



DFG-SNF Research Group FOR916

Statistical Regularization and Qualitative Constraints

Sara van de Geer

Peter Bühlmann

Shuheng Zhou

Prediction and variable selection with the adaptive Lasso

Preprint FOR916 10-05

Preprint-Series of the Research Group FOR916

Prediction and variable selection with the adaptive Lasso

February, 2010

Sara van de Geer, Peter Bühlmann, and Shuheng Zhou
Seminar for Statistics, ETH Zürich

Abstract We revisit the adaptive Lasso in a high-dimensional linear model, and provide bounds for its prediction error and for its number of false positive selections. We compare the adaptive Lasso with an “oracle” that trades off approximation error against an ℓ_0 -penalty. Considering prediction error and false positives simultaneously is a way to study variable selection performance in settings where non-zero regression coefficients can be smaller than the detection limit. We show that an appropriate choice of the tuning parameter yields a prediction error of the same order as that of the least squares refitted initial Lasso after thresholding, while the number of false positives is small, depending on the size of the trimmed harmonic mean of the oracle coefficients.

Keywords: adaptive Lasso, prediction, restricted eigenvalue, thresholding, variable selection

Running Head: Adaptive Lasso

1 Introduction

We consider the linear model

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon,$$

where $\beta \in \mathbb{R}^p$ is a vector of coefficients, \mathbf{X} is an $(n \times p)$ -design matrix, and \mathbf{Y} is an n -vector of noisy observations, ϵ being the noise term. The design matrix \mathbf{X} is treated as fixed. The Gram matrix is $\Sigma := \mathbf{X}^T \mathbf{X} / n$. We assume throughout the normalization $\Sigma_{j,j} = 1$ for all $j \in \{1, \dots, p\}$.

We examine the case $p \geq n$ (i.e., a high-dimensional situation). Regularized estimation with the ℓ_1 -norm penalty, also known as the Lasso (Tibshirani [1996]), refers to the following convex optimization problem:

$$\hat{\beta} := \arg \min_{\beta} \left\{ \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 / n + \lambda \|\beta\|_1 \right\}, \quad (1)$$

where $\lambda > 0$ is a penalization parameter.

Regularization with ℓ_1 -penalization in high-dimensional scenarios has become extremely popular. The methods are easy to use, due to recent progress in specifically tailored convex optimization (Meier et al. [2008], Friedman et al. [2010]).

Consistency results for the prediction error can be found in Greenshtein and Ritov [2004]. The prediction error of the Lasso is asymptotically oracle optimal under some conditions on the design matrix \mathbf{X} , see e.g. van de Geer [2008],

Bickel et al. [2009], Koltchinskii [2009a,b] where also estimation in terms of the ℓ_1 - or ℓ_2 -loss is considered. The “restricted eigenvalue condition” of Bickel et al. [2009] (see also Koltchinskii [2009a,b]) plays a key role here. Candès and Plan [2009] use a different approach, where assumptions on the set of true coefficients allow for a major relaxation of the incoherence condition.

For consistent variable selection with the Lasso, it is known that the so-called “neighborhood stability condition” (Meinshausen and Bühlmann [2006]) for the design matrix, which has been re-formulated in a nicer form as the “irrepresentable condition” (Zhao and Yu [2006]), is sufficient and essentially necessary. Under certain “incoherence conditions”, Wainwright [2007, 2009] analyzes the smallest sample size needed to recover a sparse signal. Because irrepresentable conditions or incoherence conditions are restrictive - they are much stronger than restricted eigenvalue conditions - (see van de Geer and Bühlmann [2009] for a comparison), we conclude that the Lasso for variable selection only works in a rather narrow range of problems, excluding cases where the design exhibits strong (empirical) correlations.

There is also a bias problem with ℓ_1 -penalization, due to the shrinking of the estimates which correspond to true signal variables. A discussion can be found in Zou [2006], Meinshausen [2007]. Regularization with the ℓ_q -“norm” with $q < 1$ would mitigate some of the bias problems but are computationally infeasible as the penalty is non-convex. As an interesting alternative, one can consider multi-step procedures where each of the steps involves a convex optimization only. A prime example is the adaptive Lasso which is a two-step algorithm and whose repeated application corresponds in some “loose” sense to a non-convex penalization scheme (Zou and Li [2008]). The adaptive Lasso was originally proposed by Zou [2006]. He analyzed the case where p is fixed. Further progress in the high-dimensional scenario has been achieved by Huang et al. [2008]. Under a rather strong mutual incoherence condition between every pair of relevant and irrelevant covariables, they prove that the adaptive Lasso recovers the correct model and has an oracle property.

Meinshausen and Yu [2009] examined the variable selection property of the Lasso followed by a thresholding procedure, when all non-zero components are large enough. Under a relaxed incoherence assumption, they show that the estimator is still consistent in the ℓ_2 -norm sense. In addition, they show it is possible to achieve variable selection consistency. Thresholding and multistage procedures are also considered in Candès et al. [2006]. In Zhou [2009, 2010], it is shown that a multi-step thresholding procedure can accurately estimate a sparse vector $\beta \in \mathbb{R}^p$ under the restricted eigenvalue condition of Bickel et al. [2009]. The two-stage procedure in Zhang [2009] applies “selective penalization” in the second stage. This procedure is studied assuming incoherence conditions. A more general framework for multi-stage variable selection was studied by Wasserman and Roeder [2009]. Their approach controls the probability of false positives (type I error) but pays a price in terms of false negatives (type II error).

A key motivation of our work is to continue the exploration of the computa-

tionally tractable adaptive Lasso for variable selection, without requiring the stringent irrerepresentable conditions or incoherence conditions on the design matrix \mathbf{X} . Furthermore, we avoid assumptions saying that the minimal non-zero coefficients β_{true} of the “true” regression are “sufficiently large” since allowing for small non-zero regression coefficients appears to be much more realistic. Consequently, it is impossible to infer the true underlying active set

$$S_{\text{true}} = \{j : \beta_{j,\text{true}} \neq 0\},$$

since co-variables j whose corresponding absolute coefficient $|\beta_{j,\text{true}}|$ is below a detection limit cannot be inferred from data (say with probability tending to 1 as $n \rightarrow \infty$). Thus, we have to tolerate some false negative selections, i.e. variables from S_{true} which are not selected by the statistical estimator.

To study and quantify the variable selection property of an estimator, we consider its prediction error together with its number of false positive selections. The prediction error is closely tied to false negative selections, as for example described in Lemma 5.2 and the remarks following it, and hence, looking at prediction error and false positive selections together translates to performance for variable selection in settings where we do not make any assumptions on the size of $\min_{j \in S_{\text{true}}} |\beta_{j,\text{true}}|$.

2 The adaptive Lasso and its target

The adaptive Lasso is

$$\hat{\beta}_{\text{adap}} := \arg \min_{\beta} \left\{ \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n + \lambda_{\text{init}} \lambda_{\text{adap}} \sum_{j=1}^p \frac{|\beta_j|}{|\hat{\beta}_{j,\text{init}}|} \right\}. \quad (2)$$

Here, $\hat{\beta}_{\text{init}}$ is the one-stage Lasso defined in (1), with initial tuning parameter $\lambda = \lambda_{\text{init}}$, and $\lambda_{\text{adap}} > 0$ is the tuning parameter for the second stage. Note that when $|\hat{\beta}_{j,\text{init}}| = 0$, we exclude variable j in the second stage.

We write $\hat{f}_{\text{init}} := \mathbf{X}\hat{\beta}_{\text{init}}$ and $\hat{f}_{\text{adap}} := \mathbf{X}\hat{\beta}_{\text{adap}}$, with active sets $\hat{S}_{\text{init}} := \{j : \hat{\beta}_{j,\text{init}} \neq 0\}$ and $\hat{S}_{\text{adap}} := \{j : \hat{\beta}_{j,\text{adap}} \neq 0\}$, respectively.

The sparse object to recover may be the “true” unknown parameter $\beta_{\text{true}} \in \mathbb{R}^p$ of the linear regression. More generally, we aim at recovering a sparse approximation of

$$\mathbf{E}\mathbf{Y} := \mathbf{f}^0,$$

when the regression \mathbf{f}^0 itself is not necessarily sparse. As we only consider estimation in a linear model and use penalized least squares loss, we can (by a projection argument) without loss of generality assume that \mathbf{f}^0 is linear, say

$$\mathbf{f}^0 = \mathbf{X}\beta_{\text{true}}.$$

Variable selection with the adaptive Lasso is generally only studied under the assumption that the true underlying signal \mathbf{f}^0 is sparse. It may however well be that many of the $|\beta_{j,\text{true}}|$ are non-zero, but very small. Thus, its active set

$$S_{\text{true}} = \{j : \beta_{j,\text{true}} \neq 0\}$$

may be quite large, and not the set we want to recover.

We believe that an extension to the case where \mathbf{f}^0 is only “approximately” sparse, better reflects the true state of nature, and will highlight the role of assumptions on the size of the coefficients. We then need to decide what we are actually targeting at. There are various sensible possibilities. Our proposal is to target at the approximation of \mathbf{f}^0 that trades off the number of non-zero coefficients against fit. It is defined as follows.

For a set S and for $\beta \in \mathbb{R}^p$, we let

$$\beta_{j,S} := \beta_j \mathbb{1}\{j \in S\}, \quad j = 1, \dots, p.$$

Given a set of indices $S \subset \{1, \dots, p\}$, the best approximation of \mathbf{f}^0 using only variables in S is

$$\mathbf{f}_S = \mathbf{X}b^S := \arg \min_{f=\mathbf{X}\beta_S} \|f - \mathbf{f}^0\|_2.$$

Thus, \mathbf{f}_S is the projection of \mathbf{f}^0 on the span of the variables in S . Our target is now the projection \mathbf{f}_{S_0} , where

$$S_0 := \arg \min_S \left\{ \|\mathbf{f}_S - \mathbf{f}^0\|_2^2/n + 7\lambda_{\text{init}}^2 |S|/\phi^2(6, S) \right\}.$$

Here, $|S|$ denotes the size of S . Moreover, $\phi^2(6, S)$ is a “restricted eigenvalue” (see Section 4 for its definition), and the constants are chosen in relation with the oracle result (see Lemma 6.1). In other words, \mathbf{f}_{S_0} is the optimal ℓ_0 -penalized approximation, albeit that it is discounted by the restricted eigenvalue $\phi^2(6, S_0)$.

We refer to \mathbf{f}_{S_0} as the “oracle”. The set S_0 is called the active set of \mathbf{f}_{S_0} , and $b^0 = b^{S_0}$ are its coefficients, i.e.,

$$\mathbf{f}_{S_0} = \mathbf{X}b^0.$$

The choice of the projection $\mathbf{f}_{S_0} = \mathbf{X}b^0$ as “target” is relatively arbitrary. It can be easily verified that in what follows, one may replace b^0 by any other vector β^0 , provided one also replaces $\|\hat{\beta}_{\text{init}} - b^0\|_q$ by $\|\hat{\beta}_{\text{init}} - \beta^0\|_q$ ($q \geq 1$). The vector b^0 is a natural candidate for β^0 (as $\hat{\beta}_{\text{init}}$ is close to b^0 , see Lemma 5.1 and Corollary 6.1). One may also argue that another natural candidate to target at is $(\beta_{\text{true}})_{S_{\text{true}}^\delta}$, where $\delta > 0$ is some threshold, and $S_{\text{true}}^\delta := \{j : |\beta_{j,\text{true}}| > \delta\}$ (see also Zhou [2010]). Moreover, since the adaptive Lasso will not take up variables that were already abandoned by the initial Lasso, one may also think of first replacing \mathbf{f}^0 by its projection $\mathbf{f}_{\hat{S}_{\text{init}}}$ on the space spanned by variables in \hat{S}_{init} , and then take the best ℓ_0 -penalized approximation of this projection. Of

course, this is then a random target. Nevertheless, with this target the adaptive Lasso will not get blamed for the variables missed by the initial Lasso.

To settle the matter, we will choose b^0 as target in what follows.

We assume that S_0 has a relatively small number $s_0 := |S_0|$ of nonzero coefficients. Inferring the sparsity pattern, i.e. variable selection, refers to the task of correctly estimating the support set S_0 , or more modestly, to have a limited number of false positives (type I errors) and false negatives (type II errors)¹. It can be verified that under reasonable conditions (e.g. i.i.d. Gaussian noise and properly chosen tuning parameter λ) the “ideal” estimator

$$\hat{\beta}_{\text{ideal}} := \arg \min_{\beta} \left\{ \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n + \lambda^2 |\{j : \beta_j \neq 0\}| \right\},$$

has $O(s_0)$ false positives (see for instance Barron et al. [1999] and van de Geer [2001]). With this in mind, we generally aim at $O(s_0)$ false positives (see also Zhou [2010]), yet keeping the prediction error as small as possible.

