



Adaptive robust estimation in sparse vector model

Laëtitia Comminges, Olivier Collier, M Ndaoud, A. Tsybakov

► To cite this version:

Laëtitia Comminges, Olivier Collier, M Ndaoud, A. Tsybakov. Adaptive robust estimation in sparse vector model. 2019. hal-01707612v2

HAL Id: hal-01707612

<https://hal.science/hal-01707612v2>

Preprint submitted on 5 Oct 2019 (v2), last revised 7 Mar 2021 (v3)

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.

ADAPTIVE ROBUST ESTIMATION IN SPARSE VECTOR MODEL

BY L. COMMINGES¹, O. COLLIER², M. NDAOUD³ AND A.B. TSYBAKOV³¹*CEREMADE, Université Paris-Dauphine and CREST*²*Modal'X, UPL, Université Paris Nanterre and CREST*³*CREST, ENSAE*

Abstract For the sparse vector model, we consider estimation of the target vector, of its ℓ_2 -norm and of the noise variance. We construct adaptive estimators and establish the optimal rates of adaptive estimation when adaptation is considered with respect to the triplet "noise level – noise distribution – sparsity". We consider classes of noise distributions with polynomially and exponentially decreasing tails as well as the case of Gaussian noise. The obtained rates turn out to be different from the minimax non-adaptive rates when the triplet is known. A crucial issue is the ignorance of the noise variance. Moreover, knowing or not knowing the noise distribution can also influence the rate. For example, the rates of estimation of the noise variance can differ depending on whether the noise is Gaussian or sub-Gaussian without a precise knowledge of the distribution. Estimation of noise variance in our setting can be viewed as an adaptive variant of robust estimation of scale in the contamination model, where instead of fixing the "nominal" distribution in advance we assume that it belongs to some class of distributions.

1. Introduction. This paper considers estimation of the unknown sparse vector, of its ℓ_2 -norm and of the noise level in the sparse sequence model. The focus is on construction of estimators that are optimally adaptive in a minimax sense with respect to the noise level, to the form of the noise distribution, and to the sparsity.

We consider the model defined as follows. Let the signal $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$ be observed with noise of unknown magnitude $\sigma > 0$:

$$(1) \quad Y_i = \theta_i + \sigma \xi_i, \quad i = 1, \dots, d.$$

The noise random variables ξ_1, \dots, ξ_d are assumed to be i.i.d. and we denote by P_ξ the unknown distribution of ξ_1 . We assume throughout that the noise is zero-mean, $\mathbf{E}(\xi_1) = 0$, and that $\mathbf{E}(\xi_1^2) = 1$, since σ needs to be identifiable. We denote by $\mathbf{P}_{\boldsymbol{\theta}, P_\xi, \sigma}$ the distribution of $\mathbf{Y} = (Y_1, \dots, Y_d)$ when the signal is $\boldsymbol{\theta}$, the noise level is σ and the distribution of the noise variables is P_ξ . We also denote by $\mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma}$ the expectation with respect to $\mathbf{P}_{\boldsymbol{\theta}, P_\xi, \sigma}$.

We assume that the signal $\boldsymbol{\theta}$ is s -sparse, *i.e.*,

$$\|\boldsymbol{\theta}\|_0 = \sum_{i=1}^d \mathbf{1}_{\theta_i \neq 0} \leq s,$$

where $s \in \{1, \dots, d\}$ is an integer. Set $\Theta_s = \{\boldsymbol{\theta} \in \mathbb{R}^d \mid \|\boldsymbol{\theta}\|_0 \leq s\}$. We consider the problems of

Keywords and phrases: variance estimation, sparsity in linear regression, functional estimation, robust estimation, adaptive estimation, minimax rate

estimating $\boldsymbol{\theta}$ under the ℓ_2 loss, estimating the variance σ^2 , and estimating the ℓ_2 -norm

$$\|\boldsymbol{\theta}\|_2 = \left(\sum_{i=1}^d \theta_i^2 \right)^{1/2}.$$

The classical Gaussian sequence model corresponds to the case where the noise ξ_i is standard Gaussian ($P_\xi = \mathcal{N}(0, 1)$) and the noise level σ is known. Then, the optimal rate of estimation of $\boldsymbol{\theta}$ under the ℓ_2 loss in a minimax sense on the class Θ_s is $\sqrt{s \log(ed/s)}$ and it is attained by thresholding estimators [10]. Also, for the Gaussian sequence model with known σ , minimax optimal estimator of the norm $\|\boldsymbol{\theta}\|_2$ as well as the corresponding minimax rate are available from [8] (see Table 1).

In this paper, we study estimation of the three objects $\boldsymbol{\theta}$, $\|\boldsymbol{\theta}\|_2$, and σ^2 in the following two settings.

- (a) *The distribution of ξ_i and the noise level σ are both unknown.* This is the main setting of our interest. For the unknown distribution of ξ_i , we consider two types of assumptions. Either P_ξ belongs to a class $\mathcal{G}_{a,\tau}$, *i.e.*, for some $a, \tau > 0$,

$$(2) \quad P_\xi \in \mathcal{G}_{a,\tau} \quad \text{iff} \quad \mathbf{E}(\xi_1) = 0, \mathbf{E}(\xi_1^2) = 1 \text{ and } \forall t \geq 2, \mathbf{P}(|\xi_1| > t) \leq 2e^{-(t/\tau)^a},$$

which includes for example sub-Gaussian distributions ($a = 2$), or to a class of distributions with polynomially decaying tails $\mathcal{P}_{a,\tau}$, *i.e.*, for some $\tau > 0$ and $a \geq 2$,

$$(3) \quad P_\xi \in \mathcal{P}_{a,\tau} \quad \text{iff} \quad \mathbf{E}(\xi_1) = 0, \mathbf{E}(\xi_1^2) = 1 \text{ and } \forall t \geq 2, \mathbf{P}(|\xi_1| > t) \leq \left(\frac{\tau}{t}\right)^a.$$

We propose estimators of $\boldsymbol{\theta}$, $\|\boldsymbol{\theta}\|_2$, and σ^2 that are optimal in non-asymptotic minimax sense on these classes of distributions and the sparsity class Θ_s . We establish the corresponding non-asymptotic minimax rates. They are given in the second and third columns of Table 1. We also provide the minimax optimal estimators.

- (b) *Gaussian noise ξ_i and unknown σ .* The results on the non-asymptotic minimax rates are summarized in the first column of Table 1. Notice an interesting effect – the rates of estimation of σ^2 and of the norm $\|\boldsymbol{\theta}\|_2$ when the noise is Gaussian are faster than the optimal rates when the noise is sub-Gaussian. This can be seen by comparing the first column of Table 1 with the particular case $a = 2$ of the second column corresponding to sub-Gaussian noise.

Some comments about Table 1 and additional details are in order.

- The difference between the minimax rates for estimation of $\boldsymbol{\theta}$ and estimation of the ℓ_2 -norm $\|\boldsymbol{\theta}\|_2$ turns out to be specific for the pure Gaussian noise model. It disappears for the classes $\mathcal{G}_{a,\tau}$ and $\mathcal{P}_{a,\tau}$. This is somewhat unexpected since $\mathcal{G}_{2,\tau}$ is the class of sub-Gaussian distributions, and it turns out that $\|\boldsymbol{\theta}\|_2$ is estimated optimally at different rates for sub-Gaussian and pure Gaussian noise. Another conclusion is that if the noise is not Gaussian and σ is unknown, the minimax rate for $\|\boldsymbol{\theta}\|_2$ does not have an elbow between the "dense" ($s > \sqrt{d}$) and the "sparse" ($s \leq \sqrt{d}$) zones.
- For the problem of estimation of variance σ^2 with *known* distribution of the noise P_ξ , we consider a more general setting than (b) mentioned above. We show that when the noise distribution is exactly known (and satisfies a rather general assumption, not necessarily Gaussian - can have polynomial tails), then the rate of estimation of σ^2 can be as fast

	Gaussian noise model	Noise in class $\mathcal{G}_{a,\tau}$,	Noise in class $\mathcal{P}_{a,\tau}$,
$\boldsymbol{\theta}$	$\sqrt{s \log(ed/s)}$ known σ [10] unknown σ [22]	$\sqrt{s \log^{\frac{1}{a}}(ed/s)}$ unknown σ	$\sqrt{s(d/s)^{\frac{1}{a}}}$ unknown σ
$\ \boldsymbol{\theta}\ _2$	$\sqrt{s \log(1 + \frac{\sqrt{d}}{s})} \wedge d^{1/4}$ known σ [8] $\sqrt{s \log(1 + \frac{\sqrt{d}}{s})} \vee \sqrt{\frac{s}{1 + \log_+(s^2/d)}}$ unknown σ	$\sqrt{s \log^{\frac{1}{a}}(ed/s)} \wedge d^{1/4}$ known σ $\sqrt{s \log^{\frac{1}{a}}(ed/s)}$ unknown σ	$\sqrt{s(d/s)^{\frac{1}{a}}} \wedge d^{1/4}$ known σ $\sqrt{s(d/s)^{\frac{1}{a}}}$ unknown σ
σ^2	$\frac{1}{\sqrt{d}} \vee \frac{s}{d(1 + \log_+(s^2/d))}$	$\frac{1}{\sqrt{d}} \vee \frac{s}{d} \log^{\frac{2}{a}}\left(\frac{ed}{s}\right)$	$\frac{1}{\sqrt{d}} \vee \left(\frac{s}{d}\right)^{1-\frac{2}{a}}$

TABLE 1
Optimal rates of convergence.

as $\max\left(\frac{1}{\sqrt{d}}, \frac{s}{d}\right)$, which is faster than the optimal rate $\max\left(\frac{1}{\sqrt{d}}, \frac{s}{d} \log\left(\frac{ed}{s}\right)\right)$ for the class of sub-Gaussian noise. In other words, the phenomenon of improved rate is not due to the Gaussian character of the noise but rather to the fact that the noise distribution is known.

- Our findings show that there is a dramatic difference between the behavior of optimal estimators of $\boldsymbol{\theta}$ in the sparse sequence model and in the sparse linear regression model with "well spread" regressors. It is known from [11, 2] that in sparse linear regression with "well spread" regressors (that is, having positive variance), the rates of estimating $\boldsymbol{\theta}$ are the same for the noise with sub-Gaussian and polynomial tails. We show that the situation is quite different in the sparse sequence model, where the optimal rates are much slower and depend on the polynomial index of the noise.
- The rates shown in Table 1 for the classes $\mathcal{G}_{a,\tau}$ and $\mathcal{P}_{a,\tau}$ are achieved on estimators that are adaptive to the sparsity index s . Thus, knowing or not knowing s does not influence the optimal rates of estimation when the distribution of ξ and the noise level are unknown.

We conclude this section by a discussion of related work. Chen, Gao and Ren [7] explore the problem of robust estimation of variance and of covariance matrix under Huber's contamination model. As explained in Section 4 below, this problem has similarities with estimation of noise level in our setting. The main difference is that instead of fixing in advance the Gaussian nominal distribution of the contamination model we assume that it belongs to a class of distributions, such as (2) or (3). Therefore, the corresponding results in Section 4 can be viewed as results on robust estimation of scale where, in contrast to the classical setting, we are interested in adaptation to the unknown nominal law. Another aspect of robust estimation of scale is analyzed by Minsker and Wei [17] who consider classes of distributions similar to $\mathcal{P}_{a,\tau}$ rather than the contamination model. The main aim in [17] is to construct estimators having sub-Gaussian deviations under weak moment assumptions. Our setting is different in that we consider the sparsity class Θ_s of vectors $\boldsymbol{\theta}$ and the rates that we obtain depend on s . Estimation of variance in sparse linear model is discussed in [20] where some upper bounds for the rates are given. We also mention the recent paper [12] that deals with estimation of

variance in linear regression in a framework that does not involve sparsity, as well as the work on estimation of signal-to-noise ratio functionals in settings involving sparsity [23, 13] and not involving sparsity [16]. Papers [9, 6] discuss estimation of other functionals than the ℓ_2 -norm $\|\boldsymbol{\theta}\|_2$ in the sparse vector model when the noise is Gaussian with unknown variance.

Notation. For $x > 0$, let $\lfloor x \rfloor$ denote the maximal integer smaller than x . For a finite set A , we denote by $|A|$ its cardinality. Let $\inf_{\hat{T}}$ denote the infimum over all estimators. The notation C, C', c, c' will be used for positive constants that can depend only a and τ and can vary from line to line.

2. Estimation of sparse vector $\boldsymbol{\theta}$. In this section, we study the problem of estimating a sparse vector $\boldsymbol{\theta}$ in ℓ_2 -norm when the variance of noise σ and the distribution of ξ_i are both unknown. We only assume that the noise distribution belongs a given class, which can be either a class of distributions with polynomial tails $\mathcal{P}_{a,\tau}$, or a class $\mathcal{G}_{a,\tau}$ with exponential decay of the tails.

First, we introduce a preliminary estimator $\tilde{\sigma}^2$ of σ^2 that will be used to define an estimator of $\boldsymbol{\theta}$. Let $\gamma \in (0, 1/2]$ be a constant that will be chosen small enough and depending only on a and τ . Divide $\{1, \dots, d\}$ into $m = \lfloor \gamma d \rfloor$ disjoint subsets B_1, \dots, B_m , each of cardinality $|B_i| \geq k := \lfloor d/m \rfloor \geq 1/\gamma - 1$. Consider the median-of-means estimator

$$(4) \quad \tilde{\sigma}^2 = \text{med}(\bar{\sigma}_1^2, \dots, \bar{\sigma}_m^2), \text{ where } \bar{\sigma}_i^2 = \frac{1}{|B_i|} \sum_{j \in B_i} Y_j^2, \quad i = 1, \dots, m.$$

Here, $\text{med}(\bar{\sigma}_1^2, \dots, \bar{\sigma}_m^2)$ denotes the median of $\bar{\sigma}_1^2, \dots, \bar{\sigma}_m^2$. The next proposition shows that the estimator $\tilde{\sigma}^2$ recovers σ^2 to within a constant factor.

PROPOSITION 1. *Let $\tau > 0, a > 2$. There exist constants $\gamma \in (0, 1/2]$, $c > 0$ and $C > 0$ depending only on a and τ such that for any integers s and d satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$ we have*

$$\inf_{P_\xi \in \mathcal{P}_{a,\tau}} \inf_{\sigma > 0} \inf_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{P}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(1/2 \leq \frac{\tilde{\sigma}^2}{\sigma^2} \leq 3/2 \right) \geq 1 - \exp(-cd),$$

$$\sup_{P_\xi \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma} |\tilde{\sigma}^2 - \sigma^2| \leq C\sigma^2,$$

and for $a > 4$,

$$\sup_{P_\xi \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma} (\tilde{\sigma}^2 - \sigma^2)^2 \leq C\sigma^4.$$

Note that the result of Proposition 1 also holds for the class $\mathcal{G}_{a,\tau}$ for all $a > 0$ and $\tau > 0$. Indeed, $\mathcal{G}_{a,\tau} \subset \mathcal{P}_{a,\tau}$ for all $a > 2$ and $\tau > 0$, while for any $0 < a \leq 2$ and $\tau > 0$, there exist $a' > 4$ and $\tau' > 0$ such that $\mathcal{G}_{a,\tau} \subset \mathcal{P}_{a',\tau'}$.

We further note that assuming $s < cd$ for some $0 < c < 1$ is natural in the context of variance estimation since σ is not identifiable when $s = d$. In what follows, all upper bounds on the risks of estimators will be obtained under this assumption.

Consider now an estimator $\hat{\boldsymbol{\theta}}$ defined as follows:

$$(5) \quad \hat{\boldsymbol{\theta}} \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left(\sum_{i=1}^d (Y_i - \theta_i)^2 + \tilde{\sigma} \|\boldsymbol{\theta}\|_* \right).$$

Here, $\|\cdot\|_*$ denotes the sorted ℓ_1 -norm:

$$(6) \quad \|\boldsymbol{\theta}\|_* = \sum_{i=1}^d \lambda_i |\theta|_{(d-i+1)},$$

where $|\theta|_{(1)} \leq \dots \leq |\theta|_{(d)}$ are the order statistics of $|\theta_1|, \dots, |\theta_d|$, and $\lambda_1 \geq \dots \geq \lambda_p > 0$ are tuning parameters.

