^r$ .

The following result states lower risk bounds for the estimation of  $f$  and thus describes the complexity of the problem.

<sup>1</sup>For a sequence  $(b_k)_{k \in \mathbb{N}}$  attaining a minimal value on  $N \subset \mathbb{N}$ , let  $\operatorname{argmin}_{n \in N} b_n := \min\{n \in N \mid b_n \leq b_k \forall k \in N\}$ .

<sup>2</sup> $b_\rho \sim c_\rho$  means that  $\lim_{\rho \rightarrow 0} b_\rho/c_\rho$  exists in  $(0, \infty)$ .

**Theorem 2.3** Suppose that we observe sequences  $Y$  and  $X$  according to the model (1.1). Consider sequences  $\omega$ ,  $\gamma$ , and  $\lambda$  satisfying Assumption 2.1. For all  $\varepsilon, \sigma \in (0, 1)$ , define

$$\rho_{k,\varepsilon} := \max\left(\frac{\omega_k}{\gamma_k}, \sum_{j=1}^k \frac{\varepsilon \omega_j}{\lambda_j}\right), \quad \chi_\varepsilon := \min_{k \in \mathbb{N}} \rho_{k,\varepsilon}, \quad k_\varepsilon^* := \operatorname{argmin}_{k \in \mathbb{N}} \rho_{k,\varepsilon}, \quad \kappa_\sigma := \max_{k \in \mathbb{N}} \left\{ \frac{\omega_k}{\gamma_k} \min\left(1, \frac{\sigma}{\lambda_k}\right) \right\}. \quad (2.1)$$

If  $\eta := \inf_{n \in \mathbb{N}} \{ \chi_\varepsilon^{-1} \min(\omega_{k_\varepsilon^*} \gamma_{k_\varepsilon^*}^{-1}, \sum_{l=1}^{k_\varepsilon^*} \varepsilon \omega_l (\lambda_l)^{-1}) \} > 0$ , then

$$\inf_{\tilde{f}} \sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \left\{ \mathbb{E} \|\tilde{f} - f\|_\omega^2 \right\} \geq \frac{1}{4d} \min(\eta, r) \min(r, 1/(2d), (1 - d^{-1/2})^2) \max(\chi_\varepsilon, \kappa_\sigma).$$

It is noteworthy that apart from the somewhat unwieldy constant, the lower bound is given by two terms ( $\chi_\varepsilon$  and  $\kappa_\sigma$ ), each of which depending only on one noise level. We show in the proof that  $\chi_\varepsilon$  is actually a lower bound when we assume the eigenvalues  $a_j$  to be known. Therefore,  $\kappa_\sigma$  shows to which extent the additional difficulty arising from the preliminary estimation of the eigenvalues  $a_j$  influences the possible estimation accuracy for  $f$ : As long as  $\chi_\varepsilon \geq \kappa_\sigma$ , the same lower bound as in the case of known eigenvalues holds. Otherwise, the lower bound increases. Notice further that the term  $\rho_{k,\varepsilon}$  above corresponds to the MISE of the orthogonal series estimator  $\hat{f}_k$  in the case of known eigenvalues  $a_j$ , and  $k_\varepsilon^*$  is its minimizer with respect to  $k$ . Under classical smoothness assumptions, the rates and  $k_\varepsilon^*$  take the following forms.

**Illustration 2.4** In the special cases defined in Illustration 2.2 above, the rates from (2.1) are

$$\begin{aligned} \text{[m]} \quad & \chi_\varepsilon \sim \varepsilon^{2(p-s)/(2p+2b+1)}, \quad k_\varepsilon^* \sim \varepsilon^{-1/(2p+2b+1)}, \quad \kappa_\sigma \sim \sigma^{((p-s) \wedge b)/b} \\ \text{[s]} \quad & \chi_\varepsilon \sim |\log \varepsilon|^{(p-s)/b}, \quad k_\varepsilon^* \sim |\log \varepsilon|^{1/(2b)}, \quad \kappa_\sigma \sim |\log \sigma|^{-(p-s)/b}. \end{aligned}$$

The following theorem shows that the orthogonal series estimator  $\hat{f}_{k_\varepsilon^*}$  with optimal parameter  $k_\varepsilon^*$  given in (2.1) actually attains the lower risk bound up to a constant and is thus minimax optimal.

**Theorem 2.5** Under the assumptions of Theorem 2.3, the estimator  $\hat{f}_{k_\varepsilon^*}$  satisfies for all  $\varepsilon, \sigma \in (0, 1)$

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \left\{ \mathbb{E} \|\hat{f}_{k_\varepsilon^*} - f\|_\omega^2 \right\} \leq 4(6d + r) \max(\chi_\varepsilon, \kappa_\sigma).$$

To conclude this section, let us summarize the resulting optimal convergence rates under the classical smoothness assumptions introduced in Illustration 2.2. In order to characterize the influence of the second noise level  $\sigma$ , we consider it as a function of the first noise level  $\varepsilon$ .

**Illustration 2.6** Let  $(\sigma_\varepsilon)_{\varepsilon \in (0,1)}$  be a noise level in  $X$  depending on the noise level  $\varepsilon$  in  $Y$ .

**[m]** Let  $p > 1/2$ ,  $b > 1$ , and  $0 \leq s \leq p$ . If  $q_1 := \lim_{\varepsilon \rightarrow 0} \varepsilon^{-2((p-s) \vee b)/(2p+2b+2)} \sigma_\varepsilon$  exists<sup>3</sup>, then

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E} \|\hat{f}_{k_\varepsilon^*}^{(s)} - f^{(s)}\|_{L^2}^2 = \begin{cases} O(\varepsilon^{2(p-s)/(2p+2b+1)}) & \text{if } q_1 < \infty \\ O(\sigma_\varepsilon^{((p-s) \wedge b)/b}) & \text{otherwise.} \end{cases}$$

**[s]** Let  $p > 1/2, b > 0$  and  $0 \leq s \leq p$ . If  $q_2 := \lim_{\varepsilon \rightarrow 0} |\log \varepsilon| |\log \sigma_\varepsilon|^{-1}$  exists, then

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E} \|\hat{f}_{k_\varepsilon^*}^{(s)} - f^{(s)}\|_{L^2}^2 = \begin{cases} O(|\log \varepsilon|^{(p-s)/b}) & \text{if } q_2 < \infty \\ O(|\log \sigma_\varepsilon|^{(p-s)/b}) & \text{otherwise.} \end{cases}$$

This illustration shows that often the same optimal rates as in the case of known eigenvalues hold even when  $\varepsilon < \sigma$ .

<sup>3</sup>The limit « $\infty$ » meaning strict divergence is authorized.

### 3. Adaptation

In this section, we construct a fully data-driven estimator of  $f$  following the procedure sketched in (1.2) and (1.3). The following Lemma will be our key tool when controlling the risk of the adaptive estimator.

**Lemma 3.1** *Let  $\text{pen}$  be an arbitrary positive sequence and  $K \in \mathbb{N}$ . Consider the sequence  $\Psi$  of contrasts  $\Psi_k := \max_{1 \leq j \leq K} \left\{ \|\widehat{f}_j - \widehat{f}_k\|_\omega^2 - \text{pen}_j \right\}$  and  $\widetilde{k} := \text{argmin}_{1 \leq j \leq K} \{\Psi_j + \text{pen}_j\}$ . Let further  $(a)_+ := (a \vee 0)$ . If  $(\text{pen}_1, \dots, \text{pen}_K)$  is non-decreasing, then we have for all  $1 \leq k \leq K$  that*

$$\|\widehat{f}_{\widetilde{k}} - f\|_\omega^2 \leq 7 \text{pen}_k + 78 \text{bias}_k^2 + 42 \max_{1 \leq j \leq K} \left( \|\widehat{f}_j - f_j\|_\omega^2 - \frac{1}{6} \text{pen}_j \right)_+, \quad (3.1)$$

where we denote by  $f_j := \sum_{k=1}^j [f]_k \varphi_k$  the projection of  $f$  on the first  $j$  basis vectors in  $H$  and by  $\text{bias}_k := \|f - f_k\|_\omega$  the bias due to the projection.

*Proof.* In view of the definition of  $\widetilde{k}$ , we have for all  $1 \leq k \leq K$  that

$$\begin{aligned} \|\widehat{f}_{\widetilde{k}} - f\|_\omega^2 &\leq 3 \left\{ \|\widehat{f}_{\widetilde{k}} - \widehat{f}_{\widetilde{k} \wedge k}\|_\omega^2 + \|\widehat{f}_{\widetilde{k} \wedge k} - \widehat{f}_k\|_\omega^2 + \|\widehat{f}_k - f\|_\omega^2 \right\} \\ &\leq 3 \left\{ \Psi_k + \text{pen}_{\widetilde{k}} + \Psi_{\widetilde{k}} + \text{pen}_k + \|\widehat{f}_k - f\|_\omega^2 \right\} \\ &\leq 6 \left\{ \Psi_k + \text{pen}_k \right\} + 3 \|\widehat{f}_k - f\|_\omega^2. \end{aligned} \quad (3.2)$$

Since  $(\text{pen}_1, \dots, \text{pen}_K)$  is non-decreasing and  $4 \text{bias}_k^2 \geq \max_{1 \leq j \leq K} \|f_k - f_j\|_\omega^2$ , we have

$$\Psi_k \leq 6 \sup_{1 \leq j \leq K} \left( \|\widehat{f}_j - f_j\|_\omega^2 - \frac{1}{6} \text{pen}_j \right)_+ + 12 \text{bias}_k^2.$$

It easily verified that for all  $1 \leq k \leq K$  we have

$$\|\widehat{f}_k - f\|_\omega^2 \leq \frac{1}{3} \text{pen}_k + 2 \text{bias}_k^2 + 2 \sup_{1 \leq j \leq K} \left( \|\widehat{f}_j - f_j\|_\omega^2 - \frac{1}{6} \text{pen}_j \right)_+.$$

The result follows combining the last estimates with (3.2).  $\square$

The Lemma being valid for any upper bound  $K$  and any monotonic sequence of penalties, we now need to specify our choice. Let us first define some auxiliary quantities needed in the construction of the penalty sequence and the upper bound  $K$ .

