



# OPTIMAL ADAPTIVE ESTIMATION OF LINEAR FUNCTIONALS UNDER SPARSITY

O Collier, Laëtitia Comminges, A B Tsybakov, Nicolas Verzelen

## ► To cite this version:

O Collier, Laëtitia Comminges, A B Tsybakov, Nicolas Verzelen. OPTIMAL ADAPTIVE ESTIMATION OF LINEAR FUNCTIONALS UNDER SPARSITY. 2017. hal-01425801v1

**HAL Id: hal-01425801**

**<https://hal.science/hal-01425801v1>**

Preprint submitted on 3 Jan 2017 (v1), last revised 6 Oct 2017 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# OPTIMAL ADAPTIVE ESTIMATION OF LINEAR FUNCTIONALS UNDER SPARSITY

O. COLLIER, L. COMMINGES, A.B. TSYBAKOV AND N. VERZELEN

**ABSTRACT.** We consider the problem of estimation of a linear functional in the Gaussian sequence model where the unknown vector  $\theta \in \mathbb{R}^d$  belongs to a class of  $s$ -sparse vectors with unknown  $s$ . We suggest an adaptive estimator achieving a non-asymptotic rate of convergence that differs from the minimax rate at most by a logarithmic factor. We also show that this optimal adaptive rate cannot be improved when  $s$  is unknown. Furthermore, we address the issue of simultaneous adaptation to  $s$  and to the variance  $\sigma^2$  of the noise. We suggest an estimator that achieves the optimal adaptive rate when both  $s$  and  $\sigma^2$  are unknown.

## 1. INTRODUCTION

We consider the model

$$(1) \quad y_j = \theta_j + \sigma \xi_j, \quad j = 1, \dots, d,$$

where  $\theta = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$  is an unknown vector of parameters,  $\xi_j$  are i.i.d. standard normal random variables, and  $\sigma > 0$  is the noise level. We study the problem of estimation of the linear functional

$$L(\theta) = \sum_{i=1}^d \theta_i,$$

based on the observations  $y = (y_1, \dots, y_d)$ .

For  $s \in \{1, \dots, d\}$ , we denote by  $\Theta_s$  the class of all  $\theta \in \mathbb{R}^d$  satisfying  $\|\theta\|_0 \leq s$ , where  $\|\theta\|_0$  denotes the number of non-zero components of  $\theta$ . We assume that  $\theta$  belongs to  $\Theta_s$  for some  $s \in \{1, \dots, d\}$ . Parameter  $s$  characterizes the sparsity of vector  $\theta$ . The problem of estimation of  $L(\theta)$  in this context arises, for example, if one wants to estimate the value of a function  $f$  at a fixed point from noisy observations of its Fourier coefficients knowing that the function admits a sparse representation with respect to the first  $d$  functions of the Fourier basis. Indeed, in this case the value  $f(0)$  is equal to the sum of Fourier coefficients of  $f$  with even indices.

As a measure of quality of an estimator  $\hat{T}$  of the functional  $L(\theta)$  based on the sample  $(y_1, \dots, y_d)$ , we consider the maximum squared risk

$$\psi_s^{\hat{T}} \triangleq \sup_{\theta \in \Theta_s} \mathbf{E}_{\theta} (\hat{T} - L(\theta))^2,$$

where  $\mathbf{E}_{\theta}$  denotes the expectation with respect to the distribution  $\mathbf{P}_{\theta}$  of  $(y_1, \dots, y_d)$  satisfying (1). For each fixed  $s \in \{1, \dots, d\}$ , the best quality of estimation is characterized by the minimax risk

$$\psi_s^* \triangleq \inf_{\hat{T}} \sup_{\theta \in \Theta_s} \mathbf{E}_{\theta} (\hat{T} - L(\theta))^2,$$

where the infimum is taken over all estimators. An estimator  $T^*$  is called rate optimal on  $\Theta_s$  if  $\psi_s^{T^*} \asymp \psi_s^*$ . Here and in the following we write  $a(d, s, \sigma) \asymp b(d, s, \sigma)$  for two functions  $a(\cdot)$  and  $b(\cdot)$  of  $d, s$  and  $\sigma$  if there exist absolute constants  $c > 0$  and  $c' > 0$  such that  $c < a(d, s, \sigma)/b(d, s, \sigma) < c'$  for all  $d$ , all  $s \in \{1, \dots, d\}$  and all  $\sigma > 0$ .

The problem of estimation of the linear functional from the minimax point of view has been analyzed in [6, 1, 2, 4, 5, 8] among others. Most of these papers study minimax estimation of linear functionals on classes of vectors  $\theta$  different from  $\Theta_s$ . Namely,  $\theta$  is considered as a vector of first  $d$  Fourier or wavelet coefficients of functions belonging to some smoothness class, such as Sobolev or Besov classes. In particular, the class of vectors  $\theta$  is assumed to be convex, which is not the case of class  $\Theta_s$ . Cai and Low [1] were the first to address the problem of constructing rate optimal estimators of  $L(\theta)$  on the sparsity class  $\Theta_s$  and evaluating the minimax risk  $\psi_s^*$ . They studied the case  $s < d^a$  for some  $a < 1/2$ , with  $\sigma = 1/\sqrt{d}$ , and established upper and lower bounds on  $\psi_s^*$  that are accurate up to a logarithmic factor in  $d$ . The sharp non-asymptotic expression for the minimax risk  $\psi_s^*$  is derived in [3] where it is shown that, for all  $d$ , all  $s \in \{1, \dots, d\}$  and all  $\sigma > 0$

$$\psi_s^* \asymp \sigma^2 s^2 \log(1 + d/s^2).$$

Furthermore, [3] proves that a simple estimator of the form

$$(2) \quad \hat{L}_s^* = \begin{cases} \sum_{j=1}^d y_j \mathbf{1}_{y_j^2 > 2\sigma^2 \log(1+d/s^2)}, & \text{if } s < \sqrt{d}, \\ \sum_{j=1}^d y_j, & \text{otherwise,} \end{cases}$$

is rate optimal. Here and in the following,  $\mathbf{1}_{\{\cdot\}}$  denotes the indicator function.

Note that the minimax risk  $\psi_s^*$  critically depends on the parameter  $s$  that in practice is usually unknown. More importantly, the rate optimal estimator  $\hat{L}_s^*$  depends on  $s$  as well, which makes it inaccessible in practice.

In this paper, we suggest adaptive estimators of  $L(\theta)$  that do not depend on  $s$  and achieve a non-asymptotic rate of convergence  $\Phi^L(\sigma, s)$  that differs from the minimax rate  $\psi_s^*$  at most by a logarithmic factor. We also show that this rate cannot be improved when  $s$  is unknown in the sense of the definition that we give in Section 2 below. Furthermore, in Section 3 we address the issue of simultaneous adaptation to  $s$  and  $\sigma$ . We suggest an estimator that achieves the best rate of adaptive estimation  $\Phi^L(\sigma, s)$  when both  $s$  and  $\sigma$  are unknown.

## 2. MAIN RESULTS

We assume throughout the paper that  $d \geq 3$ . Our aim is to show that the optimal adaptive rate of convergence is of the form

$$\Phi^L(\sigma, s) = \sigma^2 s^2 \log(1 + d(\log d)/s^2)$$

and to construct an adaptive estimator attaining this rate. Note that

$$(3) \quad \Phi^L(\sigma, s) \asymp \sigma^2 d(\log d), \quad \text{for all } \sqrt{d \log d} \leq s \leq d.$$

Indeed, since the function  $x \mapsto x \log(1 + 1/x)$  is increasing for  $x > 0$ ,

$$(4) \quad d(\log d)/2 \leq s^2 \log(1 + d(\log d)/s^2) \leq d(\log d), \quad \forall \sqrt{d \log d} \leq s \leq d, \quad d \geq 3.$$

To construct an adaptive estimator, we first consider a collection of non-adaptive estimators indexed by  $s = 1, \dots, d$ :

$$(5) \quad \hat{L}_s = \begin{cases} \sum_{j=1}^d y_j \mathbf{1}_{y_j^2 > \alpha \sigma^2 \log(1+d(\log d)/s^2)}, & \text{if } s \leq \sqrt{d \log d}, \\ \sum_{j=1}^d y_j, & \text{otherwise,} \end{cases}$$

where  $\alpha > 0$  is a constant that will be chosen large enough. Note that if in definition (5) we replace  $d(\log d)$  by  $d$ , and  $\alpha$  by 2, we obtain the estimator  $\hat{L}_s^*$  suggested in [3], cf. (2). It is proved in [3] that the estimator  $\hat{L}_s^*$  is rate optimal in the minimax non-adaptive sense. The additional  $\log d$  factor is necessary to achieve adaptivity as it will be clear from the subsequent arguments.