Set

$$(7) \quad \phi_{\text{exp}}^*(s, d) = \sqrt{s} \log^{1/a}(ed/s), \quad \phi_{\text{pol}}^*(s, d) = \sqrt{s}(d/s)^{1/a}.$$

The next theorem shows that $\hat{\boldsymbol{\theta}}$ estimates $\boldsymbol{\theta}$ with the rates $\phi_{\text{exp}}^*(s, d)$ and $\phi_{\text{pol}}^*(s, d)$ when the noise distribution belongs to the class $\mathcal{G}_{a,\tau}$ and class $\mathcal{P}_{a,\tau}$, respectively.

THEOREM 1. *Let s and d be integers satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$ where $\gamma \in (0, 1/2]$ is the tuning parameter in the definition of $\tilde{\sigma}^2$. Then for the estimator $\hat{\boldsymbol{\theta}}$ defined by (5) the following holds.*

1. *Let $\tau > 0$, $a > 0$. There exist constants $c, C > 0$ and $\gamma \in (0, 1/2]$ depending only on (a, τ) such that if $\lambda_j = c \log^{1/a}(ed/j)$, $j = 1, \dots, d$, we have*

$$\sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2 \right) \leq C \sigma^2 \left(\phi_{\text{exp}}^*(s, d) \right)^2.$$

2. *Let $\tau > 0$, $a > 2$. There exist constants $c, C > 0$ and $\gamma \in (0, 1/2]$ depending only on (a, τ) such that if $\lambda_j = c(d/j)^{1/a}$, $j = 1, \dots, d$, we have*

$$\sup_{P_\xi \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2 \right) \leq C \sigma^2 \left(\phi_{\text{pol}}^*(s, d) \right)^2.$$

Furthermore, it follows from the lower bound of Theorem 2 in Section 3 that the rates $\phi_{\text{exp}}^*(s, d)$ and $\phi_{\text{pol}}^*(s, d)$ cannot be improved in a minimax sense. Thus, the estimator $\hat{\boldsymbol{\theta}}$ defined in (5) achieves the optimal rates in a minimax sense.

From Theorem 1, we can conclude that the optimal rate ϕ_{pol}^* under polynomially decaying noise is very different from the optimal rate ϕ_{exp}^* under exponential tails, in particular, from the rate under the sub-Gaussian noise. At first sight, this phenomenon seems to contradict some results in the literature on sparse regression model. Indeed, Gautier and Tsybakov [11] consider sparse linear regression with unknown noise level σ and show that the Self-Tuned Dantzig estimator can achieve the same rate as in the case of Gaussian noise (up to a logarithmic factor) under the assumption that the noise is symmetric and has only a bounded moment of order $a > 2$. Belloni, Chernozhukov and Wang [2] show for the same model that a square-root Lasso estimator achieves analogous behavior under the assumption that the noise has a bounded moment of order $a > 2$. However, a crucial condition in [2] is that the design is "well spread", that is all components of the design vectors are random with positive variance. The same type of condition is needed in [11] to obtain a sub-Gaussian rate. This condition of "well spreadness" is not satisfied in the sparse sequence model that we are considering here. In this model viewed as a special case of linear regression, the design is deterministic, with only one non-zero component. We see that such a degenerate design turns out to be the least favorable from the point of view of the convergence rate, while the "well spread" design

is the best one. An interesting general conclusion of comparing our findings to [11] and [2] is that the optimal rate of convergence of estimators under sparsity when the noise level is unknown depends dramatically on the properties of the design. There is a whole spectrum of possibilities between the degenerate and "well spread" designs where a variety of new rates can arise depending on the properties of the design. Studying them remains an open problem.

3. Estimation of the ℓ_2 -norm. In this section, we consider the problem of estimation of the ℓ_2 -norm of a sparse vector when the variance of the noise and the form of its distribution are both unknown. We show that the rates $\phi_{\text{exp}}^*(s, d)$ and $\phi_{\text{pol}}^*(s, d)$ are optimal in a minimax sense on the classes $\mathcal{G}_{a,\tau}$ and $\mathcal{P}_{a,\tau}$, respectively. We first provide a lower bound on the risks of any estimators of the ℓ_2 -norm when the noise level σ is unknown and the unknown noise distribution P_ξ belongs either to $\mathcal{G}_{a,\tau}$ or $\mathcal{P}_{a,\tau}$. We denote by \mathcal{L} the set of all monotone non-decreasing functions $\ell : [0, \infty) \rightarrow [0, \infty)$ such that $\ell(0) = 0$ and $\ell \not\equiv 0$.

THEOREM 2. *Let s, d be integers satisfying $1 \leq s \leq d$. Let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then, for any $a > 0, \tau > 0$,*

$$(8) \quad \inf_{\hat{T}} \sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{E}_{\theta, P_\xi, \sigma} \ell \left(c(\phi_{\text{exp}}^*(s, d))^{-1} \left| \frac{\hat{T} - \|\theta\|_2}{\sigma} \right| \right) \geq c',$$

and, for any $a \geq 2, \tau > 0$,

$$(9) \quad \inf_{\hat{T}} \sup_{P_\xi \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{E}_{\theta, P_\xi, \sigma} \ell \left(\bar{c}(\phi_{\text{pol}}^*(s, d))^{-1} \left| \frac{\hat{T} - \|\theta\|_2}{\sigma} \right| \right) \geq \bar{c}'.$$

Here, $\inf_{\hat{T}}$ denotes the infimum over all estimators, and $c, \bar{c} > 0, c', \bar{c}' > 0$ are constants that can depend only on $\ell(\cdot)$, τ and a .

The lower bound (9) implies that the rate of estimation of the ℓ_2 -norm of a sparse vector deteriorates dramatically if the bounded moment assumption is imposed on the noise instead, for example, of the sub-Gaussian assumption.

Note also that (8) and (9) immediately imply lower bounds with the same rates ϕ_{exp}^* and ϕ_{pol}^* for the estimation of the s -sparse vector θ under the ℓ_2 -norm.

Given the upper bounds of Theorem 1, the lower bounds (8) and (9) are tight for the quadratic loss, and are achieved by the following plug-in estimator independent of s or σ :

$$(10) \quad \hat{N} = \|\hat{\theta}\|_2$$

where $\hat{\theta}$ is defined in (5).

In conclusion, when both P_ξ and σ are unknown the rates ϕ_{exp}^* and ϕ_{pol}^* defined in (7) are minimax optimal both for estimation of θ and of the the norm $\|\theta\|_2$.

We now compare these results with the findings in [8] regarding the (nonadaptive) estimation of $\|\theta\|_2$ when ξ_i have the standard Gaussian distribution ($P_\xi = \mathcal{N}(0, 1)$) and σ is known. It is shown in [8] that in this case the optimal rate of estimation of $\|\theta\|_2$ has the form

$$\phi_{\mathcal{N}(0,1)}(s, d) = \min \left\{ \sqrt{s \log(1 + \sqrt{d}/s)}, d^{1/4} \right\}.$$

Namely, the following proposition holds.

PROPOSITION 2 (Gaussian noise, known σ [8]). *For any $\sigma > 0$ and any integers s, d satisfying $1 \leq s \leq d$, we have*

$$c\sigma^2\phi_{\mathcal{N}(0,1)}^2(s, d) \leq \inf_{\hat{T}} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} (\hat{T} - \|\boldsymbol{\theta}\|_2)^2 \leq C\sigma^2\phi_{\mathcal{N}(0,1)}^2(s, d),$$

where $c > 0$ and $C > 0$ are absolute constants and $\inf_{\hat{T}}$ denotes the infimum over all estimators.

We have seen that, in contrast to this result, in the case of unknown P_ξ and σ the optimal rates (7) do not exhibit an elbow at $s = \sqrt{d}$ between the "sparse" and "dense" regimes. Another conclusion is that, in the "dense" zone $s > \sqrt{d}$, adaptation to P_ξ and σ is only possible with a significant deterioration of the rate. On the other hand, for the sub-Gaussian class $\mathcal{G}_{2,\tau}$, in the "sparse" zone $s \leq \sqrt{d}$ the non-adaptive rate $\sqrt{s \log(1 + \sqrt{d}/s)}$ differs only slightly from the adaptive sub-Gaussian rate $\sqrt{s \log(ed/s)}$; in fact, this difference in the rate appears only in a vicinity of $s = \sqrt{d}$.

A natural question is whether such a deterioration of the rate is caused by the ignorance of σ or by the ignorance of the distribution of ξ_i within the sub-Gaussian class $\mathcal{G}_{2,\tau}$. The answer is that both are responsible. It turns out that if only one of the two ingredients (σ or the noise distribution) is unknown, then a rate faster than the adaptive sub-Gaussian rate $\phi_{\text{exp}}^*(s, d) = \sqrt{s \log(ed/s)}$ can be achieved. This is detailed in the next two propositions.

Consider first the case of Gaussian noise and unknown σ . Set

$$\phi_{\mathcal{N}(0,1)}^*(s, d) = \max \left\{ \sqrt{s \log(1 + \sqrt{d}/s)}, \sqrt{\frac{s}{1 + \log_+(s^2/d)}} \right\},$$

where $\log_+(x) = \max(0, \log(x))$ for any $x > 0$. We divide the set $\{1, \dots, d\}$ into two disjoint subsets I_1 and I_2 with $\min(|I_1|, |I_2|) \geq \lfloor d/2 \rfloor$. Let $\hat{\sigma}^2$ be the variance estimator defined by (15), cf. Section 4.1 below, and let $\hat{\sigma}_{\text{med},1}^2, \hat{\sigma}_{\text{med},2}^2$ be the median estimators (12) corresponding to the samples $(Y_i)_{i \in I_1}$ and $(Y_i)_{i \in I_2}$, respectively. Consider the estimator

$$(11) \quad \hat{N}^* = \begin{cases} \sqrt{\left| \sum_{j=1}^d (Y_j^2 \mathbf{1}_{\{|Y_j| > \rho_j\}}) - d\alpha\hat{\sigma}^2 \right|} & \text{if } s \leq \sqrt{d}, \\ \sqrt{\left| \sum_{j=1}^d Y_j^2 - d\hat{\sigma}^2 \right|} & \text{if } s > \sqrt{d}, \end{cases}$$

where $\rho_j = 2\hat{\sigma}_{\text{med},1}\sqrt{2\log(1 + d/s^2)}$ if $j \in I_2$, $\rho_j = 2\hat{\sigma}_{\text{med},2}\sqrt{2\log(1 + d/s^2)}$ if $j \in I_1$ and $\alpha = \mathbf{E} \left(\xi_1^2 \mathbf{1}_{\{|\xi_1| > 2\sqrt{2\log(1 + d/s^2)}\}} \right)$. Note that Y_j is independent of ρ_j for every j . Note also that the estimator \hat{N}^* depends on the preliminary estimator $\tilde{\sigma}^2$ since $\hat{\sigma} > 0$ defined in (15) depends on it.

PROPOSITION 3 (Gaussian noise, unknown σ). *The following two properties hold.*

- (i) *Let s and d be integers satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$, where $\gamma \in (0, 1/2]$ is the tuning parameter in the definition of $\tilde{\sigma}^2$. There exist absolute constants $C > 0$ and $\gamma \in (0, 1/2]$ such that*

$$\sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} \left(\hat{N}^* - \|\boldsymbol{\theta}\|_2 \right)^2 \leq C\sigma^2 \left(\phi_{\mathcal{N}(0,1)}^*(s, d) \right)^2.$$

(ii) Let s and d be integers satisfying $1 \leq s \leq d$ and let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then,

$$\inf_{\hat{T}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} \ell \left(c(\phi_{\mathcal{N}(0,1)}^*(s, d))^{-1} \left| \frac{\hat{T} - \|\boldsymbol{\theta}\|_2}{\sigma} \right| \right) \geq c',$$

where $\inf_{\hat{T}}$ denotes the infimum over all estimators, and $c > 0$, $c' > 0$ are constants that can depend only on $\ell(\cdot)$.

The proof of item (ii) of Proposition 3 (the lower bound) is given in the Supplementary material.

Proposition 3 establishes the minimax optimality of the rate $\phi_{\mathcal{N}(0,1)}^*(s, d)$. It also shows that if σ is unknown, the knowledge of the Gaussian character of the noise leads to an improvement of the rate compared to the adaptive sub-Gaussian rate $\sqrt{s \log(ed/s)}$. However, the improvement is only in a logarithmic factor.

Consider now the case of unknown noise distribution in $\mathcal{G}_{a,\tau}$ and known σ . We show in the next proposition that in this case the minimax rate is of the form

$$\phi_{\text{exp}}^\circ(s, d) = \min\{\sqrt{s} \log^{\frac{1}{a}}(ed/s), d^{1/4}\}$$

and it is achieved by the estimator

$$\hat{N}_{\text{exp}}^\circ = \begin{cases} \|\hat{\boldsymbol{\theta}}\|_2 & \text{if } s \leq \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}, \\ \left| \sum_{j=1}^d Y_j^2 - d\sigma^2 \right|^{1/2} & \text{if } s > \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}, \end{cases}$$

where $\hat{\boldsymbol{\theta}}$ is defined in (5). Note $\phi_{\text{exp}}^\circ(s, d)$ can be written equivalently (up to absolute constants) as $\min\{\sqrt{s} \log^{\frac{1}{a}}(ed), d^{1/4}\}$.

PROPOSITION 4 (Unknown noise in $\mathcal{G}_{a,\tau}$, known σ). *Let $a, \tau > 0$. The following two properties hold.*

(i) Let s and d be integers satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$, where $\gamma \in (0, 1/2]$ is the tuning parameter in the definition of $\tilde{\sigma}^2$. There exist constants $c, C > 0$, and $\gamma \in (0, 1/2]$ depending only on (a, τ) such that if $\hat{\boldsymbol{\theta}}$ is the estimator defined in (5) with $\lambda_j = c \log^{\frac{1}{a}}(ed/j)$, $j = 1, \dots, d$, then

$$\sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(\hat{N}_{\text{exp}}^\circ - \|\boldsymbol{\theta}\|_2 \right)^2 \leq C\sigma^2 \left(\phi_{\text{exp}}^\circ(s, d) \right)^2.$$

(ii) Let s and d be integers satisfying $1 \leq s \leq d$ and let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then, there exist constants $c > 0$, $c' > 0$ depending only on $\ell(\cdot)$, a and τ such that

$$\inf_{\hat{T}} \sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_\xi, \sigma} \ell \left(c(\phi_{\text{exp}}^\circ(s, d))^{-1} \left| \frac{\hat{T} - \|\boldsymbol{\theta}\|_2}{\sigma} \right| \right) \geq c',$$

where $\inf_{\hat{T}}$ denotes the infimum over all estimators.

Proposition 4 establishes the minimax optimality of the rate $\phi_{\text{exp}}^{\circ}(s, d)$. It also shows that if the noise distribution is unknown and belongs to $\mathcal{G}_{a, \tau}$, the knowledge of σ leads to an improvement of the rate compared to the case when σ is unknown. In contrast to the case of Proposition 3 (Gaussian noise), the improvement here is substantial; it results not only in a logarithmic but in a polynomial factor in the dense zone $s > \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}$.

We end this section by considering the case of unknown polynomial noise and known σ . The next proposition shows that in this case the minimax rate, for a given $a > 4$, is of the form

$$\phi_{\text{pol}}^{\circ}(s, d) = \min\{\sqrt{s}(d/s)^{\frac{1}{a}}, d^{1/4}\}$$

and it is achieved by the estimator

$$\hat{N}_{\text{pol}}^{\circ} = \begin{cases} \|\hat{\boldsymbol{\theta}}\|_2 & \text{if } s \leq d^{\frac{1}{2} - \frac{1}{a-2}}, \\ \left| \sum_{j=1}^d Y_j^2 - d\sigma^2 \right|^{1/2} & \text{if } s > d^{\frac{1}{2} - \frac{1}{a-2}}, \end{cases}$$

where $\hat{\boldsymbol{\theta}}$ is defined in (5).

PROPOSITION 5 (Unknown noise in $\mathcal{P}_{a, \tau}$, known σ). *Let $\tau > 0, a > 4$. The following two properties hold.*

- (i) *Let s and d be integers satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$, where $\gamma \in (0, 1/2]$ is the tuning parameter in the definition of $\tilde{\sigma}^2$. There exist constants $c, C > 0$, and $\gamma \in (0, 1/2]$ depending only on (a, τ) such that if $\hat{\boldsymbol{\theta}}$ is the estimator defined in (5) with $\lambda_j = c(d/j)^{\frac{1}{a}}$, $j = 1, \dots, d$, then*

$$\sup_{P_{\xi} \in \mathcal{P}_{a, \tau}} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\hat{N}_{\text{pol}}^{\circ} - \|\boldsymbol{\theta}\|_2 \right)^2 \leq C\sigma^2 \left(\phi_{\text{pol}}^{\circ}(s, d) \right)^2.$$

- (ii) *Let s and d be integers satisfying $1 \leq s \leq d$ and let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then, there exist constants $c > 0$, $c' > 0$ depending only on $\ell(\cdot)$, a and τ such that*

$$\inf_{\hat{T}} \sup_{P_{\xi} \in \mathcal{P}_{a, \tau}} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \ell \left(c(\phi_{\text{pol}}^{\circ}(s, d))^{-1} \left| \frac{\hat{T} - \|\boldsymbol{\theta}\|_2}{\sigma} \right| \right) \geq c',$$

where $\inf_{\hat{T}}$ denotes the infimum over all estimators.