We show that the two-stage adaptive Lasso can combine a limited number of false positives with good prediction properties (see Theorem 3.1). Invoking these results to prove bounds for the number of false negatives is relatively straightforward.

2.1 Organization of the paper

In the next section, we present some of the main results, in a simplified formulation. Most notation and definitions are given in Section 4. In Section 5, we consider the noiseless case, i.e., the case where $\epsilon = 0$. The reason is that many of the theoretical issues involved concern the approximation properties of the two stage procedure, and not so much the fact that there is noise. By studying the noiseless case first, we separate the approximation problem from the stochastic problem. We first obtain in Subsection 5.1 some simple bounds for the initial Lasso and its thresholded version. In Subsection 5.2 we derive results for the adaptive Lasso by comparing it with a “oracle-thresholded” initial Lasso. When the trimmed harmonic mean of the squared coefficients of the target is large enough, the adaptive Lasso combines good variable selection properties with good prediction properties. Both initial and adaptive Lasso are special cases of a weighted Lasso, which we consider in Subsection 5.3. Subsection 5.4 briefly discusses the weighted irrepresentable condition. Section 6 studies the noisy case. Here, the previous results are summarized and conditions on the tuning parameters are examined. Section 7 concludes. All proofs are in Section 8.

¹For a generic estimator $\hat{\beta}$ and active set S_0 , a type I error is a non-zero estimate $\hat{\beta}_j$ when $j \notin S_0$. A type II error occurs when $\hat{\beta}_j = 0$ while $j \in S_0$.

3 A preview of the results

Clearly, by choosing the tuning parameters λ_{init} and λ_{adap} very large, one can get rid of all coefficients and hence have no false positives. However, the prediction error will then be very large. In practice, the tuning parameters are generally chosen by cross validation. With this in mind, we discuss choices of λ_{adap} which optimize our bounds for the prediction error. This means that λ_{init} and λ_{adap} cannot be chosen arbitrary large.

Let

$$\hat{\delta}_{\text{init}}^2 := \|\mathbf{X}\hat{\beta}_{\text{init}} - \mathbf{f}^0\|_2^2/n,$$

be the prediction error of the initial Lasso, and, for $q > 1$,

$$\hat{\delta}_q := \|\hat{\beta}_{\text{init}} - b^0\|_q$$

be its ℓ_q -error. Denote the prediction error of the adaptive Lasso as

$$\hat{\delta}_{\text{adap}}^2 := \|\mathbf{X}\hat{\beta}_{\text{adap}} - \mathbf{f}^0\|_2^2/n.$$

The behavior of the adaptive Lasso depends on the design, as well as on the true \mathbf{f}^0 , and actually on the interplay between the two. This means that many cases can be distinguished. If $\hat{\delta}_{\infty}$ is small, the adaptive Lasso will generally have good prediction error $\hat{\delta}_{\text{adap}}$. On the other hand, for good variable selection, it will need some assumptions on the size of b^0 . If $\hat{\delta}_{\infty}$ is large, our bounds indicate that $\hat{\delta}_{\text{adap}}$ is large. The tuning parameter λ_{adap} can then be large too, and the variable selection properties may follow without additional assumptions on b^0 .

To present a correct preview, yet keep the exposition simple, we will use order symbols. Our expressions are functions of n , p , \mathbf{X} , and \mathbf{f}^0 , and also of the tuning parameters λ_{init} and λ_{adap} . For positive functions g and h , we say that $g = O(h)$ if $\|g/h\|_{\infty}$ is bounded, and $g \asymp h$ if in addition $\|h/g\|_{\infty}$ is bounded. Our results depend on the restricted eigenvalue $\phi_0 := \phi(6, S_0, 2s_0)$, which we generally think of as being not too small, i.e., $1/\phi_0 = O(1)$. The exact definition of this constant is given in Section 4. Moreover, to simplify the expressions, we assume throughout that $\|f_{S_0} - \mathbf{f}^0\|_2^2/n = O(\lambda_{\text{init}}^2 s_0 / \phi_0^2)$, which roughly says that the oracle ‘‘squared bias’’ term is not substantially larger than the oracle ‘‘variance’’ term. We will furthermore discuss here the case where the noise ϵ is $\mathcal{N}(0, I)$ -distributed, so that we can present an explicit expression for the lower bound on λ_{init} . Extension to other distributions is straightforward, as we show in Section 6. The results may furthermore be improved when the largest eigenvalue of the Gram matrix $\mathbf{X}^T \mathbf{X} / n$ is well-behaved (e.g., when it is $O(1)$, see Section 6).

We define for $\delta > 0$, the set of thresholded coefficients

$$S_0^{\delta} := \{j : |b_j^0| > \delta\}.$$

The next theorem, in particular its result 3), contains main ingredients of the present work. Results 1) and 2) are not new (see e.g. Bunea et al. [2006, 2007a,b], Bickel et al. [2009], Koltchinskii [2009a]), albeit that we replace the perhaps non-sparse β_{true} by the sparser b^0 (see also van de Geer [2008]).

Theorem 3.1 With $\epsilon \sim \mathcal{N}(0, I)$, take $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$, where, for a given $t > 0$,

$$\lambda_{\text{noise}} = 2\sqrt{\frac{2t + 2 \log p}{n}}.$$

Then, with probability at least $1 - 2\exp[-t]$, the following statements hold.

1) There exists a bound $\delta_{\text{init}}^{\text{upper}} = O(\lambda_{\text{init}}\sqrt{s_0}/\phi_0)$ such that

$$\hat{\delta}_{\text{init}} \leq \delta_{\text{init}}^{\text{upper}}.$$

2) For $q \in \{1, 2, \infty\}$, there exists bounds δ_q^{upper} satisfying

$$\delta_1^{\text{upper}} = O(\lambda_{\text{init}}s_0/\phi_0^2), \quad \delta_2^{\text{upper}} = O(\lambda_{\text{init}}\sqrt{s_0}/\phi_0^2), \quad \delta_\infty^{\text{upper}} = O(\lambda_{\text{init}}\sqrt{s_0}/\phi_0^2),$$

such that

$$\hat{\delta}_q \leq \delta_q^{\text{upper}}, \quad q \in \{1, 2, \infty\}.$$

3) Let δ_2^{upper} and $\delta_\infty^{\text{upper}}$ be such bounds, satisfying $\delta_\infty^{\text{upper}} \geq \delta_2^{\text{upper}}/\sqrt{s_0}$, and $\delta_2^{\text{upper}} = O(\lambda_{\text{init}}\sqrt{s_0}/\phi_0^2)$. Let $|b^0|_{\text{harmonic}}^2$ be the trimmed harmonic mean

$$|b^0|_{\text{harmonic}}^2 := \left(\sum_{|b_j^0| > 2\delta_\infty^{\text{upper}}} \frac{1}{|b_j^0|^2} \right)^{-1}.$$

Suppose that

$$\lambda_{\text{adap}}^2 \asymp \left(\frac{1}{n} \left\| \mathbf{f}_{S_0^{4\delta_\infty^{\text{upper}}}} - \mathbf{f}^0 \right\|_2^2 + \lambda_{\text{init}}^2 s_0 / \phi_0^2 \right) \frac{|b^0|_{\text{harmonic}}^2}{\lambda_{\text{init}}^2 / \phi_0^2}. \quad (3)$$

Then

$$\hat{\delta}_{\text{adap}}^2 = O\left(\frac{1}{n} \left\| \mathbf{f}_{S_0^{4\delta_\infty^{\text{upper}}}} - \mathbf{f}^0 \right\|_2^2 + \lambda_{\text{init}}^2 s_0 / \phi_0^2 \right),$$

and

$$|\hat{S}_{\text{adap}} \setminus S_0| = O\left(\frac{\lambda_{\text{init}}^2 s_0}{\phi_0^6} / |b^0|_{\text{harmonic}}^2 \right).$$

Theorem 3.1 is a reformulation of part of Corollary 6.3. According to Theorem 3.1, the larger the trimmed harmonic mean $|b^0|_{\text{harmonic}}^2$, the better the variable selection properties of the adaptive Lasso are. A large value for $\delta_\infty^{\text{upper}}$ will make $|b^0|_{\text{harmonic}}^2$ large, but on the other hand can increase the bound for the prediction error $\hat{\delta}_{\text{adap}}^2$. Note that

$$|b^0|_{\text{harmonic}}^2 \geq 4(\delta_\infty^{\text{upper}})^2 / s_0.$$

This implies that when we take $\delta_\infty^{\text{upper}} \asymp \lambda_{\text{init}}\sqrt{s_0}/\phi_0^2$, then

$$|b^0|_{\text{harmonic}}^{-2} = O(\phi_0^4 / \lambda_{\text{init}}^2),$$

and hence, with large probability,

$$|\hat{S}_{\text{adap}} \setminus S_0| = O(s_0/\phi_0^2).$$

If in fact

$$|b^0|_{\text{harmonic}}^{-2} = O(\phi_0^6/(\lambda_{\text{init}}^2 s_0)), \quad (4)$$

we get that with large probability

$$|\hat{S}_{\text{adap}} \setminus S_0| = O(1).$$

We shall show in Example 5.1 (see also (13)) that without additional assumptions on the Gram matrix $\Sigma = \mathbf{X}^T \mathbf{X}/n$, condition (4) is actually necessary to guarantee the “weighted irrepresentable condition” (see Subsection 5.4 for its definition), the latter ensuring zero false positives. In that sense, our results are tight.

The assumption we make on the tuning parameter λ_{adap} is such that the prediction error $(1/n)\|f_{S_0^{4\delta_\infty^{\text{upper}}}} - \mathbf{f}^0\|_2^2$ and a penalty term are balanced. In this way, we attempt to mimic a choice for λ_{adap} based on cross validation. Thus, (3) is not to be understood as if we assume the expression on the right hand side to be known.

The bound we provide above for $\hat{\delta}_{\text{adap}}$ may be subject to improvement. In fact, we shall show that the threshold $\delta_\infty^{\text{upper}}$ can be replaced by an “oracle” threshold which minimizes (for a given λ_{adap}) bounds for the prediction error (see (8)). The choice of λ_{adap} we then advocate is the one which minimizes the prediction error obtained with the oracle threshold. This refinement is more involved and therefore postponed to Subsection 5.2 (for the noiseless case) and Section 6 (for the noisy case).

Note that Theorem 3.1 allows for a large choice of $\delta_\infty^{\text{upper}}$, larger than a tight bound for $\hat{\delta}_\infty$. However, with such a large choice, the choice (3) for the tuning parameter is also much too large. Thus, a too large threshold will not reflect in any way a choice for λ_{adap} yielding - given the procedure - an optimal prediction error, or mimic a cross validation choice for λ_{adap} . We always may take $\delta_\infty^{\text{upper}} = O(\lambda_{\text{init}}\sqrt{s_0}/\phi_0^2)$. Under incoherence conditions, one may prove that one can take as small as $\delta_\infty^{\text{upper}} = \text{constant} \times \lambda_{\text{init}}$, where the constant depends on the incoherence conditions (see Lounici [2008]). Because the adaptive Lasso inherits properties of the initial Lasso, we conjecture that the “oracle” threshold yielding optimal prediction error will not be much smaller than $\hat{\delta}_\infty$.

The situation is simplified if we assume that the minimal coefficient

$$b_{\min}^0 := \min_{j \in S_0} |b_j^0|$$

is sufficiently large. For example, when

$$b_{\min}^0 > 4\delta_\infty^{\text{upper}},$$

then thresholding at $4\delta_\infty^{\text{upper}}$ will not increase the prediction error. The bound of Theorem 3.1 then coincides with the bound for $\hat{\delta}_{\text{init}}^2$, namely

$$\hat{\delta}_{\text{adap}}^2 = O(\lambda_{\text{init}}^2 s_0 / \phi_0^2).$$

The number of false positives is again $O(s_0 / \phi_0^2)$. If b_{min} is even larger, the prediction error remains of the same order, but the number of false positives decreases, and may even vanish.

