Estimating the error variance in a high-dimensional linear model

Guo Yu* Jacob Bien[†]

Abstract

The lasso has been studied extensively as a tool for estimating the coefficient vector in the high-dimensional linear model; however, considerably less is known about estimating the error variance in this context. In this paper, we propose the natural lasso estimator for the error variance, which maximizes a penalized likelihood objective. A key aspect of the natural lasso is that the likelihood is expressed in terms of the natural parameterization of the multiparameter exponential family of a Gaussian with unknown mean and variance. The result is a remarkably simple estimator of the error variance with provably good performance in terms of mean squared error. These theoretical results do not require placing any assumptions on the design matrix or the true regression coefficients. We also propose a companion estimator, called the organic lasso, which theoretically does not require tuning of the regularization parameter. Both estimators do well empirically compared to preexisting methods, especially in settings where successful recovery of the true support of the coefficient vector is hard. Finally, we show that existing methods can do well under fewer assumptions than previously known, thus providing a fuller story about the problem of estimating the error variance in high-dimensional linear models.

1 Introduction

The linear model

$$y = X\beta^* + \varepsilon \qquad \varepsilon \sim N(0, \sigma^2 I_n),$$
 (1)

is one of the most fundamental models in statistics. It describes the relationship between a response vector $y \in \mathbb{R}^n$ and a fixed design matrix $X \in \mathbb{R}^{n \times p}$. When $p \gg n$, estimating the coefficient vector β^* is a challenging, well-studied problem. Perhaps the most common method in this setting is the lasso (Tibshirani 1996), which assumes that β^* is sparse and solves the following convex optimization problem:

$$\hat{\beta}_{\lambda} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}} \left(\frac{1}{n} \left\| y - X\beta \right\|_{2}^{2} + 2\lambda \left\| \beta \right\|_{1} \right). \tag{2}$$

Over the past decade, an extensive literature has emerged studying the properties of $\hat{\beta}_{\lambda}$ from both computational (e.g., Hastie et al. 2015) and theoretical (e.g., Bühlmann & Van De Geer 2011) perspectives.

Compared to the vast amount of work on estimating β^* , relatively little attention has been paid to the problem of estimating σ^2 , which captures the noise level or extent to which y cannot be predicted from X. Nonetheless, reliable estimation of σ^2 is important for quantifying the uncertainty in estimating β^* . A series of recent advances in high-dimensional inference (Bühlmann 2013, Zhang & Zhang 2014, Van de Geer et al. 2014, Lockhart et al. 2014, Javanmard & Montanari 2014, Lee et al. 2016, Tibshirani et al. 2016, Taylor & Tibshirani 2017, Ning &

^{*}Department of Statistics, University of Washington, Seattle, Washington, 98105, gy63@uw.edu

 $^{^\}dagger$ Data Sciences and Operations, Marshall School of Business, University of Southern California, Los Angeles, CA 90089, jbien@usc.edu

Liu 2017, etc.) may very well be the determining factor for the widespread adoption of the lasso and related methods in fields where p-values and confidence intervals are required. Thus, estimating σ^2 reliably in finite samples is crucial.

If β^* were known, then the optimal estimator for σ^2 would of course be $n^{-1}||y - X\beta^*||_2^2 = n^{-1}||\varepsilon||_2^2$. Thus, a naive estimator for σ^2 based on an estimator $\hat{\beta}$ of β^* would be

$$\hat{\sigma}_{\text{naive}}^2 = \frac{1}{n} \|y - X\hat{\beta}\|_2^2. \tag{3}$$

However, a simple calculation in the classical n > p setting shows that such an estimator is biased downward: a least-squares oracle with knowledge of the true support $S = \{j : \beta_j^* \neq 0\}$ scales this to give an unbiased estimator,

$$\hat{\sigma}_{\text{oracle}}^2 = \frac{1}{n - |S|} \|y - X_S X_S^+ y\|_2^2,\tag{4}$$

where X_S is a sub-matrix of X with columns indexed by S and X_S^+ is its pseudoinverse. Many papers in this area discuss the difficulty of estimating σ^2 and warn of the perils of underestimating it: if σ^2 is underestimated then one gets anti-conservative confidence intervals, which are highly undesirable (Tibshirani et al. 2018).