Note that here, similarly to Proposition 4, the improvement over the case of unknown σ is in a polynomial factor in the dense zone $s > d^{\frac{1}{2} - \frac{1}{a-2}}$.

4. Estimating the variance of the noise.

4.1. *Estimating σ^2 when the distribution P_{ξ} is known.* In the sparse setting when $\|\boldsymbol{\theta}\|_0$ is small, estimation of the noise level can be viewed as a problem of robust estimation of scale. Indeed, our aim is to recover the second moment of $\sigma\xi_1$ but the sample second moment cannot be used as an estimator because of the presence of a small number of outliers $\theta_i \neq 0$. Thus, the models in robustness and sparsity problems are quite similar but the questions of interest are different. When robust estimation of σ^2 is considered, the object of interest is the pure

noise component of the sparsity model while the non-zero components θ_i that are of major interest in the sparsity model play a role of nuisance.

In the context of robustness, it is known that the estimator based on sample median can be successfully applied. Recall that, when $\boldsymbol{\theta} = 0$, the median M -estimator of scale (cf. [14]) is defined as

$$(12) \quad \hat{\sigma}_{\text{med}}^2 = \frac{\hat{M}}{\beta}$$

where \hat{M} is the sample median of (Y_1^2, \dots, Y_d^2) , that is

$$\hat{M} \in \arg \min_{x \geq 0} |F_d(x) - 1/2|,$$

and β is the median of the distribution of ξ_1^2 . Here, F_d denotes the empirical c.d.f. of (Y_1^2, \dots, Y_d^2) . When F denotes the c.d.f. of ξ_1^2 , it is easy to see that

$$(13) \quad \beta = F^{-1}(1/2).$$

The following proposition specifies the rate of convergence of the estimator $\hat{\sigma}_{\text{med}}^2$.

PROPOSITION 6. *Let ξ_1^2 have a c.d.f. F with positive density, and let β be given by (13). There exist constants $\gamma \in (0, 1/8)$, $c > 0$, $c_* > 0$ and $C > 0$ depending only on F such that for any integers s and d satisfying $1 \leq s < \gamma d$ and any $t > 0$ we have*

$$\sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{P}_{\boldsymbol{\theta}, F, \sigma} \left(\left| \frac{\hat{\sigma}_{\text{med}}^2}{\sigma^2} - 1 \right| \geq c_* \left(\sqrt{\frac{t}{d}} + \frac{s}{d} \right) \right) \leq 2(e^{-t} + e^{-cd}),$$

and if $\mathbf{E}|\xi_1^2|^{2+\epsilon} < \infty$ for some $\epsilon > 0$. Then,

$$\sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \frac{\mathbf{E}_{\boldsymbol{\theta}, F, \sigma} |\hat{\sigma}_{\text{med}}^2 - \sigma^2|}{\sigma^2} \leq C \max \left(\frac{1}{\sqrt{d}}, \frac{s}{d} \right).$$

The main message of Proposition 6 is that the rate of convergence of $\hat{\sigma}_{\text{med}}^2$ in probability and in expectation is as fast as

$$(14) \quad \max \left(\frac{1}{\sqrt{d}}, \frac{s}{d} \right)$$

and it does not depend on F when F varies in a large class. The role of Proposition 6 is to contrast the subsequent results of this section dealing with unknown distribution of noise and providing slower rates. It emphasizes the fact that the knowledge of the noise distribution is crucial as it leads to an improvement of the rate of estimating the variance.

However, the rate (14) achieved by the median estimator is not necessarily optimal. As shown in the next proposition, in the case of Gaussian noise the optimal rate is even better:

$$\phi_{\mathcal{N}(0,1)}(s, d) = \max \left\{ \frac{1}{\sqrt{d}}, \frac{s}{d(1 + \log_+(s^2/d))} \right\}.$$

This rate is attained by an estimator that we are going to define now. We use the observation that, in the Gaussian case, the modulus of the empirical characteristic function $\varphi_d(t) = \frac{1}{d} \sum_{i=1}^d e^{itY_j}$ is to within a constant factor of the Gaussian characteristic function $\exp(-\frac{t^2\sigma^2}{2})$ for any t . This suggests the estimator

$$\tilde{v}^2 = -\frac{2\log(|\varphi_d(\hat{t}_1)|)}{\hat{t}_1^2},$$

with a suitable choice of $t = \hat{t}_1$ that we further set as follows:

$$\hat{t}_1 = \frac{1}{\tilde{\sigma}} \sqrt{\log(4(es/\sqrt{d} + 1))},$$

where $\tilde{\sigma}$ is the preliminary estimator (4) with some tuning parameter $\gamma \in (0, 1/2]$. The final variance estimator is defined as a truncated version of \tilde{v}^2 :

$$(15) \quad \hat{\sigma}^2 = \begin{cases} \tilde{v}^2 & \text{if } |\varphi_d(\hat{t}_1)| > (es/\sqrt{d} + 1)^{-1}/4, \\ \tilde{\sigma}^2 & \text{otherwise.} \end{cases}$$

PROPOSITION 7 (Gaussian noise). *The following two properties hold.*

- (i) *Let s and d be integers satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$, where $\gamma \in (0, 1/2]$ is the tuning parameter in the definition of $\tilde{\sigma}^2$. There exist absolute constants $C > 0$ and $\gamma \in (0, 1/2]$ such that the estimator $\hat{\sigma}^2$ defined in (15) satisfies*

$$\sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \frac{\mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} |\hat{\sigma}^2 - \sigma^2|}{\sigma^2} \leq C \phi_{\mathcal{N}(0,1)}(s, d).$$

- (ii) *Let s and d be integers satisfying $1 \leq s \leq d$ and let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then,*

$$\inf_{\hat{T}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} \ell \left(c(\phi_{\mathcal{N}(0,1)}(s, d))^{-1} \left| \frac{\hat{T}}{\sigma^2} - 1 \right| \right) \geq c',$$

where $\inf_{\hat{T}}$ denotes the infimum over all estimators, and $c > 0$, $c' > 0$ are constants that can depend only on $\ell(\cdot)$.

Estimators of variance or covariance matrix based on the empirical characteristic function have been studied in several papers [4, 5, 3, 6]. The setting in [4, 5, 3] is different from the ours as those papers deal with the model where the non-zero components of $\boldsymbol{\theta}$ are random with a smooth distribution density. The estimators in [4, 5] are also quite different. On the other hand, [3, 6] consider estimators close to \tilde{v}^2 . In particular, [6] uses a similar pilot estimator for testing in the sparse vector model where it is assumed that $\sigma \in [\sigma_-, \sigma_+]$, $0 < \sigma_- < \sigma_+ < \infty$, and the estimator depends on σ_+ . Although [6] does not provide explicitly stated result about the rate of this estimator, the proofs in [6] come close to it and we believe that it satisfies an upper bound as in item (i) of Proposition 7 with $\sup_{\sigma > 0}$ replaced by $\sup_{\sigma \in [\sigma_-, \sigma_+]}$.

4.2. *Distribution-free variance estimators.* The main drawback of the estimator $\hat{\sigma}_{\text{med}}^2$ is the dependence on the parameter β . It reflects the fact that the estimator is tailored for a given and known distribution of noise F . Furthermore, as shown below, the rate (14) cannot

be achieved if it is only known that F belongs to one of the classes of distributions that we consider in this paper.

Instead of using one particular quantile, like the median in Section 4.1, one can estimate σ^2 by an integral over all quantiles, which allows one to avoid considering distribution-dependent quantities like (13).

Indeed, with the notation $q_\alpha = G^{-1}(1 - \alpha)$ where G is the c.d.f. of $(\sigma\xi_1)^2$ and $0 < \alpha < 1$, the variance of the noise can be expressed as

$$\sigma^2 = \mathbf{E}(\sigma\xi_1)^2 = \int_0^1 q_\alpha d\alpha.$$

Discarding the higher order quantiles that are dubious in the presence of outliers and replacing q_α by the empirical quantile \hat{q}_α of level α we obtain the following estimator

$$(16) \quad \hat{\sigma}^2 = \int_0^{1-s/d} \hat{q}_\alpha d\alpha = \frac{1}{d} \sum_{k=1}^{d-s} Y_{(k)}^2,$$

where $Y_{(1)}^2 \leq \dots \leq Y_{(d)}^2$ are the ordered values of the squared observations Y_1^2, \dots, Y_d^2 . Note that $\hat{\sigma}^2$ is an L -estimator, cf. [14]. Also, up to a constant factor, $\hat{\sigma}^2$ coincides with the statistic used in Collier, Comminges and Tsybakov [8].

The following theorem provides an upper bound on the risk of the estimator $\hat{\sigma}^2$ under the assumption that the noise belongs to the class $\mathcal{G}_{a,\tau}$. Set

$$\phi_{\text{exp}}(s, d) = \max \left(\frac{1}{\sqrt{d}}, \frac{s}{d} \log^{2/a} \left(\frac{ed}{s} \right) \right).$$

THEOREM 3. *Let $\tau > 0$, $a > 0$, and let s, d be integers satisfying $1 \leq s < d/2$. Then, the estimator $\hat{\sigma}^2$ defined in (16) satisfies*

$$(17) \quad \sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \frac{\mathbf{E}_{\theta, P_\xi, \sigma} (\hat{\sigma}^2 - \sigma^2)^2}{\sigma^4} \leq C \phi_{\text{exp}}^2(s, d)$$

where $C > 0$ is a constant depending only on a and τ .

The next theorem establishes the performance of variance estimation in the case of distributions with polynomially decaying tails. Set

$$\phi_{\text{pol}}(s, d) = \max \left(\frac{1}{\sqrt{d}}, \left(\frac{s}{d} \right)^{1-\frac{2}{a}} \right).$$

THEOREM 4. *Let $\tau > 0$, $a > 4$, and let s, d be integers satisfying $1 \leq s < d/2$. Then, the estimator $\hat{\sigma}^2$ defined in (16) satisfies*

$$(18) \quad \sup_{P_\xi \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \frac{\mathbf{E}_{\theta, P_\xi, \sigma} (\hat{\sigma}^2 - \sigma^2)^2}{\sigma^4} \leq C \phi_{\text{pol}}^2(s, d),$$

where $C > 0$ is a constant depending only on a and τ .

We assume here that the noise distribution has a moment of order greater than 4, which is close to the minimum requirement since we deal with the expected squared error of a quadratic function of the observations.

We now state the lower bounds matching the results of Theorems 3 and 4.

THEOREM 5. *Let $\tau > 0$, $a > 0$, and let s, d be integers satisfying $1 \leq s \leq d$. Let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then,*

$$(19) \quad \inf_{\hat{T}} \sup_{P_{\xi} \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{E}_{\theta, P_{\xi}, \sigma} \ell \left(c(\phi_{\text{exp}}(s, d))^{-1} \left| \frac{\hat{T}}{\sigma^2} - 1 \right| \right) \geq c',$$

where $\inf_{\hat{T}}$ denotes the infimum over all estimators and $c > 0$, $c' > 0$ are constants depending only on $\ell(\cdot)$, a and τ .

Theorems 3 and 5 imply that the estimator $\hat{\sigma}^2$ is rate optimal in a minimax sense when the noise belongs to $\mathcal{G}_{a,\tau}$, in particular when it is sub-Gaussian. Interestingly, an extra logarithmic factor appears in the optimal rate when passing from the pure Gaussian distribution of ξ_i 's (cf. Proposition 7) to the class of all sub-Gaussian distributions. This factor can be seen as a price to pay for the lack of information regarding the exact form of the distribution. Also note that this logarithmic factor vanishes as $a \rightarrow \infty$.

Under polynomial tail assumption on the noise, we have the following minimax lower bound.

THEOREM 6. *Let $\tau > 0$, $a \geq 2$, and let s, d be integers satisfying $1 \leq s \leq d$. Let $\ell(\cdot)$ be any loss function in the class \mathcal{L} . Then,*

$$(20) \quad \inf_{\hat{T}} \sup_{P_{\xi} \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{E}_{\theta, P_{\xi}, \sigma} \ell \left(c(\phi_{\text{pol}}(s, d))^{-1} \left| \frac{\hat{T}}{\sigma^2} - 1 \right| \right) \geq c'$$

where $\inf_{\hat{T}}$ denotes the infimum over all estimators and $c > 0$, $c' > 0$ are constants depending only on $\ell(\cdot)$, a and τ .

This theorem shows that the rate $\phi_{\text{pol}}(s, d)$ obtained in Theorem 4 cannot be improved in a minimax sense.

A drawback of the estimator defined in (16) is in the lack of adaptivity to the sparsity parameter s . At first sight, it may seem that the estimator

$$(21) \quad \hat{\sigma}_*^2 = \frac{2}{d} \sum_{1 \leq k \leq d/2} Y_{(k)}^2$$

could be taken as its adaptive version. However, $\hat{\sigma}_*^2$ is not a good estimator of σ^2 as can be seen from the following proposition.

PROPOSITION 8. *Define $\hat{\sigma}_*^2$ as in (21). Let $\tau > 0$, $a \geq 2$, and let s, d be integers satisfying $1 \leq s \leq d$, and $d = 4k$ for an integer k . Then,*

$$\sup_{P_{\xi} \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \frac{\mathbf{E}_{\theta, P_{\xi}, \sigma} (\hat{\sigma}_*^2 - \sigma^2)^2}{\sigma^4} \geq \frac{1}{64}.$$

On the other hand, it turns out that a simple plug-in estimator

$$(22) \quad \hat{\sigma}^2 = \frac{1}{d} \|\mathbf{Y} - \hat{\boldsymbol{\theta}}\|_2^2$$

with $\hat{\boldsymbol{\theta}}$ chosen as in Section 2 achieves rate optimality adaptively to the noise distribution and to the sparsity parameter s . This is detailed in the next theorem.

THEOREM 7. *Let s and d be integers satisfying $1 \leq s < \lfloor \gamma d \rfloor / 4$, where $\gamma \in (0, 1/2]$ is the tuning parameter in the definition of $\hat{\sigma}^2$. Let $\hat{\sigma}^2$ be the estimator defined by (22) where $\hat{\boldsymbol{\theta}}$ is defined in (5). Then the following properties hold.*

1. *Let $\tau > 0, a > 0$. There exist constants $c, C > 0$ and $\gamma \in (0, 1/2]$ depending only on (a, τ) such that if $\lambda_j = c \log^{1/a}(ed/j), j = 1, \dots, d$, we have*

$$\sup_{P_{\xi} \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} |\hat{\sigma}^2 - \sigma^2| \leq C \sigma^2 \phi_{\text{exp}}(s, d).$$

2. *Let $\tau > 0, a > 4$. There exist constants $c, C > 0$ and $\gamma \in (0, 1/2]$ depending only on (a, τ) such that if $\lambda_j = c(d/j)^{1/a}, j = 1, \dots, d$, we have*

$$\sup_{P_{\xi} \in \mathcal{P}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} |\hat{\sigma}^2 - \sigma^2| \leq C \sigma^2 \phi_{\text{pol}}(s, d).$$

5. Proofs of the upper bounds.

5.1. *Proof of Proposition 1.* Fix $\boldsymbol{\theta} \in \Theta_s$ and let S be the support of $\boldsymbol{\theta}$. We will call outliers the observations Y_i with $i \in S$. There are at least $m - s$ blocks B_i that do not contain outliers. Denote by J a set of $m - s$ indices i , for which B_i contains no outliers.