We obtain an adaptive estimator via data-driven selection in the collection of estimators  $\{\hat{L}_s\}$ . The selection is based on a Lepski type scheme. For  $s = 1, \dots, d$ , consider the thresholds  $\omega_s > 0$  given by

$$\omega_s^2 = \beta \sigma^2 s^2 \log(1 + d(\log d)/s^2) = \beta \Phi^L(\sigma, s),$$

where  $\beta > 0$  is a constant that will be chosen large enough. We define the selected index  $\hat{s}$  by the relation

$$(6) \quad \hat{s} \triangleq \min \left\{ s \in \{1, \dots, \lfloor \sqrt{d \log d} \rfloor\} : |\hat{L}_s - \hat{L}_{s'}| \leq \omega_{s'} \text{ for all } s' > s \right\}$$

with the convention that  $\hat{s} = \lfloor \sqrt{d \log d} \rfloor + 1$  if the set in (6) is empty. Here,  $\lfloor \sqrt{d \log d} \rfloor$  denotes the largest integer less than  $\sqrt{d \log d}$ . Finally, we define an adaptive to  $s$  estimator of  $L$  as

$$(7) \quad \hat{L} \triangleq \hat{L}_{\hat{s}}.$$

The following theorem exhibits an upper bound on its risk.

**Theorem 1.** *There exists an absolute constant  $\alpha_0 > 0$  such that the following holds. Let  $\hat{L}$  be the estimator defined in (7) with parameters  $\alpha > \alpha_0$  and  $\beta > 37\alpha$ . Then, for all  $\sigma > 0$ , and all integers  $d \geq 3$ ,  $s \in \{1, \dots, d\}$  we have*

$$\sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\hat{L} - L(\theta))^2 \leq C \Phi^L(\sigma, s),$$

where  $C > 0$  is an absolute constant.

Observe that for small  $s$  (such that  $s \leq d^b$  for  $b < 1/2$ ), we have  $1 \leq \Phi^L(\sigma, s)/\psi_s^* \leq c'$  where  $c' > 0$  is an absolute constant. Therefore, for such  $s$  our estimator  $\hat{L}$  attains the best possible rate on  $\Theta_s$  given by the minimax risk  $\psi_s^*$  and it cannot be improved, even by estimators depending on  $s$ . Because of this, the only issue is to check that the rate  $\Phi^L(\sigma, s)$  cannot be improved if  $s$  is greater than  $d^b$  with  $b < 1/2$ . For definiteness, we consider below the case  $b = 1/4$  but with minor modifications the argument applies to any  $b < 1/2$ . Specifically, we prove that any estimator whose maximal risk over  $\Theta_s$  is smaller (within a small constant) than  $\Phi^L(\sigma, s)$  for some  $s \geq d^{1/4}$ , must have a maximal risk over  $\Theta_1$  of power order in  $d$  instead of the logarithmic order  $\Phi^L(\sigma, 1)$  corresponding to our estimator. In other words, if we find an estimator that improves upon our estimator only slightly (by a constant factor) for some  $s \geq d^{1/4}$ , then this estimator inevitably loses much more for small  $s$ , such as  $s = 1$ , since there the ratio of maximal risks of the two estimators behaves as a power of  $d$ .

**Theorem 2.** *Let  $d \geq 3$  and  $\sigma > 0$ . There exist two small absolute constants  $C_0 > 0$  and  $C_1 > 0$  such that the following holds. Any estimator  $\hat{T}$  that satisfies*

$$\sup_{\theta \in \Theta_s} \mathbf{E}_\theta [(\hat{T} - L(\theta))^2] \leq C_0 \Phi^L(\sigma, s) \quad \text{for some } s \geq d^{1/4}$$

*has a degenerate maximal risk over  $\Theta_1$ , that is*

$$\sup_{\theta \in \Theta_1} \mathbf{E}_\theta [(\hat{T} - L(\theta))^2] \geq C_1 \sigma^2 d^{1/4}.$$

The property obtained in Theorem 2 can be paraphrased in an asymptotic context to conclude that  $\Phi^L(\sigma, s)$  is the adaptive rate of convergence on the scale of classes  $\{\Theta_s, s = 1, \dots, d\}$  in the sense of the definition in [10]. Indeed, assume that  $d \rightarrow \infty$ . Following [10], we call a function  $s \mapsto \Psi_d(s)$  the *adaptive rate of convergence on the scale of classes*  $\{\Theta_s, s = 1, \dots, d\}$  if the following holds.

(i) There exists an estimator  $\hat{L}$  such that, for all  $d$ ,

$$(8) \quad \max_{s=1, \dots, d} \sup_{\theta \in \Theta_s} \mathbf{E}_\theta (\hat{L} - L(\theta))^2 / \Psi_d(s) \leq C,$$

where  $C > 0$  is a constant (clearly, such an estimator  $\hat{L}$  is adaptive since it cannot depend on  $s$ ).

(ii) If there exist another function  $s \mapsto \Psi'_d(s)$  and a constant  $C' > 0$  such that, for all  $d$ ,

$$(9) \quad \inf_{\hat{T}} \max_{s=1, \dots, d} \sup_{\theta \in \Theta_s} \mathbf{E}_\theta (\hat{T} - L(\theta))^2 / \Psi'_d(s) \leq C',$$

and

$$(10) \quad \min_{s=1, \dots, d} \frac{\Psi'_d(s)}{\Psi_d(s)} \rightarrow 0 \text{ as } d \rightarrow \infty,$$

then there exists  $\bar{s} \in \{1, \dots, d\}$  such that

$$(11) \quad \frac{\Psi'_d(\bar{s})}{\Psi_d(\bar{s})} \min_{s=1, \dots, d} \frac{\Psi'_d(s)}{\Psi_d(s)} \rightarrow \infty \text{ as } d \rightarrow \infty.$$

In words, this definition states that the adaptive rate of convergence  $\Psi_d(s)$  is such that any improvement of this rate for some  $s$  (cf. (10)) is possible only at the expense of much greater loss for another  $\bar{s}$  (cf. (11)).

**Corollary 1.** *The rate  $\Phi^L(\sigma, s)$  is the adaptive rate of convergence on the scale of classes  $\{\Theta_s, s = 1, \dots, d\}$ .*

It follows from the above results that the rate  $\Phi^L(\sigma, s)$  cannot be improved when adaptive estimation on the family of sparsity classes  $\{\Theta_s, s = 1, \dots, d\}$  is considered. The ratio between the best rate of adaptive estimation  $\Phi^L(\sigma, s)$  and the minimax rate  $\psi_s^*$  is equal to

$$\phi_s^* = \frac{\Phi^L(\sigma, s)}{\psi_s^*} = \frac{\log(1 + d(\log d)/s^2)}{\log(1 + d/s^2)}.$$

As mentioned above,  $\phi_s^* \asymp 1$  if  $s \leq d^b$  for  $b < 1/2$ . In a vicinity of  $s = \sqrt{d}$  we have  $\phi_s^* \asymp \log \log d$ , whereas for  $s \geq \sqrt{d \log d}$  the behavior of this ratio is logarithmic:  $\phi_s^* \asymp \log d$ . Thus, there are three different regimes and we see that, in all of them, rate adaptive estimation of the linear functional on the sparsity classes is impossible without loss of efficiency as compared to the minimax estimation. However, this loss is at most logarithmic in  $d$ .

3. ADAPTATION TO  $s$  WHEN  $\sigma$  IS UNKNOWN

In this section we discuss a generalization of our adaptive estimator to the case when the standard deviation  $\sigma$  of the noise is unknown.

To treat the case of unknown  $\sigma$ , we first construct an estimator  $\hat{\sigma}$  of  $\sigma$  such that, with high probability,  $\sigma \leq \hat{\sigma} \leq 10\sigma$ . Then, we consider the family of estimators defined by a relation analogous to (5):

$$(12) \quad \hat{L}'_s = \begin{cases} \sum_{j=1}^d y_j \mathbf{1}_{y_j^2 > \alpha \hat{\sigma}^2 \log(1+d(\log d)/s^2)}, & \text{if } s \leq \sqrt{d \log d}, \\ \sum_{j=1}^d y_j, & \text{otherwise,} \end{cases}$$

where  $\alpha > 0$  is a constant to be chosen large enough. The difference from (5) consists in the fact that we replace the unknown  $\sigma$  by  $\hat{\sigma}$ . Then, we define a random threshold  $\omega'_s > 0$  as

$$(\omega'_s)^2 = \beta \hat{\sigma}^2 s^2 \log(1 + d(\log d)/s^2),$$

where  $\beta > 0$  is a constant to be chosen large enough. The selected index  $\hat{s}'$  is defined by the formula analogous to (6):

$$(13) \quad \hat{s}' \triangleq \min \{s \in \{1, \dots, \lfloor \sqrt{d \log d} \rfloor\} : |\hat{L}'_s - \hat{L}'_{s'}| \leq \omega'_{s'} \text{ for all } s' > s\}.$$

Finally, the adaptive estimator when  $\sigma$  is unknown is defined as

$$\hat{L}' \triangleq \hat{L}'_{\hat{s}'},$$

The aim of this section is to show that the risk of the estimator  $\hat{L}'$  admits an upper bound with the same rate as in Theorem 1 for all  $d$  large enough. Consequently,  $\hat{L}'$  attains the best rate of adaptive estimation as follows from Section 2.