4 Notation and definitions

For the noiseless case, it is convenient to formulate the problem in $L_2(Q)$, where Q is a probability measure on some space \mathcal{X} . Let $\{\psi_j\}_{j=1}^p \subset L_2(Q)$ be a given dictionary. The ψ_j will play the role of the co-variables. The Gram matrix is

$$\Sigma := \int \psi^T \psi dQ.$$

We assume that Σ is normalized, i.e., that $\int \psi_j^2 dQ = 1$ for all j . Write a linear function of the ψ_j with coefficients $\beta \in \mathbb{R}^p$ as

$$f_\beta := \sum_{j=1}^p \psi_j \beta_j.$$

The $L_2(Q)$ -norm is denoted by $\|\cdot\|$, so that

$$\|f_\beta\|^2 = \beta^T \Sigma \beta.$$

Recall that for an arbitrary $\beta \in \mathbb{R}^p$, and an arbitrary index set S , we use the notation

$$\beta_{j,S} = \beta_j \mathbb{1}\{j \in S\}.$$

The largest eigenvalue of Σ is denoted by Λ_{max}^2 , i.e.,

$$\Lambda_{\text{max}}^2 := \max_{\|\beta\|_2=1} \beta^T \Sigma \beta.$$

We will also need the largest eigenvalue of submatrices containing the inner products of variables in S :

$$\Lambda_{\text{max}}^2(S) := \max_{\|\beta_S\|_2=1} \beta_S^T \Sigma \beta_S.$$

Its minimal eigenvalue is

$$\Lambda_{\text{min}}^2(S) := \min_{\|\beta_S\|_2=1} \beta_S^T \Sigma \beta_S.$$

4.1 Restricted eigenvalues

A restricted eigenvalue condition is a condition of similar nature as a condition on the minimal eigenvalue of Σ , but with the coefficients β restricted to subsets of \mathbb{R}^p . An overview can be found in van de Geer and Bühlmann [2009].

Define for an index set S , and for $L > 0$, the sets of restrictions

$$\mathcal{R}(L, S) := \{\beta : \|\beta_{S^c}\|_1 \leq L\sqrt{|S|}\|\beta_S\|_2\}.$$

The restricted eigenvalue condition we impose corresponds to the so-called *adaptive* version as introduced in van de Geer and Bühlmann [2009]. It differs from the restricted eigenvalue condition in Bickel et al. [2009] or Koltchinskii [2009a,b]. This is due to the fact that we want to mimic the oracle f_{S_0} , and do not choose \mathbf{f}^0 as target, so that we have to deal with a bias term $\|f_{S_0} - \mathbf{f}^0\|$.

Definition: Restricted eigenvalue. For $N \geq |S|$, we call

$$\phi^2(L, S, N) := \min \left\{ \frac{\|f_\beta\|^2}{\|\beta_{\mathcal{N}}\|_2^2} : \mathcal{N} \supset S, |\mathcal{N}| \leq N, \beta \in \mathcal{R}(L, \mathcal{N}) \right\}$$

the (L, S, N) -restricted eigenvalue. The (L, S, N) -restricted eigenvalue condition holds if $\phi(L, S, N) > 0$. For the case $N = |S|$, we write $\phi(L, S) := \phi(L, S, |S|)$.

It is easy to see that $\phi(L, S) \leq \Lambda_{\min}(S)$ for all $L > 0$. For a given S , our restricted eigenvalue condition is stronger than the one in Bickel et al. [2009] or Koltchinskii [2009a,b]. On the other hand, we apply it to the smaller set S_0 instead of to S_{true} .

Let $f_S := \arg \min_{f=f_{\beta_S}} \|f_{\beta_S} - \mathbf{f}^0\|$ be the projection of \mathbf{f}^0 on the $|S|$ -dimensional linear space spanned by the variables $\{\psi_j\}_{j \in S}$. We denote the coefficients of f_S by b^S , i.e.,

$$f_S = \sum_{j \in S} \psi_j b_j^S = f_{b^S}.$$

5 The noiseless case

Consider a fixed target $\mathbf{f}^0 = f_{\beta_{\text{true}}} \in L_2(Q)$. The initial Lasso is

$$\beta_{\text{init}} := \arg \min_{\beta} \left\{ \|f_\beta - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \|\beta\|_1 \right\}.$$

We assume that the tuning parameter λ_{init} is set at some fixed value. Of course, in the noiseless case, the optimal - in terms of prediction error - value for λ_{init} is $\lambda_{\text{init}} = 0$. However, in the noisy case, a strictly positive lower bound for λ_{init} is dictated by the noise level.

Write

$$f_{\text{init}} := f_{\beta_{\text{init}}}, S_{\text{init}} := \{j : \beta_{j,\text{init}} \neq 0\}, \delta_{\text{init}} := \|f_{\text{init}} - \mathbf{f}^0\|.$$

The adaptive Lasso is

$$\beta_{\text{adap}} := \arg \min_{\beta} \left\{ \|f_{\beta} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{adap}} \sum_{j=1}^p \frac{|\beta_j|}{|\beta_{j,\text{init}}|} \right\}.$$

The second stage tuning parameter λ_{adap} is again assumed to be strictly positive.

Write

$$f_{\text{adap}} := f_{\beta_{\text{adap}}}, \quad S_{\text{adap}} := \{j : \beta_{j,\text{adap}} \neq 0\}, \quad \delta_{\text{adap}} := \|f_{\text{adap}} - \mathbf{f}^0\|.$$

We now fix a set S_0 . This will be the set of variables we want to find. They are intuitively the ones with the largest coefficients $|\beta_{j,\text{true}}|$. However, when we delete the smallest coefficients and refit, the newly refitted coefficients are different. As explained in Section 2, we settle this by taking S_0 as the set obtained by a trading off dimension against fit, namely the set

$$S_0 := \arg \min_S \left\{ \|f_S - \mathbf{f}^0\|^2 + \frac{3\lambda_{\text{init}}^2 |S|}{\phi^2(2, S)} \right\}, \quad (5)$$

where the constants are now from Lemma 5.1. We call f_{S_0} the oracle, and S_0 the oracle active set, and we let $b^0 := b^{S_0}$.

Write

$$\phi_0 := \phi(2, S_0, 2s_0).$$

For simplicity, we assume throughout that

$$\|f_{S_0} - \mathbf{f}^0\|^2 = O(\lambda_{\text{init}}^2 s_0 / \phi_0^2),$$

which roughly says that the approximation error does not overrule the penalty term. To avoid too many details, we use here the restricted eigenvalue ϕ_0^2 instead of $\phi^2(2, S_0)$, because the smaller ϕ_0^2 will occur in the bound for the ℓ_2 -norm $\|\beta_{\text{init}} - b^0\|_2$ (see Lemma 5.1).

5.1 The initial Lasso and its thresholded version

The adaptive Lasso inherits some of its properties from the initial Lasso. In addition, we will derive theory for the adaptive Lasso via the thresholded and refitted initial Lasso. Therefore, we first consider the initial Lasso and thresholding.

Recall that

$$\delta_{\text{init}} := \|f_{\text{init}} - \mathbf{f}^0\|.$$

For $q \geq 1$, we define

$$\delta_q := \|\beta_{\text{init}} - b^0\|_q.$$

Lemma 5.1 *Let*

$$\delta_{\text{oracle}}^2 := \|\mathbf{f}_{S_0} - \mathbf{f}^0\|^2 + \frac{3\lambda_{\text{init}}^2 |S_0|}{\phi_0^2}.$$

We have

$$\delta_{\text{init}}^2 + \lambda_{\text{init}} \|(\beta_{\text{init}})_{S_0^c}\|_1 \leq 2\|\mathbf{f}_{S_0} - \mathbf{f}^0\|^2 + \frac{6\lambda_{\text{init}}^2 |S_0|}{\phi^2(2, S_0)} \leq 2\delta_{\text{oracle}}^2.$$

Moreover

$$\delta_1 \leq 3\|\mathbf{f}_{S_0} - \mathbf{f}^0\|^2 / \lambda_{\text{init}} + \frac{3\lambda_{\text{init}} |S_0|}{\phi^2(2, S_0)} \leq 3\delta_{\text{oracle}}^2 / \lambda_{\text{init}},$$

and

$$\delta_2 \leq 6\delta_{\text{oracle}}^2 / (\lambda_{\text{init}} \sqrt{s_0}).$$

We thus conclude that

$$\delta_{\text{init}}^2 = O(\lambda_{\text{init}}^2 s_0 / \phi_0^2),$$

and

$$\delta_1 = O(\lambda_{\text{init}} s_0 / \phi_0^2), \quad \delta_2 = O(\lambda_{\text{init}} \sqrt{s_0} / \phi_0^2).$$

But then also

$$\delta_{\infty} = O(\delta_2) = O(\lambda_{\text{init}} \sqrt{s_0} / \phi_0^2).$$

The latter can be improved under coherence conditions on the Gram matrix. To simplify the exposition, we will not discuss such improvements in detail (see Lounici [2008]).

We note that in Lemma 5.1, we may replace \mathbf{f}^0 by its projection $\mathbf{f}_{S_{\text{init}}}$ on the space spanned by the variables $\{\psi_j\}_{j \in S_{\text{init}}}$ that are selected by the initial Lasso. We then replace S_0 by the optimal S given in (5) with this newly defined \mathbf{f}^0 . With this replacement, the result follows by using Pythagoras' Theorem. (That is, it follows by a slight adjustment of the proof of Lemma 5.9.) This means that when S_{init} does not have the screening property (i.e., when S_{init} does not contain S_0), one can simply accept this, and see how well one can do with the remaining variables.

We now first look at false negatives and then false positives. Define

$$S_0^\delta = \{j : |b_j^0| > \delta\}.$$

For a generic estimator \tilde{S} of S_0 , we define the set of δ -false negatives as $S_0^\delta \setminus \tilde{S}$. Moreover, we define $\delta_{\infty}^{S_0} := \|(\beta_{\text{init}})_{S_0} - b^0\|_{\infty}$. The following lemma is a trivial application of the triangle inequality.

Lemma 5.2 *There are no $\delta_{\infty}^{S_0}$ -false negatives.*

Because $\delta_{\infty}^{S_0} \leq \|(\beta_{\text{init}})_{S_0} - b^0\|_2 \leq \|\mathbf{f}_{\text{init}} - \mathbf{f}^0\| / \phi_0$, the above lemma shows that once one has a good bound for the prediction error, only small values of b^0 are possibly not detected. We remark that the same reasoning applies to the adaptive Lasso. The results can be refined using the KKT conditions (see Section 5.3 for their definition), but we omit the details.

We now turn to the false positives.

Lemma 5.3 *It holds that*

$$|\mathcal{S}_{\text{init}} \setminus \mathcal{S}_0| \leq 4\Lambda_{\max}^2 \frac{\delta_{\text{init}}^2}{\lambda_{\text{init}}^2}.$$

Hence, the initial estimator has number of false positives

$$|\mathcal{S}_{\text{init}} \setminus \mathcal{S}_0| = \Lambda_{\max}^2 O(s_0/\phi_0^2).$$

In many cases, the eigenvalue Λ_{\max}^2 is quite large; it can even be almost as large as p . Therefore, from the result of Lemma 5.3 one generally cannot deduce good variable selection properties of the initial Lasso.

Next, we consider thresholding. Let for $\delta > 0$,

$$\mathcal{S}_{\text{init}}^\delta := \{j : |\beta_{j,\text{init}}| > \delta\},$$

and

$$f_{\text{init}}^\delta := f_{(\beta_{\text{init}})_{\mathcal{S}_{\text{init}}^\delta}} = \sum_{j \in \mathcal{S}_{\text{init}}^\delta} \psi_j \beta_{j,\text{init}}.$$

Observe that $f_{\mathcal{S}_{\text{init}}^\delta}$ is the refitted estimator after thresholding at δ .

The following lemma presents a bound for the prediction error of the thresholded and refitted initial estimator.

Lemma 5.4 *We have*

$$\|f_{\mathcal{S}_{\text{init}}^\delta} - \mathbf{f}^0\| \leq \|f_{\text{init}}^\delta - \mathbf{f}^0\|,$$

and moreover,

$$\begin{aligned} \|f_{\mathcal{S}_{\text{init}}^\delta} - \mathbf{f}^0\| &\leq \|f_{\mathcal{S}_0^{\delta+\delta_\infty}} - \mathbf{f}^0\| \\ &\leq \|f_{\mathcal{S}_0} - \mathbf{f}^0\| + \Lambda_{\max}(S_0 \setminus \mathcal{S}_0^{\delta+\delta_\infty}) \sqrt{|S_0 \setminus \mathcal{S}_0^{\delta+\delta_\infty}|} (\delta + \delta_\infty). \end{aligned}$$

We know that $\|f_{\mathcal{S}_0} - \mathbf{f}^0\| = O(\lambda_{\text{init}} \sqrt{s_0}/\phi_0)$ and $\delta_\infty = O(\lambda_{\text{init}} \sqrt{s_0}/\phi_0^2)$. Therefore, Lemma 5.4 with $\delta = 3\delta_\infty$ (which is the value that will be used in Corollary 5.1 ahead), gives

$$\begin{aligned} \|\mathbf{f}_{\mathcal{S}_{\text{init}}}^{3\delta_\infty} - \mathbf{f}^0\| &\leq \|f_{\mathcal{S}_0^{4\delta_\infty}} - \mathbf{f}^0\| \\ &= O(\lambda_{\text{init}} \sqrt{s_0}/\phi_0) + 4\Lambda_{\max}(S_0 \setminus \mathcal{S}_0^{4\delta_\infty}) \sqrt{|S_0 \setminus \mathcal{S}_0^{4\delta_\infty}|} \delta_\infty \\ &= O(\lambda_{\text{init}} \sqrt{s_0}/\phi_0) \left(1 + \Lambda_{\max}(S_0 \setminus \mathcal{S}_0^{4\delta_\infty}) \sqrt{|S_0 \setminus \mathcal{S}_0^{4\delta_\infty}|} / \phi_0 \right). \end{aligned} \quad (6)$$

When $b_{\min}^0 := \min_{j \in \mathcal{S}_0} |b_j^0|$ is larger than $4\delta_\infty$, we obviously have $S_0 \setminus \mathcal{S}_0^{4\delta_\infty} = \emptyset$. In that case, the prediction error after thresholding at $3\delta_\infty$ is still of the same order as the oracle bound. If b_{\min}^0 is small, the situation is less clear. The prediction error can then be worse than the oracle bound, and Lemma 5.4 does not tell us whether it improves by taking the threshold δ small, say $\delta_2/\sqrt{s_0} \leq$

$\delta < \delta_\infty$ (with the lower bound for δ being inspired by the comment following Lemma 5.5).