As $a > 2$, there exist constants $L = L(a, \tau)$ and $r = r(a, \tau) \in (1, 2]$ such that $\mathbf{E}|\xi_1^2 - 1|^r \leq L$. Using von Bahr-Esseen inequality (cf. [18]) and the fact that $|B_i| \geq k$ we get

$$\mathbf{P}\left(\left|\frac{1}{|B_i|} \sum_{j \in B_i} \xi_j^2 - 1\right| > 1/2\right) \leq \frac{2^{r+1}L}{k^{r-1}}, \quad i = 1, \dots, m.$$

Hence, there exists a constant $C_1 = C_1(a, \tau)$ such that if $k \geq C_1$ (i.e., if γ is small enough depending on a and τ), then

$$(23) \quad \mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\tilde{\sigma}_i^2 \notin I) \leq \frac{1}{4}, \quad i = 1, \dots, m,$$

where $I = [\frac{\sigma^2}{2}, \frac{3\sigma^2}{2}]$. Next, by the definition of the median, for any interval $I \subseteq \mathbb{R}$ we have

$$(24) \quad \mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\tilde{\sigma}^2 \notin I) \leq \mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}\left(\sum_{i=1}^m \mathbb{1}_{\tilde{\sigma}_i^2 \notin I} \geq \frac{m}{2}\right) \leq \mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}\left(\sum_{i \in J} \mathbb{1}_{\tilde{\sigma}_i^2 \notin I} \geq \frac{m}{2} - s\right).$$

Now, $s \leq \frac{\lfloor \gamma d \rfloor}{4} = \frac{m}{4}$, so that $\frac{m}{2} - s \geq \frac{m-s}{3}$. Set $\eta_i = \mathbb{1}_{\tilde{\sigma}_i^2 \notin I}$, $i \in J$. Due to (23) we have $\mathbf{E}(\eta_i) \leq 1/4$, and $(\eta_i, i \in J)$ are independent. Using these remarks and Hoeffding's inequality we find

$$\mathbf{P}\left(\sum_{i \in J} \eta_i \geq \frac{m}{2} - s\right) \leq \mathbf{P}\left(\sum_{i \in J} (\eta_i - \mathbf{E}(\eta_i)) \geq \frac{m-s}{12}\right) \leq \exp(-C(m-s)).$$

Note that $|J| = m - s \geq 3m/4 = 3\lfloor \gamma d \rfloor / 4$. Thus, if γ is chosen small enough depending only on a and τ then

$$\mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\tilde{\sigma}^2 \notin I) \leq \exp(-Cd).$$

This proves the desired bound in probability. To obtain the bounds in expectation, set $Z = |\tilde{\sigma}^2 - \sigma^2|$. Let first $a > 4$ and take some $r \in (1, a/4)$. Then

$$\begin{aligned} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(Z^2) &\leq \frac{\sigma^4}{4} + \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}\left(Z^2 \mathbf{1}_{Z \geq \frac{\sigma^2}{2}}\right) \\ &\leq \frac{9\sigma^4}{4} + 2(\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\tilde{\sigma}^{4r}))^{1/r} (\mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(Z \geq \sigma^2/2))^{1-1/r} \\ &\leq \frac{9\sigma^4}{4} + 2(\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\tilde{\sigma}^{4r}))^{1/r} \exp(-Cd). \end{aligned}$$

Since $m \geq 4s$, we can easily argue that $\tilde{\sigma}^{4r} \leq \sum_{i \in J} \tilde{\sigma}_i^{4r}$. It follows that

$$\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\tilde{\sigma}^{4r}) \leq C\sigma^{4r}d^2.$$

Hence $\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(Z^2) \leq C\sigma^4$. Similarly, if $a > 2$, then $\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(Z) \leq C\sigma^2$.

5.2. *Proof of Theorem 1.* Set $\mathbf{u} = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}$. It follows from Lemma A.2 in [1] that

$$2\|\mathbf{u}\|_2^2 \leq 2\sigma \sum_{i=1}^d \xi_i u_i + \tilde{\sigma} \|\boldsymbol{\theta}\|_* - \tilde{\sigma} \|\hat{\boldsymbol{\theta}}\|_*,$$

where u_i are the components of \mathbf{u} . Next, Lemma A.1 in [1] yields

$$\|\boldsymbol{\theta}\|_* - \|\hat{\boldsymbol{\theta}}\|_* \leq \left(\sum_{j=1}^s \lambda_j^2\right)^{1/2} \|\mathbf{u}\|_2 - \sum_{j=s+1}^d \lambda_j |u|_{(d-j+1)}$$

where $|u|_{(k)}$ is the k th order statistic of $|u_1|, \dots, |u_d|$. Combining these two inequalities we get

$$(25) \quad 2\|\mathbf{u}\|_2^2 \leq 2\sigma \sum_{j=1}^d \xi_j u_j + \tilde{\sigma} \left\{ \left(\sum_{j=1}^s \lambda_j^2\right)^{1/2} \|\mathbf{u}\|_2 - \sum_{j=s+1}^d \lambda_j |u|_{(d-j+1)} \right\}.$$

For some permutation $(\varphi(1), \dots, \varphi(d))$ of $(1, \dots, d)$, we have

$$(26) \quad \left| \sum_{i=1}^d \xi_i u_i \right| \leq \sum_{j=1}^d |\xi|_{(d-j+1)} |u_{\varphi(j)}| \leq \sum_{j=1}^d |\xi|_{(d-j+1)} |u|_{(d-j+1)},$$

where the last inequality is due to the fact that the sequence $|\xi|_{(d-j+1)}$ is non-increasing. Hence

$$\begin{aligned} 2\|\mathbf{u}\|_2^2 &\leq 2\sigma \sum_{j=1}^s |\xi|_{(d-j+1)} |u|_{(d-j+1)} + \tilde{\sigma} \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} \|\mathbf{u}\|_2 + \sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j) |u|_{(d-j+1)} \\ &\leq \left\{ 2\sigma \left(\sum_{j=1}^s |\xi|_{(d-j+1)}^2 \right)^{1/2} + \tilde{\sigma} \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} + \left(\sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j)_+^2 \right)^{1/2} \right\} \|\mathbf{u}\|_2. \end{aligned}$$

This implies

$$\|\mathbf{u}\|_2^2 \leq C \left\{ \sigma^2 \sum_{j=1}^s |\xi|_{(d-j+1)}^2 + \tilde{\sigma}^2 \sum_{j=1}^s \lambda_j^2 + \sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j)_+^2 \right\}.$$

From Lemmas 1 and 2 we have $\mathbf{E}(|\xi|_{(d-j+1)}^2) \leq C\lambda_j^2$. Using this and Proposition 1 we obtain

$$(27) \quad \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} (\|\mathbf{u}\|_2^2) \leq C \left(\sigma^2 \sum_{j=1}^s \lambda_j^2 + \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j)_+^2 \right) \right).$$

Define the events $\mathcal{A}_j = \{|\xi|_{(d-j+1)} \leq \lambda_j/4\} \cap \{1/2 \leq \tilde{\sigma}^2/\sigma^2 \leq 3/2\}$ for $j = s+1, \dots, d$. Then

$$\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j)_+^2 \right) \leq 4\sigma^2 \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\sum_{j=s+1}^d |\xi|_{(d-j+1)}^2 \mathbf{1}_{\mathcal{A}_j^c} \right).$$

Fixing some $1 < r < a/2$ we get

$$\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j)_+^2 \right) \leq 4\sigma^2 \sum_{j=s+1}^d \mathbf{E} \left(|\xi|_{(d-j+1)}^{2r} \right)^{1/r} \mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma} (\mathcal{A}_j^c)^{1-1/r}.$$

Lemmas 1, 2 and the definitions of parameters λ_j imply that

$$\mathbf{E} \left(|\xi|_{(d-j+1)}^{2r} \right)^{1/r} \leq C\lambda_s^2, \quad j = s+1, \dots, d.$$

Furthermore, it follows from the proofs of Lemmas 1 and 2 that if the constant c in the definition of λ_j is chosen large enough, then $\mathbf{P}(|\xi|_{(d-j+1)} > \lambda_j/4) \leq q^j$ for some $q < 1/2$ depending only on a and τ . This and Proposition 1 imply that $\mathbf{P}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(\mathcal{A}_j^c) \leq e^{-cd} + q^j$. Hence,

$$\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\sum_{j=s+1}^d (2\sigma |\xi|_{(d-j+1)} - \tilde{\sigma} \lambda_j)_+^2 \right) \leq C\sigma^2 \lambda_s^2 \sum_{j=s+1}^d (e^{-cd} + q^j)^{1-1/r} \leq C'\sigma^2 \sum_{j=1}^s \lambda_j^2.$$

Combining this inequality with (27) we obtain

$$(28) \quad \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} (\|\mathbf{u}\|_2^2) \leq C\sigma^2 \sum_{j=1}^s \lambda_j^2.$$

To complete the proof, it remains to note that $\sum_{j=1}^s \lambda_j^2 \leq C(\phi_{\text{pol}}^*(s, d))^2$ in the polynomial case and $\sum_{j=1}^s \lambda_j^2 \leq C(\phi_{\text{exp}}^*(s, d))^2$ in the exponential case, cf. Lemma 3.

5.3. *Proof of part (i) of Proposition 3.* We consider separately the "dense" zone $s > \sqrt{d}$ and the "sparse" zone $s \leq \sqrt{d}$. Let first $s > \sqrt{d}$. Then the rate $\phi_{\mathcal{N}(0,1)}^*(s, d)$ is of order $\sqrt{\frac{s}{1+\log_+(s^2/d)}}$. Thus, for $s > \sqrt{d}$ we need to prove that

$$(29) \quad \sup_{\sigma>0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} \left(\left| \frac{\hat{N}^* - \|\boldsymbol{\theta}\|_2}{\sigma} \right|^2 \right) \leq \frac{Cs}{1 + \log_+(s^2/d)}.$$

Denoting $\boldsymbol{\xi} = (\xi_1, \dots, \xi_d)$ we have

$$(30) \quad \begin{aligned} |\hat{N}^* - \|\boldsymbol{\theta}\|_2| &= \left| \left| \sum_{j=1}^d Y_j^2 - d\hat{\sigma}^2 \right|^{1/2} - \|\boldsymbol{\theta}\|_2 \right| \\ &= \left| \sqrt{\|\boldsymbol{\theta}\|_2^2 + 2\sigma(\boldsymbol{\theta}, \boldsymbol{\xi}) + \sigma^2\|\boldsymbol{\xi}\|_2^2 - d\hat{\sigma}^2} - \|\boldsymbol{\theta}\|_2 \right| \\ &\leq \left| \sqrt{\|\boldsymbol{\theta}\|_2^2 + 2\sigma(\boldsymbol{\theta}, \boldsymbol{\xi})} - \|\boldsymbol{\theta}\|_2 \right| + \sigma \sqrt{\|\boldsymbol{\xi}\|_2^2 - d} + \sqrt{d|\sigma^2 - \hat{\sigma}^2|}. \end{aligned}$$

The first term in the last line vanishes if $\boldsymbol{\theta} = 0$, while for $\boldsymbol{\theta} \neq 0$ it is bounded as follows:

$$(31) \quad \left| \sqrt{\|\boldsymbol{\theta}\|_2^2 + 2\sigma(\boldsymbol{\theta}, \boldsymbol{\xi})} - \|\boldsymbol{\theta}\|_2 \right| = \|\boldsymbol{\theta}\|_2 \left| \sqrt{1 + \frac{2\sigma(\boldsymbol{\theta}, \boldsymbol{\xi})}{\|\boldsymbol{\theta}\|_2^2}} - 1 \right| \leq \frac{2\sigma|(\boldsymbol{\theta}, \boldsymbol{\xi})|}{\|\boldsymbol{\theta}\|_2}$$

where we have used the inequality $|\sqrt{1+x}| - 1| \leq |x|$, $\forall x \in \mathbb{R}$. Since here $|(\boldsymbol{\theta}, \boldsymbol{\xi})|/\|\boldsymbol{\theta}\|_2 \sim \mathcal{N}(0, 1)$ we have, for all $\boldsymbol{\theta}$,

$$(32) \quad \mathbf{E} \left(\left| \sqrt{\|\boldsymbol{\theta}\|_2^2 + 2\sigma(\boldsymbol{\theta}, \boldsymbol{\xi})} - \|\boldsymbol{\theta}\|_2 \right|^2 \right) \leq 4\sigma^2,$$

and since $\|\boldsymbol{\xi}\|_2^2$ has a chi-square distribution with d degrees of freedom we have

$$\mathbf{E} \left(\|\boldsymbol{\xi}\|_2^2 - d \right) \leq \left(\mathbf{E} \left(\|\boldsymbol{\xi}\|_2^2 - d \right)^2 \right)^{1/2} = \sqrt{2d}.$$

Next, by Proposition 7 we have that, for $s > \sqrt{d}$,

$$(33) \quad \sup_{\sigma>0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} \left(\left| \frac{\hat{\sigma}^2}{\sigma^2} - 1 \right| \right) \leq \frac{Cs}{d(1 + \log_+(s^2/d))}$$

for some absolute constant $C > 0$. Combining (30) – (33) yields (29).

Let now $s \leq \sqrt{d}$. Then the rate $\phi_{\mathcal{N}(0,1)}^*(s, d)$ is of order $\sqrt{s \log(1 + d/s^2)}$. Thus, for $s \leq \sqrt{d}$ we need to prove that

$$(34) \quad \sup_{\sigma>0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} \left(\left| \frac{\hat{N}^* - \|\boldsymbol{\theta}\|_2}{\sigma} \right|^2 \right) \leq Cs \log(1 + d/s^2).$$

We have

$$\begin{aligned}
 (35) \quad \left| \hat{N}^* - \|\boldsymbol{\theta}\|_2 \right| &= \left| \left| \sum_{j=1}^d (Y_j^2 \mathbf{1}_{\{|Y_j| > \rho_j\}}) - d\alpha\hat{\sigma}^2 \right|^{1/2} - \|\boldsymbol{\theta}\|_2 \right| \\
 &= \left| \left| \sum_{j \in S} (Y_j^2 \mathbf{1}_{\{|Y_j| > \rho_j\}}) + \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > \rho_j\}}) - d\alpha\hat{\sigma}^2 \right|^{1/2} - \|\boldsymbol{\theta}\|_2 \right| \\
 &\leq \left| \sqrt{\sum_{j \in S} (Y_j^2 \mathbf{1}_{\{|Y_j| > \rho_j\}})} - \|\boldsymbol{\theta}\|_2 \right| + \left| \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > \rho_j\}}) - d\alpha\hat{\sigma}^2 \right|^{1/2}.
 \end{aligned}$$

Here,

$$\begin{aligned}
 (36) \quad \left| \sqrt{\sum_{j \in S} (Y_j^2 \mathbf{1}_{\{|Y_j| > \rho_j\}})} - \|\boldsymbol{\theta}\|_2 \right| &\leq \left| \sqrt{\sum_{j \in S} (Y_j \mathbf{1}_{\{|Y_j| > \rho_j\}} - \theta_j)^2} \right| \\
 &\leq \sqrt{\sum_{j \in S} \rho_j^2} + \sigma \sqrt{\sum_{j \in S} \xi_j^2}.
 \end{aligned}$$

Hence, writing for brevity $\mathbf{E}_{\boldsymbol{\theta}, \mathcal{N}(0,1), \sigma} = \mathbf{E}$, we get

$$\begin{aligned}
 \mathbf{E} \left(\left| \sqrt{\sum_{j \in S} (Y_j^2 \mathbf{1}_{\{|Y_j| > \rho_j\}})} - \|\boldsymbol{\theta}\|_2 \right|^2 \right) &\leq 16\mathbf{E} (\hat{\sigma}_{\text{med},1}^2 + \hat{\sigma}_{\text{med},2}^2) s \log(1 + d/s^2) + 2\sigma^2 s \\
 &\leq C\sigma^2 s \log(1 + d/s^2),
 \end{aligned}$$

where we have used the fact that $\mathbf{E}(|\hat{\sigma}_{\text{med},k}^2 - \sigma^2|) \leq C\sigma^2$, $k = 1, 2$, by Proposition 6. Next, we study the term $\Gamma = \left| \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > \rho_j\}}) - d\alpha\hat{\sigma}^2 \right|$. We first write

$$(37) \quad \Gamma \leq \left| \sigma^2 \sum_{j \notin S} \xi_j^2 (\mathbf{1}_{\{|\sigma|\xi_j| > \rho_j\}} - \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) \right| + \left| \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) - d\alpha\hat{\sigma}^2 \right|,$$

where $t_* = 2\sigma\sqrt{2\log(1 + d/s^2)}$. For the second summand on the right hand side of (37) we have

$$\left| \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) - d\alpha\hat{\sigma}^2 \right| \leq \sigma^2 \left| \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) - (d - |S|)\alpha \right| + |\sigma^2 - \hat{\sigma}^2| d\alpha + |S|\alpha\sigma^2,$$

where $|S|$ denotes the cardinality of S . By Proposition 7 we have $\mathbf{E}(|\hat{\sigma}^2 - \sigma^2|) \leq C/\sqrt{d}$ for $s \leq \sqrt{d}$. Hence,

$$\mathbf{E} \left| \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) - d\alpha\hat{\sigma}^2 \right| \leq \sigma^2 \sqrt{d\mathbf{E} \left(\xi_1^4 \mathbf{1}_{\{|\xi_1| > \sqrt{2\log(1 + d/s^2)}\}} \right)} + C\alpha\sigma^2 (\sqrt{d} + s).$$