Different estimators  $\hat{\sigma}$  can be used. By slightly modifying the method suggested in [3], we consider the statistic

$$(14) \quad \hat{\sigma} = 9 \left( \frac{1}{\lfloor d/2 \rfloor} \sum_{j \leq d/2} y_{(j)}^2 \right)^{1/2}$$

where  $y_{(1)}^2 \leq \dots \leq y_{(d)}^2$  are the order statistics associated to  $y_1^2, \dots, y_d^2$ . This statistic has the properties stated in the next proposition. In particular,  $\hat{\sigma}$  overestimates  $\sigma$  but it turns out to be without prejudice to the attainment of the best rate by the resulting estimator  $\hat{L}'_s$ .

**Proposition 1.** *There exists an absolute constant  $d_0 \geq 3$  such that the following holds. Let  $\hat{\sigma}$  be the estimator defined in (14). Then, for all integers  $d \geq d_0$  and  $s < d/2$  we have*

$$(15) \quad \inf_{\theta \in \Theta_s} \mathbf{P}_\theta(\sigma \leq \hat{\sigma} \leq 10\sigma) \geq 1 - d^{-5},$$

and

$$(16) \quad \sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\hat{\sigma}^4) \leq \bar{C} \sigma^4,$$

where  $\bar{C}$  is an absolute constant.

The proof of this proposition is given in Section 4. Using Proposition 1 we establish the following bound on the risk of the estimator  $\hat{L}'$ .

**Theorem 3.** *There exist large enough absolute constants  $\alpha > 0$ ,  $\beta > 0$ , and  $d_0 \geq 3$  such that the following holds. Let  $\hat{\sigma}$  be the estimator defined in (14). Then, for the estimator  $\hat{L}'$  with tuning parameters  $\alpha$  and  $\beta$ , for all  $\sigma > 0$ , and all integers  $d \geq d_0$  and  $s < d/2$  we have*

$$(17) \quad \sup_{\theta \in \Theta_s} \mathbf{E}_\theta (\hat{L}' - L(\theta))^2 \leq C \Phi^L(\sigma, s),$$

where  $C > 0$  is an absolute constant.

Thus, the estimator  $\hat{L}'$ , which is independent of both  $s$  and  $\sigma$  achieves the rate  $\Phi^L(\sigma, s)$  that is the best possible rate of adaptive estimation established in Section 2.

The condition  $s < d/2$  in this theorem can be generalized to  $s \leq \zeta d$  for some  $\zeta \in (0, 1)$ . In fact, for any  $\zeta \in (0, 1)$ , we can modify the definition of (14) by summing only over the  $(1 - \zeta)d$  smallest values of  $y_i^2$ . Then, changing the numerical constants  $\alpha$  and  $\beta$  in the definition of  $\omega'_s$ , we obtain that the corresponding estimator  $\hat{L}'$  achieves the best possible rate simultaneously for all  $s \leq \zeta d$  with a constant  $C$  in (17) that would depend on  $\zeta$ . However, we cannot set  $\zeta = 1$ . Indeed, the following proposition shows that it is not possible to construct an estimator, which is simultaneously adaptive to all  $\sigma > 0$  and to all  $s \in [1, d]$ .

**Proposition 2.** *Let  $d \geq 3$  and  $\sigma > 0$ . There exists a small absolute constant  $C_0 > 0$  such that the following holds. Any estimator  $\hat{T}$  that satisfies*

$$(18) \quad \sup_{\theta \in \Theta_1} \mathbf{E}_\theta [(\hat{T} - L(\theta))^2] \leq C_0 \sigma^2 d, \quad \forall \sigma > 0,$$

has a degenerate maximal risk over  $\Theta_d$ , that is, for any fixed  $\sigma > 0$ ,

$$(19) \quad \sup_{\theta \in \Theta_d} \mathbf{E}_\theta [(\hat{T} - L(\theta))^2] = \infty.$$

In other words, when  $\sigma$  is unknown, any estimator, for which the maximal risk over  $\Theta_d$  is finite for all  $\sigma$ , cannot achieve over  $\Theta_1$  a risk of smaller order than  $\sigma^2 d$ , and hence cannot be minimax adaptive. Indeed, as shown above, the adaptive minimax rate over  $\Theta_1$  is of the order  $\sigma^2 \log d$ .

#### 4. PROOFS OF THE UPPER BOUNDS

In the following, we will denote by  $c_1, c_2, \dots$  absolute positive constants. We will write for brevity  $L$  instead of  $L(\theta)$ .

**4.1. Proof of Theorem 1.** Let  $s \in \{1, \dots, d\}$  and assume that  $\theta$  belongs to  $\Theta_s$ . We have

$$(20) \quad \mathbf{E}_\theta (\hat{L} - L)^2 = \mathbf{E}_\theta [(\hat{L}_{\hat{s}} - L)^2 \mathbf{1}_{\hat{s} \leq s}] + \mathbf{E}_\theta [(\hat{L}_{\hat{s}} - L)^2 \mathbf{1}_{\hat{s} > s}].$$

Consider the first summand on the right hand side of (20). Set for brevity  $s_0 = \lfloor \sqrt{d \log d} \rfloor + 1$ . Using the definition of  $\hat{s}$  we obtain, on the event  $\{\hat{s} \leq s\}$ ,

$$(\hat{L}_{\hat{s}} - L)^2 \leq 2\omega_s^2 + 2(\hat{L}_s - L)^2 \text{ if } s < s_0 \text{ or } s \geq s_0, \hat{s} < s_0.$$

Thus,

$$(21) \quad \forall s < s_0 : \quad \mathbf{E}_\theta [(\hat{L}_{\hat{s}} - L)^2 \mathbf{1}_{\hat{s} \leq s}] \leq 2\beta^2 \Phi^L(\sigma, s) + 2\mathbf{E}_\theta (\hat{L}_s - L)^2,$$

$$(22) \quad \forall s \geq s_0 : \quad \begin{aligned} \mathbf{E}_\theta [(\hat{L}_{\hat{s}} - L)^2 \mathbf{1}_{\hat{s} \leq s}] &\leq \mathbf{E}_\theta [(\hat{L}_{\hat{s}} - L)^2 (\mathbf{1}_{\hat{s} \leq s, \hat{s} < s_0} + \mathbf{1}_{\hat{s} = s_0})] \\ &\leq 2\beta^2 \Phi^L(\sigma, s) + 2\mathbf{E}_\theta (\hat{L}_s - L)^2 + \mathbf{E}_\theta (\hat{L}_{s_0} - L)^2. \end{aligned}$$

By Lemma 6 proved at the end of this section, for any tuning constant  $\alpha > 0$  large enough we have

$$\sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\hat{L}_s - L)^2 \leq c_1 \Phi^L(\sigma, s), \quad s = 1, \dots, s_0 - 1.$$

Note that, in view of (3), for all  $s \in [s_0, d]$  we have

$$\Phi^L(\sigma, s_0) \leq \sigma^2 d \log d \leq 2\sigma^2 s^2 \log(1 + (d \log d)/s^2) = 2\Phi^L(\sigma, s),$$

and by definition of  $\hat{L}_s$ , for all  $s \in [s_0, d]$  and all  $\theta \in \mathbb{R}^d$ , we have  $\mathbf{E}_\theta(\hat{L}_s - L)^2 = \sigma^2 d \leq 2\Phi^L(\sigma, s)$ . Combining these remarks with (21) and (22) yields

$$(23) \quad \sup_{\theta \in \Theta_s} \mathbf{E}_\theta[(\hat{L}_s - L)^2 \mathbf{1}_{\hat{s} \leq s}] \leq c_2 \Phi^L(\sigma, s), \quad s = 1, \dots, d.$$

Consider now the second summand on the right hand side of (20). Since  $\hat{s} \leq s_0$  we obtain the following two facts. First,

$$(24) \quad \sup_{\theta \in \Theta_s} \mathbf{E}_\theta[(\hat{L}_s - L)^2 \mathbf{1}_{\hat{s} > s}] = 0, \quad \forall s \geq s_0.$$

Second, on the event  $\{\hat{s} > s\}$ ,

$$(\hat{L}_s - L)^4 \leq \sum_{s < s' \leq s_0} (\hat{L}_{s'} - L)^4.$$

Thus,

$$(25) \quad \begin{aligned} \sup_{\theta \in \Theta_s} \mathbf{E}_\theta[(\hat{L}_s - L)^2 \mathbf{1}_{\hat{s} > s}] &\leq \sup_{\theta \in \Theta_s} \left[ \sqrt{\mathbf{P}_\theta(\hat{s} > s)} (d \log d)^{1/4} \max_{s < s' \leq s_0} \sqrt{\mathbf{E}_\theta(\hat{L}_{s'} - L)^4} \right] \\ &\leq (d \log d)^{1/4} \sup_{\theta \in \Theta_s} \sqrt{\mathbf{P}_\theta(\hat{s} > s)} \max_{s' \leq s_0} \left[ \sup_{\theta \in \Theta_{s'}} \sqrt{\mathbf{E}_\theta(\hat{L}_{s'} - L)^4} \right] \end{aligned}$$

where for the second inequality we have used that  $\Theta_s \subset \Theta_{s'}$  for  $s < s'$ . To evaluate the right hand side of (25) we use the following two lemmas.