The number of false positives and false negatives of the thresholded initial Lasso is examined in the next lemma.

Lemma 5.5 *The thresholded initial estimator has number of false positives*

$$|S_{\text{init}}^\delta \setminus S_0| \leq \frac{\delta_2^2}{\delta^2}.$$

Moreover, for any $K > 0$, its number of $(K + 1)\delta$ -false negatives is

$$|S_0^{(K+1)\delta} \setminus S_{\text{init}}^\delta| \leq \frac{\delta_2^2}{K^2 \delta^2}.$$

If we take $\delta \geq \delta_2/\sqrt{s_0}$, we get from Lemma 5.5 that

$$|S_{\text{init}}^\delta \setminus S_0| \leq s_0, \tag{7}$$

i.e., then we have at most s_0 false positives after thresholding.

Clearly, the larger the threshold δ the smaller the number of false positives. On the other hand, a large δ may result in a bad prediction error. According to Lemma 5.4, the prediction error $\|\mathbf{f}_{S_{\text{init}}^\delta} - \mathbf{f}^0\|^2$ can be quite large for δ much larger than δ_∞ . With δ in the range $\delta_2/\sqrt{s_0} \leq \delta \leq 3\delta_\infty$, the prediction error is perhaps not very sensitive to the exact value of δ . Looking ahead to the noisy case, cross validation should moreover prefer a larger threshold due to the additional estimation error that occurs if one keeps too many coefficients.

5.2 The adaptive Lasso

Observe that the adaptive Lasso is somewhat more reluctant than thresholding and refitting: the latter ruthlessly disregards all coefficients with $|\beta_{j,\text{init}}| \leq \delta$ (i.e., these coefficients get penalty ∞), and puts zero penalty on coefficients with $|\beta_{j,\text{init}}| > \delta$. The adaptive Lasso gives the coefficients with $|\beta_{j,\text{init}}| \leq \delta$ a penalty of at least $\lambda_{\text{init}}(\lambda_{\text{adap}}/\delta)$ and those with $|\beta_{j,\text{init}}| > \delta$ a penalty of at most $\lambda_{\text{init}}(\lambda_{\text{adap}}/\delta)$.

Recall

$$\delta_{\text{adap}} := \|f_{\text{adap}} - \mathbf{f}^0\|.$$

Lemma 5.6 *We have, for all $\delta \geq \delta_2/\sqrt{s_0}$,*

$$\begin{aligned} & \delta_{\text{adap}}^2 + \lambda_{\text{init}} \lambda_{\text{adap}} \sum_{j \in (S_{\text{init}}^\delta)^c} \frac{|\beta_{j,\text{adap}}|}{|\beta_{j,\text{init}}|} \\ & \leq 2\|\mathbf{f}_{S_{\text{init}}^\delta} - \mathbf{f}^0\|^2 + \frac{6\lambda_{\text{init}}^2}{\phi_0^2} \lambda_{\text{adap}}^2 \sum_{j \in S_{\text{init}}^\delta} \frac{1}{\beta_{j,\text{init}}^2}. \end{aligned}$$

The above lemma is an obstructed oracle inequality, where the oracle is restricted to choose the index set as the set of variables that are left over after removing the smallest $|\beta_{j,\text{init}}|$. If λ_{adap} is chosen small enough, one sees that the prediction error $\|\mathbf{f}_{S_{\text{init}}^\delta} - \mathbf{f}^0\|^2$ of the refitted thresholded initial estimator is not overruled by the penalty term on the right hand side. This means that the prediction error of the adaptive Lasso is not of larger order than the prediction error of the refitted thresholded initial Lasso. Note that we may take $\lambda_{\text{adap}} \geq \delta$, because for $\lambda_{\text{adap}} \leq \delta$, the penalty term in the bound of Lemma 5.6 is not larger than $6\lambda_{\text{init}}^2 s_0 / \phi_0^2$ (see also Lemma 5.7), which - in order of magnitude - is the oracle bound (which cannot be improved).

Lemma 5.6 leads to defining the ‘‘oracle’’ threshold as

$$\delta_0 := \arg \min_{\delta \geq \delta_2 / \sqrt{s_0}} \left\{ \|\mathbf{f}_{S_{\text{init}}^\delta} - \mathbf{f}^0\|^2 + \frac{3\lambda_{\text{init}}^2}{\phi_0^2} \lambda_{\text{adap}}^2 \sum_{j \in S_{\text{init}}^\delta} \frac{1}{\beta_{j,\text{init}}^2} \right\}. \quad (8)$$

This oracle has active set $S_{\text{init}}^{\delta_0}$, with size $|S_{\text{init}}^{\delta_0}| = O(s_0)$. In what follows however, we will mainly choose $\delta = 3\delta_\infty$. Thus, our bounds are good when the oracle threshold δ_0 is not too different from $3\delta_\infty$. If in the range $\delta_2 / \sqrt{s_0} \leq \delta \leq 3\delta_\infty$ the prediction error $\|\mathbf{f}_{S_{\text{init}}^\delta} - \mathbf{f}^0\|$ is roughly constant in δ , the oracle threshold will at least be not much smaller than $3\delta_\infty$. When the oracle threshold is larger than this, it is straightforward to reformulate the situation. This, we have omitted to avoid too many cases.

Some further results for the prediction error δ_{adap} follow by inserting bounds for the initial Lasso.

Lemma 5.7 *It holds that*

$$\sum_{j \in S_{\text{init}}^\delta} \frac{1}{\beta_{j,\text{init}}^2} \leq \frac{1}{\delta^2} \left\{ \left| \{j : \delta - \delta_\infty < |b_j^0| \leq 2\delta_\infty\} \right| + 4\delta^2 |b^0|_{\text{harmonic}}^{-2} \right\}.$$

Moreover, for $\delta \geq \delta_2 / \sqrt{s_0}$,

$$\sum_{j \in S_{\text{init}}^\delta} \frac{1}{\beta_{j,\text{init}}^2} \leq \frac{2s_0}{\delta^2}.$$

Corollary 5.1 *With the special choice $\delta = 3\delta_\infty$, we get*

$$\sum_{j \in S_{\text{init}}^{3\delta_\infty}} \frac{1}{\beta_{j,\text{init}}^2} \leq 4|b^0|_{\text{harmonic}}^{-2}.$$

Corollary 5.2 *Using the bound of Lemma 5.7 in Lemma 5.6 gives that for all $\delta \geq \delta_2 / \sqrt{s_0}$,*

$$\begin{aligned} \delta_{\text{adap}}^2 &\leq 2\|\mathbf{f}_{S_{\text{init}}^\delta} - \mathbf{f}^0\|^2 \\ &+ \frac{6\lambda_{\text{init}}^2}{\phi_0^2} \frac{\lambda_{\text{adap}}^2}{\delta^2} \left\{ \left| \{j : \delta - \delta_\infty < |b_j^0| \leq 2\delta_\infty\} \right| + 4\delta^2 |b^0|_{\text{harmonic}}^{-2} \right\}. \end{aligned}$$

If $\delta_2/\sqrt{s_0} \leq 3\delta_\infty$, we may choose $\delta = 3\delta_\infty$ to find

$$\delta_{\text{adap}}^2 \leq 2\|f_{S_{\text{init}}^{3\delta_\infty}} - \mathbf{f}^0\|^2 + \frac{24\lambda_{\text{init}}^2}{\phi_0^2}\lambda_{\text{adap}}^2|b^0|_{\text{harmonic}}^{-2}. \quad (9)$$

According to Lemma 5.4,

$$\begin{aligned} & \|f_{S_{\text{init}}^{3\delta_\infty}} - \mathbf{f}^0\| \leq \|f_{S_0^{4\delta_\infty}} - \mathbf{f}^0\| \\ & = O(\lambda_{\text{init}}\sqrt{s_0}/\phi_0) \left(1 + \Lambda_{\max}(S_0 \setminus S_0^{4\delta_\infty}) \sqrt{|S_0 \setminus S_0^{4\delta_\infty}|} \delta_\infty \right). \end{aligned}$$

Inserting these yields

$$\begin{aligned} \delta_{\text{adap}}^2 & \leq 2\|f_{S_0^{4\delta_\infty}} - \mathbf{f}^0\|^2 + \frac{24\lambda_{\text{init}}^2}{\phi_0^2}\lambda_{\text{adap}}^2|b^0|_{\text{harmonic}}^{-2} \\ & = O\left(\frac{\lambda_{\text{init}}^2 s_0}{\phi_0^2}\right) \left(1 + \Lambda_{\max}^2(S_0 \setminus S_0^{4\delta_\infty}) |S_0 \setminus S_0^{4\delta_\infty}| \delta_\infty^2 \right) \\ & \quad + \frac{24\lambda_{\text{init}}^2}{\phi_0^2}\lambda_{\text{adap}}^2|b^0|_{\text{harmonic}}^{-2}. \end{aligned} \quad (10)$$

We proceed by considering number of false positives of the adaptive Lasso.

Lemma 5.8 *We have*

$$|S_{\text{adap}} \setminus S_0| \leq 4 \frac{\delta_{\text{adap}}^2}{\lambda_{\text{adap}}^2} \frac{\delta_2^2}{\lambda_{\text{init}}^2} \wedge 2\Lambda_{\max} \frac{\delta_{\text{adap}}}{\lambda_{\text{adap}}} \frac{\delta_2}{\lambda_{\text{init}}}.$$

We will choose $\lambda_{\text{adap}} \geq 3\delta_\infty$ in such a way that that the prediction error and the penalty term in (9) are balanced, so that

$$\frac{\delta_{\text{adap}}^2}{\lambda_{\text{adap}}^2} = O\left(\frac{\lambda_{\text{init}}^2}{\phi_0^2} \sum_{|b_j^0| > 2\delta_\infty} \frac{1}{|b_j^0|^2}\right). \quad (11)$$

Let us summarize the consequences of this choice in a corollary.

Corollary 5.3: Main result for the noiseless case

We take the choice for λ_{adap} given by (11).

a) It then holds that

$$|S_{\text{adap}} \setminus S_0| = O\left(\frac{\lambda_{\text{init}}^2 s_0}{\phi_0^4} |b_0|_{\text{harmonic}}^{-2} \wedge \Lambda_{\max} \frac{\lambda_{\text{init}} \sqrt{s_0}}{\phi_0^2} |b_0|_{\text{harmonic}}^{-1}\right).$$

- When $\lambda_{\text{init}}^2 |b_0|_{\text{harmonic}}^{-2} / \phi_0^2 = O(1)$, we get $|S_{\text{adap}} \setminus S_0| = O(s_0 / \phi_0^2)$.

- When also $\Lambda_{\max} = O(1)$, we get $|S_{\text{adap}} \setminus S_0| = O(\sqrt{s_0} / \phi_0)$.