Reid et al. (2016) carry out an extensive review and simulation study of several estimators of σ^2 (Fan et al. 2012, Sun & Zhang 2012, Dicker 2014), and they devote special attention to studying the estimator

$$\hat{\sigma}_R^2 = \frac{1}{n - \hat{s}_{\lambda}} \|y - X\hat{\beta}_{\lambda}\|_2^2,\tag{5}$$

where $\hat{\beta}_{\lambda}$ is as in (2), with λ selected using a cross-validation procedure, and \hat{s}_{λ} is the number of nonzero elements in $\hat{\beta}_{\lambda}$. They show that (5) has promising performance in a wide range of simulation settings and provide an asymptotic theoretical understanding of the estimator in the special case where X is an orthogonal matrix.

While intuition from (4) suggests that (5) is a quite reasonable estimator when S can be well recovered, it also points to the question of how well the estimator will perform when S is not well recovered by the lasso. The conditions required for the lasso to recover S are much stricter than the conditions needed for it to do well in prediction (e.g., Van de Geer & Bühlmann 2009). The scale factor $(n-\hat{s}_{\lambda})^{-1}$ used in $\hat{\sigma}_R^2$ means that this approach depends not just on the predicted values of the lasso, $X\hat{\beta}_{\lambda}$, but on the magnitude of the set of nonzero elements in $\hat{\beta}_{\lambda}$. Indeed, we find that in situations where recovering S is challenging, $\hat{\sigma}_R^2$ tends to yield less favorable empirical performance. The theoretical development in Reid et al. (2016) sidesteps this complication by working in an asymptotic regime in which $\hat{\sigma}_R^2$ behaves like the naive estimator (3). To understand the finite-sample performance of $\hat{\sigma}_R^2$ would require considering the behavior of the random variable \hat{s}_{λ} . Clearly, when $\hat{s}_{\lambda} \approx n$, even small fluctuations in \hat{s}_{λ} can lead to large fluctuations in $\hat{\sigma}_R^2$. Finally, from a practical standpoint, computing \hat{s}_{λ} is a numerically sensitive operation in that it requires the choice of a threshold size for calling a value numerically zero, and the assurance that one has solved the problem to sufficient precision.

Based on these observations, we propose in this paper a completely different approach to estimating σ^2 . The basic premise of our framework is that when both β^* and σ^2 are unknown, it is convenient to formulate the penalized log-likelihood problem in terms of

$$\phi = \frac{1}{\sigma^2}, \qquad \theta = \frac{\beta}{\sigma^2},\tag{6}$$

the natural parameters of the Gaussian multiparameter exponential family with unknown mean and variance. The negative Gaussian log-likelihood is not jointly convex in the (β, σ) parameterization. In fact, even with β fixed, it is nonconvex in σ . However, in the natural parameterization the negative log-likelihood is jointly convex in (ϕ, θ) .

We penalize this negative log-likelihood with an ℓ_1 -norm on the natural parameter θ and call this new estimator the natural lasso. We show in Section 3 that the resulting error variance estimator can in fact be very simply expressed as the minimizing value of the regular lasso problem (2):

$$\hat{\sigma}_{\lambda}^{2} = \min_{\beta \in \mathbb{R}^{p}} \left(\frac{1}{n} \|y - X\beta\|_{2}^{2} + 2\lambda \|\beta\|_{1} \right). \tag{7}$$

Observing that the first term is $\hat{\sigma}_{\text{naive}}^2$, we directly see that the natural lasso counters the naive method's downward bias through an additive correction; this is in contrast to $\hat{\sigma}_R^2$'s reliance on a multiplicative correction that sometimes may be unstable. Computing (7) is clearly no harder than solving a lasso and, unlike $\hat{\sigma}_R^2$, does not require determining a threshold for deciding which coefficient estimates are numerically zero. Furthermore, we establish finite-sample bounds on the mean squared error that hold without making any assumptions on the design matrix X. Our theoretical analysis suggests a second approach that is also based on the natural parameterization. The theory that we develop for this method, which we call the organic lasso, relies on weaker assumptions. We find that both methods have competitive empirical performance relative to $\hat{\sigma}_R^2$ and show particular strength in settings in which support recovery is known to be challenging.

