

High-Dimensional Regression with Unknown Variance

Christophe Giraud, Sylvie Huet and Nicolas Verzelen

Abstract. We review recent results for high-dimensional sparse linear regression in the practical case of unknown variance. Different sparsity settings are covered, including coordinate-sparsity, group-sparsity and variation-sparsity. The emphasis is put on nonasymptotic analyses and feasible procedures. In addition, a small numerical study compares the practical performance of three schemes for tuning the lasso estimator and some references are collected for some more general models, including multivariate regression and nonparametric regression.

Key words and phrases: Linear regression, high-dimension, unknown variance.

1. INTRODUCTION

In the present paper, we mainly focus on the linear regression model

$$(1) \quad Y = \mathbf{X}\beta_0 + \varepsilon,$$

where Y is a n -dimensional response vector, \mathbf{X} is a fixed $n \times p$ design matrix, and the vector ε is made of n i.i.d. Gaussian random variables with $\mathcal{N}(0, \sigma^2)$ distribution. In the sequel, $\mathbf{X}^{(i)}$ stands for the i th row of \mathbf{X} . Our interest is on the high-dimensional setting, where the dimension p of the unknown parameter β_0 is large, possibly larger than n .

The analysis of the high-dimensional linear regression model has attracted a lot of attention in the last decade. Nevertheless, there is a longstanding gap between the theory where the variance σ^2 is generally assumed to be known and the practice where it is often unknown. The present paper is mainly devoted to reviewing recent results on linear regression in high-dimensional settings with *unknown* variance σ^2 . A few

additional results for multivariate regression and the nonparametric regression model

$$(2) \quad Y_i = f(\mathbf{X}^{(i)}) + \varepsilon_i, \quad i = 1, \dots, n,$$

will also be mentioned.

1.1 Sparsity Assumptions

In a high-dimensional linear regression model, accurate estimation is unfeasible unless it relies on some special properties of the parameter β_0 . The most common assumption on β_0 is that it is sparse in some sense. We will consider in this paper the three following classical sparsity assumptions.

Coordinate-sparsity. Most of the coordinates of β_0 are assumed to be zero (or approximately zero). This is the most common acceptance for sparsity in linear regression.

Structured-sparsity. The pattern of zero(s) of the coordinates of β_0 is assumed to have an a priori known structure. For instance, in group-sparsity [80], the covariates are clustered into M groups, and when the coefficient $\beta_{0,i}$ corresponding to the covariate \mathbf{X}_i (the i th column of \mathbf{X}) is nonzero, then it is likely that all the coefficients $\beta_{0,j}$ with variables \mathbf{X}_j in the same cluster as \mathbf{X}_i are nonzero.

Variation-sparsity. The $p - 1$ -dimensional vector β_0^V of variation of β_0 is defined by $\beta_{0,j}^V = \beta_{0,j+1} - \beta_{0,j}$. Sparsity in variation means that most of the components of β_0^V are equal to zero (or approximately zero). When $p = n$ and $\mathbf{X} = I_n$, variation-sparse linear regression corresponds to signal segmentation.

Christophe Giraud is Professor, CMAP, UMR CNRS 7641, Ecole Polytechnique, Route de Saclay, 91128 Palaiseau Cedex, France (e-mail: christophe.giraud@polytechnique.edu). Sylvie Huet is Research Scientist, UR341 MIA, INRA, F-78350 Jouy-en-Josas, France (e-mail: sylvie.huet@jouy.inra.fr). Nicolas Verzelen is Research Scientist, UMR729 MISTEA, INRA, F-34060 Montpellier, France (e-mail: nicolas.verzelen@supagro.inra.fr).

1.2 Statistical Objectives

In the linear regression model, there are roughly two kinds of estimation objectives. In the *prediction problem*, the goal is to estimate $\mathbf{X}\beta_0$, whereas in the *inverse problem* it is to estimate β_0 . When the vector β_0 is sparse, a related objective is to estimate the *support* of β_0 (model identification problem) which is the set of the indices j corresponding to the nonzero coefficients $\beta_{0,j}$. Inverse problems and prediction problems are not equivalent in general. When the Gram matrix $\mathbf{X}\mathbf{X}^*$ is poorly conditioned, the former problems can be much more difficult than the latter. Since there are only a few results on inverse problems with unknown variance, we will focus on the prediction problem, the support estimation problem being shortly discussed in the course of the paper.

In the sequel, $\mathbb{E}_{\beta_0}[\cdot]$ stands for the expectation with respect to $Y \sim \mathcal{N}(\mathbf{X}\beta_0, \sigma^2 I_n)$, and $\|\cdot\|_2$ is the Euclidean norm. The prediction objective amounts to build estimators $\hat{\beta}$ so that the risk

$$(3) \quad \mathcal{R}[\hat{\beta}; \beta_0] := \mathbb{E}_{\beta_0}[\|\mathbf{X}(\hat{\beta} - \beta_0)\|_2^2]$$

is as small as possible.

1.3 Approaches

Most procedures that handle high-dimensional linear models [22, 26, 62, 72, 73, 81, 83, 85] rely on tuning parameters whose optimal value depends on σ . For example, the results of Bickel et al. [17] suggest to choose the tuning parameter λ of the lasso of the order of $2\sigma\sqrt{2\log(p)}$. As a consequence, all these procedures cannot be directly applied when σ^2 is unknown.

A straightforward approach is to replace σ^2 by an estimate of the variance in the optimal value of the tuning parameter(s). Nevertheless, the variance σ^2 is difficult to estimate in high-dimensional settings (see Proposition 2.3 below), so a plug-in of the variance does not necessarily yield good results. There are basically two approaches to build on this amount of work on high-dimensional estimation with known variance.

(1) *Ad-hoc estimation*. There has been some recent work [16, 68, 71] to modify procedures like the lasso in such a way that the tuning parameter does not depend anymore on σ^2 ; see Section 4.2. The challenge is to find a smart modification of the procedure, so that the resulting estimator $\hat{\beta}$ is computationally feasible and has a risk $\mathcal{R}[\hat{\beta}; \beta_0]$ as small as possible.

(2) *Estimator selection*. Given a collection $(\hat{\beta}_\lambda)_{\lambda \in \Lambda}$ of estimators, the objective of estimator selection is to pick an index $\hat{\lambda}$ such that the risk of $\hat{\beta}_{\hat{\lambda}}$ is as small as

possible, ideally as small as the risk $\mathcal{R}[\hat{\beta}_{\lambda^*}; \beta_0]$ of the so-called *oracle* estimator

$$(4) \quad \hat{\beta}_{\lambda^*} := \operatorname{argmin}_{\{\hat{\beta}_\lambda, \lambda \in \Lambda\}} \mathcal{R}[\hat{\beta}_\lambda; \beta_0].$$

Efficient estimator selection procedures can then be applied to tune the aforementioned estimation methods [22, 26, 62, 72, 73, 81, 83, 85]. Among the most famous methods for estimator selection, we mention V -fold cross-validation (Geisser [32]), AIC (Akaike [1]) and BIC (Schwarz [64]) criteria.

The objective of this survey is to describe state-of-the-art procedures for high-dimensional linear regression with unknown variance. We will review both automatic tuning methods and ad-hoc methods. There are some procedures that we will let aside. For example, Baraud [11] provides a versatile estimator selection scheme, but the procedure is computationally intractable in large dimensions. Linear or convex aggregation of estimators are also valuable alternatives to estimator selection when the goal is to perform *estimation*, but only a few theoretical works have addressed the aggregation problem when the variance is unknown [33, 34]. For these reasons, we will not review these approaches in the sequel.

1.4 Why Care about Nonasymptotic Analyses?

AIC [1], BIC [64] and V -fold cross-validation [32] are probably the most popular criteria for estimator selection. The use of these criteria relies on some classical asymptotic optimality results. These results focus on the setting where the collection of estimators $(\hat{\beta}_\lambda)_{\lambda \in \Lambda}$ and the dimension p are fixed and consider the limit behavior of the criteria when the sample size n goes to infinity. For example, under some suitable conditions, Shibata [67], Li [53] and Shao [66] prove that the risk of the estimator selected by AIC or V -fold CV (with $V = V_n \rightarrow \infty$) is asymptotically equivalent to the oracle risk $\mathcal{R}[\hat{\beta}_{\lambda^*}; \beta_0]$. Similarly, Nishii [59] shows that the BIC criterion is consistent for model selection.

All these asymptotic results can lead to misleading conclusions in modern statistical settings where the sample size remains small and the parameter's dimension becomes large. For instance it is proved in [12], Section 3.3.2, and illustrated in [12], Section 6.2, that BIC (and thus AIC) can strongly overfit and should not be used for p larger than n . Additional examples are provided in Appendix A. A nonasymptotic analysis takes into account all the characteristics of the selection problem (sample size n , parameter dimension p , number of models per dimension, design \mathbf{X} , etc.). It treats n

and p as they are, and it highlights important features hidden by the asymptotic theory. For these reasons, we will restrict this review to nonasymptotic results.

1.5 Organization of the Paper

In Section 2, we investigate how the ignorance of the variance affects the minimax risk bounds. In Section 3, some “generic” schemes for selecting estimators are presented. The coordinate-sparse setting is addressed in Section 4 where some theoretical results are collected, and a small numerical experiment compares different lasso-based procedures. The group-sparse and variation-sparse settings are reviewed in Sections 5 and 6, and Section 7 is devoted to some more general models such as multivariate regression or nonparametric regression.

In the sequel, C, C_1, \dots refer to numerical constants whose value may vary from line to line, while $\|\beta\|_0$ stands for the number of nonzero components of β and $|\mathcal{J}|$ for the cardinality of a set \mathcal{J} .

2. THEORETICAL LIMITS

The goal of this section is to address the intrinsic difficulty of a coordinate-sparse linear regression problem. We will answer the following questions: Which range of p can we reasonably consider? When the variance is unknown, can we hope to do as well as when the variance is known?

2.1 Minimax Adaptation

A classical way to assess the performance of an estimator $\hat{\beta}$ is to measure its maximal risk over a class $\mathbf{B} \subset \mathbb{R}^p$. This is the minimax point of view. As we are interested in coordinate-sparsity for β_0 , we will consider the sets $\mathbf{B}[k, p]$ of vectors that contain at most k nonzero coordinates for some $k > 0$.

Given an estimator $\hat{\beta}$, the *maximal prediction risk* of $\hat{\beta}$ over $\mathbf{B}[k, p]$ for a fixed design \mathbf{X} and a variance σ^2 is defined by $\sup_{\beta_0 \in \mathbf{B}[k, p]} \mathcal{R}[\hat{\beta}; \beta_0]$ where the risk function $\mathcal{R}[\cdot, \beta_0]$ is defined by (3). Taking the infimum of the maximal risk over all possible estimators $\hat{\beta}$, we obtain the *minimax risk*

$$(5) \quad \mathbf{R}[k, \mathbf{X}] = \inf_{\hat{\beta}} \sup_{\beta_0 \in \mathbf{B}[k, p]} \mathcal{R}[\hat{\beta}; \beta_0].$$

Minimax bounds are convenient results to assess the range of problems that are statistically feasible and the optimality of particular procedures. Below, we say that an estimator $\hat{\beta}$ is “minimax” over $\mathbf{B}[k, p]$ if its maximal prediction risk equals [up to a possible multiplicative constant $C(\mathbf{X})$] the minimax risk.

In practice, the number of nonzero coordinates of β_0 is unknown. The fact that an estimator $\hat{\beta}$ is minimax over $\mathbf{B}[k, p]$ for some specific $k > 0$ does not imply that $\hat{\beta}$ performs well when β_0 has a number of nonzero components different from k . Indeed, $\hat{\beta}$ can be strongly biased when β_0 has more than k nonzero components or the variance of $\hat{\beta}$ can be too large compared to $\mathbf{R}[k_0, \mathbf{X}]$ when β_0 has k_0 nonzero components with k_0 less than k . A good estimation procedure $\hat{\beta}$ should not require the knowledge of the sparsity k of β_0 and should perform as well as if this sparsity were known. An estimator $\hat{\beta}$ that achieves [up to a possible multiplicative constant $C(\mathbf{X})$] the minimax risk over $\mathbf{B}[k, p]$ for a range of k is said to be *adaptive* to the sparsity. Similarly, an estimator $\hat{\beta}$ is adaptive to the variance σ^2 if it does not require the knowledge of σ^2 and nearly achieves the minimax risk for all $\sigma^2 > 0$. When possible, the main challenge is to build adaptive procedures.