- With $\lambda_{\text{init}}^2 s_0 |b_0|_{\text{harmonic}}^{-2} / \phi_0^4 = O(1)$, we get $|S_{\text{adap}} \setminus S_0| = O(1)$ (or even $|S_{\text{adap}} \setminus S_0| = 0$ if the constants are small enough). This corresponds with the bound of Corollary 5.4, implying the weighted irrepresentable condition defined in Subsection 5.4, a bound which, according to Example 5.1, cannot be improved.

b) We know that $\delta_\infty = O(\delta_2) = O(\lambda_{\text{init}} \sqrt{s_0} / \phi_0^2)$. Suppose now that this cannot be improved, i.e., that also

$$\lambda_{\text{init}} \sqrt{s_0} / \phi_0^2 = O(\delta_\infty).$$

Then we get

$$|b_0|_{\text{harmonic}}^{-2} = O\left(\frac{\phi_0^4}{\lambda_{\text{init}}^2}\right),$$

Hence, when the convergence in sup-norm is slow, we can get relatively few false positives, but possibly a not-so-good prediction error. When $\delta_\infty \geq \delta_2 / \sqrt{s_0}$ is small, the bound

$$\sum_{|b_j^0| > \delta_\infty} \frac{1}{|b_j^0|^2} \leq \sum_{|b_j^0| > \delta_2 / \sqrt{s_0}} \frac{1}{|b_j^0|^2}$$

may be appropriate. Assuming this to be $O(\phi_0^4 / \lambda_{\text{init}}^2)$ amounts to assuming that “on average”, the coefficients b_j^0 are “not too small”. For example, it is allowed that $O(1)$ coefficients are as small as $\lambda_{\text{init}} / \phi_0^2$.

c) Suppose now that

$$b_{\min}^0 := \min_{j \in S_0} |b_j^0| \geq \lambda_{\text{init}} \sqrt{s_0} / \phi_0^2.$$

Then

$$\sum_{|b_j^0| > 2\delta_\infty} \frac{1}{|b_j^0|^2} \leq \phi_0^4 / \lambda_{\text{init}}^2.$$

With this (or larger) values for b_{\min}^0 , we also see that the refitted thresholded estimator $\hat{f}_{S_{\text{init}}}^{3\delta_\infty}$ has prediction error $O(\delta_{\text{init}}^2) = O(\lambda_{\text{init}}^2 s_0 / \phi_0^2)$. If

$$b_{\min}^0 \geq \lambda_{\text{init}} s_0 / \phi_0^3,$$

we in fact only have $O(1)$ false positives.

d) More generally, in view of Lemma 5.4, the prediction error of the adaptive Lasso can be bounded by

$$\delta_{\text{adap}}^2 = O\left(\frac{\lambda_{\text{init}}^2 s_0}{\phi_0^2}\right) \left(1 + \Lambda_{\max}^2(S_0 \setminus S_0^{4\delta_\infty}) |S_0 \setminus S_0^{4\delta_\infty}| \delta_\infty^2\right).$$

Thus, our theory shows similar bounds for the adaptive Lasso and thresholding, in terms of prediction error and variable selection.

5.3 The weighted Lasso

As the initial and adaptive Lasso are special cases of the weighted Lasso, some of the results in Subsections 5.1 and 5.2 are consequences of those in this subsection.

The weighted Lasso is

$$\beta_{\text{weight}} := \arg \min_{\beta} \left\{ \|f_{\beta} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j=1}^p w_j |\beta_j| \right\},$$

where the $\{w_j\}_{j=1}^p$ are non-negative weights.

We set $f_{\text{weight}} := f_{\beta_{\text{weight}}}$. Moreover, we define

$$\|w_S\|_2^2 := \sum_{j \in S} w_j^2, \quad w_{S^c}^{\min} := \min_{j \notin S} w_j.$$

By the reparametrization $\beta \mapsto \gamma := W\beta$, where $W = \text{diag}(w_1, \dots, w_p)$, one sees that the weighted Lasso is a standard Lasso with Gram matrix

$$\Sigma_{\text{weight}} := W^{-1} \Sigma W^{-1}.$$

We emphasize however that Σ_{weight} is generally not normalized, i.e., generally $\text{diag}(\Sigma_{\text{weight}}) \neq I$.

We first present a bound for the prediction error and then consider variable selection.

Lemma 5.9 *For all S satisfying $\|w_S\|_2/w_{S^c}^{\min} \leq L\sqrt{|S|}$, and all β , we have*

$$\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j,\text{weight}}| \leq 2\|f_{\beta_S} - \mathbf{f}^0\|^2 + \frac{6\lambda_{\text{init}}^2 \lambda_{\text{weight}}^2}{\phi^2(2L, S)} \|w_S\|_2^2,$$

and

$$\begin{aligned} & \lambda_{\text{init}} \lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j,\text{weight}}| \\ & \leq 3\|f_{\beta_S} - \mathbf{f}^0\|^2 + \frac{3\lambda_{\text{init}}^2 \lambda_{\text{weight}}^2}{\phi^2(2L, S)} \|w_S\|_2^2. \end{aligned}$$

Our next theme is variable selection. An important characterization of the solution β_{weight} can be derived from the *Karush-Kuhn-Tucker (KKT)* conditions (see Bertsimas and Tsitsiklis [1997]).

Weighted KKT-conditions *We have*

$$2\Sigma(\beta_{\text{weight}} - \beta_{\text{true}}) = -\lambda_w W \tau_{\text{weight}}.$$

Here, $\|\tau_{\text{weight}}\|_{\infty} \leq 1$, and moreover

$$\tau_{j,\text{weight}} \mathbb{1}\{\beta_{j,\text{weight}} \neq 0\} = \text{sign}(\beta_{j,\text{weight}}), \quad j = 1, \dots, p.$$

The KKT-conditions can be invoked to derive the next lemma, where we use the notation

$$\|(1/w)_S\|_2^2 := \sum_{j \in S} \frac{1}{w_j^2}.$$

Lemma 5.10 *We have*

$$\begin{aligned} |S_{\text{weight}} \setminus S_0| &\leq 4 \frac{\|f_{\text{weight}} - \mathbf{f}^0\|^2}{\lambda_{\text{weight}}^2} \frac{\|(1/w)_{S_{\text{weight}} \setminus S_0}\|_2^2}{\lambda_{\text{init}}^2} \\ &\wedge 2\Lambda_{\max} \frac{\|f_{\text{weight}} - \mathbf{f}^0\|}{\lambda_{\text{weight}}} \frac{\|(1/w)_{S_{\text{weight}} \setminus S_0}\|_2}{\lambda_{\text{init}}}. \end{aligned}$$

5.4 The weighted irrepresentable condition

It is known that the initial Lasso essentially needs the irrepresentable condition in order to have no false positives (Zhao and Yu [2006]). Similar statements can be made for the weighted Lasso. In Corollary 5.2, we showed that under certain conditions on the trimmed harmonic mean, the adaptive Lasso has no false positives. This subsection links this to the weighted irrepresentable condition.

For a $(p \times p)$ -matrix $\Sigma = (\sigma_{j,k})$. we define

$$\Sigma_{1,1}(S) := (\sigma_{j,k})_{j,k \in S},$$

$$\Sigma_{2,1}(S) := (\sigma_{j,k})_{j \notin S, k \in S}.$$

We let $W_S := \text{diag}(\{w_j\}_{j \in S})$.

Definition *We say that the weighted irrepresentable condition holds for S if for all vectors $\tau_S \in \mathbb{R}^{|S|}$ with $\|\tau_S\|_{\infty} \leq 1$, one has*

$$\|W_{S^c}^{-1} \Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S\|_{\infty} < 1.$$

The reparametrization $\beta \mapsto \gamma := W^{-1}\beta$ leads to the following lemma, which is the weighted variant of the first part of Lemma 6.2 in van de Geer and Bühlmann [2009]. Here, we actually take \mathbf{f}_0 as target, instead of its ℓ_0 -sparse approximation f_{S_0} . Recall

$$S_{\text{true}} := \{j : \beta_{j,\text{true}} \neq 0\}.$$

Lemma 5.11

Suppose the weighted irrepresentable condition is met for S_{true} . Then $S_{\text{weight}} \subset S_{\text{true}}$.

We now consider conditions for the weighted irrepresentable condition to hold. We recall that $w_{S^c}^{\min} := \min_{j \notin S} w_j$, and define, as in van de Geer and Bühlmann [2009], the *adaptive restricted regression*

$$\vartheta_{\text{adaptive}}(S) := \max_{\beta \in \mathcal{R}(1, S)} \frac{|(f_{\beta_{S^c}}, f_{\beta_S})|}{\|f_{\beta_S}\|^2}.$$

Here, (f, \tilde{f}) denotes the inner product between f and \tilde{f} as elements of $L_2(Q)$.

Lemma 5.12

$$\sup_{\|\tau_S\|_{\infty} \leq 1} \|W_{S^c}^{-1} \Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S\|_{\infty} \leq \frac{\|w_S\|_2}{\sqrt{|S|} w_{S^c}^{\min}} \vartheta_{\text{adaptive}}(S).$$

It is not difficult to see that $\vartheta_{\text{adaptive}}(S) \leq \sqrt{|S|}/\Lambda_{\min}(S)$. Hence, we arrive at the following corollary.

Corollary 5.4 *Suppose that*

$$\|w_S\|_2 < \Lambda_{\min}(S) w_{S^c}^{\min}. \quad (12)$$

Then the weighted irrepresentable condition holds for S .

With, for $j \in S_{\text{true}}$, the “ideal” weights

$$w_j = 1/|\beta_{j,\text{true}}|, \quad j \in S_{\text{true}},$$

and with furthermore $w_{S_{\text{true}}}^{\min} = 1/\delta_{\infty}$, where (say) $\delta_{\infty} = \lambda_{\text{init}} \sqrt{s}/\phi_0^2$, and $s_{\text{true}} = |S_{\text{true}}|$, inequality (12) reads

$$\sum_{j \in S_{\text{true}}} \frac{1}{|\beta_{j,\text{true}}|^2} \leq \frac{\Lambda_{\min}^2(S_{\text{true}}) \phi_0^4}{\lambda_{\text{init}}^2 s_{\text{true}}}. \quad (13)$$

This corresponds to the third case in Corollary 5.3 a). The next example shows that there exist Gram matrices Σ which have smallest eigenvalue $1 - \rho$, $\rho \in (0, 1)$, and where the weighted irrepresentable condition needs the separation

$$\|w_{S_{\text{true}}}\|_2 \leq w_{S_{\text{true}}}^{\min} / \rho.$$

Example 5.1 *Let $S_{\text{true}} = \{1, \dots, s\}$, with cardinality $s := |S_{\text{true}}|$, be the active set, and suppose that*

$$\Sigma := \begin{pmatrix} \Sigma_{1,1} & \Sigma_{1,2} \\ \Sigma_{2,1} & \Sigma_{2,2} \end{pmatrix},$$

where $\Sigma_{1,1} := I$ is the $(s \times s)$ -identity matrix, and

$$\Sigma_{2,1} := \rho(c_2 c_1^T),$$

with $0 \leq \rho < 1$, and with c_1 an s -vector and c_2 a $(p - s)$ -vector, satisfying $\|c_1\|_2 = \|c_2\|_2 = 1$. Moreover, suppose $\Sigma_{2,2}$ is the $((p - s) \times (p - s))$ -identity matrix. The smallest eigenvalue of Σ is $1 - \rho$. Take $c_1 = w_{S_{\text{true}}}/\|w_{S_{\text{true}}}\|_2$, and $c_2 = (0, \dots, 1, 0, \dots)^T$, where the 1 is placed at $\arg \min_{j \in S_{\text{true}}^c} w_j$. Then

$$\sup_{\|\tau_{S_{\text{true}}}\|_{\infty} \leq 1} \|W_{S_{\text{true}}}^{-1} \Sigma_{2,1} \Sigma_{1,1}^{-1} W_{S_{\text{true}}} \tau_{S_{\text{true}}}\|_{\infty} = \rho \|w_{S_{\text{true}}}\|_2 / w_{S_{\text{true}}}^{\min}.$$

6 Adding noise

Consider an n -dimensional vector of observations

$$\mathbf{Y} = \mathbf{f}^0 + \epsilon,$$

and the weighted (noisy) Lasso

$$\hat{\beta}_{\text{weight}} = \arg \min_{\beta} \left\{ \|\mathbf{Y} - f_{\beta}\|_2^2/n + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j=1}^p w_j |\beta_j| \right\}. \quad (14)$$

Here, \mathbf{f}^0 , the dictionary $\{\psi_j\}$, and $f_{\beta} := \sum \psi_j \beta_j$ are now considered as vectors in \mathbb{R}^n . The norm we use is the normalized Euclidean norm

$$\|f\| := \|f\|_n := \|f\|_2/\sqrt{n} : f \in \mathbf{R}^n,$$

induced by the inner product

$$(f, \tilde{f})_n := \frac{1}{n} \sum_{i=1}^n f_i \tilde{f}_i, \quad f, \tilde{f} \in \mathbb{R}^n.$$

We define the projections f_S , in the same way as in the previous section. The ℓ_0 -sparse projection $f_{S_0} = \sum_{j \in S_0} b_j^0$, is now defined with a larger constant (7 instead of 3) in front of the penalty term, and a larger constant ($L = 6$ instead of $L = 2$) in the restrictions $\mathcal{R}(L, S)$ of the restricted eigenvalue condition:

$$S_0 := \arg \min_S \left\{ \|f_S - \mathbf{f}^0\|_n^2 + \frac{7\lambda_{\text{init}}^2 |S|}{\phi^2(6, S)} \right\}.$$

We also change the constant ϕ_0 accordingly:

$$\phi_0 := \phi(6, S_0, 2s_0).$$

Let

$$\hat{f}_{\text{weight}} := f_{\hat{\beta}_{\text{weight}}}, \quad \hat{S}_{\text{weight}} := \{j : \hat{\beta}_{j, \text{weight}} \neq 0\}.$$

We define the estimators \hat{f}_{init} and \hat{f}_{adap} as in Section 2, with active sets \hat{S}_{init} and \hat{S}_{adap} . The unpenalized least squares estimator using the variables in S is

$$\hat{f}_S = f_{\hat{b}_S} := \arg \min_{f=f_{\beta_S}} \|\mathbf{Y} - f_{\beta_S}\|_n.$$

We define as before, for $\delta > 0$,

$$\hat{S}_{\text{init}}^{\delta} := \{j : |\hat{\beta}_{j, \text{init}}| > \delta\}, \quad S_0^{\delta} := \{j : |b_j^0| > \delta\}.$$

The refitted version after thresholding, based on the data \mathbf{Y} , is $\hat{f}_{\hat{S}_{\text{init}}^{\delta}}$.