**Lemma 1.** *For all  $s \leq s_0$  we have*

$$\sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\hat{L}_s - L)^4 \leq c_3 \sigma^4 d^4 (\log d)^2, \quad \sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\hat{L}'_s - L)^4 \leq c_3 \sigma^4 d^4 (\log d)^2.$$

**Lemma 2.** *There exist absolute constants  $c_4 > 0$  and  $\alpha_0 > 0$  such that the following holds.*

(i) *For all  $\alpha > \alpha_0$  and  $\beta > 37\alpha$  we have*

$$(26) \quad \max_{s \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_s} \mathbf{P}_\theta(\hat{s} > s) \leq c_4 d^{-5}.$$

(ii) *There exist  $\alpha > \alpha_0$  and  $\beta > \alpha_0$  such that*

$$\max_{s \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_s} \mathbf{P}_\theta(\hat{s}' > s) \leq c_4 d^{-5}.$$

From (24), (25), the first inequality in Lemma 1, and part (i) of Lemma 2 we find that

$$\sup_{\theta \in \Theta_s} \mathbf{E}_\theta[(\hat{L}_s - L)^2 \mathbf{1}_{\hat{s} > s}] \leq c_5 \sigma^2 \leq c_6 \Phi^L(\sigma, s), \quad s = 1, \dots, d.$$

Combining this inequality with (20) and (23) we obtain the theorem.



#### 4.2. Proofs of the lemmas.

*Proof of Lemma 1.* For  $s = s_0$ ,  $\hat{L}_s - L = \hat{L}'_s - L = \sigma \sum_{i=1}^d \xi_i$ . As a consequence

$$\mathbf{E}_\theta(\hat{L}_s - L)^4 = \mathbf{E}_\theta(\hat{L}'_s - L)^4 = 3\sigma^4 d^2 \leq 3\sigma^4 d^4 (\log d)^2.$$

Henceforth, we focus on the case  $s \leq \sqrt{d \log(d)}$ . We have

$$(27) \quad \hat{L}_s - L = \sigma \sum_{i=1}^d \xi_i - \sum_{i=1}^d y_i \mathbb{1}_{y_i^2 \leq \alpha \sigma^2 \log(1+d(\log d)/s^2)}.$$

Thus,

$$\mathbf{E}_\theta(\hat{L}_s - L)^4 \leq 8 \left( \sigma^4 \mathbf{E} \left( \sum_{i=1}^d \xi_i \right)^4 + d^4 \alpha^2 \sigma^4 \log^2(1 + d(\log d)/s^2) \right) \leq c_3 \sigma^4 d^4 (\log d)^2.$$

In a similar way,

$$(28) \quad \hat{L}'_s - L = \sigma \sum_{i=1}^d \xi_i - \sum_{i=1}^d y_i \mathbb{1}_{y_i^2 \leq \alpha \hat{\sigma}^2 \log(1+d(\log d)/s^2)},$$

and

$$\mathbf{E}_\theta(\hat{L}'_s - L)^4 \leq 8 \left( \sigma^4 \mathbf{E} \left( \sum_{i=1}^d \xi_i \right)^4 + d^4 \alpha^2 \mathbf{E}_\theta(\hat{\sigma}^4) \log^2(1 + d(\log d)/s^2) \right).$$

The desired bound for  $\mathbf{E}_\theta(\hat{L}'_s - L)^4$  follows from this inequality and (16).  $\square$

*Proof of Lemma 2.* We start by proving part (i) of Lemma 2. Note first that, for  $s \leq \sqrt{d \log d}$  and all  $\theta$  we have

$$(29) \quad \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_{s'}/2) \leq \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_s/3), \quad \forall s < s' \leq d.$$

Indeed, if  $s < s'$  we have  $\omega_{s'} > \omega_s$  since the function  $t \mapsto \omega_t$  is increasing for  $t > 0$ . If  $\sqrt{d \log d} \leq s' \leq d$ , we use (4), which yields  $\omega_{s'}^2 \geq \beta \sigma^2 d(\log d)/2 \geq \omega_{\sqrt{d \log d}}^2/2$  and since  $s \leq \sqrt{d \log d}$  we obtain  $\omega_{s'}^2 \geq \omega_s^2/2$  using again the fact that the function  $t \mapsto \omega_t$  is increasing.

From (29) we obtain that, for  $s \leq \sqrt{d \log d}$ , all  $s'$  such that  $s < s' \leq d$  and all  $\theta$ ,

$$\mathbf{P}_\theta(|\hat{L}_{s'} - \hat{L}_s| > \omega_{s'}) \leq \mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/2) + \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_s/3).$$

This inequality and the definition of  $\hat{s}$  imply that, for all  $s \leq \sqrt{d \log d}$  and all  $\theta$ ,

$$(30) \quad \begin{aligned} \mathbf{P}_\theta(\hat{s} > s) &\leq \sum_{s < s' \leq d} \mathbf{P}_\theta(|\hat{L}_{s'} - \hat{L}_s| > \omega_{s'}) \\ &\leq d \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_s/3) + \sum_{s < s' \leq d} \mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/2). \end{aligned}$$

Note that, for  $\sqrt{d \log d} < s' \leq d$ , we have  $\hat{L}_{s'} = \sum_{i=1}^d y_i$ , and  $\omega_{s'} \geq \sigma \sqrt{\beta d \log d}/\sqrt{2}$  due to (4). Hence, for  $\sqrt{d \log d} < s' \leq d$ , and all  $\theta$ ,

$$\mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/2) \leq \mathbf{P} \left( \left| \sum_{i=1}^d \xi_i \right| > \frac{\sqrt{\beta d \log d}}{2\sqrt{2}} \right) \leq 2d^{-\beta/16},$$

where we have used that  $\xi_i$  are i.i.d. standard Gaussian random variables. This inequality and (30) imply that, for  $s \leq \sqrt{d \log d}$ , and all  $\theta$ ,

$$(31) \quad \begin{aligned} \mathbf{P}_\theta(\hat{s} > s) &\leq \sqrt{d \log d} \max_{s < s' \leq \sqrt{d \log d}} \mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/2) \\ &\quad + d \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_s/3) + 2d^{1-\beta/16}. \end{aligned}$$

As  $\Theta_s \subset \Theta_{s'}$  for  $s < s'$ , we have

$$\max_{s < s' \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_s} \mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/2) \leq \max_{s' \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_{s'}} \mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/2).$$

Together with (31) this implies

$$\max_{s \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_s} \mathbf{P}_\theta(\hat{s} > s) \leq 2d \max_{s' \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_{s'}} \mathbf{P}_\theta(|\hat{L}_{s'} - L| > \omega_{s'}/3) + 2d^{1-\beta/16}.$$

For  $\beta > 96$  the last summand in this inequality does not exceed  $2d^{-5}$ . Thus, to prove (26) it is enough to show that

$$(32) \quad \max_{s \leq \sqrt{d \log d}} \sup_{\theta \in \Theta_s} \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_s/3) \leq c_7 d^{-6}.$$

The rest of this proof consists in demonstrating that (32) is satisfied if the tuning constants  $\alpha$  and  $\beta$  are properly chosen. Fix  $s \leq \sqrt{d \log d}$  and let  $\theta$  belong to  $\Theta_s$ . We will denote by  $S$  the support of  $\theta$  and we set for brevity

$$a \triangleq \sqrt{\log(1 + d(\log d)/s^2)}.$$

From (27) and the fact that  $y_i = \theta_i + \sigma \xi_i$  we have

$$(33) \quad \begin{aligned} |\hat{L}_s - L| &= \left| \sigma \sum_{i \in S} \xi_i - \sum_{i \in S} y_i \mathbf{1}_{y_i^2 \leq \alpha \sigma^2 a^2} + \sigma \sum_{i \notin S} \xi_i \mathbf{1}_{\xi_i^2 > \alpha a^2} \right| \\ &\leq \sigma \left| \sum_{i \in S} \xi_i \right| + \sigma \left| \sum_{i \notin S} \xi_i \mathbf{1}_{\xi_i^2 > \alpha a^2} \right| + \sqrt{\alpha} \sigma s a. \end{aligned}$$