**Definition 3.2** *For any sequence  $\alpha := (\alpha_j)_{j \in \mathbb{N}}$ , define*

$$(i) \quad \Delta_k^\alpha := \max_{1 \leq j \leq k} \omega_j \alpha_j^{-2} \quad \text{and} \quad \delta_k^\alpha := k \Delta_k^\alpha \frac{\log(\Delta_k^\alpha \vee (k+2))}{\log(k+2)};$$

$$(ii) \quad \text{given } \omega_k^+ := \max_{1 \leq j \leq k} \omega_j, \quad N_\varepsilon^\circ := \max\{1 \leq N \leq \varepsilon^{-1} \mid \omega_N^+ \leq \varepsilon^{-1}\}, \\ \text{and } v_\sigma := (8 \log(\log(\sigma^{-1} + 20)))^{-1}, \text{ let}$$

$$N_\varepsilon^\alpha := \min \left\{ 2 \leq j \leq N_\varepsilon^\circ \mid \frac{\alpha_j^2}{j \omega_j^+} \leq \varepsilon |\log \varepsilon| \right\} - 1 \quad \text{and} \quad M_\sigma^\alpha := \min \left\{ 2 \leq j \leq \sigma^{-1} \mid \alpha_j^2 \leq \sigma^{1-v_\sigma} \right\} - 1,$$

and  $K_{\varepsilon, \sigma}^\alpha := N_\varepsilon^\alpha \wedge M_\sigma^\alpha$ . If the defining set is empty, set  $N_\varepsilon^\alpha = N_\varepsilon^\circ$  or  $M_\sigma^\alpha = \lfloor \sigma^{-1} \rfloor$ , respectively.

Let us first have a closer look at the last term on the right hand side of (3.1). To this end, let us define an upper bound  $K_{\varepsilon, \sigma}^+$  and a penalty sequence.

**Definition 3.3** *Using Definition 3.2, let  $N_\varepsilon^+ := N_\varepsilon^{\sqrt{4d\lambda}}$ ,  $M_\sigma^+ := M_\sigma^{\sqrt{4d\lambda}}$ , and  $K_{\varepsilon, \sigma}^+ := K_{\varepsilon, \sigma}^{\sqrt{4d\lambda}}$ , finally  $\text{pen}_k := 60 \delta_k^\alpha \varepsilon$ .*

The following assumption is satisfied in particular under the classical smoothness assumptions considered in the illustration.

**Assumption 3.4** Suppose that  $\sigma^{-7} \lambda_{M_\sigma^+ + 1}^{-1/2} \exp\left(-\lambda_{M_\sigma^+ + 1}/(72 \sigma d)\right) \leq C(\lambda, d)$  for all  $\sigma \in (0, 1)$ .

**Proposition 3.5** There is a constant  $C > 0$  depending only on the class  $\mathcal{A}_\lambda^d$  such that

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E} \left[ \sup_{1 \leq k \leq K_{\varepsilon, \sigma}^+} \left( \|\widehat{f}_k - f_k\|_\omega^2 - \frac{1}{6} \text{pen}_k \right)_+ \right] \leq C \left\{ \varepsilon + r \kappa_\sigma + \sigma \right\}.$$

Note that we could now define an estimator using the penalty sequence  $\text{pen}$  from Definition 3.3. Combining Lemma 3.1 and Proposition 3.5, we would even obtain an upper risk bound for this estimator. Though, as  $\text{pen}$  and the upper bound  $K_{\varepsilon, \sigma}^+$  still depend on the singular values  $(a_j)$  and the operator class  $\mathcal{A}_\lambda^d$ , respectively, this would not yield an adaptive procedure. Thus, let us define randomized versions of  $\text{pen}$  and  $K$  which exclusively depend on the observations.

**Definition 3.6** Using Definition 3.2, define the sequences  $N_\varepsilon^- := N_\varepsilon^{\sqrt{\lambda/(4d)}}$ ,  $M_\sigma^- := M_\sigma^{\sqrt{\lambda/(4d)}}$ , and  $K_{\varepsilon, \sigma}^- := K_{\varepsilon, \sigma}^{\sqrt{4/(d\lambda)}}$  which are obviously element-wise smaller than the analogous sequences from Definition 3.3. Denoting by  $X$  the sequence of random variables  $(X_j)_{j \in \mathbb{N}}$ , define further the random quantities  $\widehat{N}_\varepsilon := N_\varepsilon^X$ ,  $\widehat{M}_\sigma := M_\sigma^X$ ,  $\widehat{K}_{\varepsilon, \sigma} := K_{\varepsilon, \sigma}^X$ , and  $\widehat{\text{pen}}_k := 600 \delta_k^X \varepsilon$ .

The next proposition ensures that the randomized upper bound and penalty sequence behave similarly to their deterministic counterparts with sufficiently high probability as not to deteriorate the estimation risk. This justifies the choice of the penalty.

**Proposition 3.7** For every  $\varepsilon, \sigma \in (0, 1)$ , define the event

$$\mathcal{U}_{\varepsilon, \sigma} := \left\{ \text{pen}_k \leq \widehat{\text{pen}}_k \leq 30 \text{pen}_k \quad \forall 1 \leq k \leq K_{\varepsilon, \sigma}^+ \right\} \cap \left\{ K_{\varepsilon, \sigma}^- \leq \widehat{K}_{\varepsilon, \sigma} \leq K_{\varepsilon, \sigma}^+ \right\}.$$

Then, we have that  $\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E}[\|\widehat{f}_k - f\|_\omega^2 \mathbf{1}_{\mathcal{U}_{\varepsilon, \sigma}}] \leq C \sigma \quad \forall \varepsilon, \sigma \in (0, 1)$ , where  $C > 0$  is a constant depending only on the classes  $\mathcal{F}_\gamma^r$  and  $\mathcal{A}_\lambda^d$ .

We are finally in position to state the upper risk bound of the fully data-driven estimator of  $f$ , which is the main result of this article.

**Theorem 3.8** Consider the adaptive estimator  $\widehat{f}_{\widehat{k}}$  with  $\widehat{k}$  given in (1.3). Under Assumptions 2.1 and 3.4, there is a constant  $C$  depending only on the classes  $\mathcal{F}_\gamma^r$  and  $\mathcal{A}_\lambda^d$  such that for all  $\varepsilon, \sigma \in (0, 1)$

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 \leq C \left\{ \min_{1 \leq k \leq K_{\varepsilon, \sigma}^-} \left\{ \max(\omega_k/\gamma_k, \delta_k^\lambda \varepsilon) \right\} + \kappa_\sigma + \varepsilon + \sigma \right\}.$$

*Proof.* First, decompose the risk using the event  $\mathcal{U}_{\varepsilon, \sigma}$  defined in Proposition 3.7 as

$$\mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 = \mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 \mathbf{1}_{\mathcal{U}_{\varepsilon, \sigma}} + \mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 \mathbf{1}_{\mathcal{U}_{\varepsilon, \sigma}^c}.$$

As the random sequence  $\widehat{\text{pen}}_k$  is non-decreasing in  $k$ , we may apply Lemma 3.1 and obtain for every  $1 \leq k \leq \widehat{K}_{\varepsilon, \sigma}$

$$\|\widehat{f}_k - f\|_\omega^2 \leq 7 \widehat{\text{pen}}_k + 78 \text{bias}_k^2 + 42 \max_{1 \leq j \leq \widehat{K}_{\varepsilon, \sigma}} \left( \|\widehat{f}_j - f_j\|_\omega^2 - \frac{1}{6} \widehat{\text{pen}}_j \right)_+.$$

On the event  $\mathcal{U}_{\varepsilon, \sigma}$ , this implies that for all  $1 \leq k \leq K_{\varepsilon, \sigma}^-$

$$\mathbb{E} \|\widehat{f}_k - f\|_\omega^2 \mathbf{1}_{\mathcal{U}_{\varepsilon, \sigma}} \leq 210 \text{pen}_k + 78 \text{bias}_k^2 + 42 \max_{1 \leq j \leq K_{\varepsilon, \sigma}^+} \left( \|\widehat{f}_j - f_j\|_\omega^2 - \frac{1}{6} \text{pen}_j \right)_+.$$

Thus, using  $\delta_k^a \leq d\zeta_d \delta_k^\lambda$  with  $\zeta_d = \log(3d)/\log(3)$ ,

$$\mathbb{E}\|\widehat{f}_k - f\|_\omega^2 \mathbf{1}_{\mathcal{U}_{\varepsilon,\sigma}} \leq C(d) \min_{1 \leq k \leq K_{\varepsilon,\sigma}^-} \{\max(\omega_k/\gamma_k, \delta_k^\lambda \varepsilon)\} + 42 \max_{1 \leq j \leq K_{\varepsilon,\sigma}^+} \left( \|\widehat{f}_j - f_j\|_\omega^2 - \frac{1}{6} \text{pen}_j \right)_+.$$

It remains to apply Propositions 3.5 and 3.7 to conclude.  $\square$

A comparison with the lower bound from Theorem 2.3 shows that this upper bound ensures minimax optimality of  $\widehat{f}_k$  only if

$$\chi_{\varepsilon,\sigma}^\diamond := \min_{1 \leq k \leq K_{\varepsilon,\sigma}^\lambda} \left[ \max\left(\frac{\omega_k}{\gamma_k}, \delta_k^\lambda \varepsilon\right) \right]$$

is at most of the same order as  $\max(\chi_\varepsilon, \kappa_\sigma)$ , whence the following corollary.