We let

$$\hat{\delta}_{\text{init}} := \|\hat{f}_{\text{init}} - \mathbf{f}^0\|_n, \quad \hat{\delta}_{\text{adap}} := \|\hat{f}_{\text{adap}} - \mathbf{f}^0\|_n,$$

and moreover, for $q \geq 1$,

$$\hat{\delta}_q := \|\hat{\beta}_{\text{init}} - b^0\|_q.$$

To handle the (random) noise, we define the set

$$\mathcal{T} := \left\{ \max_{1 \leq j \leq p} 2|(\epsilon, \psi_j)_n| \leq \lambda_{\text{noise}} \right\},$$

where λ_{noise} is chosen in such a way that

$$\mathbb{P}(\mathcal{T}) \geq 1 - \alpha$$

where $(1 - \alpha)$ is the confidence we want to achieve. In other words, $\lambda_{\text{noise}}/2$ is a bound for the maximal sample correlation between the noise and the variables ψ_j , which holds with large probability. Using the probability bound $\mathbb{P}(|Z| \geq \sqrt{2t}) \leq 2 \exp[-t]$ for a standard normal random variable Z , one can easily derive that when $\epsilon \sim \mathcal{N}(0, I)$, and with

$$\lambda_{\text{noise}} = 2\sqrt{\frac{2t + 2 \log p}{n}},$$

one has

$$\mathbb{P}(\mathcal{T}) \geq 1 - 2 \exp[-t].$$

We will first formulate the noisy weighted Lasso. Once this is done, results for the initial Lasso, its thresholded version, and for the adaptive Lasso, follow in the same way as in Subsections 5.1 and 5.2. Therefore, we do not repeat all derivations in detail. The main point is to take care that the tuning parameters are chosen in such a way that the noisy part due to variables in S_0^c are overruled by the penalty term. In our situation, this can be done by taking $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$, and λ_{adap} large enough. A lower bound for λ_{adap} depends on the behavior of the initial estimator (see Corollary 6.2). In Corollary 6.3, we let λ_{adap} depend on λ_{init} , s_0 and ϕ_0 , on a bound $\delta_\infty^{\text{upper}}$ for $\hat{\delta}_\infty$, on the prediction error of $f_{S_0^{A\delta_\infty^{\text{upper}}}}$, and on the trimmed harmonic mean of the $|b_j^0|^2$.

After presenting the noisy versions of Lemma 5.9 and Lemma 5.10 (their proof is a straightforward adjustment of the noiseless case, and hence omitted), we give a result for the least squares estimator using only the variables j with large enough $|\hat{\beta}_{j,\text{init}}|$. We then present the corollaries for the noisy initial and noisy adaptive Lasso, as regards prediction error and variable selection. These corollaries have “random” quantities in the bounds. We end with a corollary containing the main result for the noisy case, where the random bounds are replaced by fixed ones, and where we moreover choose a more specific lower bound for λ_{adap} .

Lemma 6.1 *Suppose we are on \mathcal{T} . Let $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$. Let S satisfy*

$$\|w_S\|_2/w_{S^c}^{\min} \leq L\sqrt{|S|}, \quad L \geq 1,$$

and

$$\lambda_{\text{weight}}\|w_S\|_2 \geq \sqrt{|S|}.$$

For all β we have

$$\begin{aligned} & \|\hat{f}_{\text{weight}} - \mathbf{f}^0\|_n^2 + \lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\hat{\beta}_{j,\text{weight}}|/2 \\ & \leq 2\|f_{\beta_S} - \mathbf{f}^0\|_n^2 + \frac{14\lambda_{\text{init}}^2\lambda_{\text{weight}}^2}{\phi^2(6L, S)}\|w_S\|_2^2, \end{aligned}$$

and

$$\begin{aligned} & \lambda_{\text{init}}\lambda_{\text{weight}}\|w_S\|_2\|(\beta_{\text{weight}})_S - \beta_S\|_2 + \lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j,\text{weight}}| \\ & \leq 5\|f_{\beta_S} - \mathbf{f}^0\|_n^2 + \frac{7\lambda_{\text{init}}^2\lambda_{\text{weight}}^2}{\phi^2(6L, S)}\|w_S\|_2^2. \end{aligned}$$

Lemma 6.2 *Suppose we are on \mathcal{T} . Assume $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$ and $\lambda_{\text{weight}}w_{S_0}^{\min} \geq 1$. We have*

$$\begin{aligned} |\hat{S}_{\text{weight}} \setminus S_0| & \leq 16 \frac{\|\hat{f}_{\text{weight}} - \mathbf{f}^0\|_n^2}{\lambda_{\text{weight}}^2} \frac{\|(1/w)\hat{S}_{\text{weight}} \setminus S_0\|_2^2}{\lambda_{\text{init}}^2} \\ & \wedge 4\Lambda_{\max} \frac{\|\hat{f}_{\text{weight}} - \mathbf{f}^0\|_n}{\lambda_{\text{weight}}} \frac{\|(1/w)\hat{S}_{\text{weight}} \setminus S_0\|_2}{\lambda_{\text{init}}}. \end{aligned}$$

The least squares estimator $\hat{f}_{\hat{S}_{\text{init}}^\delta}$ using only variables in $\hat{S}_{\text{init}}^\delta$ (i.e., the projection of $\mathbf{Y} = \mathbf{f}^0 + \epsilon$ on the linear space spanned by $\{\psi_j\}_{j \in \hat{S}_{\text{init}}^\delta}$) has similar prediction properties as $f_{\hat{S}_{\text{init}}^\delta}$ (the projection of \mathbf{f}^0 on the same linear space). This is because, as is shown in the next lemma, their difference is small.

Lemma 6.3 *Let $\delta \geq \hat{\delta}_2/\sqrt{s_0}$. Then*

$$\|\hat{f}_{\hat{S}_{\text{init}}^\delta} - f_{\hat{S}_{\text{init}}^\delta}\|_n^2 \leq \frac{2\lambda_{\text{noise}}^2 s_0}{\phi_0^2}.$$

Finally², we present two corollaries, one treating prediction error and variable selection of the initial Lasso, the other one prediction error and variable selection of the adaptive Lasso. The consequences of these two corollaries, presented in Corollary 6.3, give qualitatively the same conclusion as in the noiseless case.

²Of separate interest is a direct comparison of the noisy initial Lasso with the noisy ℓ_0 -penalized estimator. Replacing \mathbf{f}^0 by \mathbf{Y} in Lemma 5.1 gives

$$\|\mathbf{Y} - \hat{f}_{\text{init}}\|_n^2 \leq 2 \min_S \left\{ \|\mathbf{Y} - \hat{f}_S\|_n^2 + \frac{3\lambda_{\text{init}}^2|S|}{\phi^2(2, S)} \right\}.$$

Corollary 6.1 *Let*

$$\delta_{\text{oracle}}^2 := \|\mathbf{f}_{S_0} - \mathbf{f}^0\|_n^2 + \frac{7\lambda_{\text{init}}^2 |S_0|}{\phi^2(6, S_0, 2s_0)}.$$

Take $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$. We have on \mathcal{T} ,

$$\hat{\delta}_{\text{init}}^2 \leq 2\delta_{\text{oracle}}^2.$$

Moreover, on \mathcal{T} ,

$$\hat{\delta}_1 \leq 5\delta_{\text{oracle}}^2/\lambda_{\text{init}},$$

and

$$\hat{\delta}_2 \leq 10\delta_{\text{oracle}}^2/(\lambda_{\text{init}}\sqrt{s_0}).$$

Also, on \mathcal{T} ,

$$|\hat{S}_{\text{init}} \setminus S_0| \leq 16\Lambda_{\text{max}}^2 \frac{\hat{\delta}_{\text{init}}^2}{\lambda_{\text{init}}^2}.$$

Corollary 6.2 *Suppose we are on \mathcal{T} . Take $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$ and $\delta \geq \hat{\delta}_2/\sqrt{s_0}$. Let*

$$\lambda_{\text{adap}}^2 \sum_{j \in \hat{S}_{\text{init}}^\delta} \frac{1}{\hat{\beta}_{j,\text{init}}^2} \geq |\hat{S}_{\text{init}}^\delta|.$$

Then

$$\begin{aligned} & \hat{\delta}_{\text{adap}}^2 + \frac{1}{2}\lambda_{\text{init}}\lambda_{\text{adap}} \sum_{j \notin \hat{S}_{\text{init}}^\delta} \frac{|\hat{\beta}_{j,\text{adap}}|}{|\hat{\beta}_{j,\text{init}}|} \\ & \leq 2\|\mathbf{f}_{\hat{S}_{\text{init}}^\delta} - \mathbf{f}^0\|^2 + \frac{14\lambda_{\text{init}}^2}{\phi_0^2} \lambda_{\text{adap}}^2 \sum_{j \in \hat{S}_{\text{init}}^\delta} \frac{1}{\hat{\beta}_{j,\text{init}}^2}. \end{aligned}$$

If moreover

$$\lambda_{\text{adap}} \geq \|(\hat{\beta}_{\text{init}})_{S_0^c}\|_\infty,$$

then

$$|\hat{S}_{\text{adap}} \setminus S_0| \leq 16 \frac{\hat{\delta}_{\text{adap}}^2}{\lambda_{\text{adap}}^2} \frac{\hat{\delta}_2^2}{\lambda_{\text{init}}^2} \wedge 4\Lambda_{\text{max}} \frac{\hat{\delta}_{\text{adap}}}{\lambda_{\text{adap}}} \frac{\hat{\delta}_2}{\lambda_{\text{init}}}.$$

The randomness in the bounds for the adaptive Lasso can be easily handled invoking fixed bounds $\delta_2^{\text{upper}} \geq \hat{\delta}_2$ and $\delta_\infty^{\text{upper}} \geq \hat{\delta}_\infty$, that are assumed to hold on \mathcal{T} . We define

$$|b^0|_{\text{harmonic}}^2 := \left(\sum_{|b_j^0| > 2\delta_\infty^{\text{upper}}} \frac{1}{|b_j^0|^2} \right)^{-1}.$$

The special case $\delta = 3\delta_\infty^{\text{upper}}$ then gives

Corollary 6.3: Main result for the noisy case

Let $\delta_{\text{oracle}}^2 := \|\mathbf{f}_{S_0} - \mathbf{f}^0\|_n^2 + 7\lambda_{\text{init}}^2 s_0/\phi_0^2$. Let

$$\delta_2^{\text{upper}} := \frac{10\delta_{\text{oracle}}^2}{\lambda_{\text{init}}\sqrt{s_0}}.$$

Suppose we are on \mathcal{T} . Suppose $\hat{\delta}_\infty \leq \delta_\infty^{\text{upper}}$, where $3\delta_\infty^{\text{upper}} \geq \delta_2^{\text{upper}}/\sqrt{s_0}$. Let $\lambda_{\text{init}} \geq 2\lambda_{\text{noise}}$ and $\lambda_{\text{adap}} \geq 3\delta_\infty^{\text{upper}}$. Then

$$\hat{\delta}_{\text{adap}}^2 \leq 2 \left\| \mathbf{f}_{S_0^{4\delta_\infty^{\text{upper}}}} - \mathbf{f}^0 \right\|_n^2 + \frac{56\lambda_{\text{init}}^2}{\phi_0^2} \lambda_{\text{adap}}^2 |b^0|_{\text{harmonic}}^{-2}.$$