In the following subsections, we review sharp bounds on the minimax prediction risks for both known and unknown sparsity, known and unknown variance. The big picture is summed up in Figure 1. Roughly, it says that adaptation is possible as long as $2k \log(p/k) < n$. In contrast, the situation becomes more complex for the ultra-high-dimensional¹ setting where $2k \log(p/k) \geq n$. The rest of this section is devoted to explain this big picture.

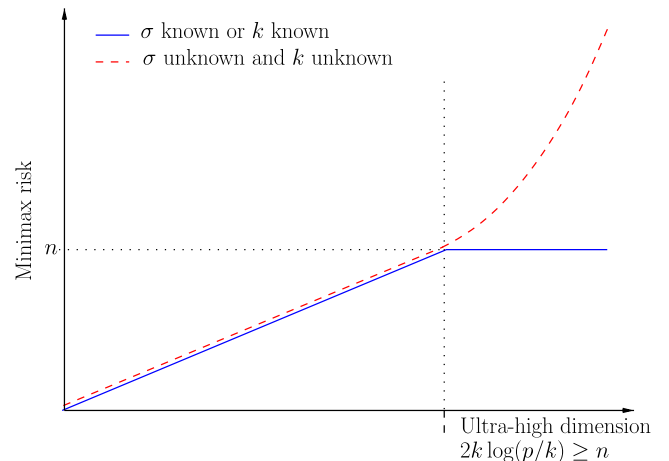


FIG. 1. *Minimal prediction risk over $\mathbf{B}[k, p]$ as a function of k .*

¹In some papers, the expression ultra-high-dimensional has been used to characterize problems such that $\log(p) = O(n^\theta)$ with $\theta < 1$. We argue here that as soon as $k \log(p)/n$ goes to 0, the case $\log(p) = O(n^\theta)$ is not intrinsically more difficult than conditions such as $p = O(n^\delta)$ with $\delta > 0$.

2.2 Minimax Risks under Known Sparsity and Known Variance

The minimax risk $\mathbf{R}[k, \mathbf{X}]$ depends on the form of the design \mathbf{X} . In order to grasp this dependency, we define for any $k > 0$, the largest and the smallest sparse eigenvalues of order k of $\mathbf{X}^* \mathbf{X}$ by

$$\Phi_{k,+}(\mathbf{X}) := \sup_{\beta \in \mathbf{B}[k,p] \setminus \{0_p\}} \frac{\|\mathbf{X}\beta\|_n^2}{\|\beta\|_p^2} \quad \text{and}$$

$$\Phi_{k,-}(\mathbf{X}) := \inf_{\beta \in \mathbf{B}[k,p] \setminus \{0_p\}} \frac{\|\mathbf{X}\beta\|_n^2}{\|\beta\|_p^2}.$$

PROPOSITION 2.1. *Assume that k and σ are known. There exist positive numerical constants C_1, C'_1, C_2 and C'_2 such that the following holds. For any (k, n, p) such that $k \leq n/2$ and any design \mathbf{X} , we have*

$$(6) \quad C_1 \frac{\Phi_{2k,-}(\mathbf{X})}{\Phi_{2k,+}(\mathbf{X})} k \log\left(\frac{p}{k}\right) \sigma^2 \leq \mathbf{R}[k, \mathbf{X}] \leq C'_1 \left[k \log\left(\frac{p}{k}\right) \wedge n \right] \sigma^2.$$

For any (k, n, p) such that $k \leq n/2$, we have

$$(7) \quad C_2 \left[k \log\left(\frac{p}{k}\right) \wedge n \right] \sigma^2 \leq \sup_{\mathbf{X}} \mathbf{R}[k, \mathbf{X}] \leq C'_2 \left[k \log\left(\frac{p}{k}\right) \wedge n \right] \sigma^2.$$

The minimax lower bound (6) has been first proved in [61, 62, 79] while (7) is stated in [77]. Let us first comment on bound (7). If the vector β_0 has k nonzero components, and if these components are a priori known, then one may build estimators that achieve a risk bound of the order k . In a (nonultra) high-dimensional setting [$2k \log(p/k) \leq n$], the minimax risk is of the order $k \log(p/k) \sigma^2$. The logarithmic term is the price to pay to cope with the fact that we do not know the position of the nonzero components in β_0 . The situation is quite different in an ultra-high-dimensional setting [$2k \log(p/k) > n$]. Indeed, the minimax risk remains of the order of $n \sigma^2$, which corresponds to the minimax risk of estimation of the vector $\mathbf{X}\beta_0$ without any sparsity assumption; see the blue curve in Figure 1. In other terms, the sparsity index k does not play a role anymore.

Dependency of $\mathbf{R}[k, \mathbf{X}]$ on the design \mathbf{X} . It follows from (6) that $\sup_{\mathbf{X}} \mathbf{R}[k, \mathbf{X}]$ is nearly achieved by designs \mathbf{X} satisfying $\Phi_{2k,-}(\mathbf{X})/\Phi_{2k,+}(\mathbf{X}) \approx 1$, when the setting is not ultra-high-dimensional. For some designs such that $\Phi_{2k,-}(\mathbf{X})/\Phi_{2k,+}(\mathbf{X})$ is small, the minimax

prediction risk $\mathbf{R}[k, \mathbf{X}]$ is possibly faster; see [77] for a discussion. In an ultra-high-dimensional setting, the form of the minimax risk ($n \sigma^2$) is related to the fact that no designs can satisfy $\Phi_{2k,-}(\mathbf{X})/\Phi_{2k,+}(\mathbf{X}) \approx 1$; see, for example, [10]. More precisely, the lower bound (6) enforces the following geometrical constrain:

$$\frac{\Phi_{2k,-}(\mathbf{X})}{\Phi_{2k,+}(\mathbf{X})} \leq C \frac{n}{k \log(p/k)}$$

for any design \mathbf{X} . The lower bound $\mathbf{R}[k, \mathbf{X}] \geq C[k \cdot \log(p/k) \wedge n] \sigma^2$ in (7) is, for instance, achieved by realizations of a standard Gaussian design, that is, designs \mathbf{X} whose components follow independent standard normal distributions. See [77] for more details.

2.3 Adaptation to the Sparsity and to the Variance

Adaptation to the sparsity when the variance is known. When σ^2 is known, there exist both model selection and aggregation procedures that achieve this $[k \log(p/k) \wedge n] \sigma^2$ risk simultaneously for all k and for all designs \mathbf{X} . Such procedures derive from the work of Birgé and Massart [18] and Leung and Barron [52]. However, these methods are intractable for large p except for specific forms of the design. We refer to the supplementary material [38] for more details.

Simultaneous adaptation to the sparsity and the variance. We first restrict to the nonultra high-dimensional setting, where the number of nonzero components k is unknown but satisfies $2k \log(p/k) < n$. In this setting, some procedures based on penalized log-likelihood [12] are simultaneous adaptive to the unknown sparsity and to the unknown variance and this for all designs \mathbf{X} . Again such procedures are intractable for large p . See the supplementary material [38] for more details. If we want to cover all k (including ultra-high-dimensional settings), the situation is different as shown in the next proposition (from [77]).

PROPOSITION 2.2 (Simultaneous adaptation is impossible). *There exist positive constants $C, C', C_1, C_2, C_3, C'_1, C'_2$ and C'_3 , such that the following holds. Consider any $p \geq n \geq C$ and $k \leq p^{1/3} \wedge n/2$ such that $k \log(p/k) \geq C'n$. There exist designs \mathbf{X} of size $n \times p$ such that for any estimator $\hat{\beta}$, we have either*

$$\sup_{\sigma^2 > 0} \frac{\mathcal{R}[\hat{\beta}; 0_p]}{\sigma^2} > C_1 n \quad \text{or}$$

$$\sup_{\beta_0 \in \mathbf{B}[k,p], \sigma^2 > 0} \frac{\mathcal{R}[\hat{\beta}; \beta_0]}{\sigma^2} > C_2 k \log\left(\frac{p}{k}\right) \exp\left[C_3 \frac{k}{n} \log\left(\frac{p}{k}\right)\right].$$

Conversely, there exist two estimators $\widehat{\beta}^{(n)}$ and $\widehat{\beta}^{BGH}$ (defined in the supplementary material [38]) that respectively satisfy

$$\begin{aligned} \sup_{\mathbf{X}} \sup_{\beta_0 \in \mathbb{R}^p, \sigma^2 > 0} \frac{\mathcal{R}[\widehat{\beta}^{(n)}; \beta_0]}{\sigma^2} &\leq C'_1 n, \\ \sup_{\mathbf{X}} \sup_{\beta_0 \in \mathbf{B}[k, p], \sigma^2 > 0} \frac{\mathcal{R}[\widehat{\beta}^{BGH}; \beta_0]}{\sigma^2} \\ &\leq C'_2 k \log\left(\frac{p}{k}\right) \exp\left[C'_3 \frac{k}{n} \log\left(\frac{p}{k}\right)\right] \end{aligned}$$

for all $1 \leq k \leq [(n-1) \wedge p]/4$.

As a consequence, simultaneous adaptation to the sparsity and to the variance is impossible in an ultra-high-dimensional setting. Indeed, any estimator $\widehat{\beta}$ that does not rely on σ^2 has to pay at least one of these two prices:

- (1) The estimator $\widehat{\beta}$ does not use the sparsity of the true parameter β_0 , and its risk for estimating $\mathbf{X}\beta_0$ is of the same order as the minimax risk over \mathbb{R}^n .
- (2) For any $1 \leq k \leq p^{1/3}$, the risk of $\widehat{\beta}$ fulfills

$$\begin{aligned} \sup_{\sigma > 0} \sup_{\beta_0 \in \mathbf{B}[k, p]} \frac{\mathcal{R}[\widehat{\beta}; \beta_0]}{\sigma^2} \\ \geq C_1 k \log(p) \exp\left[C_2 \frac{k}{n} \log(p)\right]. \end{aligned}$$

It follows that the maximal risk of $\widehat{\beta}$ increases exponentially fast in an ultra-high-dimensional setting (red curve in Figure 1), while the minimax risk is stuck to n (blue curve in Figure 1). The designs that satisfy the minimax lower bounds of Proposition 2.2 include realizations of a standard Gaussian design.

In an ultra-high-dimensional setting, the prediction problem becomes extremely difficult under unknown variance because the variance estimation itself is inconsistent as shown in the next proposition (from [77]).

PROPOSITION 2.3. *There exist positive constants C , C_1 and C_2 such that the following holds. Assume that $p \geq n \geq C$. For any $1 \leq k \leq p^{1/3}$, there exist designs \mathbf{X} such that*

$$\begin{aligned} \inf_{\widehat{\sigma}} \sup_{\sigma > 0, \beta_0 \in \mathbf{B}[k, p]} \mathbb{E}_{\beta_0} \left[\left| \frac{\sigma^2}{\widehat{\sigma}^2} - \frac{\widehat{\sigma}^2}{\sigma^2} \right| \right] \\ \geq C_1 \frac{k}{n} \log\left(\frac{p}{k}\right) \exp\left[C_2 \frac{k}{n} \log\left(\frac{p}{k}\right)\right]. \end{aligned}$$

2.4 What Should We Expect from a Good Estimation Procedure?

Let us consider an estimator $\widehat{\beta}$ that does not depend on σ^2 . Relying on the previous minimax bounds, we will say that $\widehat{\beta}$ achieves an *optimal* risk bound (with respect to the sparsity) if

$$(8) \quad \mathcal{R}[\widehat{\beta}; \beta_0] \leq C_1 \|\beta_0\|_0 \log(p) \sigma^2,$$

for any $\sigma > 0$ and any vector $\beta_0 \in \mathbb{R}^p$ such that $1 \leq \|\beta_0\|_0 \log(p) \leq C_2 n$. Such risk bounds prove that the estimator is approximately [up to a possible $\log(\|\beta_0\|_0)$ additional term] minimax adaptive to the unknown variance and the unknown sparsity. The condition $\|\beta_0\|_0 \log(p) \leq C_2 n$ ensures that the setting is not ultra-high-dimensional. As stated above, some procedures achieve (8) for all designs \mathbf{X} , but they are intractable for large p ; see [38]. One purpose of this review is to present fast procedures that achieve these kind of bounds under possible restrictive assumptions on the design matrix \mathbf{X} .

For some procedures, (8) can be improved into a bound of the form

$$(9) \quad \begin{aligned} \mathcal{R}[\widehat{\beta}; \beta_0] \\ \leq C_1 \inf_{\beta \neq 0} \{ \|\mathbf{X}(\beta - \beta_0)\|_2^2 + \|\beta\|_0 \log(p) \sigma^2 \}, \end{aligned}$$

with C_1 close to one. Again, the dimension $\|\beta_0\|_0$ is restricted to be smaller than $Cn/\log(p)$ to ensure that the setting is not ultra-high-dimensional. This kind of bound makes a clear trade-off between a bias and a variance term. For instance, when β_0 contains many components that are nearly equal to zero, the bound (9) can be much smaller than (8).