**Corollary 3.9** *Under Assumption 2.1 and if  $\sup_{\varepsilon,\sigma \in (0,1)} \{\chi_{\varepsilon,\sigma}^\diamond / \max(\chi_\varepsilon, \kappa_\sigma)\} < \infty$ , we have*

$$\mathcal{R}_\omega(\widehat{f}_k, \mathcal{F}_\gamma^r, \mathcal{A}_\lambda^d) \leq C \mathcal{R}_\omega^*(\mathcal{F}_\gamma^r, \mathcal{A}_\lambda^d) \quad \forall \varepsilon, \sigma \in (0,1).$$

We conclude this article reconsidering the framework of the preceding Illustration 2.6. Notice that the adaptive estimator is minimax optimal over a wide range of cases, even when  $\varepsilon < \sigma$ .

**Illustration 3.10** Let  $(\sigma_\varepsilon)_{\varepsilon \in (0,1)}$  be a noise level in  $X$  depending on the noise level  $\varepsilon$  in  $Y$  and suppose that the limits  $q_1$  and  $q_2$  from Illustration 2.6 exist in the respective cases. Some straightforward computations then show that the adaptive estimator attains the following rates of convergence.

**[m]** If  $p - s > b$ , the adaptive estimator  $\widehat{f}_k^{(s)}$  attains the optimal rates (cf. Illustration 2.6). In case  $p - s \leq b$ , we have, supposing that  $q_1^v := \lim_{\varepsilon \rightarrow 0} \varepsilon^{-2b/(2p+2b+1)} \sigma_\varepsilon^{1-v\sigma_\varepsilon}$  exists,

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E}\|\widehat{f}_k^{(s)} - f^{(s)}\|_{L^2}^2 = \begin{cases} O(\varepsilon^{2(p-s)/(2p+2b+1)}) & \text{if } q_1 < \infty \text{ and } q_1^v < \infty, \\ O(\sigma_\varepsilon^{(p-s)/b} \sigma_\varepsilon^{-v\sigma_\varepsilon}) & \text{otherwise.} \end{cases}$$

**[s]** The adaptive estimator attains the optimal rates.

## A. Proofs

### A.1. Minimax theory (Section 2)

#### Lower risk bound

*Proof of Theorem 2.3.* The proof consists of two steps: (A) First, we show that  $\chi_\varepsilon$  yields a lower risk bound in the case where the eigenvalues  $(a_j)$  of the operator  $A$  are known. (B) Then, we show that another lower risk bound is given by  $\kappa_\sigma$ .

*Step (A).* Given  $\zeta := \eta \min(r, 1/(2d))$  and  $\alpha_\varepsilon := \chi_\varepsilon (\sum_{j=1}^{k_\varepsilon^*} \varepsilon \omega_j / \lambda_j)^{-1}$  we consider the function  $f := (\varepsilon \zeta \alpha_\varepsilon)^{1/2} \sum_{j=1}^{k_\varepsilon^*} \lambda_j^{-1/2} \varphi_j$ . We are going to show that for any  $\theta := (\theta_j) \in \{-1, 1\}^{k_\varepsilon^*}$ , the function  $f_\theta := \sum_{j=1}^{k_\varepsilon^*} \theta_j [f]_j \varphi_j$  belongs to  $\mathcal{F}_\gamma^r$  and is hence a possible candidate for the solution.

For a fixed  $\theta$  and under the hypothesis that the solution is  $f_\theta$ , the observation  $Y_k$  is distributed according to  $\mathcal{N}(a_k [f_\theta]_k, \varepsilon)$  for any  $k \in \mathbb{N}$ . We denote by  $\mathbb{P}_\theta$  the distribution of the resulting sequence  $\{Y_k\}$  and by  $\mathbb{E}_\theta$  the expectation with respect to this distribution.

Furthermore, for  $1 \leq j \leq k_\varepsilon^*$  and each  $\theta$ , we introduce  $\theta^{(j)}$  by  $\theta_l^{(j)} = \theta_l$  for  $j \neq l$  and  $\theta_j^{(j)} = -\theta_j$ . The key argument of this proof is the following reduction scheme. If  $\tilde{f}$  denotes an estimator of  $f$  then we



conclude

$$\begin{aligned}
\sup_{f \in \mathcal{F}_\gamma^r} \mathbb{E} \|\tilde{f} - f\|_\omega^2 &\geq \sup_{\theta \in \{-1,1\}^{k_\varepsilon^*}} \mathbb{E}_\theta \|\tilde{f} - f_\theta\|_\omega^2 \geq \frac{1}{2^{k_\varepsilon^*}} \sum_{\theta \in \{-1,1\}^{2k_\varepsilon^*}} \mathbb{E}_\theta \|\tilde{f} - f_\theta\|_\omega^2 \\
&\geq \frac{1}{2^{k_\varepsilon^*}} \sum_{\theta \in \{-1,1\}^{k_\varepsilon^*}} \sum_{j=1}^{k_\varepsilon^*} \omega_j \mathbb{E}_\theta |[\tilde{f} - f_\theta]_j|^2 \\
&= \frac{1}{2^{k_\varepsilon^*}} \sum_{\theta \in \{-1,1\}^{k_\varepsilon^*}} \sum_{j=1}^{k_\varepsilon^*} \frac{\omega_j}{2} \left\{ \mathbb{E}_\theta |[\tilde{f} - f_\theta]_j|^2 + \mathbb{E}_{\theta^{(j)}} |[\tilde{f} - f_{\theta^{(j)}}]_j|^2 \right\}.
\end{aligned} \tag{A.1}$$

Below we show furthermore that for all  $\varepsilon \in (0, 1)$  we have

$$\left\{ \mathbb{E}_\theta |[\tilde{f} - f_\theta]_j|^2 + \mathbb{E}_{\theta^{(j)}} |[\tilde{f} - f_{\theta^{(j)}}]_j|^2 \right\} \geq \frac{\varepsilon \zeta \alpha_\varepsilon}{2\lambda_j}. \tag{A.2}$$

Combining the last lower bound and the reduction scheme gives

$$\sup_{f \in \mathcal{F}_\gamma^r} \mathbb{E} \|\tilde{f} - f\|_\omega^2 \geq \frac{1}{2^{k_\varepsilon^*}} \sum_{\theta \in \{-1,1\}^{k_\varepsilon^*}} \sum_{j=1}^{k_\varepsilon^*} \frac{\omega_j}{2} \frac{\varepsilon \zeta \alpha_\varepsilon}{2\lambda_j} = \frac{\zeta \alpha_\varepsilon}{4} \sum_{j=1}^{k_\varepsilon^*} \frac{\varepsilon \omega_j}{\lambda_j} = \frac{\zeta \chi_\varepsilon}{4},$$

which implies the lower bound given in the theorem by definition of  $\zeta$ .

To complete the proof, it remains to check (A.2) and  $f_\theta \in \mathcal{F}_\gamma^r$  for all  $\theta \in \{-1, 1\}^{k_\varepsilon^*}$ . The latter is easily verified if  $f \in \mathcal{F}_\gamma^r$ , which can be seen recalling that  $\omega/\gamma$  is non-increasing and noticing that the definitions of  $\zeta$ ,  $\alpha_\varepsilon$  and  $\eta$  imply  $\|f\|_\gamma^2 \leq \zeta \frac{\gamma k_\varepsilon^*}{\omega k_\varepsilon^*} \alpha_\varepsilon \left( \sum_{j=1}^{k_\varepsilon^*} \frac{\varepsilon \omega_j}{\lambda_j} \right) \leq \zeta/\eta \leq r$ .