The choice

$$\lambda_{\text{adap}}^2 = 9 \left(\left\| \mathbf{f}_{S_0^{4\delta_\infty^{\text{upper}}}} - \mathbf{f}^0 \right\|_n^2 + \lambda_{\text{init}}^2 s_0 / \phi_0^2 \right) \frac{|b^0|_{\text{harmonic}}^2}{4\lambda_{\text{init}}^2 / \phi_0^2},$$

indeed has

$$\lambda_{\text{adap}}^2 \geq s_0 |b^0|_{\text{harmonic}}^2 \geq (3\delta_{\text{inf}}^{\text{upper}})^2.$$

With this choice, we find

$$\hat{\delta}_{\text{adap}}^2 \leq 128 \left\{ \left\| \mathbf{f}_{S_0^{4\delta_\infty^{\text{upper}}}} - \mathbf{f}^0 \right\|_n^2 + \lambda_{\text{init}}^2 s_0 \phi_0^2 \right\},$$

and

$$|\hat{S}_{\text{adap}} \setminus S_0| \leq (32M)^2 \frac{\lambda_{\text{init}}^2}{\phi_0^6} |b^0|_{\text{harmonic}}^{-2} s_0 \wedge 32M \Lambda_{\text{max}} \frac{\lambda_{\text{init}}}{\phi_0^3} |b^0|_{\text{harmonic}}^{-1},$$

where

$$M := \frac{10\delta_{\text{oracle}}^2}{\lambda_{\text{init}}^2 s_0 / \phi_0^2}.$$

7 Conclusion

Estimating the support S_0 of the non-zero coefficients is a hard statistical problem. The irrepresentable condition, which is essentially a necessary condition for exact recovery of the non-zero coefficients by the one step Lasso, is much too restrictive in many cases. In this paper, our main focus is on having $O(s_0)$ false positives while achieving good prediction. This is inspired by the behavior of the “ideal” ℓ_0 -penalized estimator. As noted in Section 1, such a viewpoint describes the performance of variable selection in settings where some of the regression coefficients may be smaller than the detection limit.

When using cross validation, the best one can expect is a choice of the tuning parameters that reflects the optimal prediction error of the procedure. We have examined thresholding with least squares refitting and the adaptive Lasso, optimizing the bounds on the prediction error for choosing the tuning parameters. According to our theory (and for simplicity not exploiting the fact that the adaptive Lasso mimics thresholding and refitting using an “oracle” threshold), the two methods are comparable when the trimmed harmonic mean of the squared coefficients of the target is large enough. When the coefficients of the

initial Lasso converge rather slowly in sup-norm to the target, the condition on the trimmed harmonic mean is trivially true.

The adaptive Lasso with cross validation does fitting and variable selection in one single standard algorithm, giving the solution path for all λ_{adap} with $O(n|S_{\text{init}}| \min(n, |S_{\text{init}}|))$ essential operation counts. The tuning parameter λ_{adap} is then chosen based on the performance in the validation sets. Cross validation for thresholding and refitting amounts to removing, for each k , the k smallest estimated initial coefficients $|\hat{\beta}_{j,\text{init}}|$, and evaluating the least squares solution based on the remaining variables on the validation sets. Therefore, the two methods are also computationally comparable.

8 Proofs

8.1 Proofs for Subsection 5.1

Proof of Lemma 5.1. The first result is a special case of Lemma 5.9, taking $\beta = b^0$ and $S = S_0$. The second result then follows from this Lemma, as

$$\|\beta_{\text{init}} - b^0\|_1 \leq \sqrt{s_0} \|(\beta_{\text{init}})_{S_0} - b_{S_0}^0\|_2 + \|(\beta_{\text{init}})_{S_0^c}\|_1.$$

The third result follows from taking $\beta = b^0$ and $S = \mathcal{N}$ in Lemma 5.9, where \mathcal{N} is the set S_0 , complemented with the s_0 largest - in absolute value - coefficients of $(\beta_{\text{init}})_{S_0^c}$. Then $\|f_{b^0} - \mathbf{f}^0\| = \|f_{S_0} - \mathbf{f}^0\|$. Moreover $\phi(2, \mathcal{N}) \leq \phi(2, S_0, 2s_0) = \phi_0$. Thus, from Lemma 5.9, we get

$$\lambda_{\text{init}} \sqrt{2s_0} \|(\beta_{\text{init}})_{\mathcal{N}} - b_{\mathcal{N}}^0\|_2 + \lambda_{\text{init}} \|(\beta_{\text{init}})_{\mathcal{N}^c}\|_1 \leq 3\|f_{S_0} - \mathbf{f}^0\|^2 + \frac{6\lambda_{\text{init}}^2 s_0}{\phi_0^2}.$$

Moreover, as is shown in Lemma 2.2 in van de Geer and Bühlmann [2009] (with original reference Candès and Tao [2005], and Candès and Tao [2007]),

$$\|(\beta_{\text{init}})_{\mathcal{N}^c}\|_2 \leq \|(\beta_{\text{init}})_{S_0}\|_1 / \sqrt{s_0} \leq \frac{3\|f_{S_0} - \mathbf{f}^0\|^2 + 3\lambda_{\text{init}}^2 s_0 / \phi_0^2}{\lambda_{\text{init}} \sqrt{s_0}}.$$

So then

$$\begin{aligned} \|\beta_{\text{init}} - b_0\|_2 &\leq \|(\beta_{\text{init}})_{\mathcal{N}} - b_{\mathcal{N}}^0\|_2 + \|(\beta_{\text{init}})_{\mathcal{N}^c}\|_2 \\ &\leq \frac{6\|f_{S_0} - \mathbf{f}^0\|^2 + 9\lambda_{\text{init}}^2 s_0 / \phi_0^2}{\sqrt{s_0} \lambda_{\text{init}}} \leq \frac{6\delta_{\text{oracle}}^2}{\lambda_{\text{init}} \sqrt{s_0}}. \end{aligned}$$

□

Proof of Lemma 5.2. This follows from

$$|\beta_{j,\text{init}}| \geq |b_j^0| - |\beta_{j,\text{init}} - b_j^0|.$$

□

Proof of Lemma 5.3. This is a special case of Lemma 5.10. □

Proof of Lemma 5.4. The first inequality is trivial, as the refitted version is the projection of \mathbf{f}^0 on the space spanned by the variables in S_{init}^δ , and f_{init}^δ is in this space.

For the second result, we note that if $|\beta_{j,\text{init}}| \leq \delta$, then $|b_j^0| \leq \delta + \delta_\infty$. In other words

$$S_0^{\delta+\delta_\infty} \subset S_{\text{init}}^\delta.$$

Hence

$$\|f_{S_{\text{init}}^\delta} - \mathbf{f}^0\| \leq \|f_{S_0^{\delta+\delta_\infty}} - \mathbf{f}^0\|.$$

Moreover,

$$\begin{aligned} \|f_{b_{S_0}^0} - f_{b_{S_0^{\delta+\delta_\infty}}^0}\|^2 &= (b_{S_0 \setminus S_0^{\delta+\delta_\infty}}^0)^T \Sigma (b_{S_0 \setminus S_0^{\delta+\delta_\infty}}^0) \\ &\leq \Lambda_{\max}^2(S_0 \setminus S_0^{\delta+\delta_\infty}) \|b_{S_0 \setminus S_0^{\delta+\delta_\infty}}^0\|_2^2 \\ &\leq \Lambda_{\max}^2(S_0 \setminus S_0^{\delta+\delta_\infty}) |S_0 \setminus S_0^{\delta+\delta_\infty}| (\delta + \delta_\infty)^2. \end{aligned}$$

But then

$$\begin{aligned} \|f_{S_0^{\delta+\delta_\infty}} - \mathbf{f}^0\| &= \min_{\beta = \beta_{S_0^{\delta+\delta_\infty}}} \|f_\beta - \mathbf{f}^0\| \\ &\leq \|f_{b_{S_0^{\delta+\delta_\infty}}^0} - \mathbf{f}^0\| \leq \|f_{S_0} - \mathbf{f}^0\| + \|f_{b_{S_0}^0} - f_{b_{S_0^{\delta+\delta_\infty}}^0}\| \leq \\ &\|f_{S_0} - \mathbf{f}^0\| + \Lambda_{\max}(S_0 \setminus S_0^{\delta+\delta_\infty}) \sqrt{|S_0 \setminus S_0^{\delta+\delta_\infty}|} (\delta + \delta_\infty). \end{aligned}$$

□

Proof of Lemma 5.5. We clearly have

$$\delta^2 |S_{\text{init}}^\delta \setminus S_0| \leq \sum_{j \in S_{\text{init}}^\delta \setminus S_0} |\beta_{j,\text{init}}|^2 \leq \|\beta_{\text{init}} - b^0\|_2^2 \leq \delta_2^2.$$

Whence the first result.

Moreover,

$$K^2 \delta^2 |S_0^{(K+1)\delta} \setminus S_{\text{init}}^\delta| \leq \sum_{j \in S_0^{(K+1)\delta} \setminus S_{\text{init}}^\delta} |\beta_{j,\text{init}} - b_j^0|^2 \leq \|\beta_{\text{init}} - b^0\|_2^2 \leq \delta_2^2.$$

This gives the second result.

□

8.2 Proofs for Subsection 5.2

Proof of Lemma 5.6. This follows from applying Lemma 5.9, with $L = 1$. We only have to show that

$$\phi(2, S_{\text{init}}) \geq \phi_0.$$

Because (see (7)),

$$|S_{\text{init}}^\delta \setminus S_0| \leq s_0.$$

indeed

$$\phi(2, S_{\text{init}}^\delta) \geq \phi(2, S_{\text{init}}^\delta \cup S_0) \geq \phi(2, S_0, 2s_0) = \phi_0.$$

□

Proof of Lemma 5.7. We use that if $|\beta_{j,\text{init}}| > \delta$, then $|b_j^0| > \delta - \delta_\infty$. Moreover, if $|b_j^0| > 2\delta_\infty$, then $|\beta_{j,\text{init}}| \geq |b_j^0|/2$. Hence,

$$\begin{aligned} \sum_{j \in S_{\text{init}}^\delta} \frac{1}{\beta_{j,\text{init}}^2} &= \sum_{|\beta_{j,\text{init}}| > \delta, |b_j^0| \leq 2\delta_\infty} \frac{1}{\beta_{j,\text{init}}^2} + \sum_{|\beta_{j,\text{init}}| > \delta, |b_j^0| > 2\delta_\infty} \frac{1}{\beta_{j,\text{init}}^2} \\ &\leq \frac{1}{\delta^2} \left\{ \left| \{j : \delta - \delta_\infty < |b_j^0| \leq 2\delta_\infty\} \right| + 4\delta^2 \sum_{|b_j^0| > 2\delta_\infty} \frac{1}{|b_j^0|^2} \right\}. \end{aligned}$$

The second result follows from

$$\sum_{j \in S_{\text{init}}^\delta} \frac{1}{\beta_{j,\text{init}}^2} \leq \frac{1}{\delta^2} |S_{\text{init}}^\delta|,$$

and, invoking (7),

$$|S_{\text{init}}^\delta| \leq |S_{\text{init}}^\delta \setminus S_0| + |S_0| \leq 2s_0.$$

□

Proof of Lemma 5.8. This is a special case of Lemma 5.10. □

8.3 Proofs for Subsection 5.3

Proof of Lemma 5.9. We have

$$\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j=1}^p w_j |\beta_{j,\text{weight}}| \leq \|f_{\beta_S} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \in S} w_j |\beta_j|,$$

and hence

$$\begin{aligned} &\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j,\text{weight}}| \\ &\leq \|f_{\beta_S} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \in S} w_j |\beta_{j,\text{weight}} - \beta_j| \\ &\leq \|f_{\beta_S} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2. \end{aligned}$$

Case i). If

$$\|f_{\beta_S} - \mathbf{f}^0\|^2 \leq \lambda_{\text{init}} \lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2,$$

we get

$$\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j,\text{weight}}| \leq 2\lambda_{\text{init}} \lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2. \quad (15)$$

It follows that

$$\|(\beta_{\text{weight}})_{S^c}\|_1 \leq 2L\sqrt{|S|}\|(\beta_{\text{weight}})_S - (\beta)_S\|_2.$$

But then

$$\begin{aligned} \|(\beta_{\text{weight}})_S - \beta_S\|_2 &\leq \|f_{\text{weight}} - f_{\beta_S}\|/\phi(2L, S) \\ &\leq \|f_{\text{weight}} - \mathbf{f}^0\|/\phi(2L, S) + \|f_{\beta_S} - \mathbf{f}^0\|/\phi(2L, S). \end{aligned}$$