2.5 Other Statistical Problems in an Ultra-High-Dimensional Setting

We have seen that adaptation becomes impossible for the prediction problem in an ultra-high-dimensional setting. For other statistical problems, including the prediction problem with random design, the inverse problem (estimation of β_0), the variable selection problem (estimation of the support of β_0), the dimension reduction problem [47, 77, 78], the minimax risks increase exponentially fast in an ultra-high-dimensional setting. This kind of phase transition has been observed in a wide range of random geometry problems [28], suggesting some universality in this limitation. In practice, the sparsity index k is not known, but given (n, p) , we can compute $k^* := \max\{k : 2k \log(p/k) \geq n\}$. One may interpret that the problem is still reasonably difficult as long as $k \leq k^*$. This gives a simple rule of

thumb to know what we can hope from a given regression problem. For example, setting $p = 5000$ and $n = 50$ leads to $k^* = 3$, implying that the prediction problem becomes extremely difficult when there are more than 4 relevant covariates; see the simulations in [77].

3. SOME GENERIC SELECTION SCHEMES

Among the selection schemes not requiring the knowledge of the variance σ^2 , some are very specific to a particular algorithm, while some others are more generic. We describe in this section three versatile selection principles and refer to the examples for the more specific schemes.

3.1 Cross-Validation Procedures

The cross-validation schemes are nearly universal in the sense that they can be implemented in most statistical frameworks and for most estimation procedures. The principle of the cross-validation schemes is to split the data into a *training* set and a *validation* set: the estimators are built on the *training* set, and the *validation* set is used for estimating their prediction risk. This training/validation splitting is eventually repeated several times. The most popular cross-validation schemes are:

- *Hold-out* [27, 57] which is based on a single split of the data for *training* and *validation*.
- *V-fold CV* [32]. The data is split into V subsamples. Each subsample is successively removed for *validation*, the remaining data being used for *training*.
- *Leave-one-out* [69] which corresponds to n -fold CV.
- *Leave- q -out* (also called *delete- q -CV*) [65] where every possible subset of cardinality q of the data is removed for *validation*, the remaining data being used for *training*.

We refer to Arlot and Céliste [5] for a review of the cross-validation schemes and their theoretical properties.

3.2 Penalized Empirical Loss

Penalized empirical loss criteria form another class of versatile selection schemes, yet less universal than CV procedures. The principle is to select among a family $(\hat{\beta}_\lambda)_{\lambda \in \Lambda}$ of estimators by minimizing a criterion of the generic form

$$(10) \quad \text{Crit}(\lambda) = \mathcal{L}_X(Y, \hat{\beta}_\lambda) + \text{pen}(\lambda),$$

where $\mathcal{L}_X(Y, \hat{\beta}_\lambda)$ is a measure of the distance between Y and $X\hat{\beta}_\lambda$, and pen is a function from Λ to \mathbb{R}^+ . The penalty function sometimes depends on data.

Penalized log-likelihood. The most famous criteria of the form (10) are AIC and BIC. They have been designed to select among estimators $\hat{\beta}_\lambda$ obtained by maximizing the likelihood of (β, σ) with the constraint that β lies on a linear space S_λ (called *model*). In the Gaussian case, these estimators are given by $X\hat{\beta}_\lambda = \Pi_{S_\lambda} Y$, where Π_{S_λ} denotes the orthogonal projector onto the model S_λ . For AIC and BIC, the function \mathcal{L}_X corresponds to twice the negative log-likelihood $\mathcal{L}_X(Y, \hat{\beta}_\lambda) = n \log(\|Y - X\hat{\beta}_\lambda\|_2^2)$ and the penalties are $\text{pen}(\lambda) = 2 \dim(S_\lambda)$ and $\text{pen}(\lambda) = \dim(S_\lambda) \log(n)$, respectively. We recall that these two criteria can perform very poorly in a high-dimensional setting.

In the same setting, Baraud et al. [12] propose alternative penalties built from a nonasymptotic perspective. The resulting criterion can handle the high-dimensional setting where p is possibly larger than n , and the risk of the selection procedure is controlled by a bound of the form (9); see Theorem 2 in [12].

Plug-in criteria. Many other penalized-empirical-loss criteria have been developed in the last decades. Several selection criteria [14, 18] have been designed from a nonasymptotic point of view to handle the case where the variance is known. These criteria usually involve the residual least-square $\mathcal{L}_X(Y, \hat{\beta}_\lambda) = \|Y - X\hat{\beta}_\lambda\|_2^2$ and a penalty $\text{pen}(\lambda)$ depending on the variance σ^2 . A common practice is then to plug in the penalty an estimate $\hat{\sigma}^2$ of the variance in place of the variance. For linear regression, when the design matrix X has a rank less than n , a classical choice for $\hat{\sigma}^2$ is

$$\hat{\sigma}^2 = \frac{\|Y - \Pi_X Y\|_2^2}{n - \text{rank}(X)},$$

with Π_X the orthogonal projector onto the range of X . In the Gaussian setting, this estimator $\hat{\sigma}^2$ has the nice feature to be independent of $\Pi_X Y$ on which usually rely the estimators $\hat{\beta}_\lambda$. Nevertheless, the variance of $\hat{\sigma}^2$ is of order $\sigma^4 / (n - \text{rank}(X))$ which is small only when the sample size n is quite large in front of the rank of X . This situation is unfortunately not likely to happen in a high-dimensional setting where p can be larger than n .

**3.3 Approximation Versus Complexity
Penalization: LinSelect**

The criterion proposed by Baraud et al. [12] can handle high-dimensional settings, but it suffers from two rigidities. First, it can only handle *fixed* collections of models $(S_\lambda)_{\lambda \in \Lambda}$. In some situations, the size of Λ is huge (e.g., for complete variable selection) and the estimation procedure can then be computationally intractable. In this case, we may want to work with a

subcollection of models $(S_\lambda)_{\lambda \in \widehat{\Lambda}}$, where $\widehat{\Lambda} \subset \Lambda$ may depend on data. For example, for complete variable selection, the subset $\widehat{\Lambda}$ could be generated by efficient algorithms like Lars [30]. The second rigidity of the procedure of Baraud et al. [12] is that it can only handle constrained-maximum-likelihood estimators. This procedure then does not help for selecting among arbitrary estimators such as the lasso or elastic-net.

These two rigidities have been addressed recently by Baraud et al. [13]. They propose a selection procedure, LinSelect, which can handle both data-dependent collections of models and arbitrary estimators $\widehat{\beta}_\lambda$. The procedure is based on a collection \mathbb{S} of linear spaces which gives a collection of possible ‘‘approximative’’ supports for the estimators $(\widehat{\beta}_\lambda)_{\lambda \in \Lambda}$. A measure of complexity on \mathbb{S} is provided by a weight function $\Delta: \mathbb{S} \rightarrow \mathbb{R}^+$. We refer to Sections 4.1 and 5 for examples of collection \mathbb{S} and weight Δ in the context of coordinate-sparse and group-sparse regression. We present below a simplified version of the LinSelect procedure. For a suitable, possibly data-dependent, subset $\widehat{\mathbb{S}} \subset \mathbb{S}$ (depending on the statistical problem), the estimator $\widehat{\beta}_{\widehat{\lambda}}$ is selected by minimizing the criterion

$$(11) \quad \text{Crit}(\widehat{\beta}_\lambda) = \inf_{S \in \widehat{\mathbb{S}}} \left[\|Y - \Pi_S \mathbf{X} \widehat{\beta}_\lambda\|_2^2 + \frac{1}{2} \|\mathbf{X} \widehat{\beta}_\lambda - \Pi_S \mathbf{X} \widehat{\beta}_\lambda\|_2^2 + \text{pen}(S) \widehat{\sigma}_S^2 \right],$$

where Π_S is the orthogonal projector onto S ,

$$\widehat{\sigma}_S^2 = \frac{\|Y - \Pi_S Y\|_2^2}{n - \dim(S)}$$

and $\text{pen}(S)$ is a penalty depending on Δ . In the cases we will consider here, the penalty $\text{pen}(S)$ is roughly of the order of $\Delta(S)$, and therefore it penalizes S according to its complexity. We refer to Appendix B for a precise definition of this penalty and more details on its characteristics. We emphasize that criterion (11) and the family of estimators $\{\widehat{\beta}_\lambda, \lambda \in \Lambda\}$ are based on the *same* data Y and \mathbf{X} . In other words, there is no data-splitting occurring in the LinSelect procedure. The first term in (11) quantifies the fit of the projected estimator to the data, the second term evaluates the approximation quality of the space S and the last term penalizes S according to its complexity. We refer to Proposition B.1 in Appendix B and Theorem 1 in [12] for risk bounds on the selected estimator. Instantiations of the procedure and more specific risks bounds are given in Sections 4 and 5 in the context of coordinate-sparsity and group-sparsity.

From a computational point of view, the algorithmic complexity of LinSelect is at most proportional to $|\Lambda| \times |\widehat{\mathbb{S}}|$, and in many cases there is no need to scan the whole set $\widehat{\mathbb{S}}$ for each $\lambda \in \Lambda$ to minimize (11). In the examples of Sections 4 and 5, the whole procedure is computationally less intensive than V -fold CV; see Table 3. Finally, we mention that for the constrained least-square estimators $\mathbf{X} \widehat{\beta}_\lambda = \Pi_{S_\lambda} Y$, the LinSelect procedure with $\widehat{\mathbb{S}} = \{S_\lambda : \lambda \in \Lambda\}$ simply coincides with the procedure of Baraud et al. [12].

4. COORDINATE-SPARSITY

In this section, we focus on the high-dimensional linear regression model $Y = \mathbf{X} \beta_0 + \varepsilon$ where the vector β_0 itself is assumed to be sparse. This setting has attracted a lot of attention in the last decade, and many estimation procedures have been developed. Most of them require the choice of tuning parameters which depend on the unknown variance σ^2 . This is, for instance, the case for the lasso [24, 72], Dantzig Selector [22], Elastic Net [85], MC+ [81], aggregation techniques [21, 26], etc.

We first discuss how the generic schemes introduced in the previous section can be instantiated for tuning these procedures and for selecting among them. Then, we pay a special attention to the calibration of the lasso. Finally, we discuss the problem of support estimation and present a small numerical study.

4.1 Automatic Tuning Methods

Cross-validation. Arguably, V -fold cross-validation is the most popular technique for tuning the above-mentioned procedures. To our knowledge, there are no theoretical results for V -fold CV in large dimensional settings.

In practice, V -fold CV seems to give rather good results. The problem of choosing the best V has not yet been solved [5], Section 10, but it is often reported that a good choice for V is between 5 and 10. Indeed, the statistical performance does not increase for larger values of V , and averaging over 10 splits remains computationally feasible [42], Section 7.10.

LinSelect. The procedure LinSelect can be used for selecting among a collection $(\widehat{\beta}_\lambda)_{\lambda \in \Lambda}$ of sparse regressors as follows. For $\mathcal{J} \subset \{1, \dots, p\}$, we define $\mathbf{X}_{\mathcal{J}}$ as the matrix $[\mathbf{X}_{ij}]_{i=1, \dots, n, j \in \mathcal{J}}$ obtained by only keeping the columns of \mathbf{X} with index in \mathcal{J} . We recall that the collection \mathbb{S} gives some possible ‘‘approximative’’ supports for the estimators $(\widehat{\beta}_\lambda)_{\lambda \in \Lambda}$. For sparse linear re-

gression, a possible collection \mathbb{S} and measure of complexity Δ are

$$\mathbb{S} = \{S = \text{range}(\mathbf{X}_{\mathcal{J}}), \mathcal{J} \subset \{1, \dots, p\}, \\ 1 \leq |\mathcal{J}| \leq n/(3 \log p)\}$$

and

$$\Delta(S) = \log \binom{p}{\dim(S)} + \log(\dim(S)).$$

Let us introduce the spaces $\widehat{S}_\lambda = \text{range}(\mathbf{X}_{\text{supp}(\widehat{\beta}_\lambda)})$ and the subcollection of \mathbb{S}

$$\widehat{\mathbb{S}} = \{\widehat{S}_\lambda, \lambda \in \widehat{\Lambda}\} \quad \text{where } \widehat{\Lambda} = \{\lambda \in \Lambda : \widehat{S}_\lambda \in \mathbb{S}\}.$$

The following proposition gives a risk bound when selecting $\widehat{\lambda}$ with LinSelect with the above choice of $\widehat{\mathbb{S}}$ and Δ .