It remains to show (A.2). Consider the Hellinger affinity  $\rho(\mathbb{P}_1, \mathbb{P}_{-1}) = \int \sqrt{d\mathbb{P}_1 d\mathbb{P}_{-1}}$ , then we obtain for any estimator  $\tilde{f}$  of  $f$  that

$$\begin{aligned}
\rho(\mathbb{P}_1, \mathbb{P}_{-1}) &\leq \int \frac{|[\tilde{f} - f_{\theta^{(j)}}]_j|}{|[f_\theta - f_{\theta^{(j)}}]_j|} \sqrt{d\mathbb{P}_1 d\mathbb{P}_{-1}} + \int \frac{|[\tilde{f} - f_\theta]_j|}{|[f_\theta - f_{\theta^{(j)}}]_j|} \sqrt{d\mathbb{P}_1 d\mathbb{P}_{-1}} \\
&\leq \left( \int \frac{|[\tilde{f} - f_{\theta^{(j)}}]_j|^2}{|[f_\theta - f_{\theta^{(j)}}]_j|^2} d\mathbb{P}_1 \right)^{1/2} + \left( \int \frac{|[\tilde{f} - f_\theta]_j|^2}{|[f_\theta - f_{\theta^{(j)}}]_j|^2} d\mathbb{P}_{-1} \right)^{1/2}.
\end{aligned}$$

Rewriting the last estimate we obtain

$$\left\{ \mathbb{E}_\theta |[\tilde{f} - f_\theta]_j|^2 + \mathbb{E}_{\theta^{(j)}} |[\tilde{f} - f_{\theta^{(j)}}]_j|^2 \right\} \geq \frac{1}{2} |[f_\theta - f_{\theta^{(j)}}]_j|^2 \rho^2(\mathbb{P}_1, \mathbb{P}_{-1}). \tag{A.3}$$

Next, we bound the Hellinger affinity  $\rho(\mathbb{P}_1, \mathbb{P}_{-1})$  from below. Consider the Kullback-Leibler divergence of these two distributions first. The components of the two sequences corresponding to the distributions  $\mathbb{P}_1$  and  $\mathbb{P}_{-1}$  are pairwise equally distributed except for the  $j$ -th component. Thus, we have  $\log(d\mathbb{P}_\theta/d\mathbb{P}_{\theta^{(j)}}) = (2y_j a_j \theta_j [f]_j / \varepsilon)$ , and taking the integral over  $y_j$  with respect to  $\mathbb{P}_\theta$ , we find

$$KL(\mathbb{P}_1, \mathbb{P}_{-1}) = \frac{2}{\varepsilon} a_j^2 [f]_j^2 \leq \frac{2d}{\varepsilon} [f]_j^2 \lambda_j = 2d\zeta \alpha_\varepsilon \leq 1,$$

Using the well-known relationship  $\rho(\mathbb{P}_1, \mathbb{P}_{-1}) \geq 1 - (1/2)KL(\mathbb{P}_1, \mathbb{P}_{-1})$  between the Kullback-Leibler divergence and the Hellinger affinity, we obtain that  $\rho(\mathbb{P}_1, \mathbb{P}_{-1}) \geq 1/2$ . Using this estimate, (A.3) becomes  $\left\{ \mathbb{E}_\theta |[\tilde{f} - f_\theta]_j|^2 + \mathbb{E}_{\theta^{(j)}} |[\tilde{f} - f_{\theta^{(j)}}]_j|^2 \right\} \geq \frac{1}{2} [f]_j^2$ , and combining this with (A.1) implies the result by construction of the solution  $f$ .

*Step (B).* First, we construct two solutions  $f_\theta \in \mathcal{F}_\gamma^r$  and operators  $A_\theta \in \mathcal{A}_\lambda^d$  (with  $\theta \in \{-1, 1\}$ ) such that the resulting images  $g_\theta$  satisfy  $g_{-1} = g_1$ . To this end, we define  $k_\sigma^* := \operatorname{argmax}_{j \in \mathbb{N}} \{\omega_j \gamma_j^{-1} \min(1, \sigma \lambda_j^{-1})\}$  and  $\alpha_\sigma := \zeta \min(1, \sigma^{1/2} \lambda_{k_\sigma^*}^{-1/2})$  with  $\zeta := \min(2^{-1}, (1 - d^{-1/2}))$ . Observe that  $1 \geq (1 - \alpha_\sigma)^2 \geq (1 - (1 - 1/d^{1/2}))^2 \geq 1/d$  and  $1 \leq (1 + \alpha_\sigma)^2 \leq (1 + (1 - 1/d^{1/2}))^2 = (2 - 1/d^{1/2})^2 \leq d$ , which implies  $1/d \leq (1 + \alpha_\sigma)^2 \leq d$ . These inequalities will be used below without further reference. We show below that for each  $\theta$  the function  $f_\theta := (1 - \theta \alpha_\sigma) \frac{r}{d} \gamma_{k_\sigma^*}^{-1/2} \varphi_{k_\sigma^*}$  belongs to  $\mathcal{F}_\gamma^r$  and that the operator  $A_\theta$  with the singular values  $a_k^\theta = [1 + \theta \alpha_\sigma \mathbf{1}\{k = k_\sigma^*\}] \sqrt{\lambda_k}$  is an element of  $\mathcal{A}_\lambda^d$ . We obviously have that  $A_1 f_f = (1 - \alpha_\sigma^2) (\lambda_{k_\sigma^*} / \gamma_{k_\sigma^*})^{1/2} (r/d) \psi_{k_\sigma^*} = A_{-1} f_{-1}$ . For  $\theta \in \{\pm 1\}$ , denote by  $\mathbb{P}_\theta$  the joint distribution of the two sequences  $(X_1, X_2, \dots)$  and  $(Y_1, Y_2, \dots)$ , and let  $\mathbb{E}_\theta$  denote the expectation with respect to  $\mathbb{P}_\theta$ .

Applying a reduction scheme as under Step (A) above, we deduce that for each estimator  $\tilde{f}$  of  $f$

$$\sup_{f \in \mathcal{F}_\gamma^r} \sup_{A \in \mathcal{A}_\lambda^d} \mathbb{E} \|\tilde{f} - f\|_\omega^2 \geq \max_{\theta \in \{-1, 1\}} \mathbb{E}_\theta \|\tilde{f} - f_\theta\|_\omega^2 \geq \frac{1}{2} \left\{ \mathbb{E}_1 \|\tilde{f} - f_1\|_\omega^2 + \mathbb{E}_{-1} \|\tilde{f} - f_{-1}\|_\omega^2 \right\}.$$

Below we show furthermore that

$$\mathbb{E}_1 \|\tilde{f} - f_1\|_\omega^2 + \mathbb{E}_{-1} \|\tilde{f} - f_{-1}\|_\omega^2 \geq \frac{1}{8} \|f_1 - f_{-1}\|_\omega^2. \quad (\text{A.4})$$

Moreover, we have  $\|f_1 - f_{-1}\|_\omega^2 = 4\alpha_\sigma^2 (r/d) \omega_{k_\sigma^*} \gamma_{k_\sigma^*}^{-1} = 4\zeta^2 (r/d) \omega_{k_\sigma^*} \gamma_{k_\sigma^*}^{-1} \min\left(1, \frac{\sigma}{\lambda_{k_\sigma^*}}\right)$ . Combining the last lower bound with the reduction scheme and the definition of  $k_\sigma^*$  implies the result of the theorem. To conclude the proof, it remains to check (A.4),  $f_\theta \in \mathcal{F}_\gamma^r$  and  $A_\theta \in \mathcal{A}_\lambda^d$  for both  $\theta$ . In order to show  $f_\theta \in \mathcal{F}_\gamma^r$  observe that  $\|f_\theta\|_\gamma^2 = \gamma_{k_\sigma^*} |[f_\theta]_{k_\sigma^*}|^2 \leq \gamma_{k_\sigma^*} |(1 - \theta \alpha_\sigma) (r/d) \gamma_{k_\sigma^*}^{-1/2}|^2 \leq r$ .

To check that  $A_\theta \in \mathcal{A}_\lambda^d$ , it remains to show that  $1/d \leq (a_j^\theta)^2 / \lambda_j \leq d$  for all  $j \geq 1$ . These inequalities are obviously satisfied for all  $j \neq k_\sigma^*$ , and as well for  $j = k_\sigma^*$  by construction of the operator  $A$ . Finally consider (A.4). As in Step (A) above by employing the Hellinger affinity  $\rho(\mathbb{P}_1, \mathbb{P}_{-1})$  we obtain for any estimator  $\tilde{f}$  of  $f$  that

$$\mathbb{E}_1 \|\tilde{f} - f_1\|_\omega^2 + \mathbb{E}_{-1} \|\tilde{f} - f_{-1}\|_\omega^2 \geq \frac{1}{2} \|f_1 - f_{-1}\|_\omega^2 \rho^2(\mathbb{P}_1, \mathbb{P}_{-1}).$$

Next, we bound the Hellinger affinity  $\rho(\mathbb{P}_1, \mathbb{P}_{-1})$  from below for all  $\sigma \in (0, 1)$ , which proves (A.4). Notice that by construction of  $f_\theta$  and  $A_\theta$ , the distribution of  $X_i$  and  $Y_i$  does not depend on  $\theta$ , except for  $X_{k_\sigma^*}^\theta$ . It is thus easily seen that the Kullback-Leibler divergence can be controlled as follows,

$$KL(\mathbb{P}_1, \mathbb{P}_{-1}) = \frac{(a_{k_\sigma^*}^1 - a_{k_\sigma^*}^{-1})^2}{2\sigma} = \frac{2\alpha_\sigma^2}{\sigma} \lambda_{k_\sigma^*} \leq 1$$

Using  $\rho(\mathbb{P}_1, \mathbb{P}_{-1}) \geq 1 - (1/2)KL(\mathbb{P}_1, \mathbb{P}_{-1})$  again, (A.4) is shown and so is the theorem.  $\square$

## Upper risk bound

The following proof uses Lemma A.1 from the auxiliary results section A.3 below.