It is not hard to check (cf., e.g., [8, Lemma 4]) that, for $s \leq \sqrt{d}$,

$$\alpha \leq C(\log(1 + d/s^2))^{1/2} \frac{s^2}{d},$$

and

$$\mathbf{E} \left(\xi_1^4 \mathbf{1}_{\{|\xi_1| > \sqrt{2 \log(1 + d/s^2)}\}} \right) \leq C(\log(1 + d/s^2))^{3/2} \frac{s^2}{d},$$

so that

$$\mathbf{E} \left| \sigma^2 \sum_{j \notin S} (\xi_j^2 \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) - d\alpha \hat{\sigma}^2 \right| \leq C\sigma^2 s \log(1 + d/s^2).$$

Thus, to complete the proof it remains to show that

$$(38) \quad \sigma^2 \sum_{j \notin S} \mathbf{E} \left| \xi_j^2 (\mathbf{1}_{\{|\sigma|\xi_j| > \rho_j\}} - \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) \right| \leq C\sigma^2 s \log(1 + d/s^2).$$

Recall that ρ_j is independent from ξ_j . Hence, conditioning on ρ_j we obtain

$$(39) \quad \sigma^2 \mathbf{E} \left(\left| \xi_j^2 (\mathbf{1}_{\{|\sigma|\xi_j| > \rho_j\}} - \mathbf{1}_{\{|\sigma|\xi_j| > t_*\}}) \right| \rho_j \right) \leq |\rho_j^2 - t_*^2| e^{-t_*^2/(8\sigma^2)} + \sigma^2 \mathbf{1}_{\{\rho_j < t_*/2\}},$$

where we have used the fact that, for $b > a > 0$,

$$\int_a^b x^2 e^{-x^2/2} dx \leq \int_a^b x e^{-x^2/4} dx \leq |b^2 - a^2| e^{-\min(a^2, b^2)/4} / 2.$$

Using Proposition 6 and definitions of ρ_j and t_* , we get that, for $s \leq \sqrt{d}$,

$$(40) \quad \mathbf{E} (|\rho_j^2 - t_*^2|) e^{-t_*^2/(8\sigma^2)} \leq 8 \max_{k=1,2} \mathbf{E} (|\hat{\sigma}_{\text{med},k}^2 - \sigma^2|) \frac{s^2}{d} \log(1 + d/s^2) \\ \leq C\sigma^2 \frac{s}{d} \log(1 + d/s^2).$$

Next, it follows from Proposition 6 that there exists $\gamma \in (0, 1/8)$ small enough such that for $s \leq \gamma d$ we have $\max_{k=1,2} \mathbf{P}(\hat{\sigma}_{\text{med},k}^2 < \sigma^2/2) \leq 2e^{-c_\gamma d}$ where $c_\gamma > 0$ is a constant. Thus, $\sigma^2 \mathbf{P}(\rho_j < t_*/2) \leq 2\sigma^2 e^{-c_\gamma d} \leq C\sigma^2 (s/d) \log(1 + d/s^2)$. Combining this with (39) and (40) proves (38).

5.4. *Proof of part (i) of Proposition 4 and part (i) of Proposition 5.* We only prove Proposition 4 since the proof of Proposition 5 is similar taking into account that $\mathbf{E}(\xi_1^4) < \infty$. We consider separately the "dense" zone $s > \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}$ and the "sparse" zone $s \leq \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}$. Let first $s > \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}$. Then the rate $\phi_{\text{exp}}^\circ(s, d)$ is of order $d^{1/4}$ and thus we need to prove that

$$\sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\|\theta\|_0 \leq s} \mathbf{E}_{\theta, P_\xi, \sigma} (|\hat{N}_{\text{exp}}^\circ - \|\theta\|_2|^2) \leq C\sigma^2 \sqrt{d}.$$

Since σ is known, arguing similarly to (30) - (31) we find

$$|\hat{N}_{\text{exp}}^\circ - \|\theta\|_2| \leq \left| \frac{2\sigma |(\theta, \xi)|}{\|\theta\|_2} \right| \mathbf{1}_{\theta \neq 0} + \sigma \sqrt{\|\xi\|_2^2 - d}.$$

As $\mathbf{E}(\xi_1^4) < \infty$, this implies

$$\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(|\hat{N}_{\text{exp}}^{\circ} - \|\boldsymbol{\theta}\|_2|^2) \leq 8\sigma^2 + C\sigma^2\sqrt{d},$$

which proves the result in the dense case. Next, in the sparse case $s \leq \frac{\sqrt{d}}{\log^{\frac{2}{a}}(ed)}$, we need to prove that

$$\sup_{P_{\xi} \in \mathcal{G}_{a, \tau}} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}(|\hat{N}_{\text{exp}}^{\circ} - \|\boldsymbol{\theta}\|_2|^2) \leq C\sigma^2 s \log^{\frac{2}{a}}(ed).$$

This is immediate by Theorem 1 and the fact that $|\hat{N}_{\text{exp}}^{\circ} - \|\boldsymbol{\theta}\|_2|^2 \leq \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2$ for the plug-in estimator $\hat{N}_{\text{exp}}^{\circ} = \|\hat{\boldsymbol{\theta}}\|_2$.

5.5. Proof of Proposition 6. Denote by G the cdf of $(\sigma\xi_1)^2$ and by G_d the empirical cdf of $((\sigma\xi_i)^2 : i \notin S)$, where S is the support of $\boldsymbol{\theta}$. Let M be the median of G , that is $G(M) = 1/2$. By the definition of \hat{M} ,

$$(41) \quad |F_d(\hat{M}) - 1/2| \leq |F_d(M) - 1/2|.$$

It is easy to check that $|F_d(x) - G_d(x)| \leq s/d$ for all $x > 0$. Therefore,

$$(42) \quad |G_d(\hat{M}) - 1/2| \leq |G_d(M) - 1/2| + 2s/d.$$

The DKW inequality [24, page 99], yields that $\mathbf{P}(\sup_{x \in \mathbb{R}} |G_d(x) - G(x)| \geq u) \leq 2e^{-2u^2(d-s)}$ for all $u > 0$. Fix $t > 0$ such that $\sqrt{\frac{t}{d}} + \frac{s}{d} \leq 1/8$, and consider the event

$$\mathcal{A} := \left\{ \sup_{x \in \mathbb{R}} |G_d(x) - G(x)| \leq \sqrt{\frac{t}{2(d-s)}} \right\}.$$

Then, $\mathbf{P}(\mathcal{A}) \geq 1 - 2e^{-t}$. On the event \mathcal{A} , we have

$$(43) \quad |G(\hat{M}) - 1/2| \leq |G(M) - 1/2| + 2 \left(\sqrt{\frac{t}{2(d-s)}} + \frac{s}{d} \right) \leq 2 \left(\sqrt{\frac{t}{d}} + \frac{s}{d} \right) \leq \frac{1}{4},$$

where the last two inequalities are due to the fact that $G(M) = 1/2$ and to the assumption about t . Notice that

$$(44) \quad |G(\hat{M}) - 1/2| = |G(\hat{M}) - G(M)| = |F(\hat{M}/\sigma^2) - F(M/\sigma^2)|.$$

Using (43), (44) and the fact that $M = \sigma^2 F^{-1}(1/2)$ we obtain that, on the event \mathcal{A} ,

$$(45) \quad F^{-1}(1/4) \leq \hat{M}/\sigma^2 \leq F^{-1}(3/4).$$

This and (44) imply

$$(46) \quad |G(\hat{M}) - 1/2| \geq c_{**} |\hat{M}/\sigma^2 - M/\sigma^2| = c_{**} \beta |\hat{\sigma}_{\text{med}}^2/\sigma^2 - 1|.$$

where $c_{**} = \min_{x \in [F^{-1}(1/4), F^{-1}(3/4)]} F'(x) > 0$, and $\beta = F^{-1}(1/2)$. Combining the last inequality with (43) we get that, on the event \mathcal{A} ,

$$|\hat{\sigma}_{\text{med}}^2/\sigma^2 - 1| \leq c_{**}^{-1} \beta \left(\sqrt{\frac{t}{d}} + \frac{s}{d} \right).$$

Recall that we assumed that $\sqrt{\frac{t}{d}} + \frac{s}{d} \leq 1/8$. Thus, there exists a constant $c_* > 0$ depending only on F such that for $t > 0$ and integers s, d satisfying $\sqrt{\frac{t}{d}} + \frac{s}{d} \leq 1/8$ we have

$$(47) \quad \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{P}_{\theta, F, \sigma} \left(\left| \frac{\hat{\sigma}_{\text{med}}^2}{\sigma^2} - 1 \right| \geq c_* \left(\sqrt{\frac{t}{d}} + \frac{s}{d} \right) \right) \leq 2e^{-t}.$$

This and the assumption that $\frac{s}{d} \leq \gamma < 1/8$ imply the result of the proposition in probability. We now prove the result in expectation. Set $Z = |\hat{\sigma}_{\text{med}}^2 - \sigma^2|/\sigma^2$. We have

$$\mathbf{E}_{\theta, F, \sigma}(Z) \leq c_* s/d + \int_{c_* s/d}^{c_*/8} \mathbf{P}_{\theta, F, \sigma}(Z > u) du + \mathbf{E}_{\theta, F, \sigma}(Z \mathbf{1}_{Z \geq c_*/8}).$$

Using (47), we get

$$\int_{c_* s/d}^{c_*/8} \mathbf{P}_{\theta, F, \sigma}(Z > u) du \leq \frac{2c_*}{\sqrt{d}} \int_0^\infty e^{-t^2} dt \leq \frac{C}{\sqrt{d}}.$$

As $s < d/2$, one may check that $\hat{\sigma}_{\text{med}}^{2+\epsilon} \leq (\max_{i \notin S} (\sigma \xi_i)^2 / \beta)^{1+\epsilon/2} \leq (\sigma^2 / \beta)^{1+\epsilon/2} \sum_{i=1}^d |\xi_i|^{2+\epsilon}$. Since $\mathbf{E}|\xi_1|^{2+\epsilon} < \infty$ this yields $\mathbf{E}_{\theta, F, \sigma}(Z^{1+\epsilon}) \leq Cd$. It follows that

$$\mathbf{E}_{\theta, F, \sigma}(Z \mathbf{1}_{Z \geq c_*/8}) \leq (\mathbf{E}_{\theta, F, \sigma}(Z^{1+\epsilon}))^{1/(1+\epsilon)} \mathbf{P}_{\theta, F, \sigma}(Z \geq c_*/8)^{\epsilon/(1+\epsilon)} \leq Cde^{-d/C}.$$

Combining the last three displays yields the desired bound in expectation.

5.6. *Proof of part (i) of Proposition 7.* In this proof, we write for brevity $\mathbf{E} = \mathbf{E}_{\theta, \sigma, \mathcal{N}(0,1)}$ and $\mathbf{P} = \mathbf{P}_{\theta, \sigma, \mathcal{N}(0,1)}$. Set

$$\varphi_d(t) = \frac{1}{d} \sum_{i=1}^d e^{itY_j}, \quad \varphi(t) = \mathbf{E}(\varphi_d(t)), \quad \varphi_0(t) = e^{-\frac{t^2 \sigma^2}{2}}.$$

Since $s/d < 1/8$ and $\varphi(t) = \varphi_0(t)(1 - \frac{|S|}{d} + \frac{1}{d} \sum_{j \in S} \exp(i\theta_j t))$, we have

$$(48) \quad \frac{3}{4} \varphi_0(t) \leq \left(1 - \frac{2s}{d}\right) \varphi_0(t) \leq |\varphi(t)| \leq \varphi_0(t).$$

Consider the events

$$\mathcal{B}_1 = \left\{ \sigma^2/2 \leq \tilde{\sigma}^2 \leq 3\sigma^2/2 \right\} \quad \text{and} \quad \mathcal{A}_u = \left\{ \sup_{v \in \mathbb{R}} |\varphi_d(v) - \varphi(v)| \leq \sqrt{\frac{u}{d}} \right\}, \quad u > 0.$$

By Proposition 1, \mathcal{B}_1 holds with probability at least $1 - e^{-cd}$ if the tuning parameter γ in the definition of $\tilde{\sigma}^2$ is small enough. Using Hoeffding's inequality, it is not hard to check that \mathcal{A}_u holds with probability at least $1 - 4e^{-u}$. Moreover,

$$(49) \quad \mathbf{E} \left(\sqrt{d} \sup_{v \in \mathbb{R}} |\varphi_d(v) - \varphi(v)| \right) \leq C.$$

Notice that on the event $\mathcal{D} = \{|\varphi_d(\hat{t}_1)| > (es/\sqrt{d} + 1)^{-1}/4\}$ we have $\hat{\sigma}^2 = \tilde{v}^2 \leq 2\tilde{\sigma}^2$. First, we bound the risk restricted to $\mathcal{D} \cap \mathcal{B}_1^c$. We have

$$\mathbf{E}(|\hat{\sigma}^2 - \sigma^2| \mathbf{1}_{\mathcal{D} \cap \mathcal{B}_1^c}) \leq \mathbf{E}(|2\tilde{\sigma}^2 + \sigma^2| \mathbf{1}_{\mathcal{B}_1^c}).$$

Thus, using the Cauchy-Schwarz inequality and Proposition 1 we find

$$(50) \quad \mathbf{E}(|\hat{\sigma}^2 - \sigma^2| \mathbf{1}_{\mathcal{D} \cap \mathcal{B}_1^c}) \leq C\sigma^2 e^{-d/C} \leq \frac{C'\sigma^2}{\sqrt{d}}.$$

Next, we bound the risk restricted to \mathcal{D}^c . It will be useful to note that $\mathcal{A}_{\log d} \cap \mathcal{B}_1 \subset \mathcal{D}$. Indeed, on $\mathcal{A}_{\log d} \cap \mathcal{B}_1$, using the assumption $s < d/8$ we have

$$|\varphi_d(\hat{t}_1)| \geq \frac{3}{4}\varphi_0(\hat{t}_1) - \sqrt{\frac{\log d}{d}} \geq \frac{3}{4(es/\sqrt{d} + 1)^{1/3}} - \sqrt{\frac{\log d}{d}} > \frac{1}{4(es/\sqrt{d} + 1)}.$$

Thus, applying again the Cauchy-Schwarz inequality and Proposition 1 we find

$$(51) \quad \begin{aligned} \mathbf{E}(|\hat{\sigma}^2 - \sigma^2| \mathbf{1}_{\mathcal{D}^c}) &= \mathbf{E}(|\tilde{\sigma}^2 - \sigma^2| \mathbf{1}_{\mathcal{D}^c}) \leq (\mathbf{E}(|\tilde{\sigma}^2 - \sigma^2|^2))^{1/2} (\mathbf{P}(\mathcal{D}^c))^{1/2} \\ &\leq C\sigma^2 \sqrt{\mathbf{P}(\mathcal{A}_{\log d}^c) + \mathbf{P}(\mathcal{B}_1^c)} \leq C\sigma^2 \sqrt{\frac{4}{d} + e^{-cd}} \leq \frac{C'\sigma^2}{\sqrt{d}}. \end{aligned}$$

To complete the proof, it remains to handle the risk restricted to the event $\mathcal{C} = \mathcal{D} \cap \mathcal{B}_1$. We will use the following decomposition

$$(52) \quad |\hat{\sigma}^2 - \sigma^2| \leq \left| \frac{2\log(|\varphi_d(\hat{t}_1)|)}{\hat{t}_1^2} - \frac{2\log(|\varphi(\hat{t}_1)|)}{\hat{t}_1^2} \right| + \left| -\frac{2\log(|\varphi(\hat{t}_1)|)}{\hat{t}_1^2} - \sigma^2 \right|.$$

Since $-2\log(|\varphi_0(\hat{t}_1)|)/\hat{t}_1^2 = \sigma^2$, it follows from (48) that

$$\left| -\frac{2\log(|\varphi(\hat{t}_1)|)}{\hat{t}_1^2} - \sigma^2 \right| \leq \frac{Cs}{d\hat{t}_1^2} = \frac{Cs\tilde{\sigma}^2}{d\log(4(es/\sqrt{d} + 1))}.$$

Therefore,

$$(53) \quad \mathbf{E}\left(\left| -\frac{2\log(|\varphi(\hat{t}_1)|)}{\hat{t}_1^2} - \sigma^2 \right| \mathbf{1}_{\mathcal{C}}\right) \leq \frac{Cs\sigma^2}{d\log(es/\sqrt{d} + 1)}.$$

Next, using the inequality

$$|\log(|\varphi_d(t)|) - \log(|\varphi(t)|)| \leq \frac{|\varphi_d(t) - \varphi(t)|}{|\varphi(t)| \wedge |\varphi_d(t)|}, \quad \forall t \in \mathbb{R},$$

we find

$$\begin{aligned} \left| \frac{\log(|\varphi_d(\hat{t}_1)|)}{\hat{t}_1^2} - \frac{\log(|\varphi(\hat{t}_1)|)}{\hat{t}_1^2} \right| \mathbf{1}_{\mathcal{C}} &\leq \frac{\sup_{v \in \mathbb{R}} |\varphi_d(v) - \varphi(v)|}{\hat{t}_1^2 |\varphi(\hat{t}_1)| \wedge |\varphi_d(\hat{t}_1)|} \mathbf{1}_{\mathcal{C}} \\ &\leq \frac{C\sigma^2 U}{\sqrt{d} \log(es/\sqrt{d} + 1)} \left(\frac{es}{\sqrt{d}} + 1 \right), \end{aligned}$$

where $U = \sqrt{d} \sup_{v \in \mathbb{R}} |\varphi_d(v) - \varphi(v)|$. Bounding $\mathbf{E}(U)$ by (49) we finally get

$$(54) \quad \mathbf{E}\left[\left| \frac{\log(|\varphi_d(\hat{t}_1)|)}{\hat{t}_1^2} - \frac{\log(|\varphi(\hat{t}_1)|)}{\hat{t}_1^2} \right| \mathbf{1}_{\mathcal{C}}\right] \leq C\sigma^2 \max\left(\frac{1}{\sqrt{d}}, \frac{s}{d\log(es/\sqrt{d} + 1)}\right).$$

We conclude by combining inequalities (50) - (54).