PROPOSITION 4.1. *There exists a numerical constant $C > 1$ such that for any minimizer $\widehat{\lambda}$ of the criterion (11), we have*

$$\begin{aligned} & \mathcal{R}[\widehat{\beta}_{\widehat{\lambda}}; \beta_0] \\ & \leq C \mathbb{E} \left[\inf_{\lambda \in \widehat{\Lambda}} \left\{ \|\mathbf{X}\widehat{\beta}_\lambda - \mathbf{X}\beta_0\|_2^2 \right. \right. \\ & \quad \left. \left. + \inf_{S \in \widehat{\mathbb{S}}} \left\{ \|\mathbf{X}\widehat{\beta}_\lambda - \Pi_S \mathbf{X}\widehat{\beta}_\lambda\|_2^2 \right. \right. \right. \\ & \quad \left. \left. \left. + \dim(S) \log(p) \sigma^2 \right\} \right\} \right] \\ & \leq C \mathbb{E} \left[\inf_{\lambda \in \widehat{\Lambda}} \left\{ \|\mathbf{X}\widehat{\beta}_\lambda - \mathbf{X}\beta_0\|_2^2 \right. \right. \\ & \quad \left. \left. + \|\widehat{\beta}_\lambda\|_0 \log(p) \sigma^2 \right\} \right]. \end{aligned} \tag{12}$$

Proposition 4.1 is a simple corollary of Proposition B.1 in Appendix B. The first bound involves three terms: the loss of the estimator $\widehat{\beta}_\lambda$, an approximation loss, and a variance term. Hence, LinSelect chooses an estimator $\widehat{\beta}_\lambda$ that achieves a trade-off between the loss of $\widehat{\beta}_\lambda$ and the closeness of $\mathbf{X}\widehat{\beta}_\lambda$ to some small dimensional subspace S . Bound (12) cannot be formulated in the form (9) due to the random nature of the set $\widehat{\Lambda}$. Nevertheless, a bound similar to (8) can be deduced from (12) when the estimators $\widehat{\beta}_\lambda$ are least-squares estimators; see Corollary 4 in [13]. Furthermore, we note that increasing the size of Λ leads to a better risk bound for $\widehat{\beta}_{\widehat{\lambda}}$. It is then advisable to consider a family of candidate estimators $\{\widehat{\beta}_\lambda, \lambda \in \Lambda\}$ as large as possible. Proposition 4.1 is valid for any family of estimators $\{\widehat{\beta}_\lambda, \lambda \in \Lambda\}$. For the specific family of lasso estimators $\{\widehat{\beta}_\lambda^L, \lambda > 0\}$ we provide a refined bound in Proposition 4.3, Section 4.3.

4.2 Lasso-Type Estimation under Unknown Variance

The lasso is certainly one of the most popular methods for variable selection in a high-dimensional setting. Given $\lambda > 0$, the lasso estimator $\widehat{\beta}_\lambda^L$ is defined by $\widehat{\beta}_\lambda^L := \text{argmin}_{\beta \in \mathbb{R}^p} \|Y - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$. A sensible choice of λ must be homogeneous with the square-root of the variance σ^2 . As explained above, when the variance σ^2 is unknown, one may apply V -fold CV or LinSelect to select λ . Some alternative approaches have also been developed for tuning the lasso. Their common idea is to modify the ℓ_1 criterion so that the tuning parameter becomes pivotal with respect to σ^2 . This means that the method remains valid for any $\sigma > 0$ and that the choice of the tuning parameter does not depend on σ . For the sake of simplicity, we assume throughout this subsection and the next one that the columns of \mathbf{X} are normalized to one.

ℓ_1 -Penalized log-likelihood. In low-dimensional regression, it is classical to consider a penalized log-likelihood criterion instead of a penalized least-square criterion to handle the unknown variance. Following this principle, Städler et al. [68] propose to minimize the ℓ_1 -penalized log-likelihood criterion

$$\begin{aligned} & \widehat{\beta}_\lambda^{LL}, \widehat{\sigma}_\lambda^{LL} \\ & := \text{argmin}_{\beta \in \mathbb{R}^p, \sigma' > 0} \left[n \log(\sigma') \right. \\ & \quad \left. + \frac{\|Y - \mathbf{X}\beta\|_2^2}{2\sigma'^2} + \lambda \frac{\|\beta\|_1}{\sigma'} \right]. \end{aligned} \tag{13}$$

By reparametrizing (β, σ) , Städler et al. [68] obtain a convex criterion that can be efficiently minimized. Interestingly, the penalty level λ is pivotal with respect to σ . Under suitable conditions on the design matrix \mathbf{X} , Sun and Zhang [70] show that the choice $\lambda = c\sqrt{2 \log p}$, with $c > 1$ yields optimal risk bounds in the sense of (8).

Square-root lasso and scaled lasso. Sun and Zhang [71], following an idea of Antoniadis [3], propose to minimize a penalized Huber loss [45], page 179,

$$\begin{aligned} & \widehat{\beta}_\lambda^{SR}, \widehat{\sigma}_\lambda^{SR} \\ & := \text{argmin}_{\beta \in \mathbb{R}^p, \sigma' > 0} \left[\frac{n\sigma'}{2} + \frac{\|Y - \mathbf{X}\beta\|_2^2}{2\sigma'} \right. \\ & \quad \left. + \lambda \|\beta\|_1 \right]. \end{aligned} \tag{14}$$

This convex criterion can be minimized with roughly the same computational complexity as a Lars-lasso

path [30]. Interestingly, their procedure (called the scaled lasso in [71]) is equivalent to the square-root lasso estimator previously introduced by Belloni et al. [16]. The square-root lasso of Belloni et al. is obtained by replacing the residual sum of squares in the lasso criterion by its square-root

$$(15) \quad \hat{\beta}_\lambda^{SR} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left[\sqrt{\|Y - \mathbf{X}\beta\|_2^2} + \frac{\lambda}{\sqrt{n}} \|\beta\|_1 \right].$$

The equivalence between the two definitions follows from the minimization of the criterion in (14) with respect to σ' . In (14) and (15), the penalty level λ is again pivotal with respect to σ . Sun and Zhang [71] state sharp oracle inequalities for the estimator $\hat{\beta}_\lambda^{SR}$ with $\lambda = c\sqrt{2\log(p)}$, with $c > 1$; see Proposition 4.2 below. Their empirical results suggest that criterion (15) provides slightly better results than the ℓ^1 -penalized log-likelihood. In the sequel, we shall refer to $\hat{\beta}_\lambda^{SR}$ as the square-root lasso estimator.

Bayesian lasso. The Bayesian paradigm allows us to put prior distributions on the variance σ^2 and the tuning parameter λ , as in the Bayesian lasso [60]. Bayesian procedures straightforwardly handle the case of unknown variance, but no frequentist analysis of these procedures is so far available.

4.3 Risk Bounds for Square-Root Lasso and Lasso-LinSelect

Let us state a bound on the prediction error for the square-root lasso (also called scaled lasso). For the sake of conciseness, we only present a simplified version of Theorem 1 in [71]. Consider some number $\xi > 0$ and some subset $T \subset \{1, \dots, p\}$. The compatibility constant $\kappa[\xi, T]$ is defined by

$$\kappa[\xi, T] = \min_{u \in \mathcal{C}(\xi, T)} \left\{ \frac{|T|^{1/2} \|\mathbf{X}u\|_2}{\|u_T\|_1} \right\}$$

where $\mathcal{C}(\xi, T) = \{u : \|u_{T^c}\|_1 < \xi \|u_T\|_1\}$.

PROPOSITION 4.2. *There exist positive numerical constants C_1, C_2 and C_3 such that the following holds. Let us consider the square-root lasso with the tuning parameter $\lambda = 2\sqrt{2\log(p)}$. If we assume that:*

- (1) $p \geq C_1$;
- (2) $\|\beta_0\|_0 \leq C_2 \kappa^2[4, \operatorname{supp}(\beta_0)] \frac{n}{\log(p)}$

then, with high probability,

$$\begin{aligned} & \|\mathbf{X}(\hat{\beta}^{SR} - \beta_0)\|_2^2 \\ & \leq \inf_{\beta \neq 0} \left\{ \|\mathbf{X}(\beta_0 - \beta)\|_2^2 + C_3 \frac{\|\beta\|_0 \log(p)}{\kappa^2[4, \operatorname{supp}(\beta)]} \sigma^2 \right\}. \end{aligned}$$

This bound is comparable to the general objective (9) stated in Section 2.4. Interestingly, the constant before the bias term $\|\mathbf{X}(\beta_0 - \beta)\|_2^2$ equals one. If $\|\beta_0\|_0 = k$, the square-root lasso achieves the minimax loss $k \log(p) \sigma^2$ as long as $k \log(p)/n$ is small, and $\kappa[4, \operatorname{supp}(\beta_0)]$ is away from zero. This last condition ensures that the design \mathbf{X} is not too far from orthogonality on the cone $\mathcal{C}(4, \operatorname{supp}(\beta_0))$. State of the art results for the classical lasso with known variance [17, 49, 74] all involve this condition.

In what follows, we call lasso-LinSelect the lasso estimator of β_0 obtained by choosing the parameter λ in $\Lambda = \mathbb{R}^+$ with LinSelect. We next state a risk bound for this procedure. For $\mathcal{J} \subset \{1, \dots, p\}$, we define $\phi_{\mathcal{J}}$ as the largest eigenvalue of $X_{\mathcal{J}}^T X_{\mathcal{J}}$. The following proposition involves the restricted eigenvalue $\phi_* = \max\{\phi_{\mathcal{J}} : \operatorname{Card}(\mathcal{J}) \leq n/(3 \log p)\}$.

PROPOSITION 4.3. *There exist positive numerical constants C, C_1, C_2 and C_3 such that the following holds. Take $\Lambda = \mathbb{R}^+$, and assume that*

$$\|\beta_0\|_0 \leq C \frac{\kappa^2[5, \operatorname{supp}(\beta_0)]}{\phi_*} \times \frac{n}{\log(p)}.$$

Then, with probability at least $1 - C_1 p^{-C_2}$, the lasso estimator $\hat{\beta}_\lambda^L$ selected according to the LinSelect procedure described in Section 4.1 fulfills

$$(16) \quad \begin{aligned} & \|\mathbf{X}(\beta_0 - \hat{\beta}_\lambda^L)\|_2^2 \\ & \leq C_3 \inf_{\beta \neq 0} \left\{ \|\mathbf{X}(\beta_0 - \beta)\|_2^2 \right. \\ & \quad \left. + \frac{\phi_* \|\beta\|_0 \log(p)}{\kappa^2[5, \operatorname{supp}(\beta)]} \sigma^2 \right\}. \end{aligned}$$

The bound (16) is similar to the bound stated above for the square-root lasso, the most notable differences being the constant larger than 1 in front of the bias term and the quantity ϕ_* in front of the variance term. We refer to the supplementary material [38] for a proof of Proposition 4.3.

4.4 Support Estimation and Inverse Problem

Until now, we only discussed estimation methods that perform well in prediction. Little is known when the objective is to infer β_0 or its support under unknown variance.

Inverse problem. The square-root lasso [16, 71] is proved to achieve near optimal risk bound for the inverse problems under suitable assumptions on the design \mathbf{X} .

Support estimation. Up to our knowledge, there are no nonasymptotic results on support estimation for the aforementioned procedures in the unknown variance setting. Nevertheless, some related results and heuristics have been developed for the cross-validation scheme. If the tuning parameter λ is chosen to minimize the prediction error [i.e., take $\lambda = \lambda^*$ as defined in (4)], the lasso is not consistent for support estimation; see [51, 56] for results in a random design setting. One idea to overcome this problem, is to choose the parameter λ that minimizes the risk of the so-called Gauss-lasso estimator $\hat{\beta}_\lambda^{GL}$ which is the least square estimator over the support of the lasso estimator $\hat{\beta}_\lambda^L$

$$(17) \quad \hat{\beta}_\lambda^{GL} := \underset{\beta \in \mathbb{R}^p: \text{supp}(\beta) \subset \text{supp}(\hat{\beta}_\lambda^L)}{\text{argmin}} \quad \|Y - \mathbf{X}\beta\|_2^2.$$

When the objective is support estimation, some numerical simulations [62] suggest that it may be more advisable not to apply the selection schemes based on prediction risk (such as V -fold CV or LinSelect) to the lasso estimators, but rather to the Gauss-lasso estimators. Similar remarks also apply for the Dantzig Selector [22].

4.5 Numerical Experiments

We present two numerical experiments to illustrate the behavior of some of the above mentioned procedures for high-dimensional sparse linear regression. The first one concerns the problem of tuning the parameter λ of the lasso algorithm for estimating $\mathbf{X}\beta_0$. The procedures will be compared on the basis of the prediction risk. The second one concerns the problem of support estimation with lasso-type estimators. We will focus on the false discovery rates (FDR) and the proportion of true discoveries (Power).