This gives

$$\begin{aligned} &\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j, \text{weight}}| \\ &\leq 2\lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|f_{\text{weight}} - \mathbf{f}^0\|/\phi(2L, S) \\ &\quad + 2\lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|f_{\beta_S} - \mathbf{f}^0\|/\phi(2L, S) \\ &\leq \frac{1}{2} \|f_{\text{weight}} - \mathbf{f}^0\|^2 + \|f_{\beta_S} - \mathbf{f}^0\|^2 + \frac{3\lambda_{\text{init}}^2 \lambda_{\text{weight}}^2 \|w_S\|_2^2}{\phi^2(2L, S)}. \end{aligned}$$

Hence,

$$\|f_{\text{weight}} - \mathbf{f}^0\|^2 + 2\lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j, \text{weight}}| \leq 2\|f_{\beta_S} - \mathbf{f}^0\|^2 + \frac{6\lambda_{\text{init}}^2 \lambda_{\text{weight}}^2 \|w_S\|_2^2}{\phi^2(2L, S)}.$$

Case ii) If

$$\|f_{\beta_S} - \mathbf{f}^0\|^2 > \lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2,$$

we get

$$\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j, \text{weight}}| \leq 2\|f_{\beta_S} - \mathbf{f}^0\|^2.$$

For the second result, we add in Case i), $\lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2$ to the left and right hand side of (15):

$$\begin{aligned} &\|f_{\text{weight}} - \mathbf{f}^0\|^2 + \lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2 + \lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j, \text{weight}}| \\ &\leq 3\lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2. \end{aligned}$$

The same arguments now give

$$\begin{aligned} &3\lambda_{\text{init}}\lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2 \leq \\ &\|f_{\text{weight}} - \mathbf{f}^0\|^2 + 3\|f_{\beta_S} - \mathbf{f}^0\|^2 + \frac{3\lambda_{\text{init}}^2 \lambda_{\text{weight}}^2 \|w_S\|_2^2}{\phi^2(2L, S)}. \end{aligned}$$

In Case ii), we have

$$\lambda_{\text{init}}\lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j, \text{weight}}| \leq 2\|f_{\beta_S} - \mathbf{f}^0\|^2,$$

and also

$$\lambda_{\text{init}} \lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2 < \|f_{\beta_S} - \mathbf{f}^0\|^2.$$

So then

$$\begin{aligned} \lambda_{\text{init}} \lambda_{\text{weight}} \|w_S\|_2 \|(\beta_{\text{weight}})_S - \beta_S\|_2 + \lambda_{\text{init}} \lambda_{\text{weight}} \sum_{j \notin S} w_j |\beta_{j, \text{weight}}| \\ < 3 \|f_{\beta_S} - \mathbf{f}^0\|^2. \end{aligned}$$

□

Proof of Lemma 5.10. By the weighted KKT conditions, for all j

$$2(\psi_j, f_{\text{weight}} - \mathbf{f}^0) = -\lambda_{\text{init}} \lambda_{\text{weight}} w_j \tau_{j, \text{weight}}.$$

Hence,

$$\begin{aligned} \sum_{j \in S_{\text{weight}} \setminus S_0} 2|(\psi_j, f_{\text{weight}} - \mathbf{f}^0)|^2 &\geq \lambda_{\text{init}}^2 \lambda_{\text{weight}}^2 \|w_{S_{\text{weight}} \setminus S_0}\|_2^2 \\ &\geq \lambda_{\text{init}}^2 \lambda_{\text{weight}}^2 |S_{\text{weight}} \setminus S_0|^2 / \|(1/w)_{S_{\text{weight}} \setminus S_0}\|_2^2. \end{aligned}$$

On the other hand

$$\sum_{j \in S_{\text{weight}} \setminus S_0} |(\psi_j, f_{\text{weight}} - \mathbf{f}^0)|^2 \leq \Lambda_{\max}^2(S_{\text{weight}} \setminus S_0) \|f_{\text{weight}} - \mathbf{f}^0\|^2.$$

Clearly,

$$\Lambda_{\max}^2(S_{\text{weight}} \setminus S_0) \leq \Lambda_{\max}^2 \wedge |S_{\text{weight}} \setminus S_0|.$$

□

8.4 Proofs for Subsection 5.4

Proof of Lemma 5.12. Clearly,

$$\|W_{S^c}^{-1} \Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S\|_{\infty} \leq \|\Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S\|_{\infty} / w_{S^c}^{\min}.$$

Define

$$\beta_S := \Sigma_{1,1}^{-1}(S) W_S \tau_S.$$

Then

$$\begin{aligned} \|W_{S^c}^{-1} \Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S\|_{\infty} &= \sup_{\|\gamma_{S^c}\|_1 \leq 1} |\gamma_{S^c}^T W_{S^c}^{-1} \Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S| \\ &= \sup_{\|W_{S^c} \beta_{S^c}\|_1 \leq 1} |\beta_S^T \Sigma_{2,1}(S) \beta_S| = \sup_{\|W_{S^c} \beta_{S^c}\|_1 \leq 1} |(f_{\beta_{S^c}}, f_{\beta_S})| \\ &\leq \sup_{\|\beta_{S^c}\|_1 \leq 1/w_{S^c}^{\min}} |(f_{\beta_{S^c}}, f_{\beta_S})| \end{aligned}$$

$$\begin{aligned}
&= \sup_{\|\beta_{S^c}\|_1 \leq \|w_S\|_2 \|\beta_S\|_2 / w_{S^c}^{\min}} \frac{|(f_{\beta_{S^c}}, f_{\beta_S})|}{\|w_S\|_2 \|\beta_S\|_2} \\
&= \sup_{\|\beta_{S^c}\|_1 \leq \|w_S\|_2 \|\beta_S\|_2 / w_{S^c}^{\min}} \frac{|(f_{\beta_{S^c}}, f_{\beta_S})|}{\|f_{\beta_S}\|^2} \frac{\|f_{\beta_S}\|^2}{\|w_S\|_2 \|\beta_S\|_2}.
\end{aligned}$$

But

$$\frac{\|f_{\beta_S}\|^2}{\|w_S\|_2 \|\beta_S\|_2} = \frac{\tau_S^T W_S \Sigma_{1,1}^{-1}(S) W_S \tau_S}{\sqrt{\tau_S^T W_S^2 \tau_S} \sqrt{\tau_S W_S \Sigma_{1,1}^{-2}(S) W_S \tau_S}} \frac{\|W_S \tau_S\|_2}{\|w_S\|_2} \leq 1.$$

We conclude that

$$\begin{aligned}
\|W_{S^c}^{-1} \Sigma_{2,1}(S) \Sigma_{1,1}^{-1}(S) W_S \tau_S\|_\infty &\leq \sup_{\|\beta_{S^c}\|_1 \leq \|w_S\|_2 \|\beta_S\|_2 / w_{S^c}^{\min}} \frac{|(f_{\beta_{S^c}}, f_{\beta_S})|}{\|f_{\beta_S}\|^2} \\
&= \frac{\|w_S\|_2}{\sqrt{|S|} w_{S^c}^{\min}} \vartheta_{\text{adaptive}}(S).
\end{aligned}$$

□

8.5 Proofs for Section 6

The proofs of Lemma 6.1 and 6.2 are a straightforward extension of their noiseless versions, and therefore omitted.

Proof of Lemma 6.3. This follows from

$$\|\hat{f}_{\hat{S}_{\text{init}}^\delta} - f_{\hat{S}_{\text{init}}^\delta}\|_n^2 \leq 2(\epsilon, \hat{f}_{\hat{S}_{\text{init}}^\delta} - f_{\hat{S}_{\text{init}}^\delta})_n,$$

and

$$\begin{aligned}
2(\epsilon, \hat{f}_{\hat{S}_{\text{init}}^\delta} - f_{\hat{S}_{\text{init}}^\delta})_n &\leq \lambda_{\text{noise}} \|\hat{b}^{\hat{S}_{\text{init}}^\delta} - b^{\hat{S}_{\text{init}}^\delta}\|_1 \\
&\leq \sqrt{2s_0} \|\hat{b}^{\hat{S}_{\text{init}}^\delta} - b^{\hat{S}_{\text{init}}^\delta}\|_2 \leq \sqrt{2s_0} \|\hat{f}_{\hat{S}_{\text{init}}^\delta} - f_{\hat{S}_{\text{init}}^\delta}\|_n / \phi_0^2.
\end{aligned}$$

□

References

- A. Barron, L. Birge, and P. Massart. Risk bounds for model selection via penalization. *Probability Theory and Related Fields*, 113:301–413, 1999.
- D. Bertsimas and J. Tsitsiklis. *Introduction to linear optimization*. Athena Scientific Belmont, MA, 1997.
- P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37:1705–1732, 2009.

- F. Bunea, A.B. Tsybakov, and M.H. Wegkamp. Aggregation and sparsity via ℓ_1 -penalized least squares. In *Proceedings of 19th Annual Conference on Learning Theory, COLT 2006. Lecture Notes in Artificial Intelligence 4005*, pages 379–391, Heidelberg, 2006. Springer Verlag.
- F. Bunea, A.B. Tsybakov, and M.H. Wegkamp. Aggregation for Gaussian regression. *Annals of Statistics*, 35:1674–1697, 2007a.
- F. Bunea, A. Tsybakov, and M.H. Wegkamp. Sparsity oracle inequalities for the Lasso. *Electronic Journal of Statistics*, 1:169–194, 2007b.
- E. Candès and Y. Plan. Near-ideal model selection by ℓ_1 minimization. *Annals of Statistics*, 37:2145–2177, 2009.
- E. Candès and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51:4203–4215, 2005.
- E. Candès and T. Tao. The Dantzig selector: statistical estimation when p is much larger than n . *Annals of Statistics*, 35:2313–2351, 2007.
- E.J. Candès, J.K. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics*, 59:1207–1223, 2006.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularized paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33, 2010.
- E. Greenshtein and Y. Ritov. Persistency in high dimensional linear predictor-selection and the virtue of over-parametrization. *Bernoulli*, 10:971–988, 2004.
- J. Huang, S. Ma, and C.-H. Zhang. Adaptive Lasso for sparse high-dimensional regression models. *Statistica Sinica*, 18:1603–1618, 2008.
- V. Koltchinskii. Sparsity in penalized empirical risk minimization. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, 45:7–57, 2009a.
- V. Koltchinskii. The Dantzig selector and sparsity oracle inequalities. *Bernoulli*, 15:799–828, 2009b.
- K. Lounici. Sup-norm convergence rate and sign concentration property of Lasso and Dantzig estimators. *Electronic Journal of Statistics*, 2:90–102, 2008.
- L. Meier, S. van de Geer, and P. Bühlmann. The group Lasso for logistic regression. *Journal of the Royal Statistical Society Series B*, 70:53–71, 2008.
- N. Meinshausen. Relaxed Lasso. *Computational Statistics and Data Analysis*, 52:374–393, 2007.
- N. Meinshausen and P. Bühlmann. High dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics*, 37:246–270, 2009.

- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society Series B*, 58:267–288, 1996.
- S. van de Geer. High-dimensional generalized linear models and the Lasso. *Annals of Statistics*, 36:614–645, 2008.
- S. van de Geer. Least squares estimation with complexity penalties. *Mathematical Methods of Statistics*, pages 355–374, 2001.
- S. van de Geer and P. Bühlmann. On the conditions used to prove oracle results for the Lasso. *Electronic Journal of Statistics*, pages 1360–1392, 2009.
- M. Wainwright. Information-theoretic limitations on sparsity recovery in the high-dimensional and noisy setting. *IEEE Transactions on Information Theory*, 55:5728–5741, 2007.
- M. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Transactions on Information Theory*, 55:2183–2202, 2009.
- L. Wasserman and K. Roeder. High dimensional variable selection. *Annals of statistics*, 37:2178–2201, 2009.
- T. Zhang. Some sharp performance bounds for least squares regression with ℓ_1 regularization. *Annals of Statistics*, 37:2109–2144, 2009.
- P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2567, 2006.
- S. Zhou. Thresholding procedures for high dimensional variable selection and statistical estimation. In *Advances in Neural Information Processing Systems 22*. MIT Press, 2009.
- S. Zhou. Thresholded lasso for high dimensional variable selection and statistical estimation, 2010. arXiv:1002.1583v2, shorter version in *Advances in Neural Information Processing Systems 22 (NIPS 2009)*.
- H. Zou. The adaptive Lasso and its oracle properties. *Journal of the American Statistical Association*, 101:1418–1429, 2006.
- H. Zou and R. Li. One-step sparse estimates in nonconcave penalized likelihood models (with discussion). *Annals of Statistics*, 36:1509–1566, 2008.