5.7. *Proof of Theorems 3 and 4.* Let $\|\boldsymbol{\theta}\|_0 \leq s$ and denote by S the support of $\boldsymbol{\theta}$. Note first that, by the definition of $\hat{\sigma}^2$,

$$(55) \quad \frac{\sigma^2}{d} \sum_{i=1}^{d-2s} \xi_{(i)}^2 \leq \hat{\sigma}^2 \leq \frac{\sigma^2}{d} \sum_{i \in S^c} \xi_i^2,$$

where $\xi_{(1)}^2 \leq \dots \leq \xi_{(d)}^2$ are the ordered values of ξ_1^2, \dots, ξ_d^2 . Indeed, the right hand inequality in (55) follows from the relations

$$\sum_{k=1}^{d-s} Y_{(k)}^2 = \min_{J: |J|=d-s} \sum_{i \in J} Y_{(i)}^2 \leq \sum_{i \in S^c} Y_{(i)}^2 = \sum_{i \in S^c} \sigma^2 \xi_i^2.$$

To show the left hand inequality in (55), notice that at least $d-2s$ among the $d-s$ order statistics $Y_{(1)}^2, \dots, Y_{(d-s)}^2$ correspond to observations Y_k of pure noise, *i.e.*, $Y_k = \sigma \xi_k$. The sum of squares of such observations is bounded from below by the sum of the smallest $d-2s$ values $\sigma^2 \xi_{(1)}^2, \dots, \sigma^2 \xi_{(d-2s)}^2$ among $\sigma^2 \xi_1^2, \dots, \sigma^2 \xi_d^2$.

Using (55) we get

$$\left(\hat{\sigma}^2 - \frac{\sigma^2}{d} \sum_{i=1}^d \xi_i^2 \right)^2 \leq \frac{\sigma^4}{d^2} \left(\sum_{i=d-2s+1}^d \xi_{(i)}^2 \right)^2,$$

so that

$$\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\hat{\sigma}^2 - \frac{\sigma^2}{d} \sum_{i=1}^d \xi_i^2 \right)^2 \leq \frac{\sigma^4}{d^2} \left(\sum_{i=1}^{2s} \sqrt{\mathbf{E} \xi_{(d-i+1)}^4} \right)^2.$$

Then

$$\begin{aligned} \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} (\hat{\sigma}^2 - \sigma^2)^2 &\leq 2\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\hat{\sigma}^2 - \frac{\sigma^2}{d} \sum_{i=1}^d \xi_i^2 \right)^2 + 2\mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma} \left(\frac{\sigma^2}{d} \sum_{i=1}^d \xi_i^2 - \sigma^2 \right)^2 \\ &\leq \frac{2\sigma^4}{d^2} \left(\sum_{i=1}^{2s} \sqrt{\mathbf{E} \xi_{(d-i+1)}^4} \right)^2 + \frac{2\sigma^4 \mathbf{E}(\xi_1^4)}{d}. \end{aligned}$$

Note that under assumption (2) we have $\mathbf{E}(\xi_1^4) < \infty$ and Lemmas 1 and 3 yield

$$\sum_{i=1}^{2s} \sqrt{\mathbf{E} \xi_{(d-i+1)}^4} \leq \sqrt{C} \sum_{i=1}^{2s} \log^{2/a} (ed/i) \leq C' \sqrt{C} s \log^{2/a} \left(\frac{ed}{2s} \right).$$

This proves Theorem 3. To prove Theorem 4, we act analogously by using Lemma 2 and the fact that $\mathbf{E}(\xi_1^4) < \infty$ under assumption (3) with $a > 4$.

5.8. *Proof of Theorem 7.* With the same notation as in the proof of Theorem 1, we have

$$(56) \quad \hat{\sigma}^2 - \sigma^2 = \frac{\sigma^2}{d} (\|\boldsymbol{\xi}\|_2^2 - d) + \frac{1}{d} (\|\mathbf{u}\|_2^2 - 2\sigma \mathbf{u}^T \boldsymbol{\xi}).$$

It follows from (25) that

$$\|\mathbf{u}\|_2^2 + 2\sigma |\mathbf{u}^T \boldsymbol{\xi}| \leq 3\sigma |\mathbf{u}^T \boldsymbol{\xi}| + \frac{\tilde{\sigma}}{2} \left\{ \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} \|\mathbf{u}\|_2 - \sum_{j=s+1}^d \lambda_j |u|_{(d-j+1)} \right\}.$$

Arguing as in the proof of Theorem 1, we obtain

$$\|\mathbf{u}\|_2^2 + 2\sigma|\mathbf{u}^T \boldsymbol{\xi}| \leq \left(U_1 + \frac{\tilde{\sigma}}{2} \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} + U_2 \right) \|\mathbf{u}\|_2,$$

where

$$U_1 = 3\sigma \left(\sum_{j=1}^s |\xi|_{(d-j+1)}^2 \right)^{1/2}, \quad U_2 = \left(\sum_{j=s+1}^d \left(3\sigma |\xi|_{(d-j+1)} - \frac{\tilde{\sigma}}{2} \lambda_j \right)_+^2 \right)^{1/2}$$

Using the Cauchy-Schwarz inequality, Proposition 1 and (28) and writing for brevity $\mathbf{E} = \mathbf{E}_{\boldsymbol{\theta}, P_{\xi}, \sigma}$ we find

$$\mathbf{E} \left(\tilde{\sigma} \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} \|\mathbf{u}\|_2 \right) \leq \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} \sqrt{\mathbf{E}(\tilde{\sigma}^2)} \sqrt{\mathbf{E}(\|\mathbf{u}\|_2^2)} \leq C\sigma^2 \sum_{j=1}^s \lambda_j^2.$$

Since $\mathbf{E}(\xi_1^4) < \infty$ we also have $\mathbf{E} \left| \|\boldsymbol{\xi}\|_2^2 - d \right| \leq C\sqrt{d}$. Finally, using again (28) we get, for $k = 1, 2$,

$$\mathbf{E}(U_k \|\mathbf{u}\|_2) \leq \sqrt{\mathbf{E}(\|\mathbf{u}\|_2^2)} \sqrt{\mathbf{E}(U_k^2)} \leq \sigma \left(\sum_{j=1}^s \lambda_j^2 \right)^{1/2} \sqrt{\mathbf{E}(U_k^2)} \leq C\sigma^2 \sum_{j=1}^s \lambda_j^2,$$

where the last inequality follows from the same argument as in the proof of Theorem 1. These remarks together with (56) imply

$$\mathbf{E} (|\hat{\sigma}^2 - \sigma^2|) \leq \frac{C}{d} \left(\sigma^2 \sqrt{d} + \sigma^2 \sum_{j=1}^s \lambda_j^2 \right).$$

We conclude the proof by bounding $\sum_{j=1}^s \lambda_j^2$ in the same way as in the end of the proof of Theorem 1.

6. Proofs of the lower bounds.

6.1. *Proof of Theorems 5 and 6 and part (ii) of Proposition 7.* Since we have $\ell(t) \geq \ell(A) \mathbf{1}_{t \geq A}$ for any $A > 0$, it is enough to prove the theorems for the indicator loss $\ell(t) = \mathbf{1}_{t \geq 1}$. This remark is valid for all the proofs of this section and will not be further repeated.

(i) We first prove the lower bounds with the rate $1/\sqrt{d}$ in Theorems 5 and 6. Let $f_0 : \mathbb{R} \rightarrow [0, \infty)$ be a probability density with the following properties: f_0 is continuously differentiable, symmetric about 0, supported on $[-3/2, 3/2]$, with variance 1 and finite Fisher information $I_{f_0} = \int (f_0'(x))^2 (f_0(x))^{-1} dx$. The existence of such f_0 is shown in Lemma 7. Denote by F_0 the probability distribution corresponding to f_0 . Since F_0 is zero-mean, with variance 1 and supported on $[-3/2, 3/2]$ it belongs to $\mathcal{G}_{a, \tau}$ with any $\tau > 0$, $a > 0$, and to $\mathcal{P}_{a, \tau}$ with any $\tau > 0$, $a \geq 2$. Define $\mathbf{P}_0 = \mathbf{P}_{0, F_0, 1}$, $\mathbf{P}_1 = \mathbf{P}_{0, F_0, \sigma_1}$ where $\sigma_1^2 = 1 + c_0/\sqrt{d}$ and $c_0 > 0$ is a small constant to be fixed later. Denote by $H(\mathbf{P}_1, \mathbf{P}_0)$ the Hellinger distance between \mathbf{P}_1 and \mathbf{P}_0 . We have

$$(57) \quad H^2(\mathbf{P}_1, \mathbf{P}_0) = 2(1 - (1 - h^2/2)^d)$$

where $h^2 = \int (\sqrt{f_0(x)} - \sqrt{f_0(x/\sigma_1)/\sigma_1})^2 dx$. By Theorem 7.6. in Ibragimov and Hasminskii [15],

$$h^2 \leq \frac{(1 - \sigma_1)^2}{4} \sup_{t \in [1, \sigma_1]} I(t)$$

where $I(t)$ is the Fisher information corresponding to the density $f_0(x/t)/t$, that is $I(t) = t^{-2}I_{f_0}$. It follows that $h^2 \leq \bar{c}c_0^2/d$ where $\bar{c} > 0$ is a constant. This and (57) imply that for c_0 small enough we have $H(\mathbf{P}_1, \mathbf{P}_0) \leq 1/2$. Finally, choosing such a small c_0 and using Theorem 2.2(ii) in Tsybakov [21] we obtain

$$\begin{aligned} & \inf_{\hat{T}} \max \left\{ \mathbf{P}_0 \left(\left| \hat{T} - 1 \right| \geq \frac{c_0}{2(1+c_0)\sqrt{d}} \right), \mathbf{P}_1 \left(\left| \frac{\hat{T}}{\sigma_1^2} - 1 \right| \geq \frac{c_0}{2(1+c_0)\sqrt{d}} \right) \right\} \\ & \geq \inf_{\hat{T}} \max \left\{ \mathbf{P}_0 \left(\left| \hat{T} - 1 \right| \geq \frac{c_0}{2\sqrt{d}} \right), \mathbf{P}_1 \left(\left| \hat{T} - \sigma_1^2 \right| \geq \frac{c_0}{2\sqrt{d}} \right) \right\} \geq \frac{1 - H(\mathbf{P}_1, \mathbf{P}_0)}{2} \geq \frac{1}{4}. \end{aligned}$$

(ii) We now prove the lower bound with the rate $\frac{s}{d} \log^{2/a}(ed/s)$ in Theorem 5. It is enough to conduct the proof for $s \geq s_0$ where $s_0 > 0$ is an arbitrary absolute constant. Indeed, for $s \leq s_0$ we have $\frac{s}{d} \log^{2/a}(ed/s) \leq C/\sqrt{d}$ where $C > 0$ is an absolute constant and thus Theorem 5 follows already from the lower bound with the rate $1/\sqrt{d}$ proved in item (i). Therefore, in the rest of this proof we assume without loss of generality that $s \geq 32$.

We take $P_\xi = U$ where U is the Rademacher distribution, that is the uniform distribution on $\{-1, 1\}$. Clearly, $U \in \mathcal{G}_{a,\tau}$. Let $\delta_1, \dots, \delta_d$ be i.i.d. Bernoulli random variables with probability of success $\mathbf{P}(\delta_1 = 1) = \frac{s}{2d}$, and let $\epsilon_1, \dots, \epsilon_d$ be i.i.d. Rademacher random variables that are independent of $(\delta_1, \dots, \delta_d)$. Denote by μ the distribution of $(\alpha\delta_1\epsilon_1, \dots, \alpha\delta_d\epsilon_d)$ where $\alpha = (\tau/2) \log^{1/a}(ed/s)$. Note that μ is not necessarily supported on $\Theta_s = \{\boldsymbol{\theta} \in \mathbb{R}^d \mid \|\boldsymbol{\theta}\|_0 \leq s\}$ as the number of nonzero components of a vector drawn from μ can be larger than s . Therefore, we consider a restricted to Θ_s version of μ defined by

$$(58) \quad \bar{\mu}(A) = \frac{\mu(A \cap \Theta_s)}{\mu(\Theta_s)}$$

for all Borel subsets A of \mathbb{R}^d . Finally, we introduce two mixture probability measures

$$(59) \quad \mathbb{P}_\mu = \int \mathbf{P}_{\boldsymbol{\theta}, U, 1} \mu(d\boldsymbol{\theta}) \quad \text{and} \quad \mathbb{P}_{\bar{\mu}} = \int \mathbf{P}_{\boldsymbol{\theta}, U, 1} \bar{\mu}(d\boldsymbol{\theta}).$$

Notice that there exists a probability measure $\tilde{P} \in \mathcal{G}_{a,\tau}$ such that

$$(60) \quad \mathbb{P}_\mu = \mathbf{P}_{0, \tilde{P}, \sigma_0}$$

where $\sigma_0 > 0$ is defined by

$$(61) \quad \sigma_0^2 = 1 + \frac{\tau^2 s}{8d} \log^{2/a}(ed/s) \leq 1 + \frac{\tau^2}{8}.$$

Indeed, $\sigma_0^2 = 1 + \frac{\alpha^2 s}{2d}$ is the variance of zero-mean random variable $\alpha\delta\epsilon + \xi$, where $\xi \sim U$, $\epsilon \sim U$, $\delta \sim \mathcal{B}(\frac{s}{2d})$ and ϵ, ξ, δ are jointly independent. Thus, to prove (60) it is enough to show that, for all $t \geq 2$,

$$(62) \quad \mathbf{P}((\tau/2) \log^{1/a}(ed/s) \delta\epsilon + \xi > t\sigma_0) \leq e^{-(t/\tau)^a}.$$

But this inequality immediately follows from the fact that for $t \geq 2$ the probability in (62) is smaller than

$$(63) \quad \mathbf{P}(\epsilon = 1, \delta = 1) \mathbf{1}_{(\tau/2) \log^{1/a}(ed/s) > t-1} \leq \frac{s}{4d} \mathbf{1}_{\tau \log^{1/a}(ed/s) > t} \leq e^{-(t/\tau)^a}.$$

Now, for any estimator \hat{T} and any $u > 0$ we have

$$\begin{aligned}
& \sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{P}_{\theta, P_\xi, \sigma} \left(\left| \frac{\hat{T}}{\sigma^2} - 1 \right| \geq u \right) \\
& \geq \max \left\{ \mathbf{P}_{0, \tilde{P}, \sigma_0} (|\hat{T} - \sigma_0^2| \geq \sigma_0^2 u), \int \mathbf{P}_{\theta, U, 1} (|\hat{T} - 1| \geq u) \bar{\mu}(d\theta) \right\} \\
(64) \quad & \geq \max \left\{ \mathbb{P}_\mu (|\hat{T} - \sigma_0^2| \geq \sigma_0^2 u), \mathbb{P}_{\bar{\mu}} (|\hat{T} - 1| \geq \sigma_0^2 u) \right\}
\end{aligned}$$

where the last inequality uses (60). Write $\sigma_0^2 = 1 + 2\phi$ where $\phi = \frac{\tau^2 s}{16d} \log^{2/a}(ed/s)$ and choose $u = \phi/\sigma_0^2 \geq \phi/(1 + \tau^2/8)$. Then, the expression in (64) is bounded from below by the probability of error in the problem of distinguishing between two simple hypotheses \mathbb{P}_μ and $\mathbb{P}_{\bar{\mu}}$, for which Theorem 2.2 in Tsybakov [21] yields

$$(65) \quad \max \left\{ \mathbb{P}_\mu (|\hat{T} - \sigma_0^2| \geq \phi), \mathbb{P}_{\bar{\mu}} (|\hat{T} - 1| \geq \phi) \right\} \geq \frac{1 - V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}})}{2}$$

where $V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}})$ is the total variation distance between \mathbb{P}_μ and $\mathbb{P}_{\bar{\mu}}$. The desired lower bound follows from (65) and Lemma 5 for any $s \geq 32$.