*Proof of Theorem 2.5.* Define  $\tilde{f} := \sum_{j=1}^{k_\sigma^*} [f]_j \mathbf{1}\{X_j^2 \geq \sigma\} e_j$  and decompose the risk into two terms,

$$\mathbb{E} \|\hat{f} - f\|_\omega^2 = \mathbb{E} \|\hat{f} - \tilde{f}\|_\omega^2 + \mathbb{E} \|\tilde{f} - f\|_\omega^2 =: A + B, \quad (\text{A.5})$$

which we bound separately. Consider first  $A$  which we decompose further,

$$\begin{aligned} \mathbb{E}\|\widehat{f} - \widetilde{f}\|_\omega^2 &= \sum_{j=1}^{k_\varepsilon^*} \omega_j \mathbb{E} \left[ \frac{(Y_j - \mathbb{E}Y_j)^2}{X_j^2} \mathbf{1}\{X_j^2 \geq \sigma\} \right] \\ &\quad + \sum_{j=1}^{k_\varepsilon^*} \omega_j |[f]_j|^2 \mathbb{E} \left[ \frac{(X_j - \mathbb{E}X_j)^2}{X_j^2} \mathbf{1}\{X_j^2 \geq \sigma\} \right] =: A_1 + A_2. \end{aligned}$$

As far as  $A_1$  is considered, we use Lemma A.1 (iii) from Section A.3 below and write

$$A_1 = \sum_{j=1}^{k_\varepsilon^*} \frac{\omega_j \varepsilon}{\mathbb{E}[X_j]^2} \mathbb{E} \left[ \left( \frac{\mathbb{E}[X_j]}{X_j} \right)^2 \mathbf{1}\{X_j^2 \geq \sigma\} \right] \leq 4d \sum_{j=1}^{k_\varepsilon^*} \frac{\omega_j \varepsilon}{\lambda_j} \leq 4d\chi_\varepsilon.$$

As for  $A_2$ , we apply Lemma A.1 (i) and obtain

$$A_2 \leq 8d \sum_{j=1}^{k_\varepsilon^*} \omega_j |[f]_j|^2 \min \left( 1, \frac{\sigma}{\lambda_j} \right) \leq 8d\kappa_\sigma$$

Consider now  $B$  which we decompose further into

$$\begin{aligned} \mathbb{E}\|\widetilde{f} - f\|_\omega^2 &= \sum_{j \in \mathbb{N}} \omega_j |[f]_j|^2 \mathbb{E}[(1 - \mathbf{1}\{1 \leq j \leq k_\varepsilon^*\} \mathbf{1}\{X_j^2 \geq \sigma\})^2] \\ &= \sum_{j > k_\varepsilon^*} \omega_j |[f]_j|^2 + \sum_{j=1}^{k_\varepsilon^*} \omega_j |[f]_j|^2 \mathbf{P}(X_j^2 < \sigma) =: B_1 + B_2, \end{aligned}$$

where  $B_1 \leq \|f\|_\gamma^2 \omega_{k_\varepsilon^*} \gamma_{k_\varepsilon^*}^{-1} \leq r\chi_\varepsilon$  because  $f \in \mathcal{F}_\gamma$ . Moreover,  $B_2 \leq 4dr\kappa_\sigma$  using Lemma A.1 (ii). The result of the theorem follows now by combination of the decomposition (A.5) and the estimates of  $A_1, A_2, B_1$  and  $B_2$ .  $\square$

## A.2. Adaptive estimation (Section 3)

The proofs in this section use the Lemmas A.3– A.6 from the auxiliary results section A.3 below.

*Proof of Proposition 3.5.* Using the model equation  $Y_j = [g]_j + \sqrt{\varepsilon} \xi_j$ , we have for all  $t \in \mathcal{S}_k$  that

$$[\widehat{f}_k - f_k]_j = \frac{\sqrt{\varepsilon} \xi_j}{a_j} + \left( \frac{1}{X_j} \mathbf{1}_{[X_j^2 \geq \sigma]} - \frac{1}{a_j} \right) \sqrt{\varepsilon} \xi_j + \left( \frac{1}{X_j} \mathbf{1}_{[X_j^2 \geq \sigma]} - \frac{1}{a_j} \right) [g]_j.$$

Thus, we may decompose the norm  $\|\widehat{f}_k - f_k\|_\omega^2$  in three terms according to

$$\begin{aligned} \|\widehat{f}_k - f_k\|_\omega^2 &\leq 3 \sum_{j=1}^k \frac{\omega_j}{a_j} \varepsilon \xi_j^2 + 3 \sum_{j=1}^k \omega_j \left( \frac{1}{X_j} \mathbf{1}_{[X_j^2 \geq \sigma]} - \frac{1}{a_j} \right)^2 \varepsilon \xi_j^2 + 3 \sum_{j=1}^k \omega_j \left( \frac{1}{X_j} \mathbf{1}_{[X_j^2 \geq \sigma]} - \frac{1}{a_j} \right)^2 [g]_j^2 \\ &=: 3 \{T_k^{(1)} + T_k^{(2)} + T_k^{(3)}\}. \end{aligned}$$

Define the event

$$\Omega_\sigma := \left\{ \forall 0 < j \leq M_\sigma^+ \left| \frac{1}{X_j} - \frac{1}{a_j} \right| \leq \frac{1}{2a_j} \wedge X_j^2 \geq \sigma \right\}.$$

Since  $\mathbf{1}\{X_j^2 \geq \sigma\} \mathbf{1}\{\Omega_\sigma\} = \mathbf{1}\{\Omega_\sigma\}$ , it follows that for all  $1 \leq j \leq K_{\varepsilon, \sigma}^+$  we have

$$\left( \frac{a_j}{X_j} \mathbf{1}\{X_j^2 \geq \sigma\} - 1 \right)^2 \mathbf{1}\{\Omega_\sigma\} = a_j^2 \mathbf{1}\{\Omega_\sigma\} \left| \frac{1}{X_j} - \frac{1}{a_j} \right|^2 \leq \frac{1}{4}.$$

Hence,  $T_k^{(2)} \mathbf{1}_{\Omega_\sigma} \leq \frac{1}{4} T_k^{(1)}$  for all  $1 \leq k \leq K_{\varepsilon, \sigma}^+$ , and thus

$$\begin{aligned} \sup_{1 \leq k \leq K_{\varepsilon, \sigma}^+} \left( \|\widehat{f}_k - f_k\|_\omega^2 - \frac{1}{6} \text{pen}_k \right)_+ &\leq 4 \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \left( \sum_{j=1}^k \frac{\omega_j}{a_j} \varepsilon \xi_j^2 - 2\delta_k \varepsilon \right)_+ \\ &\quad + 3 \sup_{1 \leq k \leq K_{\varepsilon, \sigma}^+} T_k^{(2)} \mathbf{1}_{\Omega_\sigma^c} + 3 \sup_{1 \leq k \leq K_{\varepsilon, \sigma}^+} T_k^{(3)}. \end{aligned}$$

Note that  $\mathbf{P}[\Omega_\sigma^c] \leq C(d)\sigma^2$  by virtue of Lemma A.6. The result follows immediately using Lemmas A.3, A.4, and A.5 below.  $\square$

*Proof of Proposition 3.7.* Let  $\check{f}_k := \sum_{1 \leq j \leq k} [f]_j \mathbf{1}\{X_j^2 \geq \sigma\} e_j$ . It is easy to see that  $\|\widehat{f}_k - \check{f}_k\|^2 \leq \|\widehat{f}_{k'} - \check{f}_{k'}\|^2$  for all  $k' \leq k$  and  $\|\check{f}_k - f\|^2 \leq \|f\|^2$  for all  $k \geq 1$ . Thus, using that  $1 \leq \widehat{k} \leq (N_\varepsilon^\circ \wedge \sigma^{-1})$ , we can write

$$\begin{aligned} \mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} &\leq 2 \left\{ \mathbb{E} \|\widehat{f}_{\widehat{k}} - \check{f}_{\widehat{k}}\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} + \mathbb{E} \|\check{f}_{\widehat{k}} - f\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} \right\} \\ &\leq 2 \left\{ \mathbb{E} \|\widehat{f}_{(N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)} - \check{f}_{(N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)}\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} + \|f\|_\omega^2 \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c] \right\}. \end{aligned}$$

Moreover, using the Cauchy-Schwarz inequality, we conclude

$$\begin{aligned} &\mathbb{E} \|\widehat{f}_{(N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)} - \check{f}_{(N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)}\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} \\ &\leq 2\sigma^{-1} \sum_{1 \leq j \leq (N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)} \omega_j \left\{ \mathbb{E} (Y_j - a_j [f]_j)^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} + \mathbb{E} (a_j [f]_j - X_j [f]_j)^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} \right\} \\ &\leq 2\sigma^{-1} \left\{ \sum_{1 \leq j \leq (N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)} \omega_j \left[ \mathbb{E} (Y_j - [g]_j)^4 \right]^{1/2} \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c]^{1/2} \right. \\ &\quad \left. + \sum_{1 \leq j \leq (N_\varepsilon^\circ \wedge \lfloor \sigma^{-1} \rfloor)} \omega_j [f]_j^2 \left[ \mathbb{E} (X_j - a_j)^4 \right]^{1/2} \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c]^{1/2} \right\} \\ &\leq 2\sqrt{3}\sigma^{-1} \left\{ (\sigma^{-1} \max_{1 \leq j \leq N_\varepsilon^\circ} \omega_j) \varepsilon + \sigma \|f\|_\omega^2 \right\} \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c]^{1/2}, \end{aligned}$$

which implies

$$\mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} \leq C \left\{ (\sigma^{-2} + \|f\|_\omega^2) \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c]^{1/2} + \|f\|_\omega^2 \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c] \right\}.$$