In the following we assume that  $\beta > 37\alpha$ . Using this and recalling that  $\omega_s = \sqrt{\beta} \sigma s a$  we find

$$(34) \quad \begin{aligned} \mathbf{P}_\theta(|\hat{L}_s - L| > \omega_s/3) &\leq \mathbf{P}\left(\left| \sum_{i \notin S} \xi_i \mathbf{1}_{\xi_i^2 > \alpha a^2} \right| > \sqrt{\beta} s a / 6\right) + \mathbf{P}\left(\left| \sum_{i \in S} \xi_i \right| > (\sqrt{\beta}/6 - \sqrt{\alpha}) s a\right) \\ &\leq \mathbf{P}\left(\left| \sum_{i \notin S} \xi_i \mathbf{1}_{\xi_i^2 > \alpha a^2} \right| > \sqrt{\alpha} s a\right) + \mathbf{P}\left(\left| \sum_{i \in S} \xi_i \right| > c_8 \sqrt{\alpha} s a\right). \end{aligned}$$

Since  $\xi_i$  are i.i.d.  $\mathcal{N}(0, 1)$  random variables, we have

$$(35) \quad \mathbf{P}\left(\left| \sum_{i \in S} \xi_i \right| > c_8 \sqrt{\alpha} s a\right) \leq 2 \exp\left(-\frac{c_8^2}{2} \alpha s a^2\right).$$

We now use the relation

$$(36) \quad s a^2 = s \log(1 + d(\log d)/s^2) \geq (\log d)/2 \quad \text{for all } s \in [1, \sqrt{d \log d}].$$

Indeed, for  $s \in [1, \sqrt{d \log d}/3]$  the left hand side of (36) is monotone increasing in  $s$ , while for  $s \in [\sqrt{d \log d}/3, \sqrt{d \log d}]$  the inequality in (36) is trivial. It follows from (35) and (36) that,

for all  $\beta > 37\alpha$  and all suitably large  $\alpha$ ,

$$(37) \quad \mathbf{P}\left(\left|\sum_{i \in S} \xi_i\right| > c_8 \sqrt{\alpha} s a\right) \leq d^{-6}.$$

Next, consider the first probability on the right hand side of (34). To bound it from above, we invoke the following lemma.

**Lemma 3.** *For any absolute constant  $\alpha > 0$  large enough, for all  $s \leq \sqrt{d \log d}$ , and all  $U \subseteq \{1, \dots, d\}$ ,*

$$\mathbf{P}\left(\sup_{t \in [1, 10]} \left|\sum_{i \in U} \xi_i \mathbf{1}_{|\xi_i| > \sqrt{\alpha} a t}\right| > \sqrt{\alpha} s a\right) \leq c_9 d^{-6}.$$

Combining (34), (37) and Lemma 3 we obtain (32). Thus, part (i) of Lemma 2 follows.

We now proceed to the proof of part (ii) of Lemma 2. Proposition 1 implies that, for  $s \leq \sqrt{d \log d}$  and  $\theta \in \Theta_s$ ,

$$\mathbf{P}_\theta(\hat{s}' > s) \leq \mathbf{P}_\theta(\hat{s}' > s, \hat{\sigma} \in [\sigma, 10\sigma]) + d^{-5}.$$

On the event  $\{\hat{\sigma} \in [\sigma, 10\sigma]\}$ , we can replace  $\hat{\sigma}$  in the definition of  $\hat{s}'$  either by  $\sigma$  or by  $10\sigma$  according to cases, thus making the analysis of  $\mathbf{P}_\theta(\hat{s}' > s, \hat{\sigma} \in [\sigma, 10\sigma])$  equivalent, up to the values of numerical constants, to the analysis of  $\mathbf{P}_\theta(\hat{s} > s)$  given below. The only non-trivial difference consists in the fact that the analog of (33) when  $\hat{L}_s$  is replaced by  $\hat{L}'_s$  contains the term  $\sigma \left| \sum_{i \notin S} \xi_i \mathbf{1}_{\xi_i^2 > \alpha \hat{\sigma}^2 a^2 / \sigma^2} \right|$  instead of  $\sigma \left| \sum_{i \notin S} \xi_i \mathbf{1}_{\xi_i^2 > \alpha a^2} \right|$  while  $\hat{\sigma}$  depends on  $\xi_1, \dots, \xi_d$ . This term is evaluated using Lemma 3 and the fact that

$$\mathbf{P}\left(\left|\sum_{i \notin S} \xi_i \mathbf{1}_{|\xi_i| > \sqrt{\alpha} \hat{\sigma} a / \sigma}\right| > \sqrt{\alpha} s a, \hat{\sigma} \in [\sigma, 10\sigma]\right) \leq \mathbf{P}\left(\sup_{t \in [1, 10]} \left|\sum_{i \notin S} \xi_i \mathbf{1}_{|\xi_i| > \sqrt{\alpha} a t}\right| > \sqrt{\alpha} s a\right).$$

We omit further details that are straightforward from inspection of the proof of part (i) of Lemma 2 given above. Thus, part (ii) of Lemma 2 follows.  $\square$

For the proof of Lemma 3, recall the following fact about the tails of the standard Gaussian distribution.

**Lemma 4.** *Let  $X \sim \mathcal{N}(0, 1)$ ,  $x > 1$  and  $q \in \mathbb{N}$ . There exists a constant  $C_q^*$  depending only on  $q$  such that*

$$\mathbf{E}\left[X^{2q} \mathbf{1}_{|X| > x}\right] \leq C_q^* x^{2q-1} e^{-x^2/2}.$$

We will also use the Fuk-Nagaev inequality [9, page 78] that we state here for reader's convenience.

**Lemma 5** (Fuk-Nagaev inequality). *Let  $p > 2$  and  $v > 0$ . Assume that  $X_1, \dots, X_n$  are independent random variables with  $\mathbf{E}(X_i) = 0$  and  $\mathbf{E}|X_i|^p < \infty$ ,  $i = 1, \dots, n$ . Then,*

$$\mathbf{P}\left(\sum_{i=1}^n X_i > v\right) \leq (1 + 2/p)^p \sum_{i=1}^n \mathbf{E}|X_i|^p v^{-p} + \exp\left(-\frac{2v^2}{(p+2)^2 e^p \sum_{i=1}^n \mathbf{E}X_i^2}\right).$$

*Proof of Lemma 3.* We have

$$\begin{aligned} p_0 &\triangleq \mathbf{P}\left(\sup_{t \in [1, 10]} \left|\sum_{i \in U} \xi_i \mathbf{1}_{|\xi_i| > \sqrt{\alpha} a t}\right| > \sqrt{\alpha} s a\right) \\ &= \mathbf{E}\left[\mathbf{P}\left(\sup_{t \in [1, 10]} \left|\sum_{i \in U} \epsilon_i |\xi_i| \mathbf{1}_{|\xi_i| > \sqrt{\alpha} a t}\right| > \sqrt{\alpha} s a \mid |\xi_i|, i \in U\right)\right] \end{aligned}$$

where  $\epsilon_i$  denotes the sign of  $\xi_i$ . Consider the function  $g(x) = \sup_{t \in [1,10]} \left| \sum_{i \in U} x_i |\xi_i| \mathbb{1}_{|\xi_i| > \sqrt{\alpha} a t} \right|$  where  $x = (x_i, i \in U)$  with  $x_i \in \{-1, 1\}$ . For any  $i_0 \in U$ , let  $g_{i_0, u}(x)$  denote the value of this function when we replace  $x_{i_0}$  by  $u \in \{-1, 1\}$ . Note that, for any fixed  $(|\xi_i|, i \in U)$ , we have the bounded differences condition:

$$\sup_x |g(x) - g_{i_0, u}(x)| \leq 2|\xi_{i_0}| \mathbb{1}_{|\xi_{i_0}| > \sqrt{\alpha} a} \triangleq 2Z_{i_0} \quad \forall u \in \{-1, 1\}, i_0 \in U.$$

The vector of Rademacher random variables  $(\epsilon_1, \dots, \epsilon_d)$  is independent from  $(|\xi_1|, \dots, |\xi_d|)$ . Thus, for any fixed  $(|\xi_i|, i \in U)$  we can use the bounded differences inequality, which yields

$$(38) \quad p_0 \leq \mathbf{E} \left[ \exp \left( - \frac{\alpha s^2 a^2}{2 \sum_{i \in U} Z_i^2} \right) \right] \leq \exp \left( - \frac{\alpha s^2 a^2}{2\Delta} \right) + \mathbf{P} \left( \sum_{i \in U} Z_i^2 > \Delta \right), \quad \forall \Delta > 0.$$

We now set  $\Delta = \sum_{i \in U} \mathbf{E} Z_i^2 + d \exp(-\alpha a^2/(2p))$  for some integer  $p > 2$  that will be chosen large enough.