Simulation design. The simulation design is the same as the one described in Sections 6.1 and 8.2 of [13], except that we restrict to the case $n = p = 100$. Therefore, 165 examples are simulated. They are inspired by examples found in [43, 72, 84, 85] and cover a large variety of situations. The simulation were carried out with R www.r-project.org, using the library `elasticnet`.

Experiment 1: Tuning the lasso for prediction. In the first experiment, we compare 10-fold CV [32], LinSelect [13] and the square-root lasso [16, 71] (also called scaled lasso) for tuning the lasso. The lasso-LinSelect requires one to minimize criterion (11) over $\Lambda = R^+$. In fact, the minimum of (11) is achieved at some λ corresponding to a Lars-lasso [30] step. Consequently,

only the steps of the regularization path of the lasso have to be considered so that the lasso-LinSelect estimator can be computed with roughly the same computational complexity as a Lars-lasso path [30]. The square-root lasso is computed using the algorithm described in Sun and Zhang [71]. We set $\lambda = \sqrt{2\log(p)}$ which corresponds to the value recommended by the authors.²

For each tuning procedure $\ell \in \{10\text{-fold CV, LinSelect, square-root lasso}\}$, we focus on the prediction risk $\mathcal{R}[\hat{\beta}_{\lambda_\ell}^L; \beta_0]$ of the selected lasso estimator $\hat{\beta}_{\lambda_\ell}^L$.

For each simulated example $e = 1, \dots, 165$, we estimate on the basis of 400 runs:

- the risk of the oracle (4): $\mathcal{R}_e = \mathcal{R}[\hat{\beta}_{\lambda^*}; \beta_0]$;
- the risk when selecting λ with procedure ℓ : $\mathcal{R}_{\ell,e} = \mathcal{R}[\hat{\beta}_{\lambda_\ell}; \beta_0]$.

The comparison between the procedures is based on the comparison of the means, standard deviations and quantiles of the risk ratios $\mathcal{R}_{\ell,e}/\mathcal{R}_e$ computed over all the simulated examples $e = 1, \dots, 165$. The results are displayed in Table 1.

For 10-fold CV and LinSelect, the risk ratios are close to one. For 90% of the examples, the risk of the lasso-LinSelect is smaller than the risk of the lasso-CV, but there are a few examples where the risk of the lasso-LinSelect is significantly larger than the risk of the lasso-CV. For the square-root lasso procedure, the risk ratios are clearly larger than for the two others. An inspection of the results reveals that the square-root lasso selects estimators with supports of small size. This feature can be interpreted as follows. Due to the bias of the lasso-estimator, the residual variance tends to overestimate the variance, leading the square-root lasso to

TABLE 1

For each procedure ℓ , mean, standard-error and quantiles of the ratios $\{\mathcal{R}_{\ell,e}/\mathcal{R}_e, e = 1, \dots, 165\}$

Procedure	Mean	Std-err	Quantiles				
			0%	50%	75%	90%	95%
Lasso 10-fold CV	1.13	0.08	1.03	1.11	1.15	1.19	1.24
Lasso LinSelect	1.19	0.48	0.97	1.03	1.06	1.19	2.52
Square-root lasso	5.15	6.74	1.32	2.61	3.37	11.2	17

²More precisely, Sun and Zhang advocate $\lambda_0 = \sqrt{2\log(p)/n}$ in their paper. This choice is equivalent to $\lambda = \sqrt{2\log(p)}$ in (14) because the normalizations of \mathbf{X} and of (14) are different from [71].

TABLE 2
For each procedure ℓ , mean, standard-error and quantiles of FDR and Power values

Procedure	Mean	Std-err	Quantiles				
			0%	25%	50%	75%	90%
False discovery rate							
Gauss-lasso 10-fold CV	0.28	0.26	0	0.08	0.22	0.35	0.74
Gauss-lasso LinSelect	0.12	0.25	0	0.002	0.02	0.13	0.33
Square-root lasso	0.13	0.26	0	0.009	0.023	0.07	0.32
Power							
Gauss-lasso 10-fold CV	0.67	0.18	0.4	0.52	0.65	0.71	1
Gauss-lasso LinSelect	0.56	0.33	0.002	0.23	0.56	0.93	1
Square-root lasso	0.59	0.28	0.013	0.41	0.57	0.80	1

select a lasso estimator $\widehat{\beta}_\lambda^L$ with large λ . Consequently the risk is high.

Experiment 2: Variable selection with Gauss-lasso and square-root lasso. We consider now the problem of support estimation, sometimes referred as the problem of variable selection. We implement three procedures. The Gauss-lasso procedure tuned by either 10-fold CV or LinSelect and the square-root lasso. The support of β_0 is estimated by the support of the selected estimator.

For each simulated example, the FDR and the Power are estimated on the basis of 400 runs. The results are given on Table 2.

It appears that the Gauss-lasso CV procedure gives greater values of the FDR than the two others. The Gauss-lasso LinSelect and the square-root lasso behave similarly for the FDR, but the values of the power are more variable for the LinSelect procedure.

Computation time. Let us conclude this numerical section with the comparison of the computation times between the methods. For all methods the computation time depends on the maximum number of steps in the lasso algorithm, and for the LinSelect method, it depends on the cardinality of \mathbb{S} or equivalently on the maximum number of nonzero components of $\widehat{\beta}$. The results are shown at Table 3. The square-root lasso is the less time consuming method, closely followed by the lasso LinSelect method. The V -fold CV carried out with the function `cv.enet` of the R package `elasticnet`, pays the price of several calls to the lasso algorithm.

5. GROUP-SPARSITY

In the previous section, we have made no prior assumptions on the form of β_0 . In some applications,

there are some known structures between the covariates. As an example, we treat the now classical case of group sparsity. The covariates are assumed to be clustered into M groups, and when the coefficient $\beta_{0,i}$ corresponding to the covariate \mathbf{X}_i is nonzero, then it is likely that all the coefficients $\beta_{0,j}$ with variables \mathbf{X}_j in the same group as \mathbf{X}_i are nonzero. We refer to the introduction of [8] for practical examples of this so-called group-sparsity assumption. Let G_1, \dots, G_M form a given partition of $\{1, \dots, p\}$. For $\lambda = (\lambda_1, \dots, \lambda_M)$, the group-lasso estimator $\widehat{\beta}_\lambda$ is defined as the minimizer of the convex optimization criterion

$$(18) \quad \|Y - \mathbf{X}\beta\|_2^2 + \sum_{k=1}^M \lambda_k \|\beta^{G_k}\|_2,$$

where $\beta^{G_k} = (\beta_j)_{j \in G_k}$. Criterion (18) promotes solutions where all the coordinates of β^{G_k} are either zero or nonzero, leading to group selection [80]. Under some assumptions on \mathbf{X} , Huang and Zhang [44] or Lounici et al. [54] provide a suitable choice of $\lambda = (\lambda_1, \dots, \lambda_M)$ that leads to near optimal prediction bounds. As expected, this choice of $\lambda = (\lambda_1, \dots, \lambda_M)$ is proportional to σ .

As for the lasso, V -fold CV is widely used in practice to tune the penalty parameter $\lambda = (\lambda_1, \dots, \lambda_M)$. To our knowledge, there is not yet any extension of the procedures described in Section 4.2 to the group lasso. An alternative to cross-validation is to use LinSelect.

Tuning the Group-Lasso with LinSelect

For any $\mathcal{K} \subset \{1, \dots, M\}$, we define the submatrix $\mathbf{X}_{(\mathcal{K})}$ of \mathbf{X} by only keeping the columns of \mathbf{X} with index in $\bigcup_{k \in \mathcal{K}} G_k$. We also write \mathbf{X}_{G_k} for the submatrix of \mathbf{X} , built from the columns with index in G_k . The

TABLE 3

For each procedure computation time for different values of n and p . The maximum number of steps in the lasso algorithm is taken as $\text{max.steps} = \min\{n, p\}$. For the LinSelect procedure, the maximum number of nonzero components of $\hat{\beta}$, denoted k_{max} is taken as $k_{\text{max}} = \min\{p, n/\log(p)\}$

n	p	max.steps	k_{max}	Lasso 10-fold CV	Lasso LinSelect	Square-root lasso
100	100	100	21	4 s	0.21 s	0.18 s
100	500	100	16	4.8 s	0.43 s	0.4 s
500	500	500	80	300 s	11 s	6.3 s

collection \mathbb{S} and the function Δ are given by

$$\mathbb{S} = \left\{ \text{range}(\mathbf{X}_{(\mathcal{K})}) : 1 \leq |\mathcal{K}| \leq n/(3 \log(M)) \right.$$

$$\left. \text{and } \sum_{k \in \mathcal{K}} |G_k| \leq n/2 - 1 \right\}$$

and $\Delta(\text{range}(\mathbf{X}_{(\mathcal{K})})) = \log[|\mathcal{K}| \binom{M}{|\mathcal{K}|}]$. For a given $\Lambda \subset \mathbb{R}_+^M$, similarly to Section 4.1, we define $\hat{\mathcal{K}}_\lambda = \{k : \|\hat{\beta}_\lambda^{G_k}\|_2 \neq 0\}$ and

$$\hat{\mathbb{S}} = \{ \text{range}(\mathbf{X}_{(\hat{\mathcal{K}}_\lambda)}), \lambda \in \hat{\Lambda} \},$$

with $\hat{\Lambda} = \{ \lambda \in \Lambda, \text{range}(\mathbf{X}_{(\hat{\mathcal{K}}_\lambda)}) \in \mathbb{S} \}$.

Proposition B.1 in Appendix B ensures that we have for some constant $C > 1$,

$$\mathcal{R}[\hat{\beta}_\lambda; \beta_0] \leq C \mathbb{E} \left[\inf_{\lambda \in \hat{\Lambda}} \{ \|\mathbf{X}\hat{\beta}_\lambda - \mathbf{X}\beta_0\|_2^2 + (\|\hat{\beta}_\lambda\|_0 \vee |\hat{\mathcal{K}}_\lambda| \log(M)) \sigma^2 \} \right].$$

In the following, we provide a more explicit bound. For simplicity, we restrict to the specific case where each group G_k has the same cardinality T . For $\mathcal{K} \subset \{1, \dots, M\}$, we define $\phi_{(\mathcal{K})}$ as the largest eigenvalue of $\mathbf{X}_{(\mathcal{K})}^T \mathbf{X}_{(\mathcal{K})}$, and we set

$$(19) \quad \phi_* = \max \left\{ \phi_{(\mathcal{K})} : 1 \leq |\mathcal{K}| \leq \frac{n-2}{2T \vee 3 \log(M)} \right\}.$$

We assume that all the columns of \mathbf{X} are normalized to 1, and following Lounici et al. [54], we introduce for $1 \leq s \leq M$,

$$(20) \quad \kappa_G[\xi, s] = \min_{1 \leq |\mathcal{K}| \leq s} \min_{u \in \Gamma(\xi, \mathcal{K})} \frac{\|\mathbf{X}u\|_2}{\|u_{(\mathcal{K})}\|_2},$$

where $\Gamma(\xi, \mathcal{K})$ is the cone of vectors $u \in \mathbb{R}^M \setminus \{0\}$ such that $\sum_{k \in \mathcal{K}^c} \lambda_k \|u^{G_k}\|_2 \leq \xi \sum_{k \in \mathcal{K}} \lambda_k \|u^{G_k}\|_2$. In the sequel, \mathcal{K}_0 stands for the set of groups containing nonzero components of β_0 .

PROPOSITION 5.1. *There exist positive numerical constants C, C_1, C_2 and C_3 such that the following holds. Assume that Λ contains $\bigcup_{\lambda \in \mathbb{R}_+} \{(\lambda, \dots, \lambda)\}$, that $T \leq (n-2)/4$ and that*

$$1 \leq |\mathcal{K}_0| \leq C \frac{\kappa_G^2[3, |\mathcal{K}_0|]}{\phi_*} \times \frac{n-2}{\log(M) \vee T}.$$

Then, with probability larger than $1 - C_1 M^{-C_2}$, we have

$$\|\mathbf{X}\hat{\beta}_\lambda - \mathbf{X}\beta_0\|_2^2 \leq C_3 \frac{\phi_*}{\kappa_G^2[3, |\mathcal{K}_0|]} |\mathcal{K}_0| (T \vee \log(M)).$$

This proposition provides a bound comparable to the bounds of Lounici et al. [54], without requiring the knowledge of the variance. Its proof can be found in the supplementary material [38].