(iii) Finally, we prove the lower bound with the rate $\tau^2(s/d)^{1-2/a}$ in Theorem 6. Again, we do not consider the case $s \leq 32$ since in this case the rate $1/\sqrt{d}$ is dominating and Theorem 6 follows from item (i) above. For $s \geq 32$, the proof uses the same argument as in item (ii) above but we choose $\alpha = (\tau/2)(d/s)^{1/a}$. Then the variance of $\alpha\delta\epsilon + \xi$ is equal to

$$\sigma_0^2 = 1 + \frac{\tau^2(s/d)^{1-2/a}}{8}.$$

Furthermore, with this definition of σ_0^2 there exists $\tilde{P} \in \mathcal{P}_{a,\tau}$ such that (60) holds. Indeed, analogously to (62) we now have, for all $t \geq 2$,

$$(66) \quad \mathbf{P}(\alpha\delta\epsilon + \xi > t\sigma_0) \leq \mathbf{P}(\epsilon = 1, \delta = 1) \mathbf{1}_{(\tau/2)(d/s)^{1/a} > t-1} \leq \frac{s}{4d} \mathbf{1}_{\tau(d/s)^{1/a} > t} \leq (t/\tau)^a.$$

To finish the proof, it remains to repeat the argument of (64) and (65) with $\phi = \frac{\tau^2(s/d)^{1-2/a}}{16}$.

6.2. Proof of Theorem 2. We argue similarly to the proof of Theorems 5 and 6, in particular, we set $\alpha = (\tau/2) \log^{1/a}(ed/s)$ when proving the bound on the class $\mathcal{G}_{a,\tau}$, and $\alpha = (\tau/2)(d/s)^{1/a}$ when proving the bound on $\mathcal{P}_{a,\tau}$. In what follows, we only deal with the class $\mathcal{G}_{a,\tau}$ since the proof for $\mathcal{P}_{a,\tau}$ is analogous. Consider the measures $\mu, \bar{\mu}, \mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}$ and \tilde{P} defined in Section 6.1. Similarly to (64), for any estimator \hat{T} and any $u > 0$ we have

$$\begin{aligned}
& \sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{P}_{\theta, P_\xi, \sigma} (|\hat{T} - \|\theta\|_2| \geq \sigma u) \\
& \geq \max \left\{ \mathbf{P}_{0, \tilde{P}, \sigma_0} (|\hat{T}| \geq \sigma_0 u), \int \mathbf{P}_{\theta, U, 1} (|\hat{T} - \|\theta\|_2| \geq u) \bar{\mu}(d\theta) \right\} \\
& \geq \max \left\{ \mathbb{P}_\mu (|\hat{T}| \geq \sigma_0 u), \mathbb{P}_{\bar{\mu}} (|\hat{T} - \|\theta\|_2| \geq \sigma_0 u) \right\} \\
& \geq \max \left\{ \mathbb{P}_\mu (|\hat{T}| \geq \sigma_0 u), \mathbb{P}_{\bar{\mu}} (|\hat{T}| < \sigma_0 u, \|\theta\|_2 \geq 2\sigma_0 u) \right\} \\
& \geq \min_B \max \left\{ \mathbb{P}_\mu(B), \mathbb{P}_{\bar{\mu}}(B^c) - \bar{\mu}(\|\theta\|_2 < 2\sigma_0 u) \right\} \\
(67) \quad & \geq \min_B \frac{\mathbb{P}_\mu(B) + \mathbb{P}_{\bar{\mu}}(B^c)}{2} - \frac{\bar{\mu}(\|\theta\|_2 < 2\sigma_0 u)}{2}
\end{aligned}$$

where σ_0 is defined in (61), U denotes the Rademacher law and \min_B is the minimum over all Borel sets. The third line in the last display is due to (60) and to the inequality $\sigma_0 \geq 1$. Since $\min_B \{\mathbb{P}_\mu(B) + \mathbb{P}_{\bar{\mu}}(B^c)\} = 1 - V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}})$, we get

$$(68) \quad \sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\sigma > 0} \sup_{\|\theta\|_0 \leq s} \mathbf{P}_{\theta, P_\xi, \sigma}(|\hat{T} - \|\theta\|_2|/\sigma \geq u) \geq \frac{1 - V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) - \bar{\mu}(\|\theta\|_2 < 2\sigma_0 u)}{2}.$$

Consider first the case $s \geq 32$. Set $u = \frac{\alpha\sqrt{s}}{4\sigma_0}$. Then (77) and (80) imply that

$$V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) \leq e^{-\frac{3s}{16}}, \quad \bar{\mu}(\|\theta\|_2 < 2\sigma_0 u) \leq 2e^{-\frac{s}{16}},$$

which, together with (68) and the fact that $s \geq 32$ yields the result.

Let now $s < 32$. Then we set $u = \frac{\alpha\sqrt{s}}{8\sqrt{2}\sigma_0}$. It follows from (78) and (81) that

$$1 - V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) - \bar{\mu}(\|\theta\|_2 < 2\sigma_0 u) \geq \mathbf{P}\left(\mathcal{B}(d, \frac{s}{2d}) = 1\right) = \frac{s}{2}\left(1 - \frac{s}{2d}\right)^{d-1}.$$

It is not hard to check that the minimum of the last expression over all integers s, d such that $1 \leq s < 32$, $s \leq d$, is bounded from below by a positive number independent of d . We conclude by combining these remarks with (68).

6.3. Proof of part (ii) of Proposition 4 and part (ii) of Proposition 5. We argue similarly to the proof of Theorems 5 and 6, in particular, we set $\alpha = (\tau/2) \log^{1/a}(ed/s)$ when proving the bound on the class $\mathcal{G}_{a,\tau}$, and $\alpha = (\tau/2)(d/s)^{1/a}$ when proving the bound on $\mathcal{P}_{a,\tau}$. In what follows, we only deal with the class $\mathcal{G}_{a,\tau}$ since the proof for $\mathcal{P}_{a,\tau}$ is analogous. Without loss of generality we assume that $\sigma = 1$.

To prove the lower bound with the rate $\phi_{\text{exp}}^\circ(s, d)$, we only need to prove it for s such that $(\phi_{\text{exp}}^\circ(s, d))^2 \leq c_0 \sqrt{d} / \log^{2/a}(ed)$ with any small absolute constant $c_0 > 0$, since the rate is increasing with s .

Consider the measures $\mu, \bar{\mu}, \mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}$ defined in Section 6.1 with $\sigma_0 = 1$. Let ξ_1 be distributed with c.d.f. F_0 defined in item (i) of the proof of Theorems 5 and 6. Using the notation as in the proof of Theorems 5 and 6, we define \tilde{P} as the distribution of $\tilde{\xi}_1 = \sigma_1 \xi_1 + \alpha \delta_1 \epsilon_1$ with $\sigma_1^2 = (1 + \alpha^2 s / (2d))^{-1}$ where now δ_1 is the Bernoulli random variable with $\mathbf{P}(\delta_1 = 1) = \frac{s}{2d}(1 + \alpha^2 s / (2d))^{-1}$. By construction, $\mathbf{E}\tilde{\xi}_1 = 0$ and $\mathbf{E}\tilde{\xi}_1^2 = 1$. Since the support of F_0 is in $[-3/2, 3/2]$ one can check as in item (ii) of the proof of Theorems 5 and 6 that $\tilde{P} \in \mathcal{G}_{a,\tau}$. Next, analogously to (67) - (68) we obtain that, for any $u > 0$,

$$\sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\|\theta\|_0 \leq s} \mathbf{P}_{\theta, P_\xi, 1}(|\hat{T} - \|\theta\|_2| \geq u) \geq \frac{1 - V(\mathbb{P}_{\bar{\mu}}, P_{0, \tilde{P}, 1}) - \bar{\mu}(\|\theta\|_2 < 2u)}{2}.$$

Let \mathbf{P}_0 and \mathbf{P}_1 denote the distributions of (ξ_1, \dots, ξ_d) and of $(\sigma_1 \xi_1, \dots, \sigma_1 \xi_d)$, respectively. Acting as in item (i) of the proof of Theorems 5 and 6 and using the bound

$$|1 - \sigma_1| \leq \alpha^2 s / d = \frac{\tau^2}{4} \frac{s}{d} \log^{2/a}(ed/s) \leq C c_0 / \sqrt{d}$$

we find that $V(\mathbf{P}_0, \mathbf{P}_1) \leq H(\mathbf{P}_0, \mathbf{P}_1) \leq 2\kappa c_0^2$ for some $\kappa > 0$. Therefore, $V(\mathbb{P}_\mu, P_{0, \tilde{P}, 1}) = V(\mathbf{P}_0 * \mathbf{Q}, \mathbf{P}_1 * \mathbf{Q}) \leq V(\mathbf{P}_0, \mathbf{P}_1) \leq 2\kappa c_0^2$ where \mathbf{Q} denotes the distribution of $(\alpha \delta_1 \epsilon_1, \dots, \alpha \delta_d \epsilon_d)$. This bound and the fact that $V(\mathbb{P}_{\bar{\mu}}, P_{0, \tilde{P}, 1}) \leq V(\mathbb{P}_{\bar{\mu}}, \mathbb{P}_\mu) + V(\mathbb{P}_\mu, P_{0, \tilde{P}, 1})$ imply

$$\sup_{P_\xi \in \mathcal{G}_{a,\tau}} \sup_{\|\theta\|_0 \leq s} \mathbf{P}_{\theta, P_\xi, 1}(|\hat{T} - \|\theta\|_2| \geq u) \geq \frac{1 - V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) - \bar{\mu}(\|\theta\|_2 < 2u)}{2} - \kappa c_0^2.$$

We conclude by repeating the argument after (68) in the proof of Theorem 2 and choosing $c_0 > 0$ small enough to guarantee that the right hand side of the last display is positive.

6.4. *Proof of part (ii) of Proposition 7.* The lower bound with the rate $1/\sqrt{d}$ follows from the argument as in item (i) of the proof of Theorems 5 and 6 if we replace there F_0 by the standard Gaussian distribution. The lower bound with the rate $\frac{s}{d(1+\log_+(s^2/d))}$ follows from Lemma 8 and the lower bound for estimation of $\|\boldsymbol{\theta}\|_2$ in Proposition 3.

6.5. *Proof of Proposition 8.* Assume that $\boldsymbol{\theta} = 0$, $\sigma = 1$ and set

$$\xi_i = \sqrt{3}\epsilon_i u_i,$$

where the ϵ_i 's and the u_i are independent, with Rademacher and uniform distribution on $[0, 1]$ respectively. Then note that

$$(69) \quad \mathbf{E}_{0, P_{\xi}, 1}(\hat{\sigma}_*^2 - 1)^2 \geq (\mathbf{E}_{0, P_{\xi}, 1}(\hat{\sigma}_*^2) - 1)^2 = \left(\mathbf{E}_{0, P_{\xi}, 1} \left\{ \hat{\sigma}_*^2 - \frac{3}{d} \sum_{i=1}^d u_i^2 \right\} \right)^2,$$

since $\mathbf{E}(u_i^2) = 1/3$. Note also that $\hat{\sigma}_*^2 = \frac{3}{d/2} \sum_{i=1}^{d/2} u_{(i)}^2$. Now,

$$\begin{aligned} \frac{1}{d/2} \sum_{i=1}^{d/2} u_{(i)}^2 - \frac{1}{d} \sum_{i=1}^d u_i^2 &= \frac{1}{d} \sum_{i=1}^{d/2} u_{(i)}^2 - \frac{1}{d} \sum_{i=d/2+1}^d u_{(i)}^2 \\ &\leq \frac{1}{d} \sum_{i=1}^{d/4} u_{(i)}^2 - \frac{1}{d} \sum_{i=3d/4+1}^d u_{(i)}^2 \\ &\leq \frac{1}{4} (u_{(d/4)}^2 - u_{(3d/4)}^2). \end{aligned}$$

Since $u_{(i)}$ follows a Beta distribution with parameters $(i, d-i+1)$ we have $\mathbf{E}(u_{(i)}^2) = \frac{i(i+1)}{(d+1)(d+2)}$, and

$$\mathbf{E}_{0, P_{\xi}, 1} \left(\frac{1}{d/2} \sum_{i=1}^{d/2} u_{(i)}^2 - \frac{1}{d} \sum_{i=1}^d u_i^2 \right) \leq \frac{1}{4} \mathbf{E}_{0, P_{\xi}, 1} (u_{(d/4)}^2 - u_{(3d/4)}^2) = -\frac{d}{8(d+2)} \leq -\frac{1}{24}.$$

This and (69) prove the proposition.

7. Lemmas.

7.1. Lemmas for the upper bounds.

LEMMA 1. Let $z_1, \dots, z_d \stackrel{iid}{\sim} P$ with $P \in \mathcal{G}_{a, \tau}$ for some $a, \tau > 0$ and let $z_{(1)} \leq \dots \leq z_{(d)}$ be the order statistics of $|z_1|, \dots, |z_d|$. Then for $u > 2^{1/a} \tau \vee 2$, we have

$$(70) \quad \mathbf{P} \left(z_{(d-j+1)} \leq u \log^{1/a} (ed/j), \forall j = 1, \dots, d \right) \geq 1 - 4e^{-u^a/2},$$

and, for any $r > 0$,

$$(71) \quad \mathbf{E}(z_{(d-j+1)}^r) \leq C \log^{r/a} (ed/j), \quad j = 1, \dots, d,$$

where $C > 0$ is a constant depending only on τ , a and r .

PROOF. Using the definition of $\mathcal{G}_{a,\tau}$ we get that, for any $t \geq 2$,

$$\mathbf{P}(z_{(d-j+1)} \geq t) \leq \binom{d}{j} \mathbf{P}^j(|z_1| \geq t) \leq 2 \left(\frac{ed}{j}\right)^j e^{-j(t/\tau)^a}, \quad j = 1, \dots, d.$$

Thus, for $v \geq 2^{1/a} \vee (2/\tau)$ we have

$$(72) \quad \mathbf{P}(z_{(d-j+1)} \geq v\tau \log^{1/a}(ed/j)) \leq 2 \left(\frac{ed}{j}\right)^{j(1-v^a)} \leq 2e^{-jv^a/2}, \quad j = 1, \dots, d,$$

and

$$\mathbf{P}\left(\exists j \in \{1, \dots, d\} : z_{(d-j+1)} \geq v\tau \log^{1/a}(ed/j)\right) \leq 2 \sum_{j=1}^d e^{-jv^a/2} \leq 4e^{-v^a/2}$$

implying (70). Finally, (71) follows by integrating (72). \square

LEMMA 2. Let $z_1, \dots, z_d \stackrel{iid}{\sim} P$ with $P \in \mathcal{P}_{a,\tau}$ for some $a, \tau > 0$ and let $z_{(1)} \leq \dots \leq z_{(d)}$ be the order statistics of $|z_1|, \dots, |z_d|$. Then for $u > (2e)^{1/a}\tau \vee 2$, we have

$$(73) \quad \mathbf{P}\left(z_{(d-j+1)} \leq u \left(\frac{d}{j}\right)^{1/a}, \forall j = 1, \dots, d\right) \geq 1 - \frac{2e\tau^a}{u^a}$$

and, for any $r \in (0, a)$,

$$(74) \quad \mathbf{E}(z_{(d-j+1)}^r) \leq C \left(\frac{d}{j}\right)^{r/a}, \quad j = 1, \dots, d,$$

where $C > 0$ is a constant depending only on τ, a and r .