Lemma A.6 below yields, for some  $C(d) > 0$  depending only on  $d$ ,

$$\mathbb{E} \|\widehat{f}_{\widehat{k}} - f\|_\omega^2 \mathbf{1}\{\mathcal{U}_{\varepsilon, \sigma}^c\} \leq C(d) \left\{ \sigma + \|f\|_\omega^2 \sigma^6 + \|f\|_\omega^2 \sigma^{12} \right\}$$

which completes the proof due to  $f \in \mathcal{F}_\gamma^r$ .  $\square$

### A.3. Auxiliary results

**Lemma A.1** For every  $j \in \mathbb{N}$ ,

$$(i) R_j^I := \mathbb{E} \left[ \left( \frac{a_j}{X_j} - 1 \right)^2 \mathbf{1}\{X_j^2 \geq \sigma\} \right] \leq \min \left\{ 1, \frac{8\sigma}{a_j^2} \right\}$$

$$(ii) R_j^{II} := \mathbf{P}[X_j^2 < \sigma] \leq \min \left\{ 1, \frac{4\sigma}{a_j^2} \right\}$$

$$(iii) \mathbb{E} \left[ \left( \frac{\mathbb{E}[X_j]}{X_j} \right)^2 \mathbf{1}\{X_j^2 \geq \sigma\} \right] \leq 4$$

*Proof.* (i) It is easy to see that

$$R_j^I = \mathbb{E} \left[ \frac{|X_j - a_j|^2}{X_j^2} \mathbf{1}\{X_j^2 \geq \sigma\} \right] \leq \sigma^{-1} \mathbb{V}\text{ar}(X_j) = 1. \quad (\text{A.6})$$

On the other hand, using that  $\mathbb{E}[(X_j - a_j)^4] = 3\sigma^2$ , we obtain

$$\begin{aligned} R_j^I &\leq \mathbb{E} \left[ \frac{(X_j - a_j)^2}{X_j^2} \mathbf{1}\{X_j^2 \geq \sigma\} 2 \left\{ \frac{(X_j - a_j)^2}{a_j^2} + \frac{X_j^2}{a_j^2} \right\} \right] \\ &\leq \frac{2 \mathbb{E}[(X_j - a_j)^4]}{\sigma a_j^2} + \frac{2 \mathbb{V}\text{ar}(X_j)}{a_j^2} = \frac{8\sigma}{a_j^2}. \end{aligned}$$

Combining with (A.6) gives  $R_j^I \leq \min \left\{ 1, \frac{8\sigma}{a_j^2} \right\}$ , which completes the proof of (i).

(ii) Trivially,  $R_j^{II} \leq 1$ . If  $1 \leq 4\sigma/a_j^2$ , then obviously  $R_j^{II} \leq \min \left\{ 1, \frac{4\sigma}{a_j^2} \right\}$ . Otherwise, we have  $\sigma < a_j^2/4$  and hence, using Tchebychev's inequality,

$$R_j^{II} \leq \mathbf{P}[|X_j - a_j| > |a_j|/2] \leq \frac{4 \mathbb{V}\text{ar}(X_j)}{a_j^2} \leq \min \left\{ 1, \frac{4\sigma}{a_j^2} \right\},$$

where we have used that  $\mathbb{V}\text{ar}(X_j) = \sigma$  for all  $j$ .

$$(iii) \mathbb{E} \left[ \left( \frac{\mathbb{E}[X_j]}{X_j} \right)^2 \mathbf{1}\{X_j^2 \geq \sigma\} \right] \leq 2 \mathbb{E} \left[ \left( \frac{X_j - \mathbb{E}[X_j]}{X_j} \right)^2 \mathbf{1}\{X_j^2 \geq \sigma\} + \mathbf{1}\{X_j^2 < \sigma\} \right] \leq 4. \quad \square$$

**Lemma A.2** Under Assumption 2.1, we have that

$$(i) \varepsilon \delta_{N_\varepsilon^+} \leq 32 d^2 \text{ for all } \varepsilon \in (0, 1),$$

and for  $\sigma^{-1} \geq \exp(512 \log(3d)^2)$  that

$$(ii) \min_{1 \leq j \leq M_\sigma^+} a_j^2 \geq 2\sigma.$$

*Proof.* (i) For  $N_\varepsilon^+ = 0$ , we have  $\delta_{N_\varepsilon^+} = 0$  and there is nothing to show. If  $0 < N_\varepsilon^+ \leq n$ , one can show that  $\omega_{N_\varepsilon^+}^+ / \lambda_{N_\varepsilon^+} \leq 4d / (\varepsilon N_\varepsilon^+ |\log \varepsilon|)$ , which we use in the following computation:

$$\begin{aligned} \delta_{N_\varepsilon^+} &= N_\varepsilon^+ \frac{\omega_{N_\varepsilon^+}^+}{\lambda_{N_\varepsilon^+}} \frac{\log((\omega_{N_\varepsilon^+}^+ / \lambda_{N_\varepsilon^+}) \vee (N_\varepsilon^+ + 2))}{\log(N_\varepsilon^+ + 2)} \leq \frac{4d}{\varepsilon |\log \varepsilon|} \frac{\log \left( \frac{4d}{N_\varepsilon^+ \varepsilon |\log \varepsilon|} \vee (N_\varepsilon^+ + 2) \right)}{\log(N_\varepsilon^+ + 2)} \\ &\leq \varepsilon^{-1} \begin{cases} 4d & (\log(\varepsilon^{-1} + 2) \geq 4d) \\ 4d(4d + \log(4d)) / (\log(\varepsilon^{-1} + 2)) & (\text{otherwise}), \end{cases} \end{aligned}$$

which implies  $\varepsilon \delta_{N_\varepsilon^+} \leq 4d(4d + \log(4d)) \leq 32d^2$  for all  $\varepsilon \in (0, 1)$ . (iii) We have that

$$\min_{1 \leq j \leq M_\sigma^+} a_j^2 \geq \min_{1 \leq j \leq M_\sigma^+} \frac{\lambda_j}{d} \geq \frac{\sigma^{1-v_\sigma}}{4d^2} \geq 2\sigma,$$

where the last step holds for  $\sigma^{-1} \geq \exp(128 \log(8d^2)^2)$  as some algebra shows.  $\square$

**Lemma A.3** *We have that*

$$\sum_{k=1}^{K_{\varepsilon, \sigma}^+} \mathbb{E} \left( \sum_{j=1}^k \frac{\omega_j}{a_j} \varepsilon \xi_j^2 - 2\delta_k^a \varepsilon \right)_+ \leq 6720 \varepsilon.$$

*Proof.* Representing the expectation of the positive random variable by the integral over its tail probabilities and using  $\delta_k^a \geq \sum_{j=1}^k (\omega_j/a_j^2)$ , we may write

$$\begin{aligned} \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \mathbb{E} \left( \sum_{j=1}^k \frac{\omega_j}{a_j} \varepsilon \xi_j^2 - 2\delta_k^a \varepsilon \right)_+ &\leq \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \int_0^\infty \mathbf{P} \left[ \sum_{j=1}^k \frac{\varepsilon \omega_j}{a_j^2} (\xi_j^2 - 1) \geq x + 2\varepsilon \delta_k^a - \varepsilon \sum_{j=1}^k \frac{\omega_j}{a_j^2} \right] dx \\ &\leq \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \int_0^\infty \mathbf{P} \left[ \sum_{j=1}^k \frac{\varepsilon \omega_j}{a_j^2} (\xi_j^2 - 1) \geq x + \varepsilon \delta_k^a \right] dx \end{aligned}$$

Define  $\rho_k := (\varepsilon \omega_k)/a_k^2$ ,  $H_k := 4\varepsilon \Delta_k^a$ , and  $B_k := 2\varepsilon^2 \sum_{j=1}^k \omega_j^2/a_j^4$ . It can be shown (see proof of Proposition A.1 in Dahlhaus and Polonik (2006)) that for all  $1 \leq k' \leq k$  and  $m \geq 2$ , we have

$$\left| \mathbb{E} \left[ \left( \frac{\varepsilon \omega_{k'}}{a_{k'}^2} (\xi_{k'}^2 - 1) \right)^m \right] \right| \leq m! \rho_{k'}^2 H_k^{m-2}.$$

Hence, the assumption of Theorem 2.8 from Petrov (1995) is satisfied and splitting up the integral, get the following bound:

$$\begin{aligned} \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \mathbb{E} \left( \sum_{j=1}^k \frac{\omega_j}{a_j} \varepsilon \xi_j^2 - 2\delta_k^a \varepsilon \right)_+ &\leq \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \int_0^{B_k/H_k - \varepsilon \delta_k^a} \exp \left( -\frac{(x + \varepsilon \delta_k^a)^2}{4B_k} \right) dx + \int_{B_k/H_k - \varepsilon \delta_k^a}^\infty \exp \left( -\frac{x + \varepsilon \delta_k^a}{4H_k} \right) dx \end{aligned}$$