Our final contribution is to show that existing methods can also attain high-dimensional consistency under no assumptions on X. In particular, we provide finite-sample bounds for $\hat{\sigma}_{\text{naive}}^2$, with $\hat{\beta}$ in (3) taken to be the standard lasso or the square-root/scaled lasso estimator (Belloni et al. 2011, Sun & Zhang 2012). Previous results about $\hat{\sigma}_{\text{naive}}^2$ have placed strong assumptions on X. Thus, our work provides a fuller story about the problem of estimating the error variance in high-dimensional linear models.

2 Natural parameterization

The negative log-likelihood function in (1) is, up to a constant,

$$L(\beta, \sigma^2 | X, y) = \frac{n}{2} \log \sigma^2 + \frac{\|y - X\beta\|_2^2}{2\sigma^2}.$$

When σ^2 is known, the σ dependence can be ignored, leading to the standard least-squares criterion; however, when σ is unknown, performing a full minimization of the penalized negative log-likelihood amounts to solving a nonconvex optimization problem even with a convex penalty.

The nonconvexity of the Gaussian negative log-likelihood in its variance, or more generally, covariance matrix, is a well-known difficulty (Bien & Tibshirani 2011). In this context, working instead with the inverse covariance matrix is common (Yuan & Lin 2007, Banerjee et al. 2008, Friedman et al. 2008). We take an analogous approach when estimation of σ^2 is of interest, considering the natural parameterization (6) of the Gaussian multiparameter exponential family with unknown variance,

$$L\left(\phi^{-1}\theta,\phi^{-1}|X,y\right) = -\frac{n}{2}\log\phi + \frac{1}{2}\phi\left\|y - X\frac{\theta}{\phi}\right\|_2^2 = -\frac{n}{2}\log\phi + \phi\frac{\|y\|_2^2}{2} - y^TX\theta + \frac{\|X\theta\|_2^2}{2\phi}.$$

Observing that attaining sparsity in θ is equivalent to attaining sparsity in β , we propose the following penalized maximum log-likelihood estimator:

$$\left(\hat{\theta}_{\lambda}, \hat{\phi}_{\lambda}\right) \in \operatorname*{arg\,min}_{\phi > 0, \; \theta} \left\{ -\frac{1}{2} \log \phi + \phi \frac{\|y\|_2^2}{2n} - \frac{1}{n} y^T X \theta + \frac{\|X\theta\|_2^2}{2n\phi} + \lambda \Omega(\theta, \phi) \right\} \tag{8}$$

for a convex penalty $\Omega(\theta, \phi)$ that induces sparsity in θ . We will focus on $\Omega(\theta, \phi) = \|\theta\|_1$ in Section 3 and $\Omega(\theta, \phi) = \phi^{-1} \|\theta\|_1^2$ in Section 4. This problem is jointly convex in (θ, ϕ) . While

this is a general property of exponential families due to the convexity of the cumulant generating function, we can see it in this special case because of the convexity of $-\log$ and the convexity of the quadratic-over-linear function (Boyd & Vandenberghe 2004, Rockafellar 2015). Given a solution to (8), we can reverse (6) to get estimators for σ^2 and β^* :

$$\tilde{\sigma}_{\lambda}^{2} = \frac{1}{\hat{\phi}_{\lambda}}, \qquad \tilde{\beta}_{\lambda} = \frac{\hat{\theta}_{\lambda}}{\hat{\phi}_{\lambda}}.$$
 (9)

Before proceeding with an analysis of the estimator (9) with specific choices of $\Omega(\theta, \phi)$, we point out a similarity between our method and that of Städler et al. (2010), who consider a different convexifying reparameterization of the Gaussian log-likelihood, using $\rho = \sigma^{-1}$ and $\gamma = \sigma^{-1}\beta$. They put an ℓ_1 -norm penalty on γ , which has the same sparsity pattern as β , and solve

$$\min_{\rho > 0, \gamma} \left(-\log \rho + \frac{1}{2n} \|\rho y - X\gamma\|_{2}^{2} + \lambda \|\gamma\|_{1} \right). \tag{10}$$