6. VARIATION-SPARSITY

We focus in this section on the *variation-sparse* regression. We recall that the vector $\beta^V \in \mathbb{R}^{p-1}$ of the variations of β has for coordinates $\beta_j^V = \beta_{j+1} - \beta_j$ and that the variation-sparse setting corresponds to the setting where the vector of variations β_0^V is coordinate-sparse. In the following, we restrict to the case where $n = p$, and \mathbf{X} is the identity matrix. In this case, the problem of variation-sparse regression coincides with the problem of segmentation of the mean of the vector $Y = \beta_0 + \varepsilon$.

For any subset $\mathcal{I} \subset \{1, \dots, n-1\}$, we define $S_{\mathcal{I}} = \{\beta \in \mathbb{R}^n : \text{supp}(\beta^V) \subset \mathcal{I}\}$ and $\hat{\beta}_{\mathcal{I}} = \Pi_{S_{\mathcal{I}}} Y$. For any integer $q \in \{0, \dots, n-1\}$, we define also the “best” subset of size q by

$$\hat{\mathcal{I}}_q = \underset{|\mathcal{I}|=q}{\text{argmin}} \|Y - \hat{\beta}_{\mathcal{I}}\|_2^2.$$

Though the number of subsets $\mathcal{I} \subset \{1, \dots, n-1\}$ of cardinality q is of order n^{q+1} , this minimization can be performed using dynamic programming with a complexity of order n^2 [40]. To select $\hat{\mathcal{I}} = \hat{\mathcal{I}}_{\hat{q}}$ with \hat{q} in $\{0, \dots, n-1\}$, any of the generic selection schemes of Section 3 can be applied. Below, we instantiate these schemes and present some alternatives.

6.1 Penalized Empirical Loss

When the variance σ^2 is known, penalized log-likelihood model selection amounts to select a subset $\widehat{\mathcal{I}}$ which minimizes a criterion of the form $\|Y - \widehat{\beta}_{\mathcal{I}}\|_2^2 + \text{pen}(\text{Card}(\mathcal{I}))$. This is equivalent to select $\widehat{\mathcal{I}} = \widehat{\mathcal{I}}_{\widehat{q}}$ with \widehat{q} minimizing

$$(21) \quad \text{Crit}(q) = \|Y - \widehat{\beta}_{\mathcal{I}_q}\|_2^2 + \text{pen}(q).$$

Following the work of Birgé and Massart [18], Lebarbier [50] considers the penalty

$$\text{pen}(q) = (q + 1)(c_1 \log(n/(q + 1)) + c_2)\sigma^2$$

and determines the constants $c_1 = 2, c_2 = 5$ by extensive numerical experiments (see also Comte and Rozenholc [25] for a similar approach in a more general setting). With this choice of the penalty, the procedure satisfies a bound of the form

$$(22) \quad \begin{aligned} & \mathcal{R}[\widehat{\beta}_{\widehat{\mathcal{I}}}, \beta_0] \\ & \leq C \inf_{\mathcal{I} \subset \{1, \dots, n-1\}} \{ \|\widehat{\beta}_{\mathcal{I}} - \beta_0\|_2^2 \\ & \quad + (1 + |\mathcal{I}|) \log(n/(1 + |\mathcal{I}|))\sigma^2 \}. \end{aligned}$$

When σ^2 is unknown, several approaches have been proposed.

Plug-in estimator. The idea is to replace σ^2 in $\text{pen}(q)$ by an estimator of the variance such as $\widehat{\sigma}^2 = \sum_{i=1}^{n/2} (Y_{2i} - Y_{2i-1})^2/n$, or one of the estimators proposed by Hall et al. [41]. No theoretical results are proved in a nonasymptotic framework.

Estimating the variance by the residual least-squares. Baraud et al. [12], in Section 5.4.2, propose to select q by minimizing a penalized log-likelihood criterion. This criterion can be written in the form $\text{Crit}(q) = \|Y - \widehat{\beta}_{\mathcal{I}_q}\|_2^2(1 + K \text{pen}(q))$, with $K > 1$ and the penalty $\text{pen}(q)$ solving

$$\mathbb{E}[(U - \text{pen}(q)V)_+] = \frac{1}{(q + 1) \binom{n-1}{q}},$$

where $(\cdot)_+ = \max(\cdot, 0)$, and U, V are two independent χ^2 variables with respectively $q + 2$ and $n - q - 2$ degrees of freedom. The resulting estimator $\widehat{\beta}_{\widehat{\mathcal{I}}}$, with $\widehat{\mathcal{I}} = \widehat{\mathcal{I}}_{\widehat{q}}$, satisfies a nonasymptotic risk bound similar to (22) for all $K > 1$. The choice $K = 1.1$ is suggested for the practice.

Slope heuristic. Lebarbier [50] implements the slope heuristic introduced by Birgé and Massart [19] for handling the unknown variance σ^2 . The method consists in calibrating the penalty directly, without estimating $\widehat{\sigma}^2$. It is based on the following principle. First, there exists a so-called *minimal* penalty $\text{pen}_{\min}(q)$ such that choosing $\text{pen}(q) = K \text{pen}_{\min}(q)$ in (21) with $K < 1$ can lead to a strong overfit, whereas for $K > 1$, the bound (22) is met. Second, it can be shown that there exists a *dimension jump* around the minimal penalty, allowing one to estimate $\text{pen}_{\min}(q)$ from the data. The slope heuristic then proposes to select q by minimizing the criterion $\text{Crit}(q) = \|Y - \widehat{\beta}_{\mathcal{I}_q}\|_2^2 + 2\widehat{\text{pen}}_{\min}(q)$. Arlot and Massart [7] provide a nonasymptotic risk bound for this procedure. Their results are proved in a general regression model with heteroscedastic and non-Gaussian errors, but with a constraint on the number of models per dimension which is not met for the family of models $(S_{\mathcal{I}})_{\mathcal{I} \subset \{1, \dots, n-1\}}$. Nevertheless, the authors indicate how to generalize their results for the problem of signal segmentation.

Finally, for practical issues, different procedures for estimating the minimal penalty are compared and implemented in Baudry et al. [15].

6.2 CV Procedure

In a recent paper, Arlot and Céliste [6] consider the problem of signal segmentation using cross-validation. Their results apply in the heteroscedastic case. They consider several CV-methods, the leave-one-out, leave- p -out and V -fold CV for estimating the quadratic loss. They propose two cross-validation schemes. The first one, denoted *Procedure 5*, aims to estimate directly $\mathbb{E}[\|\beta_0 - \beta_{\widehat{\mathcal{I}}_q}\|_2^2]$, while the second one, denoted *Procedure 6*, relies on two steps, where the cross-validation is used first for choosing the best partition of dimension q , then the best dimension q . They show that the leave- p -out CV method can be implemented with a complexity of order n^2 , and they give a control of the expected CV risk. The use of CV leads to some restrictions on the subsets \mathcal{I} that compete for estimating β_0 . This problem is discussed in [6], Section 3 of the supplementary material.

6.3 Alternative for Very High-Dimensional Settings

When n is very large, the dynamic programming optimization can become computationally too intensive. An attractive alternative is based on the fused lasso proposed by Tibshirani et al. [73]. The estimator $\widehat{\beta}_{\lambda}^{TV}$ is

defined by minimizing the convex criterion

$$\|Y - \beta\|_2^2 + \lambda \sum_{j=1}^{n-1} |\beta_{j+1} - \beta_j|,$$

where the total-variation norm $\sum_j |\beta_{j+1} - \beta_j|$ promotes solutions which are variation-sparse. The family $(\hat{\beta}_\lambda^{TV})_{\lambda \geq 0}$ can be computed very efficiently with the Lars-algorithm; see Vert and Bleakley [75]. A sensible choice of the parameter λ must be proportional to σ . When the variance σ^2 is unknown, the parameter λ can be selected either by V -fold CV or by LinSelect (see Section 5.1 in [13] for details).

7. EXTENSIONS

7.1 Gaussian Design and Graphical Models

Assume that the design \mathbf{X} is now random and that the n rows $\mathbf{X}^{(i)}$ are independent observations of a Gaussian vector with mean 0_p and unknown covariance matrix Σ . This setting is mainly motivated by applications in compressed sensing [29] and in Gaussian graphical modeling. Indeed, Meinshausen and Bühlmann [56] have proved that it is possible to estimate the graph of a Gaussian graphical model by studying linear regression with Gaussian design and unknown variance. If we work conditionally on the observed \mathbf{X} design, then all the results and methodologies described in this survey still apply. Nevertheless, these prediction results do not really take into account the fact that the design is random. In this setting, it is more natural to consider the integrated prediction risk $\mathbb{E}[\|\Sigma^{1/2}(\hat{\beta} - \beta_0)\|_2^2]$ rather than the risk (3). Some procedures [35, 76] have been proved to achieve optimal risk bounds with respect to this risk, but they are computationally intractable in a high-dimensional setting. In the context of Gaussian graphical modeling, the procedure GGMselect [39] is designed to select among any collection of graph estimators, and it is proved to achieve near optimal risk bounds in terms of the integrated prediction risk.

7.2 Non-Gaussian Noise

A few results do not require that the noise ε follows a Gaussian distribution. The lasso-type procedures, such as the square-root lasso [16, 71], do not require the normality of the noise and extend to other distributions. In practice, it seems that cross-validation procedures still work well for other distributions of the noise.

7.3 Multivariate Regression

Multivariate regression deals with T simultaneous linear regression models $y_k = \mathbf{X}\beta_k + \varepsilon_k, k = 1, \dots, T$. Stacking the y_k 's in a $n \times T$ matrix Y , we obtain the model $Y = \mathbf{X}B_0 + E$, where B_0 is a $p \times T$ matrix with columns given by β_k and E is a $n \times T$ matrix with i.i.d. entries. The classical structural assumptions on B_0 are either that most rows of B_0 are identically zero, or the rank of B_0 is small. The first case is a simple case of group sparsity and can be handled by the group-lasso as in Section 5. The second case, first considered by Anderson [2] and Izenman [46], is much more nonlinear. Writing $\|\cdot\|_F$ for the Frobenius (or Hilbert–Schmidt) norm, the problem of selecting among the estimators

$$\hat{B}_r = \underset{B: \text{rank}(B) \leq r}{\text{argmin}} \|Y - \mathbf{X}B\|_F^2,$$

$$r \in \{1, \dots, \min(T, \text{rank}(\mathbf{X}))\}$$

has been investigated recently from a nonasymptotic point of view by Bunea et al. [20] and Giraud [36]. The prediction risk of \hat{B}_r is of order of

$$\mathbb{E}[\|\mathbf{X}\hat{B}_r - \mathbf{X}B_0\|_F^2] \asymp \sum_{k \geq r} s_k^2(\mathbf{X}B_0) + r(n + \text{rank}(X))\sigma^2,$$

where $s_k(M)$ denotes the k th largest singular value of the matrix M . Therefore, a sensible choice of r depends on σ^2 . The first selection criterion introduced by Bunea et al. [20] requires the knowledge of the variance σ^2 . To handle the case of unknown variance, Bunea et al. [20] propose to plug an estimate of the variance in their selection criterion [which works when $\text{rank}(\mathbf{X}) < n$], whereas Giraud [36] introduces a penalized log-likelihood criterion independent of the variance. Both papers provide oracle risk bounds for the resulting estimators showing rate-minimax adaptation.

Several recent papers [9, 20, 49, 58, 63] have investigated another strategy for the low-rank setting. For a positive λ , the matrix B_0 is estimated by

$$\hat{B}_\lambda \in \underset{B \in \mathbb{R}^{p \times T}}{\text{argmin}} \left\{ \|Y - \mathbf{X}B\|_F^2 + \lambda \sum_k s_k(B) \right\}.$$

Translating the work on trace regression of Koltchinskii et al. [49] into the setting of multivariate regression provides (under some conditions on \mathbf{X}) an oracle bound on the risk of \hat{B}_{λ^*} with $\lambda^* = 3s_1(X)(\sqrt{T} + \sqrt{\text{rank}(X)})\sigma$. We also refer to Giraud [37] for a slight variation of this result requiring no condition on the design \mathbf{X} . Again, the value of λ^* is proportional to σ .

To handle the case of unknown variance, Klopp [48] adapts the concept of the square-root lasso [16] to this setting and provides an oracle risk bound for the resulting procedure.

7.4 Nonparametric Regression

In the nonparametric regression model (2), classical estimation procedures include local-polynomial estimators, kernel estimators, basis-projection estimators, k -nearest neighbors, etc. All these procedures depend on one (or several) tuning parameter(s), whose optimal value(s) scales with the variance σ^2 . V -fold CV is widely used in practice for choosing these parameters, but little is known on its theoretical performance.