The second integral is equal to  $4H_k \exp(-B_k/(4H_k^2))$ . Some computation shows that the first one is bounded from above by  $4H_k [\exp(-\varepsilon^2(\delta_k^a)^2/(4B_k)) - \exp(-B_k/(4H_k^2))]$ . Thus, the two identical terms cancel, and we get

$$\sum_{k=1}^{K_{\varepsilon, \sigma}^+} \mathbb{E} \left( \sum_{j=1}^k \frac{\omega_j}{a_j} \varepsilon \xi_j^2 - 2\delta_k^a \varepsilon \right)_+ \leq 16 \varepsilon \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \Delta_k^a \exp \left( -\frac{(\delta_k^a)^2}{8k(\Delta_k^a)^2} \right).$$

To complete the proof, we bound the sum on the right hand side as follows,

$$\begin{aligned} \sum_{k=1}^{K_{\varepsilon, \sigma}^+} \Delta_k^a \exp \left( -\frac{(\delta_k^a)^2}{8k(\Delta_k^a)^2} \right) &\leq \sum_{k=1}^\infty \exp \left( -\log(\Delta_k^a \vee (k+2)) \left[ \frac{k}{8 \log(k+2)} - 1 \right] \right) \\ &\leq e \sum_{k=1}^\infty \exp \left( -\frac{k}{8 \log(k+2)} \right) \leq e \sum_{k=1}^\infty \exp \left( -\frac{\sqrt{k}}{8 \log(3)} \right) \\ &\leq e \int_0^\infty \exp \left( -\frac{\sqrt{x}}{8 \log(3)} \right) dx = 128 \log^2(3) e, \end{aligned}$$

where we have used  $\log(k+2) \leq \log(3)\sqrt{k}$  for all  $k \geq 1$ .  $\square$

**Lemma A.4** For every  $k \in \mathbb{N}$  and  $\sigma \in (0, 1)$ ,

$$\mathbb{E} \left[ \sum_{j=1}^k \omega_j [g]_j^2 \left( \frac{1}{X_j} \mathbf{1}_{[X_j \geq \sigma]} - \frac{1}{a_j} \right)^2 \right] \leq 8 d r \kappa_\sigma(\gamma, \lambda, \omega).$$

*Proof.* Firstly, as  $f \in \mathcal{F}_\gamma^r$ , it is easily seen that

$$\mathbb{E} \left[ \sum_{j=1}^k \omega_j [g]_j^2 \left( \frac{1}{X_j} \mathbf{1}_{[X_j \geq \sigma]} - \frac{1}{a_j} \right)^2 \right] \leq r \sup_{1 \leq j \leq k} \frac{\omega_j}{\gamma_j} \mathbb{E}[|R_j|^2],$$

where  $R_j$  is defined as

$$R_j := \left( \frac{a_j}{X_j} \mathbf{1}_{\{X_j^2 \geq \sigma^2\}} - 1 \right). \quad (\text{A.7})$$

In view of the definition of  $\kappa_\sigma$  in Theorem 2.3, the result follows from  $\mathbb{E}[|R_j|^2] \leq d \min \left\{ 1, \frac{8\sigma}{\lambda_j} \right\}$ , which is a consequence of the decomposition

$$\mathbb{E}|R_j|^2 = \mathbb{E} \left[ \left( \frac{a_j}{X_j} - 1 \right)^2 \mathbf{1}_{\{X_j^2 \geq \sigma^2\}} \right] + \mathbf{P}[X_j^2 < \sigma^2] \quad (\text{A.8})$$

and Lemma A.1.  $\square$

**Lemma A.5** We have that

$$\mathbb{E} \left[ \sum_{j=1}^{K_{\varepsilon, \sigma}^+} \omega_j \left( \frac{1}{X_j} \mathbf{1}_{[X_j \geq \sigma]} - \frac{1}{a_j} \right)^2 \varepsilon \xi_j^2 \mathbf{1}_{\Omega_\sigma^c} \right] \leq 64 d^3 (\mathbf{P}[\Omega_\sigma^c])^{1/2}.$$

*Proof.* Given  $R_j$  from (A.7), we begin our proof observing that

$$\mathbb{E} \left[ \sum_{j=1}^{K_{\varepsilon, \sigma}^+} \omega_j \left( \frac{1}{X_j} \mathbf{1}_{[X_j \geq \sigma]} - \frac{1}{a_j} \right)^2 \sqrt{\varepsilon} \xi_j^2 \mathbf{1}_{\Omega_\sigma^c} \right] \leq \varepsilon \sum_{j=1}^{K_{\varepsilon, \sigma}^+} \frac{\omega_j}{a_j^2} \mathbb{E}[|R_j|^2 \mathbf{1}_{\Omega_\sigma^c}],$$

where we have used the independence of  $X$  and  $Y$  and  $\text{Var}(Y_j) = \varepsilon$ . Since  $d\delta_k^\lambda \geq \sum_{j=1}^k \frac{\omega_j}{a_j^2}$  for all  $A \in \mathcal{A}_\lambda^d$ , the Cauchy-Schwarz inequality yields

$$\mathbb{E} \left[ \sum_{j=1}^{K_{\varepsilon, \sigma}^+} \omega_j \left( \frac{1}{X_j} \mathbf{1}_{[X_j \geq \sigma]} - \frac{1}{a_j} \right)^2 \varepsilon \xi_j^2 \mathbf{1}_{\Omega_\sigma^c} \right] \leq d (\mathbf{P}[\Omega_\sigma^c])^{1/2} \varepsilon \delta_{N_\varepsilon^+}^\lambda \max_{0 < j \leq N_\varepsilon^+} (\mathbb{E}[|R_j|^4])^{1/2}.$$

Proceeding analogously to (A.6) and (A.8), one can show that that  $\mathbb{E}[|R_j|^4] \leq 4$ . The result follows then by definition of  $N_\varepsilon^+$ .  $\square$

**Lemma A.6** For  $k \in \mathbb{N}$ , define the events

$$\tilde{\Omega}_k := \left\{ \left| \frac{X_j}{a_j} - 1 \right| \leq \frac{1}{3} \quad \forall 1 \leq j \leq k \right\}$$

and suppose that Assumptions 2.1 and 3.4 hold. For all  $\varepsilon, \sigma \in (0, 1)$ , we have

$$(i) \Omega_\sigma \subseteq \{\text{pen}_k \leq \widehat{\text{pen}}_k \leq 30 \text{pen}_k \quad \forall 1 \leq k \leq K_{\varepsilon, \sigma}^+\},$$

$$(ii) \widetilde{\Omega}_{M_\sigma^+ + 1} \subseteq \{K_{\varepsilon, \sigma}^- \leq \widehat{K}_{\varepsilon, \sigma} \leq K_{\varepsilon, \sigma}^+\},$$

$$(iii) \mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c] \leq C(\lambda, d)\sigma^6.$$

*Proof.* Consider (i). Notice first that  $\delta_k^a \leq \delta_k^\lambda d \zeta_d$  for all  $k \geq 1$  with  $\zeta_d := (\log(3d))/(\log 3)$ . Observe that on  $\Omega_\sigma$  we have  $(1/2)\Delta_k^a \leq \Delta_k^X \leq (3/2)\Delta_k^a$  for all  $1 \leq k \leq \widetilde{M}_\sigma$  and hence  $(1/2)[\Delta_k^a \vee (k+2)] \leq [\Delta_k^X \vee (k+2)] \leq (3/2)[\Delta_k^a \vee (k+2)]$ , which implies

$$\begin{aligned} & (1/2)k\Delta_k^a \left( \frac{\log[\Delta_k^a \vee (k+2)]}{\log(k+2)} \right) \left( 1 - \frac{\log 2}{\log(k+2)} \frac{\log(k+2)}{\log(\Delta_k^a \vee [k+2])} \right) \\ & \leq \delta_k^X \leq (3/2)k\Delta_k^a \left( \frac{\log(\Delta_k^a \vee [k+2])}{\log(k+2)} \right) \left( 1 + \frac{\log 3/2}{\log(k+2)} \frac{\log(k+2)}{\log(\Delta_k^a \vee [k+2])} \right). \end{aligned}$$

Using  $\log(\Delta_k^a \vee (k+2))/\log(k+2) \geq 1$ , we conclude from the last estimate that

$$\begin{aligned} \delta_k^a/10 & \leq (\log 3/2)/(2 \log 3) \delta_k^a \leq (1/2) \delta_k^a [1 - (\log 2)/\log(k+2)] \leq \delta_k^X \\ & \leq (3/2) \delta_k^a [1 + (\log 3/2)/\log(k+2)] \leq 3\delta_k^a. \end{aligned}$$

It follows that on  $\Omega_\sigma$  we have  $\text{pen}_k \leq \widehat{\text{pen}}_k \leq 30 \text{pen}_k$  for all  $1 \leq k \leq M_\sigma^+$  as desired.