Sun & Zhang (2010) give an asymptotic analysis of the solution to (10) under a compatibility condition. A modification of this problem (Antoniadis 2010) gives the scaled lasso (Sun & Zhang 2012), which is known to be equivalent to the square–root lasso (Belloni et al. 2011):

$$\tilde{\beta}_{\text{SQRT}} = \operatorname*{arg\,min}_{\beta \in \mathbb{R}^p} \left(\frac{1}{\sqrt{n}} \left\| y - X\beta \right\|_2 + \lambda \left\| \beta \right\|_1 \right), \qquad \tilde{\sigma}_{\text{SQRT}}^2 = \frac{1}{n} \left\| y - X\tilde{\beta}_{\text{SQRT}} \right\|_2^2. \tag{11}$$

With the same parameterization (ρ, γ) , Dalalyan & Chen (2012) propose the scaled Dantzig selector under the assumption of fused sparsity. Under the restricted eigenvalue condition, they establish the same rate of convergence in estimating the error variance as the fast prediction error rate of the standard lasso (Hebiri & Lederer 2013, Lederer et al. 2016, Dalalyan et al. 2017).

3 The natural lasso estimator of error variance

We first propose the natural lasso, which is the solution to (8) with $\Omega(\theta, \phi) = \|\theta\|_1$. One might think that solving the natural lasso would involve a specialized algorithm. The following proposition shows, remarkably, that this is not the case.

Proposition 1. The natural lasso estimator $(\tilde{\beta}_{\lambda}, \tilde{\sigma}_{\lambda}^2)$ defined in (9), where $(\hat{\theta}_{\lambda}, \hat{\phi}_{\lambda})$ is a solution to (8) with $\Omega(\theta, \phi) = \|\theta\|_1$, satisfies the following properties:

- 1. $\tilde{\beta}_{\lambda} = \hat{\beta}_{\lambda}$, a solution to the standard lasso (2);
- 2. $\tilde{\sigma}_{\lambda}^2 = \hat{\sigma}_{\lambda}^2$, the standard lasso's optimal value (7).

Furthermore, $\hat{\sigma}_{\lambda}^2 = n^{-1}(\|y\|_2^2 - \|X\hat{\beta}_{\lambda}\|_2^2).$

The proof of this proposition and all theoretical results that follow can be found in the Appendices. Thus, to get the natural lasso estimator of (β^*, σ^2) , one simply solves the standard lasso (2) and returns a solution and the minimal value.

An attractive property of the natural lasso estimator $\hat{\sigma}_{\lambda}^2$ is the relative ease with which one can prove bounds about its performance. Since $\hat{\sigma}_{\lambda}^2$ is the optimal value of the lasso problem, the objective value at any vector β provides an upper bound on $\hat{\sigma}_{\lambda}^2$. Likewise, any dual feasible vector provides a lower bound on $\hat{\sigma}_{\lambda}^2$. These considerations are used to prove the following lemma, which shows that for a suitably chosen λ , the natural lasso variance estimator gets close to the oracle estimator of σ^2 .

Lemma 2. If
$$\lambda \ge n^{-1} \|X^T \varepsilon\|_{\infty}$$
, then $|\hat{\sigma}_{\lambda}^2 - n^{-1} \|\varepsilon\|_2^2 | \le 2\lambda \|\beta^*\|_1$.

The result above is deterministic in that it does not rely on any statistical assumptions or arguments. The next result adds such considerations to give a mean squared error bound for the natural lasso.

Theorem 3. Suppose that each column X_j of the matrix $X \in \mathbb{R}^{n \times p}$ has been scaled so that $\|X_j\|_2^2 = n$ for all $j = 1, \ldots, p$, and assume that $\varepsilon \sim N\left(0, \sigma^2 I_n\right)$. Then, for any constant M > 1, the natural lasso estimator (7) with $\lambda = \sigma(2Mn^{-1}\log p)^{1/2}$ satisfies the following relative mean squared error bound:

$$E\left\{ \left(\frac{\hat{\sigma}_{\lambda}^{2}}{\sigma^{2}} - 1 \right)^{2} \right\} \leq \left\{ \left(8M + 8 \frac{p^{1 - 8M}}{\log p} \right)^{1/2} \frac{\|\beta^{*}\|_{1}}{\sigma} \left(\frac{\log p}{n} \right)^{1/2} + \left(\frac{2}{n} \right)^{1/2} \right\}^{2}.$$

Corollary 4.

$$E\left|\frac{\hat{\sigma}_{\lambda}^{2}}{\sigma^{2}} - 1\right| = O\left\{\frac{\|\beta^{*}\|_{1}}{\sigma} \left(\frac{\log p}{n}\right)^{1/2}\right\}. \tag{12}$$

Proof. This follows from Jensen's inequality.