To finish the proof, it remains to show that each of the two summands on the right hand side of (38) does not exceed  $c_9 d^{-6}/2$  if  $p$  and  $\alpha$  are large enough. To bound from above the probability  $\mathbf{P} \left( \sum_{i \in U} Z_i^2 > \Delta \right)$  we apply Lemma 5 with  $X_i = Z_i^2 - \mathbf{E}(Z_i^2)$  and  $v = d \exp(-\alpha a^2/(2p))$ . The random variables  $X_i$  are centered and satisfy, in view of Lemma 4,

$$(39) \quad \mathbf{E}|X_i|^p \leq 2^{p-1} \mathbf{E}|Z_i|^{2p} \leq 2^{p-1} C_p^* (\sqrt{\alpha} a)^{2p-1} e^{-\alpha a^2/2}$$

for any  $p \in \mathbb{N}^*$ . Thus, Lemma 5 yields

$$\mathbf{P} \left( \sum_{i \in U} Z_i^2 > \Delta \right) \leq C_p^* 2^{p-1} (1 + 2/p)^p \frac{(\sqrt{\alpha} a)^{2p-1}}{d^p} + \exp \left( - \frac{d \exp(\alpha a^2(1/2 - 1/p))}{(p+2)^2 e^p C_2^* (\sqrt{\alpha} a)^3} \right).$$

The expression in the last display can be rendered smaller than  $c_9 d^{-6}/2$  for some absolute constant  $c_9 > 0$  and all  $d \geq 3$  by choosing  $p$  large enough.

Finally, using (39) we find

$$\frac{\alpha s^2 a^2}{2\Delta} \geq \frac{\alpha s^2 a^2}{2d(C_1^* \sqrt{\alpha} a \exp(-\alpha a^2/2) + \exp(-\alpha a^2/(2p)))} \geq \frac{c_{10} \sqrt{\alpha} a s^2 \exp(\alpha a^2/(2p))}{d},$$

whereas

$$\frac{s^2 \exp(\alpha a^2/(2p))}{d} = \frac{s^2}{d} \left( 1 + \frac{d \log d}{s^2} \right)^{\alpha/(2p)} \geq \log d \left( \frac{d \log d}{s^2} \right)^{\alpha/(2p)-1} \geq \log d$$

for any  $\alpha \geq 2p$  and any  $s \leq \sqrt{d \log d}$ . Hence, for such  $\alpha$  and  $s$ ,

$$\exp \left( - \frac{\alpha s^2 a^2}{2\Delta} \right) \leq \exp(-c_{10} \sqrt{\alpha} a \log d)$$

that is smaller than  $c_9 d^{-6}/2$  for all  $d \geq 3$  provided that the absolute constant  $\alpha$  is large enough. Thus, Lemma 3 follows.  $\square$

**Lemma 6.** *There exists an absolute constant  $c_1 > 0$  such that, for all  $\alpha > 0$  large enough,*

$$\sup_{\theta \in \Theta_s} \mathbf{E}_\theta (\hat{L}_s - L)^2 \leq c_1 \Phi^L(\sigma, s), \quad \sup_{\theta \in \Theta_s} \mathbf{E}_\theta (\hat{L}'_s - L)^2 \leq c_1 \Phi^L(\sigma, s), \quad \forall s \leq \sqrt{d \log d}.$$

*Proof.* We easily deduce from (33) that

$$\mathbf{E}_\theta(\hat{L}_s - L)^2 \leq 3\sigma^2 \left( s + d\mathbf{E} \left[ X^2 \mathbf{1}_{X^2 > \alpha a^2} \right] + \alpha s^2 a^2 \right),$$

where  $X \sim \mathcal{N}(0, 1)$ . By (39),

$$d\mathbf{E} \left[ X^2 \mathbf{1}_{X^2 > 2a^2} \right] \leq 2C_1^* ad \exp(-a^2) = \frac{2C_1^* ad s^2}{s^2 + d \log d} \leq \frac{2C_1^* s^2 a}{\log d},$$

which implies that the desired bound for  $\mathbf{E}_\theta(\hat{L}_s - L)^2$  holds whenever  $\alpha \geq 2$ . Next, we prove the bound of the lemma for  $\mathbf{E}_\theta(\hat{L}'_s - L)^2$ . Similarly to (33),

$$\hat{L}'_s - L = \sigma \sum_{i \in S} \xi_i - \sum_{i \in S} y_i \mathbf{1}_{y_i^2 \leq \alpha \hat{\sigma}^2 a^2} + \sigma \sum_{i \notin S} \xi_i \mathbf{1}_{\sigma^2 \xi_i^2 > \alpha \hat{\sigma}^2 a^2}.$$

This implies

$$\begin{aligned} (40) \quad \mathbf{E}_\theta \left[ (\hat{L}'_s - L)^2 \mathbf{1}_{\hat{\sigma} \in [\sigma, 10\sigma]} \right] &\leq \mathbf{E}_\theta \left( \sigma \left| \sum_{i \in S} \xi_i \right| + \sqrt{\alpha} \hat{\sigma} s a + \sigma W \right)^2 \\ &\leq 3 \left( \sigma^2 s + \alpha \mathbf{E}_\theta(\hat{\sigma}^2) a^2 s^2 + \sigma^2 \mathbf{E}(W^2) \right), \end{aligned}$$

where  $W \triangleq \sup_{t \in [1, 10]} \left| \sum_{i \notin S} \xi_i \mathbf{1}_{|\xi_i| > \sqrt{\alpha} a t} \right|$ . Using Lemma 3 we find that, for all  $\alpha > 0$  large enough,

$$\begin{aligned} \mathbf{E}(W^2) &\leq (\sqrt{\alpha} s a)^2 + \mathbf{E} \left( \sum_{i \notin S} |\xi_i| \right)^2 \mathbf{1}_{W > \sqrt{\alpha} s a} \\ &\leq \alpha s^2 a^2 + \left[ \mathbf{E} \left( \sum_{i \notin S} |\xi_i| \right)^4 \right]^{1/2} c_9 d^{-3} \leq \alpha s^2 a^2 + c_9 \sqrt{3} d^{-1}. \end{aligned}$$

Plugging this bound in (40) and using (16) we get

$$\mathbf{E}_\theta \left[ (\hat{L}'_s - L)^2 \mathbf{1}_{\hat{\sigma} \in [\sigma, 10\sigma]} \right] \leq c_{11} \Phi^L(\sigma, s).$$

On the other hand, by virtue of Lemma 1 and (15),

$$\mathbf{E}_\theta \left[ (\hat{L}'_s - L)^2 \mathbf{1}_{\hat{\sigma} \notin [\sigma, 10\sigma]} \right] \leq \sqrt{\mathbf{P}_\theta(\hat{\sigma} \notin [\sigma, 10\sigma])} \sqrt{\mathbf{E}_\theta(\hat{L}'_s - L)^4} \leq \frac{\sqrt{c_3} \sigma^2 \log d}{d^{1/2}} \leq c_{12} \Phi^L(\sigma, s).$$

The desired bound for  $\mathbf{E}_\theta(\hat{L}'_s - L)^2$  follows from the last two displays.  $\square$

#### 4.3. Proofs of Proposition 1 and of Theorem 3.

*Proof of Proposition 1.* Since  $s \leq d/2$ , there exists a subset  $T$  of size  $\lfloor d/2 \rfloor$  such that  $T \cap S = \emptyset$ . By Definition of  $\hat{\sigma}^2$ , we obtain that

$$\hat{\sigma}^2 \leq \frac{81\sigma^2}{\lfloor d/2 \rfloor} \sum_{i \in T} \xi_i^2.$$

This immediately implies (16). To prove (15), note that the Gaussian concentration inequality (cf. [7]) yields

$$\mathbf{P} \left( \left( \sum_{i \in T} \xi_i^2 \right)^{1/2} > \sqrt{100 \lfloor d/2 \rfloor / 81} \right) \leq \exp(-cd),$$

for a positive constant  $c$ . Therefore,

$$(41) \quad \mathbf{P}_\theta(\hat{\sigma} \leq 10\sigma) \geq 1 - \exp(-cd).$$

Next, let  $\mathcal{G}$  be the collection of all subsets of  $\{1, \dots, d\}$  of cardinality  $\lfloor d/2 \rfloor$ . We now establish a bound on the deviations of random variables  $Z_G = \frac{1}{\sigma^2} \sum_{i \in G} y_i^2$  uniformly over all  $G \in \mathcal{G}$ . Fix any  $G \in \mathcal{G}$ . The random variable  $Z_G$  has a chi-square distribution with  $\lfloor d/2 \rfloor$  degrees of freedom and non-centrality parameter  $\sum_{i \in G} \theta_i^2$ . In particular, this distribution is stochastically larger than a central chi-square distribution with  $d' = \lfloor d/2 \rfloor$  degrees of freedom. Let  $Z$  be a random variable with this central chi-square distribution. By Lemma 11.1 in [12],

$$\mathbf{P}\left(Z \leq \frac{d'}{e} x^{2/d'}\right) \leq x, \quad \forall x > 0.$$