Proof of (ii). Define the events  $\Omega_I := \{K_{\varepsilon, \sigma}^- > \widehat{K}_{\varepsilon, \sigma}\}$  and  $\Omega_{II} := \{\widehat{K}_{\varepsilon, \sigma} > K_{\varepsilon, \sigma}^+\}$ . Then we have  $\{K_{\varepsilon, \sigma}^- \leq \widehat{K}_{\varepsilon, \sigma} \leq K_{\varepsilon, \sigma}^+\}^c = \Omega_I \cup \Omega_{II}$ . Consider  $\Omega_I = \{\widehat{N}_\varepsilon < K_{\varepsilon, \sigma}^-\} \cup \{\widehat{M}_\sigma < K_{\varepsilon, \sigma}^-\}$  first. By definition of  $N_\varepsilon^-$ , we have that  $\min_{1 \leq j \leq N_\varepsilon^-} \frac{a_j^2}{j \omega_j^+} \geq 4\varepsilon |\log \varepsilon|$ , which implies, keeping in mind that  $K_{\varepsilon, \sigma}^- \leq N_{\varepsilon, \sigma}^-$ ,

$$\begin{aligned} \{\widehat{N}_\varepsilon < K_{\varepsilon, \sigma}^-\} & \subseteq \left\{ \exists 1 \leq j \leq K_{\varepsilon, \sigma}^- : \frac{X_j^2}{j \omega_j^+} < \varepsilon |\log \varepsilon| \right\} \\ & \subseteq \bigcup_{1 \leq j \leq K_{\varepsilon, \sigma}^-} \left\{ \frac{X_j}{a_j} \leq \frac{1}{2} \right\} \subseteq \bigcup_{1 \leq j \leq K_{\varepsilon, \sigma}^-} \left\{ \left| \frac{X_j}{a_j} - 1 \right| \geq \frac{1}{2} \right\}. \end{aligned}$$

One can see that from  $\min_{1 \leq j \leq M_\sigma^-} a_j^2 \geq 4\sigma^{1-v_\sigma}$  it follows in the same way that

$$\{\widehat{M}_\sigma < K_{\varepsilon, \sigma}^-\} \subseteq \bigcup_{1 \leq j \leq K_{\varepsilon, \sigma}^-} \left\{ \left| \frac{X_j}{a_j} - 1 \right| \geq \frac{1}{2} \right\}.$$

Therefore,  $\Omega_I \subseteq \bigcup_{1 \leq j \leq M_\sigma^+} \left\{ |X_j/a_j - 1| \geq 1/2 \right\} \subseteq \widetilde{\Omega}_{M_\sigma^+ + 1}^c$ , since  $M_\sigma^- \leq M_\sigma^+$ .

Consider  $\Omega_{II} = \{\widehat{N}_\varepsilon > K_{\varepsilon, \sigma}^+\} \cap \{\widehat{M}_\sigma > K_{\varepsilon, \sigma}^+\}$ . In case  $K_{\varepsilon, \sigma}^+ = N_{\varepsilon, \sigma}^+$ , note that by definition of  $N_{\varepsilon, \sigma}^+$ , we have  $\varepsilon |\log \varepsilon|/4 \geq \frac{a_{N_{\varepsilon, \sigma}^+ + 1}^2}{(N_{\varepsilon, \sigma}^+ + 1) \omega_{N_{\varepsilon, \sigma}^+ + 1}^+}$ , such that

$$\begin{aligned} \Omega_{II} & \subseteq \{\widehat{N}_\varepsilon > N_{\varepsilon, \sigma}^+\} \subseteq \left\{ \forall 1 \leq j \leq N_{\varepsilon, \sigma}^+ + 1 : \frac{X_j^2}{j \omega_j^+} \geq \varepsilon |\log \varepsilon| \right\} \\ & \subseteq \left\{ \frac{X_{N_{\varepsilon, \sigma}^+ + 1}}{a_{N_{\varepsilon, \sigma}^+ + 1}} \geq 2 \right\} \subseteq \left\{ \left| \frac{X_{N_{\varepsilon, \sigma}^+ + 1}}{a_{N_{\varepsilon, \sigma}^+ + 1}} - 1 \right| \geq 1 \right\}. \end{aligned}$$

In case  $K_{\varepsilon, \sigma}^+ = M_\sigma^+$ , it follows analogously from  $\sigma^{1-v_\sigma} \geq 4 \max_{j \geq M_\sigma^+ + 1} a_j^2$  that

$$\Omega_{II} \subseteq \{\widehat{M}_\sigma > M_\sigma^+\} \subseteq \left\{ |X_{M_\sigma^+ + 1}/a_{M_\sigma^+ + 1} - 1| \geq 1 \right\}.$$



Therefore, we have  $\Omega_{II} \subseteq \left\{ |X_{K_{\varepsilon, \sigma}^+ + 1} / a_{K_{\varepsilon, \sigma}^+ + 1} - 1| \geq 1 \right\} \subseteq \tilde{\Omega}_{M_{\sigma}^+ + 1}^c$  and (ii) is shown.

Proof of (iii). We distinguish the cases  $\sigma \leq \sigma_0 := \exp(-512 \log(3d)^2)$  and  $\sigma > \sigma_0$ . The assertion is trivial for  $\sigma > \sigma_0$  (keeping in mind that  $\mathbf{P}[\mathcal{U}_{\varepsilon, \sigma}^c] \leq \sigma_0^{-6} \sigma^6$ ). Consider the case  $\sigma \leq \sigma_0$ , where  $a_j^2 \geq 2\sigma$  for all  $1 \leq j \leq M_{\sigma}^+$  due to Lemma A.2 (ii). This yields for the complement of  $\Omega_{\sigma}$

$$\Omega_{\sigma}^c = \left\{ \exists 1 \leq j \leq M_{\sigma}^+ : \left| \frac{a_j}{X_j} - 1 \right| > \frac{1}{2} \vee |X_j|^2 < \sigma \right\} \subseteq \left\{ \exists 1 \leq j \leq M_{\sigma}^+ : \left| \frac{X_j}{a_j} - 1 \right| > \frac{1}{3} \right\} = \tilde{\Omega}_{M_{\sigma}^+}^c,$$

and thus  $\tilde{\Omega}_{M_{\sigma}^+ + 1}^c \subseteq \Omega_{\sigma}$  since trivially  $\tilde{\Omega}_{M_{\sigma}^+ + 1}^c \subseteq \tilde{\Omega}_{M_{\sigma}^+}^c$ . It follows with assertion (ii) that  $\mathcal{U}_{\varepsilon, \sigma}^c \subseteq \tilde{\Omega}_{M_{\sigma}^+ + 1}^c$  for all  $\sigma \leq \sigma_0$ . For  $Z \sim \mathcal{N}(0, 1)$  and  $z \geq 0$ , one has  $\mathbf{P}[Z > z] \leq (2\pi z^2)^{-1/2} \exp(-z^2/2)$ . Hence, there is a constant  $C(d)$  depending on  $d$  such that for every  $1 \leq j \leq M_{\sigma}^+ + 1$ ,

$$\mathbf{P}[|X_j/a_j - 1| > 1/3] \leq C(d) \left( \frac{\sigma}{\lambda_{M_{\sigma}^+ + 1}} \right)^{1/2} \exp\left( -\frac{\lambda_{M_{\sigma}^+ + 1}}{18\sigma d} \right).$$

Consequently, as  $M_{\sigma}^+ \leq \sigma^{-1}$ ,

$$\mathbf{P}[\tilde{\Omega}_{M_{\sigma}^+ + 1}^c] \leq C(d) (\sigma \lambda_{M_{\sigma}^+ + 1})^{-1/2} \exp\left( -\frac{\lambda_{M_{\sigma}^+ + 1}}{18\sigma d} \right)$$

which implies the assertion (iii) by virtue of Assumption 3.4.  $\square$

## References

- Barron, A., Birgé, L., and Massart, P. (1999). Risk bounds for model selection via penalization. *Probability Theory and Related Fields*, 113:301–413.
- Cavalier, L., Golubev, G., Picard, D., and Tsybakov, A. (2002). Oracle inequalities for inverse problems. *Ann. Stat.*, 30:843–874.
- Cavalier, L. and Hengartner, N. W. (2005). Adaptive estimation for inverse problems with noisy operators. *Inverse Problems*, 21:1345–1361.
- Dahlhaus, R. and Polonik, W. (2006). Nonparametric quasi-maximum likelihood estimation for Gaussian locally stationary processes. *Ann. Stat.*, 34:2790–2824.
- Efromovich, S. (1997). Density estimation for the case of supersmooth measurement error. *Journal of the American Statistical Association*, 92:526–535.
- Ermakov, M. (1990). On optimal solutions of the deconvolution problem. *Inverse Probl.*, 6(5):863–872.
- Fan, J. (1991). On the optimal rates of convergence for nonparametric deconvolution problems. *The Annals of Statistics*, 19:1257–1272.
- Goldenshluger, A. and Lepski, O. (2011). Bandwidth selection in kernel density estimation: oracle inequalities and adaptive minimax optimality. *Ann. Stat.*, 39(3):1608–1632.
- Johnstone, I. M. and Silverman, B. W. (1990). Speed of estimation in positron emission tomography and related inverse problems. *Ann. Stat.*, 18(1):251–280.
- Mair, B. A. and Ruyngaert, F. H. (1996). Statistical inverse estimation in Hilbert scales. *SIAM Journal on Applied Mathematics*, 56(5):1424–1444.

- Mathé, P. and Pereverzev, S. V. (2001). Optimal discretization of inverse problems in Hilbert scales. Regularization and self-regularization of projection methods. *SIAM J. Numer. Anal.*, 38(6):1999–2021.
- Neumann, M. H. (1997). On the effect of estimating the error density in nonparametric deconvolution. *Journal of Nonparametric Statistics*, 7:307–330.
- Petrov, V. V. (1995). *Limit theorems of probability theory. Sequences of independent random variables.* Oxford Studies in Probability. Clarendon Press., Oxford, 4. edition.
- Stefanski, L. and Carroll, R. J. (1990). Deconvoluting kernel density estimators. *Statistics*, 21:169–184.