Remark 5. Theorem 3 can be easily generalized to the case where the independently and identically distributed zero-mean error ε_i with variance σ^2 is sub-Gaussian or sub-exponential. A high probability bound can be obtained for ε_i with bounded polynomial moments. In particular, for any $m \geq 3$, if $\mathrm{E}(|\varepsilon_i|^m) \leq (m!)^{-1} 2K^{m-2}$ for some K > 0, and if each column X_j is scaled so that $\sum_{i=1}^n X_{ij}^m = n$ for $j = 1, \ldots, p$, then with $\lambda = 4K\sigma n^{-1/2}(\log p)^{1/2}$ we have that

$$\left| \hat{\sigma}_{\lambda}^2 - \frac{\|\varepsilon\|_2^2}{n} \right| = O\left\{ \sigma \, \|\beta^*\|_1 \left(\frac{\log p}{n} \right)^{1/2} \right\}$$

holds with probability greater than $1 - p^{-1}$.

To put Theorem 3 in context, we devote the remainder of this section to considering what bounds are available for other methods for estimating σ^2 . Bayati et al. (2013) propose an estimator of σ^2 based on estimating the mean squared error of the lasso. They show that their estimator of σ^2 is asymptotically consistent with fixed p as $n \to \infty$. In contrast, we provide finite sample results and these include the $p \gg n$ case. Also, the consistency result in Bayati et al. (2013) is based on the assumption of independent Gaussian features, and in extending this to the case of correlated Gaussian features, the authors invoke a conjecture. In comparison, (12) is essentially free of assumptions on the design matrix.

The natural lasso also compares favorably to the method-of-moments-based estimator of Dicker (2014) in terms of mean squared error bounds. In particular, Dicker (2014) establishes a $O_P[(\sigma^{-2}\tau^2+1)\{n^{-2}(p+n)\}^{1/2}]$ relative mean squared error rate, where $\tau^2=\|\Sigma^{-1/2}\beta^*\|_2^2$ and Σ is the covariance of features X. This rate can be much slower for large p.

Notably, the mean squared error bound in Theorem 3 does not put any assumption on X, β^* , or σ^2 . In this sense, the result is analogous to a slow rate bound (Rigollet & Tsybakov 2011, Dalalyan et al. 2017), which appears in the lasso prediction consistency context. While it is well known (Sun & Zhang 2012) or can be easily verified that under stronger conditions, i.e., compatibility or restricted eigenvalue conditions, the naive estimator (3) based on the lasso and $\tilde{\sigma}^2_{\text{SQRT}}$ in (11) attain a faster rate, $O(|S|n^{-1}\log p)$, it is natural to ask whether these two estimators also attain a rate bound as in (12) when the conditions on X are not assumed. The following two results give an affirmative answer to this question.

Proposition 6. Under the conditions of Theorem 3, the naive estimator (3) based on the lasso estimator $\hat{\beta}_{\lambda}$ with $\lambda = 4\sigma(n^{-1}\log p)^{1/2}$ has the following bound with probability greater than $1 - p^{-1}$:

$$\left|\hat{\sigma}_{\text{naive}}^2 - \frac{\|\varepsilon\|_2^2}{n}\right| \le 16\sigma \|\beta^*\|_1 \left(\frac{\log p}{n}\right)^{1/2}.$$
 (13)

Relatedly, Chatterjee & Jafarov (2015) also consider a setting with no assumptions on X and derive an error bound $O\{\|\beta^*\|_1^{1/2}(n^{-1}\log p)^{1/4}\}$ for (3) for a lasso estimator $\hat{\beta}_{\lambda}$ with λ in (2) selected using a cross-validation procedure.