Take  $x = \left(\frac{d}{d'}\right)^{-1} e^{-d'/2}$ . Using the bound  $\log\left(\frac{d}{d'}\right) \leq d' \log(ed/d')$  it follows that  $\log(1/x) \leq d'(\frac{3}{2} + \log(\frac{d}{d'})) \leq d'(\frac{3}{2} + \log 2) + 1$ . Taking the union bound over all  $G \in \mathcal{G}$  we conclude that

$$\mathbf{P}\left(\inf_{G \in \mathcal{G}} Z_G \leq \frac{d'}{4e^3} \left(1 - \frac{2}{d'}\right)\right) \leq e^{-d'/2} < d^{-5}/2$$

for all  $d$  large enough. Since  $\hat{\sigma}^2 = \sigma^2 \frac{81}{d'} \inf_{G \in \mathcal{G}} Z_G^2$ , we obtain that  $\hat{\sigma}^2 \geq \sigma^2$  with probability at least  $1 - d^{-5}/2$  for all  $d$  large enough. Combining this with (41), we get (15) for all  $d$  large enough.  $\square$

*Proof of Theorem 3.* We repeat the proof of Theorem 1 replacing there  $\hat{L}_s$  by  $\hat{L}'_s$  and  $\hat{s}$  by  $\hat{s}'$ . The difference is that, in view of (16), the relation (21) now holds with  $c_{15}\beta^2\Phi^L(\sigma, s)$  instead of  $\beta^2\Phi^L(\sigma, s)$ , and we use the results of Lemmas 1, 2 and 6 related to  $\hat{L}'_s$  rather than to  $\hat{L}_s$ .  $\square$

## 5. PROOFS OF THE LOWER BOUNDS

**5.1. Proof of Theorem 2.** Theorem 2 is an immediate consequence of the following lemma.

**Lemma 7.** *For all  $d \geq 3$  and all  $s \geq d^{1/4}$ ,*

$$(42) \quad R(s) \triangleq \inf_{\tilde{L}} \left\{ \sup_{\theta \in \Theta_1} \mathbf{E}_\theta(\tilde{L} - L)^2 \sigma^{-2} d^{-1/4} + \sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\tilde{L} - L)^2 (\Phi^L(\sigma, s))^{-1} \right\} \geq \frac{1}{160}.$$

*Proof.* We first introduce some notation. For a probability measure  $\mu$  on  $\Theta_s$ , we denote by  $\mathbb{P}_\mu$  the mixture probability measure  $\mathbb{P}_\mu = \int_{\Theta_s} \mathbf{P}_\theta \mu(d\theta)$ . Let  $\mathcal{S}(s, d)$  denote the set of all subsets of  $\{1, \dots, d\}$  of size  $s$ , and let  $S$  be a set-valued random variable uniformly distributed on  $\mathcal{S}(s, d)$ . For any  $\rho > 0$ , denote by  $\mu_\rho$  the distribution of the random variable  $\sigma\rho \sum_{j \in S} e_j$  where  $e_j$  is the  $j$ th canonical basis vector in  $\mathbb{R}^d$ . Next, let  $\chi^2(Q, P) = \int (dQ/dP)^2 dP - 1$  denote the chi-square divergence between two probability measures  $Q$  and  $P$  such that  $Q \ll P$ , and  $\chi^2(Q, P) = +\infty$  if  $Q \not\ll P$ .

Take any  $s \geq d^{1/4}$  and set

$$\rho \triangleq \sqrt{\log(1 + d(\log d)/s^2)}/2 = (\Phi^L(\sigma, s))^{1/2}/(2s\sigma).$$

Consider the mixture distribution  $\mathbb{P}_{\mu_\rho}$  with this value of  $\rho$ . For any estimator  $\tilde{L}$ , we have  $\sup_{\theta \in \Theta_s} \mathbf{E}_\theta(\tilde{L} - L)^2 \geq \mathbb{E}_{\mu_\rho}(\tilde{L} - L)^2 \geq \mathbb{E}_{\mu_\rho}(\tilde{L} - \mathbb{E}_{\mu_\rho}(L))^2 = \mathbb{E}_{\mu_\rho}(\tilde{L} - \sigma s \rho)^2$ . Therefore,

$$\begin{aligned}
 R(s) &\geq \inf_{\tilde{L}} \left\{ \mathbf{E}_0(\tilde{L}^2) \sigma^{-2} d^{-1/4} + \mathbb{E}_{\mu_\rho}(\tilde{L} - \sigma s \rho)^2 (\Phi^L(\sigma, s))^{-1} \right\} \\
 &\geq \frac{1}{16} \inf_{\tilde{L}} \left\{ \mathbf{P}_0(\tilde{L} > \sigma s \rho / 2) \sigma^{-2} d^{-1/4} \Phi^L(\sigma, s) + \mathbb{P}_{\mu_\rho}(\tilde{L} < \sigma s \rho / 2) \right\} \\
 (43) \quad &\geq \frac{1}{16} \inf_{\mathcal{A}} \left\{ \mathbf{P}_0(\mathcal{A}) \sigma^{-2} d^{-1/4} \Phi^L(\sigma, s) + \mathbb{P}_{\mu_\rho}(\mathcal{A}^c) \right\},
 \end{aligned}$$

where  $\inf_{\mathcal{A}}$  denotes the infimum over all measurable events  $\mathcal{A}$ , and  $\mathcal{A}^c$  denotes the complement of  $\mathcal{A}$ . It remains to prove that the expression in (43) is not smaller than  $1/160$ . This will be deduced from the following lemma, the proof of which is given at the end of this section.

**Lemma 8.** *Let  $P$  and  $Q$  be two probability measures on a measurable space  $(X, \mathcal{U})$ . Then, for any  $q > 0$ ,*

$$\inf_{\mathcal{A} \in \mathcal{U}} \{P(\mathcal{A})q + Q(\mathcal{A}^c)\} \geq \max_{0 < \tau < 1} \left[ \frac{q\tau}{1 + q\tau} (1 - \tau(\chi^2(Q, P) + 1)) \right].$$

We now apply Lemma 8 with  $P = \mathbf{P}_0$ ,  $Q = \mathbb{P}_{\mu_\rho}$ , and

$$(44) \quad q = \sigma^{-2} d^{-1/4} \Phi^L(\sigma, s) = s^2 d^{-1/4} \log \left( 1 + \frac{d(\log d)}{s^2} \right).$$

By Lemma 1 in [3], the chi-square divergence  $\chi^2(\mathbb{P}_{\mu_\rho}, \mathbf{P}_0)$  satisfies

$$\chi^2(\mathbb{P}_{\mu_\rho}, \mathbf{P}_0) \leq \left( 1 - \frac{s}{d} + \frac{s}{d} e^{\rho^2} \right)^s - 1 \leq \left( 1 + \frac{s}{d} (e^{\rho^2} - 1) \right)^s.$$

Since  $\rho^2 = \left( \log \left( 1 + \frac{d(\log d)}{s^2} \right) \right) / 4$ , we find

$$\begin{aligned}
 (45) \quad \chi^2(\mathbb{P}_{\mu_\rho}, \mathbf{P}_0) &\leq \exp \left[ s \log \left[ 1 + \frac{s}{d} \left( \left( 1 + \frac{d(\log d)}{s^2} \right)^{1/4} - 1 \right) \right] \right] \\
 &\leq \exp \left[ s \log \left( 1 + \frac{\log d}{4s} \right) \right] \leq d^{1/4},
 \end{aligned}$$

where we have used that  $(1 + x)^{1/4} \leq 1 + x/4$  for  $x > 0$ . Take

$$(46) \quad \tau = (d^{1/4} + 1)^{-1}/2.$$

Then, using (44) and the inequality  $s \geq d^{1/4}$  we find

$$(47) \quad q\tau = \frac{s^2 \log \left( 1 + \frac{d(\log d)}{s^2} \right)}{2d^{1/4}(d^{1/4} + 1)} \geq \frac{d^{1/2} \log(1 + d^{1/2}(\log d))}{2d^{1/4}(d^{1/4} + 1)} > \frac{1}{4}, \quad \forall d \geq 3.$$