PROOF. Using the definition of $\mathcal{P}_{a,\tau}$ we get that, for any $t \geq 2$,

$$\mathbf{P}(z_{(d-j+1)} \geq t) \leq \left(\frac{ed}{j}\right)^j \left(\frac{\tau}{t}\right)^{ja}.$$

Set $t_j = u \left(\frac{d}{j}\right)^{1/a}$ and $q = e(\tau/u)^a$. The assumption on u yields that $q < 1/2$, so that

$$\mathbf{P}\left(\exists j \in \{1, \dots, d\} : z_{(d-j+1)} \geq u \left(\frac{d}{j}\right)^{1/a}\right) \leq \sum_{j=1}^d \left(\frac{ed}{j}\right)^j \left(\frac{\tau}{t_j}\right)^{ja} = \sum_{j=1}^d q^j \leq 2q.$$

This proves (73). The proof of (74) is analogous to that of (71). \square

LEMMA 3. For all $a > 0$ and all integers $1 \leq s \leq d$,

$$\sum_{i=1}^s \log^{2/a}(ed/i) \leq Cs \log^{2/a}\left(\frac{ed}{s}\right)$$

where $C > 0$ depends only on a .

The proof is simple and we omit it.

7.2. *Lemmas for the lower bounds.* For two probability measures P_1 and P_2 on a measurable space (Ω, \mathcal{U}) , we denote by $V(P_1, P_2)$ the total variation distance between P_1 and P_2 :

$$V(P_1, P_2) = \sup_{B \in \mathcal{U}} |P_1(B) - P_2(B)|.$$

LEMMA 4 (Deviations of the binomial distribution). *Let $\mathcal{B}(d, p)$ denote the binomial random variable with parameters d and $p \in (0, 1)$. Then, for any $\lambda > 0$,*

$$(75) \quad \mathbf{P}(\mathcal{B}(d, p) \geq \lambda\sqrt{d} + dp) \leq \exp\left(-\frac{\lambda^2}{2p(1-p)\left(1 + \frac{\lambda}{3p\sqrt{d}}\right)}\right),$$

$$(76) \quad \mathbf{P}(\mathcal{B}(d, p) \leq -\lambda\sqrt{d} + dp) \leq \exp\left(-\frac{\lambda^2}{2p(1-p)}\right).$$

Inequality (75) is a combination of formulas (3) and (10) on pages 440–441 in [19]. Inequality (76) is formula (6) on page 440 in [19].

LEMMA 5. *Let \mathbb{P}_μ and $\mathbb{P}_{\bar{\mu}}$ be the probability measures defined in (59). The total variation distance between these two measures satisfies*

$$(77) \quad V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) \leq \mathbf{P}\left(\mathcal{B}\left(d, \frac{s}{2d}\right) > s\right) \leq e^{-\frac{3s}{16}},$$

and

$$(78) \quad V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) \leq 1 - \mathbf{P}\left(\mathcal{B}\left(d, \frac{s}{2d}\right) = 0\right) - \mathbf{P}\left(\mathcal{B}\left(d, \frac{s}{2d}\right) = 1\right).$$

PROOF. We have

$$V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) = \sup_B \left| \int \mathbf{P}_{\boldsymbol{\theta}, U, 1}(B) d\mu(\boldsymbol{\theta}) - \int \mathbf{P}_{\boldsymbol{\theta}, U, 1}(B) d\bar{\mu}(\boldsymbol{\theta}) \right| \leq \sup_{|f| \leq 1} \left| \int f d\mu - \int f d\bar{\mu} \right| = V(\mu, \bar{\mu}).$$

Furthermore, $V(\mu, \bar{\mu}) \leq \mu(\Theta_s^c)$ since for any Borel subset B of \mathbb{R}^d we have $|\mu(B) - \bar{\mu}(B)| \leq \mu(B \cap \Theta_s^c)$. Indeed,

$$\mu(B) - \bar{\mu}(B) \leq \mu(B) - \mu(B \cap \Theta) = \mu(B \cap \Theta^c)$$

and

$$\bar{\mu}(B) - \mu(B) = \frac{\mu(B \cap \Theta)}{\mu(\Theta)} - \mu(B \cap \Theta) - \mu(B \cap \Theta^c) \geq -\mu(B \cap \Theta^c).$$

Thus,

$$(79) \quad V(\mathbb{P}_\mu, \mathbb{P}_{\bar{\mu}}) \leq \mu(\Theta_s^c) = \mathbf{P}\left(\mathcal{B}\left(d, \frac{s}{2d}\right) > s\right).$$

Combining this inequality with (75) we obtain (77). To prove (78), we use again (79) and notice that $\mathbf{P}\left(\mathcal{B}\left(d, \frac{s}{2d}\right) > s\right) \leq \mathbf{P}\left(\mathcal{B}\left(d, \frac{s}{2d}\right) \geq 2\right)$ for any integer $s \geq 1$. \square

LEMMA 6. *Let $\bar{\mu}$ be defined in (58) with some $\alpha > 0$. Then*

$$(80) \quad \bar{\mu}\left(\|\boldsymbol{\theta}\|_2 < \frac{\alpha}{2}\sqrt{s}\right) \leq 2e^{-\frac{s}{16}},$$

and, for all $s \leq 32$,

$$(81) \quad \bar{\mu}\left(\|\boldsymbol{\theta}\|_2 < \frac{\alpha\sqrt{s}}{4\sqrt{2}}\right) = \mathbf{P}\left(\mathcal{B}(d, \frac{s}{2d}) = 0\right).$$

PROOF. First, note that

$$(82) \quad \mu\left(\|\boldsymbol{\theta}\|_2 < \frac{\alpha}{2}\sqrt{s}\right) = \mathbf{P}\left(\mathcal{B}(d, \frac{s}{2d}) < \frac{s}{4}\right) \leq e^{-\frac{s}{16}}$$

where the last inequality follows from (76). Next, inspection of the proof of Lemma 5 yields that $\bar{\mu}(B) \leq \mu(B) + e^{-\frac{3s}{16}}$ for any Borel set B . Taking here $B = \{\|\boldsymbol{\theta}\|_2 \leq \alpha\sqrt{s}/2\}$ and using (82) proves (80). To prove (81), it suffices to note that $\mu\left(\|\boldsymbol{\theta}\|_2 < \frac{\alpha\sqrt{s}}{4\sqrt{2}}\right) = \mathbf{P}\left(\mathcal{B}(d, \frac{s}{2d}) < \frac{s}{32}\right)$. \square

LEMMA 7. *There exists a probability density $f_0 : \mathbb{R} \rightarrow [0, \infty)$ with the following properties: f_0 is continuously differentiable, symmetric about 0, supported on $[-3/2, 3/2]$, with variance 1 and finite Fisher information $I_{f_0} = \int (f'_0(x))^2 (f_0(x))^{-1} dx$.*

PROOF. Let $K : \mathbb{R} \rightarrow [0, \infty)$ be any probability density, which is continuously differentiable, symmetric about 0, supported on $[-1, 1]$, and has finite Fisher information I_K , for example, the density $K(x) = \cos^2(\pi x/2) \mathbf{1}_{|x| \leq 1}$. Define $f_0(x) = [K_h(x + (1 - \varepsilon)) + K_h(x - (1 - \varepsilon))]/2$ where $h > 0$ and $\varepsilon \in (0, 1)$ are constants to be chosen, and $K_h(u) = K(u/h)/h$. Clearly, we have $I_{f_0} < \infty$ since $I_K < \infty$. It is straightforward to check that the variance of f_0 satisfies $\int x^2 f_0(x) dx = (1 - \varepsilon)^2 + h^2 \sigma_K^2$ where $\sigma_K^2 = \int u^2 K(u) du$. Choosing $h = \sqrt{2\varepsilon - \varepsilon^2}/\sigma_K$ and $\varepsilon \leq \sigma_K^2/8$ guarantees that $\int x^2 f_0(x) dx = 1$ and the support of f_0 is contained in $[-3/2, 3/2]$. \square

LEMMA 8. *Let $\tau > 0$, $a > 4$ and let s, d be integers satisfying $1 \leq s \leq d$. Let \mathcal{P} be any subset of $\mathcal{P}_{a, \tau}$. Assume that for some function $\phi(s, d)$ of s and d and for some positive constants c_1, c_2, c'_1, c'_2 we have*

$$(83) \quad \inf_{\hat{T}} \sup_{P_\xi \in \mathcal{P}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{P}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(\left| \frac{\hat{T}}{\sigma^2} - 1 \right| \geq \frac{c_1}{\sqrt{d}} \right) \geq c'_1,$$

and

$$(84) \quad \inf_{\hat{T}} \sup_{P_\xi \in \mathcal{P}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{P}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(\left| \frac{\hat{T} - \|\boldsymbol{\theta}\|_2}{\sigma} \right| \geq c_2 \phi(s, d) \right) \geq c'_2.$$

Then

$$\inf_{\hat{T}} \sup_{P_\xi \in \mathcal{P}} \sup_{\sigma > 0} \sup_{\|\boldsymbol{\theta}\|_0 \leq s} \mathbf{P}_{\boldsymbol{\theta}, P_\xi, \sigma} \left(\left| \frac{\hat{T}}{\sigma^2} - 1 \right| \geq c_3 \max \left(\frac{1}{\sqrt{d}}, \frac{\phi^2(s, d)}{d} \right) \right) \geq c'_3$$

for some constants $c_3, c'_3 > 0$.

PROOF. Let $\hat{\sigma}^2$ be an arbitrary estimator of σ^2 . Based on $\hat{\sigma}^2$, we can construct an estimator $\hat{T} = \hat{N}^*$ of $\|\theta\|_2$ defined by formula (11), case $s > \sqrt{d}$. It follows from (30), (31) and (84) that

$$\begin{aligned} c'_2 \leq & \mathbf{P} \left(2|(\theta, \xi)| \geq c_2 \|\theta\|_2 \phi(s, d)/3 \right) + \mathbf{P} \left(\sqrt{\|\xi\|_2^2 - d} \geq c_2 \phi(s, d)/3 \right) \\ & + \mathbf{P} \left(\sqrt{d \left| \frac{\hat{\sigma}^2}{\sigma^2} - 1 \right|} \geq c_2 \phi(s, d)/3 \right), \end{aligned}$$

where we write for brevity $\mathbf{P} = \mathbf{P}_{\theta, P_\xi, \sigma}$. Hence

$$\mathbf{P} \left(\left| \frac{\hat{\sigma}^2}{\sigma^2} - 1 \right| \geq c_2^2 \phi^2(s, d)/(9d) \right) \geq c'_2 - c^* \max \left(\frac{d}{\phi^4(s, d)}, \frac{1}{\phi^2(s, d)} \right)$$

for some constant $c^* > 0$ depending only on a and τ . If $\phi^2(s, d) > \max \left(\sqrt{\frac{2c^*d}{c_2^2}}, \frac{2c^*}{c_2^2} \right)$, then

$$\mathbf{P} \left(\left| \frac{\hat{\sigma}^2}{\sigma^2} - 1 \right| \geq C \max \left(\frac{1}{\sqrt{d}}, \frac{\phi^2(s, d)}{d} \right) \right) \geq c'_2/2.$$

If $\phi^2(s, d) \leq \max \left(\sqrt{\frac{2c^*d}{c_2^2}}, \frac{2c^*}{c_2^2} \right)$, then $\max \left(\frac{1}{\sqrt{d}}, \frac{\phi^2(s, d)}{d} \right)$ is of order $\frac{1}{\sqrt{d}}$ and the result follows from (83). \square

8. Acknowledgements. The work of O. Collier has been conducted as part of the project Labex MME-DII (ANR11-LBX-0023-01). The work of M.Ndaoud and 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).

References.

- [1] P. BELLEC, G. LECUÉ AND A.B. TSYBAKOV. Slope meets Lasso: improved oracle bounds and optimality. *Annals of Statistics*, **46**, 3603–3642, 2018.
- [2] A. BELLONI, V. CHERNOZHUKOV AND L. WANG. Pivotal estimation via Square-Root Lasso in nonparametric regression. *Annals of Statistics* **42** 757–788, 2014.
- [3] D. BELOMESTNY, M. TRABS AND A.B. TSYBAKOV. Sparse covariance matrix estimation in high-dimensional deconvolution. arXiv preprint, arXiv:1710.10870. To appear in *Bernoulli*.
- [4] C. BUTUCEA AND C. MATIAS. Minimax estimation of the noise level and of the deconvolution density in a semiparametric convolution model. *Bernoulli*, **11**, 309–340, 2005.
- [5] T.T. CAI AND J. JIN. Optimal rates of convergence for estimating the null and proportion of non-null effects in large-scale multiple testing. *Annals of Statistics*, **38**, 100–145, 2010.
- [6] A. CARPENTIER AND N. VERZELEN. Adaptive estimation of the sparsity in the Gaussian vector model. *Annals of Statistics*, **47**, 93–126, 2019.
- [7] M. CHEN, C. GAO AND Z. REN. Robust covariance and scatter matrix estimation under Huber’s contamination model. *Annals of Statistics*, **46**, 1932–1960, 2018.
- [8] O. COLLIER, L. COMMINGES AND A.B. TSYBAKOV. Minimax estimation of linear and quadratic functionals under sparsity constraints. *Annals of Statistics*, **45**, 923–958, 2017.
- [9] O. COLLIER, L. COMMINGES, A.B. TSYBAKOV AND N. VERZELEN. Optimal adaptive estimation of linear functionals under sparsity. *Annals of Statistics*, **46**, 3130–3150, 2018.
- [10] D.L. DONOHO, I.M. JOHNSTONE, J.C. HOCH AND A.S. STERN. Maximum entropy and the nearly black object. *Journal of the Royal Statistical Society. Series B*, **54**, 41–81, 1992.
- [11] E. GAUTIER AND A.B. TSYBAKOV. Pivotal estimation in high-dimensional regression via linear programming. In: Empirical Inference. Festschrift in Honor of Vladimir N. Vapnik, B.Schölkopf, Z. Luo, V. Vovk eds., 195 - 204. Springer, New York e.a., 2013.

- [12] Y. GOLUBEV AND E. KRYMOVA. On estimation of the noise variance in high-dimensional linear models. arXiv preprint, arXiv:1711.09208.
- [13] Z. GUO, W. WANG, T.T. CAI AND H. LI. Optimal estimation of genetic relatedness in high-dimensional linear models. *J. of the American Statist. Assoc.*, published online, 2018.
- [14] P.J. HUBER. *Robust Statistics*, J. Wiley, 1981.
- [15] I.A.IBRAGIMOV AND R.Z.HASMINSKII. *Statistical Estimation. Asymptotic Theory*, Springer, New York, 1981.
- [16] L.JANSON, R. FOYGEL BARBER AND E.CANDES. EigenPrism: inference for high dimensional signal-to-noise ratios. *Journal of the Royal Statistical Society. Series B* , **79**, 1037–1065, 2017.
- [17] S. MINSKER AND X. WEI. Estimation of the covariance structure of heavy-tailed distributions. arXiv preprint, arXiv:1708.00502.
- [18] V.V. PETROV. *Limit Theorems of Probability Theory*. Clarendon Press, Oxford, 1995.
- [19] G. SHORACK AND J. WELLNER. *Empirical Processes with Applications to Statistics*, John Wiley, New York, 1986.
- [20] T. SUN AND C.-H. ZHANG. Scaled sparse linear regression. *Biometrika*, **99**, 879–898, 2012.
- [21] A.B. TSYBAKOV. *Introduction to Nonparametric Estimation*. Springer, New York, 2009.
- [22] N. VERZELEN. Minimax risks for sparse regressions: Ultra-high dimensional phenomenons. *Electronic J. of Statist.* **6** 38–90, 2012.
- [23] N. VERZELEN AND E. GASSIAT. Adaptive estimation of high-dimensional signal-to-noise ratios. *Bernoulli* **24**, 3683–3710, 2018.
- [24] L. WASSERMAN. *All of Statistics*. Springer, New York, 2005.

CEREMADE, UNIVERSITÉ PARIS-DAUPHINE
 PSL RESEARCH UNIVERSITY
 75016 PARIS, FRANCE
 E-MAIL: laetitia.comminges@dauphine.fr

MODAL'X, UPL, UNIVERSITÉ PARIS NANTERRE
 92000 NANTERRE FRANCE
 E-MAIL: olivier.collier@parisnanterre.fr

CREST (UMR CNRS 9194), ENSAE
 5, AV. HENRY LE CHATELIER, 91764 PALAISEAU, FRANCE
 E-MAIL: ndaoudm@gmail.com; alexandre.tsybakov@ensae.fr