The class of linear estimators (including spline smoothing, Nadaraya estimators, k -nearest neighbors, low-pass filters, kernel ridge regression, etc.) has attracted some attention in the last years. Some papers have investigated the tuning of some specific family of estimators. For example, Cao and Golubev [23] provides a tuning procedure for spline smoothing while Zhang [82] provides a sharp analysis of kernel ridge regression. Recently, two papers have focused on the tuning of arbitrary linear estimators when the variance σ^2 is unknown. Arlot and Bach [4] generalize the slope heuristic to symmetric linear estimators with spectrum in $[0, 1]$ and prove an oracle bound for the resulting estimator. Baraud et al. [13], in Section 4, shows that Lin-Select can be used for selecting among a (almost) completely arbitrary collection of linear estimators (possibly nonsymmetric and/or singular). An oracle bound for the selected estimator is given by Corollary 2 in [13].

APPENDIX A: A NOTE ON BIC TYPE CRITERIA

The BIC criterion has been initially introduced [64] to select an estimator among a collection of constrained maximum likelihood estimators. Nevertheless, modified versions of this criterion are often used for tuning more general estimation procedures. The purpose of this appendix is to illustrate why we advise against this approach in a high-dimensional setting.

DEFINITION A.1 (A modified BIC criterion). Suppose we are given a collection $(\hat{\beta}_\lambda)_{\lambda \in \Lambda}$ of estimators depending on a tuning parameter $\lambda \in \Lambda$. For any $\lambda \in \Lambda$, we consider $\hat{\sigma}_\lambda^2 = \|Y - \mathbf{X}\hat{\beta}_\lambda\|_2^2/n$, and define the modified BIC

$$(A.1) \quad \hat{\lambda} \in \underset{\lambda \in \hat{\Lambda}}{\operatorname{argmin}} \{-2\mathbf{L}_n(\hat{\beta}_\lambda, \hat{\sigma}_\lambda) + \log(n)\|\hat{\beta}_\lambda\|_0\},$$

where \mathbf{L}_n is the log-likelihood and $\hat{\Lambda} = \{\lambda \in \Lambda : \|\hat{\beta}_\lambda\|_0 \leq n/2\}$.

Sometimes, the $\log(n)$ term is replaced by $\log(p)$. Replacing Λ by $\hat{\Lambda}$ allows us to avoid trivial estimators. First, we would like to emphasize that there is *no* theoretical warranty that the selected estimator does not overfit in a high-dimensional setting. In practice, using this criterion often leads to overfitting. Let us illustrate this with a simple experiment.

Setting

We consider the model

$$(A.2) \quad Y_i = \beta_{0,i} + \varepsilon_i, \quad i = 1, \dots, n,$$

with $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ so that $p = n$ and $\mathbf{X} = I_n$. Here, we fix $n = 10,000$, $\sigma = 1$ and $\beta_0 = 0_n$.

Methods

We apply the modified BIC criterion to tune the lasso [72], SCAD [31] and the hard thresholding estimator. The hard thresholding estimator $\hat{\beta}_\lambda^{HT}$ is defined for any $\lambda > 0$ by $[\hat{\beta}_\lambda^{HT}]_i = Y_i \mathbf{1}_{|Y_i| \geq \lambda}$. Given $\lambda > 0$ and $a > 2$, the SCAD estimator $\hat{\beta}_{\lambda,a}^{SCAD}$ is defined as the minimizer of the penalized criterion $\|Y - \mathbf{X}\beta\|_2^2 + \sum_{i=1}^n p_\lambda(|\beta_i|)$, where for $x > 0$,

$$p'_\lambda(x) = \lambda \mathbf{1}_{x \leq \lambda} + (a\lambda - x) \mathbf{1}_{x > \lambda} / (a - 1).$$

For the sake of simplicity we fix $a = 3$. We note $\hat{\beta}^{L;\text{BIC}}$, $\hat{\beta}_a^{\text{SCAD};\text{BIC}}$, and $\hat{\beta}^{\text{HT};\text{BIC}}$ for the lasso, hard thresholding, and SCAD estimators selected by the modified BIC criterion.

Results

We have realized $N = 200$ experiments. For each of these experiments, the estimator $\hat{\beta}^{L;\text{BIC}}$, $\hat{\beta}_a^{\text{SCAD};\text{BIC}}$ and $\hat{\beta}^{\text{HT};\text{BIC}}$ are computed. The mean number of nonzero components and the estimated risk $\mathcal{R}[\hat{\beta}^{*;\text{BIC}}; 0_n]$ are reported in Table 4.

Obviously, the SCAD and hard thresholding methods select too many irrelevant variables when they are tuned with BIC. Moreover, their risks are quite high. Intuitively, this is due to the fact that the $\log(n)$ [or $\log(p)$] term in the BIC penalty is too small in this high-dimensional setting ($n = p$).

For the lasso estimator, a very specific phenomenon occurs due to the soft thresholding effect. In the discussion of [30], Loubes and Massart advocate that

TABLE 4
Estimated risk and Estimated number of nonzero components for $\hat{\beta}^{L;\text{BIC}}$, $\hat{\beta}_a^{\text{SCAD};\text{BIC}}$ and $\hat{\beta}^{\text{HT};\text{BIC}}$

	Lasso	SCAD	Hard thres.
$\mathcal{R}[\hat{\beta}^{*;\text{BIC}}; 0_p]$	4.6×10^{-2}	1.6×10^1	3.0×10^2
Mean of $\ \hat{\beta}^{*;\text{BIC}}\ _0$	0.025	86.9	28.2

soft thresholding estimators penalized by Mallows' C_p [55] penalties should yield good results, while hard thresholding estimators penalized by Mallows's C_p are known to highly overfit. This strange behavior is due to the bias of the soft thresholding estimator. Nevertheless, Loubes and Massart's arguments have been developed for an orthogonal design. In fact, there is no nonasymptotic justification that the lasso tuned by BIC or AIC performs well for general designs \mathbf{X} .

Conclusion

The use of the modified BIC criterion to tune estimation procedures in a high-dimensional setting is not supported by theoretical results. It is proved to overfit in the case of thresholding estimators [12], Section 3.2.2. Empirically, BIC seems to overfit except for the lasso. We advise the practitioner to avoid BIC (and AIC) when p is at least of the same order as n . For instance, LinSelect is supported by nonasymptotic arguments and by empirical results [13] in contrast to BIC.

APPENDIX B: COMPLEMENTS ON LINSELECT

B.1 More Details on the Selection Procedure

The penalty $\text{pen}(S)$ involved in the LinSelect criterion (11) is defined by $\text{pen}(S) = 1.1 \text{pen}_\Delta(S)$ where $\text{pen}_\Delta(S)$ is the unique solution of

$$\mathbb{E} \left[\left(U - \frac{\text{pen}_\Delta(S)}{n - \dim(S)} V \right)_+ \right] = e^{-\Delta(S)},$$

where U and V are two independent chi-square random variables with $\dim(S) + 1$ and $n - \dim(S) - 1$ degrees of freedom respectively. It is also the solution in x of

$$e^{-\Delta(S)} = (D + 1) \mathbb{P} \left(F_{D+3, N-1} \geq x \frac{N-1}{N(D+3)} \right) - x \frac{N-1}{N} \mathbb{P} \left(F_{D+1, N+1} \geq x \frac{N+1}{N(D+1)} \right),$$

where $D = \dim(S)$, $N = n - \dim(S)$ and $F_{d,r}$ is a Fisher random variable with d and r degrees of freedom.

Proposition 4 in [12] ensures the following upper-bound on $\text{pen}_\Delta(S)$. For any $0 < \kappa < 1$, there exists a constant $C_\kappa > 1$ such that for any $S \in \mathbb{S}$ fulfilling $1 \leq \dim(S) \vee \Delta(S) \leq \kappa n$ we have

$$\text{pen}_\Delta(S) \leq C_\kappa (\dim(S) \vee \Delta(S)).$$

Conversely, Lemma 2.3 in the supplement [38] ensures that $\text{pen}_\Delta(S) \geq 2\Delta(S) + \dim(S) - C$ for some constant $C \geq 0$.

B.2 A General Risk Bound for LinSelect

We set

$$(B.1) \quad \Sigma = \sigma^2 \sum_{S \in \mathbb{S}} e^{-\Delta(S)}.$$

The following proposition gives a risk bound when selecting $\hat{\lambda}$ by minimizing (11).

PROPOSITION B.1. *Assume that $1 \leq \dim(S) \leq n/2 - 1$ and $\Delta(S) \leq 2n/3$ for all $S \in \mathbb{S}$. Then, there exists a constant $C > 1$ such that for any minimizer $\hat{\lambda}$ of criterion (11), we have*

$$(B.2) \quad \begin{aligned} & C^{-1} \mathcal{R}[\hat{\beta}_{\hat{\lambda}}; \beta_0] \\ & \leq \mathbb{E} \left[\inf_{\lambda \in \Lambda} \left\{ \|\mathbf{X}\beta_\lambda - \mathbf{X}\beta_0\|_2^2 \right. \right. \\ & \quad \left. \left. + \inf_{S \in \mathbb{S}} \left\{ \|\mathbf{X}\hat{\beta}_\lambda - \Pi_S \mathbf{X}\hat{\beta}_\lambda\|_2^2 \right. \right. \right. \\ & \quad \left. \left. \left. + [\Delta(S) \vee \dim(S)] \sigma^2 \right\} \right\} \right] + \Sigma. \end{aligned}$$

Furthermore, with probability larger than $1 - e^{-C_0 n} - C_1 \sum_{S \in \mathbb{S}} e^{-C_2 [\Delta(S) \wedge n]} e^{-\Delta(S)}$, we have for some $C > 1$

$$\begin{aligned} & C^{-1} \|\mathbf{X}\beta_0 - \mathbf{X}\hat{\beta}_{\hat{\lambda}}\|_2^2 \\ & \leq \inf_{\lambda \in \Lambda} \left\{ \|\mathbf{X}\hat{\beta}_\lambda - \mathbf{X}\beta_0\|_2^2 \right. \\ & \quad \left. + \inf_{S \in \mathbb{S}} \left\{ \|\mathbf{X}\hat{\beta}_\lambda - \Pi_S \mathbf{X}\hat{\beta}_\lambda\|_2^2 \right. \right. \\ & \quad \left. \left. + [\Delta(S) \vee \dim(S)] \sigma^2 \right\} \right\}. \end{aligned}$$

The first part of Proposition B.1 is a slight variation of Theorem 1 in [13]. We refer to the supplementary material [38] for a proof of these two results.

SUPPLEMENTARY MATERIAL

Supplement to “High-dimensional regression with unknown variance” (DOI: [10.1214/12-STS398SUPP](https://doi.org/10.1214/12-STS398SUPP); .pdf). This supplement contains a description of estimation procedures that are minimax adaptive to the sparsity for all designs \mathbf{X} . It also contains the proofs of the risk bounds for LinSelect and lasso-LinSelect.