Lemma 8 and inequalities (45) – (47) imply

$$\inf_{\mathcal{A}} \left\{ \mathbf{P}_0(\mathcal{A}) \sigma^{-2} d^{-1/4} \Phi^L(\sigma, s) + \mathbb{P}_{\mu_\rho}(\mathcal{A}^c) \right\} \geq \frac{q\tau}{2(1 + q\tau)} \geq \frac{1}{10}.$$

□

*Proof of Lemma 8.* We follow the same lines as in the proof of Proposition 2.4 in [11]. Thus, for any  $\tau \in (0, 1)$ ,

$$P(\mathcal{A}) \geq \tau(Q(\mathcal{A}) - v), \quad \text{where } v = Q\left(\frac{dP}{dQ} < \tau\right) \leq \tau(\chi^2(Q, P) + 1).$$

Then,

$$\begin{aligned} \inf_{\mathcal{A}} \{P(\mathcal{A})q + Q(\mathcal{A}^c)\} &\geq \inf_{\mathcal{A}} \{q\tau(Q(\mathcal{A}) - v) + Q(\mathcal{A}^c)\} \\ &\geq \min_{0 \leq t \leq 1} \max(q\tau(t - v), 1 - t) = \frac{q\tau(1 - v)}{1 + q\tau}. \end{aligned}$$

□

**5.2. Proof of Corollary 1.** First, note that condition (8) with  $\Psi_d(s) = C\Phi^L(\sigma, s)$  is satisfied due to Theorem 1. Next, the minimum in condition (10) with  $\Psi_d(s) = C\Phi^L(\sigma, s)$  can be only attained for  $s \geq d^{1/4}$ , since for  $s < d^{1/4}$  we have  $\Phi^L(\sigma, s) \asymp \psi_s^*$  where  $\psi_s^*$  is the minimax rate on  $\Theta_s$ . Thus, it is not possible to achieve a faster rate than  $\Phi^L(\sigma, s)$  for  $s < d^{1/4}$ , and therefore (10) is equivalent to the condition

$$\min_{s \geq d^{1/4}} \frac{\Psi'_d(s)}{\Phi^L(\sigma, s)} \rightarrow 0,$$

and

$$\min_{s=1, \dots, d} \frac{\Psi'_d(s)}{\Phi^L(\sigma, s)} \asymp \min_{s \geq d^{1/4}} \frac{\Psi'_d(s)}{\Phi^L(\sigma, s)}.$$

Obviously,  $\Psi'_d(s)$  cannot be of smaller order than the minimax rate  $\psi_s^*$ , which implies that

$$\min_{s \geq d^{1/4}} \frac{\Psi'_d(s)}{\Phi^L(\sigma, s)} \geq \min_{s \geq d^{1/4}} \frac{c\psi_s^*}{\Phi^L(\sigma, s)} = \min_{s \geq d^{1/4}} \frac{c \log(1 + d/s^2)}{\log(1 + d(\log d)/s^2)} \geq \frac{c'}{\log d}$$

where  $c, c' > 0$  are absolute constants. On the other hand, Theorem 2 yields

$$\frac{C'\Psi'_d(1)}{\Phi^L(\sigma, 1)} \geq \frac{C'C_1\sigma^2 d^{1/4}}{\Phi^L(\sigma, 1)} = \frac{C'C_1 d^{1/4}}{\log(1 + d(\log d))}.$$

Combining the last three displays, we find

$$\frac{\Psi'_d(1)}{\Phi^L(\sigma, 1)} \min_{s=1, \dots, d} \frac{\Psi'_d(s)}{\Phi^L(\sigma, s)} \geq \frac{c'C'C_1 d^{1/4}}{(\log d) \log(1 + d(\log d))} \rightarrow \infty,$$

as  $d \rightarrow \infty$ , thus proving (11) with  $\bar{s} = 1$ .

**5.3. Proof of Proposition 2.** Since in this proof we consider different values of  $\sigma$ , we denote the probability distribution of  $(y_1, \dots, y_d)$  satisfying (1) by  $\mathbf{P}_{\theta, \sigma^2}$ . Let  $\mathbf{E}_{\theta, \sigma^2}$  be the corresponding expectation. Assume that  $\hat{T}$  satisfies (18) with  $C_0 = 1/512$ . We will prove that (19) holds for  $\sigma = 1$ . The extension to arbitrary  $\sigma > 0$  is straightforward and is therefore omitted.

Let  $a > 1$  be a positive number and let  $\mu$  be the  $d$ -dimensional normal distribution with zero mean and covariance matrix  $a^2 \mathbf{I}_d$  where  $\mathbf{I}_d$  is the identity matrix. In what follows, we consider the mixture probability measure  $\mathbb{P}_\mu = \int_{\Theta_d} \mathbf{P}_{\theta, 1} \mu(d\theta)$ . Observe that  $\mathbb{P}_\mu = \mathbf{P}_{0, 1+a^2}$ .



Fixing  $\theta = 0$  and  $\sigma^2 = 1 + a^2$  in (18), we get  $\mathbf{E}_{0,1+a^2}[\widehat{T}^2] \leq 2C_0a^2d$  and therefore  $\mathbf{P}_{0,1+a^2}(|\widehat{T}| \geq \frac{1}{8}a\sqrt{d}) \leq \frac{1}{4}$ . Since  $\mathbb{P}_\mu = \mathbf{P}_{0,1+a^2}$ , this implies

$$(48) \quad \mathbb{P}_\mu(|\widehat{T}| < \frac{1}{8}a\sqrt{d}) > \frac{3}{4}.$$

For  $\theta$  distributed according to  $\mu$ ,  $L(\theta)$  has a normal distribution with mean 0 and variance  $a^2d$ . Hence, using the table of standard normal distribution, we find

$$\mu\left(|L(\theta)| \leq \frac{a}{4}\sqrt{d}\right) < \frac{1}{4}.$$

Combining this with (48), we conclude that, with  $\mathbb{P}_\mu$ -probability greater than  $1/2$ , we have simultaneously  $|L(\theta)| > a\sqrt{d}/4$  and  $|\widehat{T}| < a\sqrt{d}/8$ . Hence,

$$\sup_{\theta \in \Theta_d} \mathbf{E}_{\theta,1}[(\widehat{T} - L(\theta))^2] \geq \mathbb{E}_\mu[(\widehat{T} - L(\theta))^2] \geq \frac{1}{128}a^2d$$

where  $\mathbb{E}_\mu$  denotes the expectation with respect to  $\mathbb{P}_\mu$ . The result now follows by letting  $a$  tend to infinity.

**Acknowledgement.** The work of A.B.Tsybakov was supported by GENES and by the French National Research Agency (ANR) under the grants IPANEMA (ANR-13-BSH1-0004-02) and Labex Ecodec (ANR-11-LABEX-0047). It was also supported by the "Chaire Economie et Gestion des Nouvelles Données", under the auspices of Institut Louis Bachelier, Havas-Media and Paris-Dauphine.

#### REFERENCES

- [1] CAI, T. T. AND LOW, M.L. (2004). Minimax estimation of linear functionals over nonconvex parameter spaces. *Ann. Statist.* **32** 552–576.
- [2] CAI, T. T. AND LOW, M.L. (2005). On adaptive estimation of linear functionals. *Ann. Statist.* **33** 2311–2343.
- [3] COLLIER, O., COMMINGES, L., AND TSYBAKOV, A.B. (2016). Minimax estimation of linear and quadratic functionals under sparsity constraints. *Ann. Statist.*, to appear.
- [4] GOLUBEV, G.K. (2004). The method of risk envelopes in the estimation of linear functionals. *Problemy Peredachi Informatsii* **40** 58–72.
- [5] GOLUBEV, Y. AND LEVIT, B. (2004). An oracle approach to adaptive estimation of linear functionals in a Gaussian model. *Math. Methods Statist.* **13** 392–408.
- [6] IBRAGIMOV, I.A. AND HASMINSKII, R.Z. *Nonparametric estimation of the value of a linear functional in Gaussian white noise*. Theory Probab. Appl., 29, 18-32, 1984.
- [7] LEDOUX, M. AND TALAGRAND, M. (1991) *Probability in Banach Spaces*. Springer, Berlin, Heidelberg.
- [8] LAURENT, B., LUDENA, C. AND PRIEUR, C. (2008). Adaptive estimation of linear functionals by model selection. *Electron. J. Stat.* **2** 993–1020.
- [9] PETROV, V.V. (1995). *Limit Theorems of Probability Theory*. Clarendon Press, Oxford.
- [10] TSYBAKOV A.B. (1998). Pointwise and sup-norm sharp adaptive estimation of functions on the Sobolev classes. *Ann. Statist.* **26** 2420–2469.
- [11] TSYBAKOV A.B. (2009). *Introduction to Nonparametric Estimation*. Springer, New York.
- [12] VERZELEN, N. (2012). Minimax risks for sparse regressions: Ultra-high dimensional phenomena. *Electron. J. Stat.* **6** 38–90.