REFERENCES

[1] AKAIKE, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory (Tsahkadsor, 1971)* 267–281. Akadémiai Kiadó, Budapest. MR0483125

- [2] ANDERSON, T. W. (1951). Estimating linear restrictions on regression coefficients for multivariate normal distributions. *Ann. Math. Statistics* **22** 327–351. [MR0042664](#)
- [3] ANTONIADIS, A. (2010). Comments on: ℓ_1 -penalization for mixture regression models. *TEST* **19** 257–258. [MR2677723](#)
- [4] ARLLOT, S. and BACH, F. (2009). Data-driven calibration of linear estimators with minimal penalties. In *Advances in Neural Information Processing Systems* **22** (Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams and A. Culotta, eds.) 46–54. Curran Associates, New York.
- [5] ARLLOT, S. and CELISSE, A. (2010). A survey of cross-validation procedures for model selection. *Stat. Surv.* **4** 40–79. [MR2602303](#)
- [6] ARLLOT, S. and CELISSE, A. (2011). Segmentation of the mean of heteroscedastic data via cross-validation. *Stat. Comput.* **21** 613–632. [MR2826696](#)
- [7] ARLLOT, S. and MASSART, P. (2010). Data-driven calibration of penalties for least-squares regression. *J. Mach. Learn. Res.* **10** 245–279.
- [8] BACH, F. R. (2008). Consistency of the group lasso and multiple kernel learning. *J. Mach. Learn. Res.* **9** 1179–1225. [MR2417268](#)
- [9] BACH, F. R. (2008). Consistency of trace norm minimization. *J. Mach. Learn. Res.* **9** 1019–1048. [MR2417263](#)
- [10] BARANIUK, R., DAVENPORT, M., DEVORE, R. and WAKIN, M. (2008). A simple proof of the restricted isometry property for random matrices. *Constr. Approx.* **28** 253–263. [MR2453366](#)
- [11] BARAUD, Y. (2011). Estimator selection with respect to Hellinger-type risks. *Probab. Theory Related Fields* **151** 353–401. [MR2834722](#)
- [12] BARAUD, Y., GIRAUD, C. and HUET, S. (2009). Gaussian model selection with an unknown variance. *Ann. Statist.* **37** 630–672. [MR2502646](#)
- [13] BARAUD, Y., GIRAUD, C. and HUET, S. (2010). Estimator selection in the Gaussian setting. Available at [arXiv:1007.2096v2](#).
- [14] BARRON, A., BIRGÉ, L. and MASSART, P. (1999). Risk bounds for model selection via penalization. *Probab. Theory Related Fields* **113** 301–413. [MR1679028](#)
- [15] BAUDRY, J.-P., MAUGIS, C. and MICHEL, B. (2012). Slope heuristics: Overview and implementation. *Statist. Comput.* **22** 455–470.
- [16] BELLONI, A., CHERNOZHUKOV, V. and WANG, L. (2011). Square-root lasso: Pivotal recovery of sparse signals via conic programming. *Biometrika* **98** 791–806. [MR2860324](#)
- [17] BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Ann. Statist.* **37** 1705–1732. [MR2533469](#)
- [18] BIRGÉ, L. and MASSART, P. (2001). Gaussian model selection. *J. Eur. Math. Soc. (JEMS)* **3** 203–268. [MR1848946](#)
- [19] BIRGÉ, L. and MASSART, P. (2007). Minimal penalties for Gaussian model selection. *Probab. Theory Related Fields* **138** 33–73. [MR2288064](#)
- [20] BUNEA, F., SHE, Y. and WEGKAMP, M. H. (2011). Optimal selection of reduced rank estimators of high-dimensional matrices. *Ann. Statist.* **39** 1282–1309. [MR2816355](#)
- [21] BUNEA, F., TSYBAKOV, A. B. and WEGKAMP, M. H. (2007). Aggregation for Gaussian regression. *Ann. Statist.* **35** 1674–1697. [MR2351101](#)
- [22] CANDES, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *Ann. Statist.* **35** 2313–2351. [MR2382644](#)
- [23] CAO, Y. and GOLUBEV, Y. (2006). On oracle inequalities related to smoothing splines. *Math. Methods Statist.* **15** 398–414. [MR2301659](#)
- [24] CHEN, S. S., DONOHO, D. L. and SAUNDERS, M. A. (1998). Atomic decomposition by basis pursuit. *SIAM J. Sci. Comput.* **20** 33–61. [MR1639094](#)
- [25] COMTE, F. and ROZENHOLC, Y. (2004). A new algorithm for fixed design regression and denoising. *Ann. Inst. Statist. Math.* **56** 449–473. [MR2095013](#)
- [26] DALALYAN, A. and TSYBAKOV, A. (2008). Aggregation by exponential weighting, sharp oracle inequalities and sparsity. *Machine Learning* **72** 39–61.
- [27] DEVROYE, L. P. and WAGNER, T. J. (1979). The L_1 convergence of kernel density estimates. *Ann. Statist.* **7** 1136–1139. [MR0536515](#)
- [28] DONOHO, D. and TANNER, J. (2009). Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philos. Trans. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci.* **367** 4273–4293. [MR2546388](#)
- [29] DONOHO, D. L. (2006). Compressed sensing. *IEEE Trans. Inform. Theory* **52** 1289–1306. [MR2241189](#)
- [30] EFRON, B., HASTIE, T., JOHNSTONE, I. and TIBSHIRANI, R. (2004). Least angle regression. *Ann. Statist.* **32** 407–499. [MR2060166](#)
- [31] FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.* **96** 1348–1360. [MR1946581](#)
- [32] GEISSER, S. (1975). The predictive sample reuse method with applications. *J. Amer. Statist. Assoc.* **70** 320–328.
- [33] GERCHINOVITZ, S. (2011). Sparsity regret bounds for individual sequences in online linear regression. In *Proceedings of COLT 2011*. Microtome Publishing, Brookline, MA.
- [34] GIRAUD, C. (2008). Mixing least-squares estimators when the variance is unknown. *Bernoulli* **14** 1089–1107. [MR2543587](#)
- [35] GIRAUD, C. (2008). Estimation of Gaussian graphs by model selection. *Electron. J. Stat.* **2** 542–563. [MR2417393](#)
- [36] GIRAUD, C. (2011). Low rank multivariate regression. *Electron. J. Stat.* **5** 775–799. [MR2824816](#)
- [37] GIRAUD, C. (2011). A pseudo-RIP for multivariate regression. Available at [arXiv:1106.5599v1](#).
- [38] GIRAUD, C., HUET, S. and VERZELEN, N. (2012). Supplement to “High-dimensional regression with unknown variance.” DOI:[10.1214/12-STS398SUPP](#).
- [39] GIRAUD, C., HUET, S. and VERZELEN, N. (2012). Graph selection with GGMselect. *Stat. Appl. Genet. Mol. Biol.* **11** 1–50.
- [40] GUTHERY, S. B. (1974). A transformation theorem for one-dimensional F -expansions. *J. Number Theory* **6** 201–210. [MR0342484](#)
- [41] HALL, P., KAY, J. W. and TITTERINGTON, D. M. (1990). Asymptotically optimal difference-based estimation of variance in nonparametric regression. *Biometrika* **77** 521–528. [MR1087842](#)
- [42] HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. (2009). *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*, 2nd ed. Springer, New York. [MR2722294](#)

- [43] HUANG, J., MA, S. and ZHANG, C.-H. (2008). Adaptive Lasso for sparse high-dimensional regression models. *Statist. Sinica* **18** 1603–1618. [MR2469326](#)
- [44] HUANG, J. and ZHANG, T. (2010). The benefit of group sparsity. *Ann. Statist.* **38** 1978–2004. [MR2676881](#)
- [45] HUBER, P. J. (1981). *Robust Statistics*. Wiley, New York. [MR0606374](#)
- [46] IZENMAN, A. J. (1975). Reduced-rank regression for the multivariate linear model. *J. Multivariate Anal.* **5** 248–264. [MR0373179](#)
- [47] JI, P. and JIN, J. (2010). UPS delivers optimal phase diagram in high dimensional variable selection. Available at <http://arxiv.org/abs/1010.5028>.
- [48] KLOPP, O. (2011). High dimensional matrix estimation with unknown variance of the noise. Available at [arXiv:1112.3055v1](http://arxiv.org/abs/1112.3055v1).
- [49] KOLTCHINSKII, V., LOUNICI, K. and TSYBAKOV, A. B. (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.* **39** 2302–2329. [MR2906869](#)
- [50] LEBARBIER, E. (2005). Detecting multiple change-points in the mean of Gaussian process by model selection. *Signal Processing* **85** 717–736.
- [51] LENG, C., LIN, Y. and WAHBA, G. (2006). A note on the lasso and related procedures in model selection. *Statist. Sinica* **16** 1273–1284. [MR2327490](#)
- [52] LEUNG, G. and BARRON, A. R. (2006). Information theory and mixing least-squares regressions. *IEEE Trans. Inform. Theory* **52** 3396–3410. [MR2242356](#)
- [53] LI, K.-C. (1987). Asymptotic optimality for C_p , C_L , cross-validation and generalized cross-validation: Discrete index set. *Ann. Statist.* **15** 958–975. [MR0902239](#)
- [54] LOUNICI, K., PONTIL, M., VAN DE GEER, S. and TSYBAKOV, A. B. (2011). Oracle inequalities and optimal inference under group sparsity. *Ann. Statist.* **39** 2164–2204. [MR2893865](#)
- [55] MALLOWS, C. L. (1973). Some comments on C_p . *Technometrics* **15** 661–675.
- [56] MEINSHAUSEN, N. and BÜHLMANN, P. (2006). High-dimensional graphs and variable selection with the lasso. *Ann. Statist.* **34** 1436–1462. [MR2278363](#)
- [57] MOSTELLER, F. and TUKEY, J. W. (1968). Data analysis, including statistics. In *Handbook of Social Psychology, Vol. 2* (G. Lindzey and E. Aronson, eds.). Addison-Wesley, Reading, MA.
- [58] NEGAHBAN, S. and WAINWRIGHT, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Ann. Statist.* **39** 1069–1097. [MR2816348](#)
- [59] NISHII, R. (1984). Asymptotic properties of criteria for selection of variables in multiple regression. *Ann. Statist.* **12** 758–765. [MR0740928](#)
- [60] PARK, T. and CASELLA, G. (2008). The Bayesian lasso. *J. Amer. Statist. Assoc.* **103** 681–686. [MR2524001](#)
- [61] RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2011). Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE Trans. Inform. Theory* **57** 6976–6994. [MR2882274](#)
- [62] RIGOLLET, P. and TSYBAKOV, A. (2011). Exponential screening and optimal rates of sparse estimation. *Ann. Statist.* **39** 731–771. [MR2816337](#)
- [63] ROHDE, A. and TSYBAKOV, A. B. (2011). Estimation of high-dimensional low-rank matrices. *Ann. Statist.* **39** 887–930. [MR2816342](#)
- [64] SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461–464. [MR0468014](#)
- [65] SHAO, J. (1993). Linear model selection by cross-validation. *J. Amer. Statist. Assoc.* **88** 486–494. [MR1224373](#)
- [66] SHAO, J. (1997). An asymptotic theory for linear model selection. *Statist. Sinica* **7** 221–264. With comments and a rejoinder by the author. [MR1466682](#)
- [67] SHIBATA, R. (1981). An optimal selection of regression variables. *Biometrika* **68** 45–54. [MR0614940](#)
- [68] STÄDLER, N., BÜHLMANN, P. and VAN DE GEER, S. (2010). ℓ_1 -penalization for mixture regression models. *TEST* **19** 209–256. [MR2677722](#)
- [69] STONE, M. (1974). Cross-validated choice and assessment of statistical predictions. *J. Roy. Statist. Soc. Ser. B* **36** 111–147. [MR0356377](#)
- [70] SUN, T. and ZHANG, C.-H. (2010). Comments on: ℓ_1 -penalization for mixture regression models. *TEST* **19** 270–275. [MR2677726](#)
- [71] SUN, T. and ZHANG, C. H. (2011). Scaled sparse linear regression. Available at [arXiv:1104.4595](http://arxiv.org/abs/1104.4595).
- [72] TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. [MR1379242](#)
- [73] TIBSHIRANI, R., SAUNDERS, M., ROSSET, S., ZHU, J. and KNIGHT, K. (2005). Sparsity and smoothness via the fused lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **67** 91–108. [MR2136641](#)
- [74] VAN DE GEER, S. A. and BÜHLMANN, P. (2009). On the conditions used to prove oracle results for the Lasso. *Electron. J. Stat.* **3** 1360–1392. [MR2576316](#)
- [75] VERT, J. P. and BLEAKLEY, K. (2010). Fast detection of multiple change-points shared by many signals using group LARS. In *Advances in Neural Information Processing Systems* **23** (J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel and A. Culotta, eds.) 2343–2351. Curran Associates, New York.
- [76] VERZELEN, N. (2010). High-dimensional Gaussian model selection on a Gaussian design. *Ann. Inst. H. Poincaré Probab. Stat.* **46** 480–524. [MR2667707](#)
- [77] VERZELEN, N. (2012). Minimax risks for sparse regressions: Ultra-high-dimensional phenomena. *Electron. J. Stat.* **6** 38–90.
- [78] WAINWRIGHT, M. J. (2009). Information-theoretic limits on sparsity recovery in the high-dimensional and noisy setting. *IEEE Trans. Inform. Theory* **55** 5728–5741. [MR2597190](#)
- [79] YE, F. and ZHANG, C.-H. (2010). Rate minimaxity of the Lasso and Dantzig selector for the ℓ_q loss in ℓ_r balls. *J. Mach. Learn. Res.* **11** 3519–3540. [MR2756192](#)
- [80] YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **68** 49–67. [MR2212574](#)
- [81] ZHANG, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.* **38** 894–942. [MR2604701](#)
- [82] ZHANG, T. (2005). Learning bounds for kernel regression using effective data dimensionality. *Neural Comput.* **17** 2077–2098. [MR2175849](#)

- [83] ZHANG, T. (2011). Adaptive forward–backward greedy algorithm for learning sparse representations. *IEEE Trans. Inform. Theory* **57** 4689–4708. [MR2840485](#)
- [84] ZOU, H. (2006). The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* **101** 1418–1429. [MR2279469](#)
- [85] ZOU, H. and HASTIE, T. (2005). Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **67** 301–320. [MR2137327](#)