Lederer et al. (2016) derive a slow rate bound for the prediction error of the square root lasso. They show, in Lemma 2.1, that there exists a value of λ for which $\lambda = 3n^{-1/2}\|X^T\varepsilon\|_{\infty}\|y-X\tilde{\beta}_{\text{SQRT}}\|_2^{-1}$ and bound $\|X\tilde{\beta}_{\text{SQRT}}-X\beta^*\|_2^2$ at this value. The following result establishes the high-dimensional consistency of $\tilde{\sigma}_{\text{SQRT}}^2$ under no assumptions on X.

Proposition 7. Under the conditions of Theorem 3, the square-root/scaled lasso estimator $\tilde{\sigma}_{SQRT}^2$ in (11) based on $\tilde{\beta}_{SQRT}$ with $\lambda = 3n^{-1/2} \|X^T \varepsilon\|_{\infty} \|y - X \tilde{\beta}_{SQRT}\|_2^{-1}$ has the following bound with probability greater than $1 - p^{-1}$:

$$\left| \tilde{\sigma}_{SQRT}^2 - \frac{\|\varepsilon\|_2^2}{n} \right| \le 12\sigma \|\beta^*\|_1 \left(\frac{\log p}{n} \right)^{1/2}. \tag{14}$$

We see the rate of the natural lasso in (12) matches, up to a constant factor, the rates (13) and (14). The values of λ used in Propositions 6 and 7 are larger than would be necessary for standard prediction error bounds; we learned of this technique from Irina Gaynanova (Gaynanova 2018), and it is key to the proofs of the two propositions. Although the same rate is obtained in Theorem 3, Proposition 6, and Proposition 7, we have not established that this is the best possible rate obtainable in this setting that makes no assumption on X.

4 The organic lasso estimate of error variance

4.1 Method formulation

In practice, the value of the regularization parameter λ in (7) may be chosen via cross-validation; however, Theorem 3 has a regrettable theoretical shortcoming: it requires using a value of λ that itself depends on σ , the very quantity that we are trying to estimate! This is a well-known theoretical limitation of the lasso and related methods that motivated the square-root/scaled lasso. In this section, we propose a second new method, which retains the natural parameterization, but remedies the natural lasso's theoretical shortcoming by using a modified penalty. We define the organic lasso as a solution to (8) with $\Omega(\theta, \phi) = \phi^{-1} \|\theta\|_1^2$, i.e.,

$$(\check{\theta}_{\lambda}, \check{\phi}_{\lambda}) = \underset{\phi > 0, \ \theta}{\min} \left(-\frac{1}{2} \log \phi + \phi \frac{\|y\|_{2}^{2}}{2n} - \frac{1}{n} y^{T} X \theta + \frac{\|X\theta\|_{2}^{2}}{2n\phi} + \lambda \frac{\|\theta\|_{1}^{2}}{\phi} \right). \tag{15}$$

We observe that the penalty $\phi^{-1}\|\theta\|_1^2$ is jointly convex in (ϕ, θ) since it can be expressed as $g(h(\theta), \phi)$ where $h(\theta) = \|\theta\|_1$ is convex and $g(x, \phi) = \phi^{-1}x^2$ is a jointly convex function that is strictly increasing in x for $x \ge 0$ (Boyd & Vandenberghe 2004, Rockafellar 2015).

Given a solution to the above problem, we can reverse (6) to give the organic lasso estimators of (β^*, σ^2) , i.e., $\check{\beta}_{\lambda} = \check{\phi}_{\lambda}^{-1}\check{\theta}_{\lambda}$, $\check{\sigma}_{\lambda}^2 = \check{\phi}_{\lambda}^{-1}$. Furthermore, $\phi^{-1}\|\theta\|_1^2$ still induces sparsity in θ , and thus the final estimate $\check{\beta}_{\lambda}$ is sparse. In direct analogy to the natural lasso, the following proposition shows that we can find $\check{\sigma}_{\lambda}^2$ and $\check{\beta}_{\lambda}$ without actually solving (15).

Proposition 8. The organic lasso estimators $(\check{\beta}_{\lambda}, \check{\sigma}_{\lambda}^2)$ correspond to the solution and minimal value of an ℓ_1^2 -penalized least-squares problem:

$$\check{\beta}_{\lambda} = \underset{\beta \in \mathbb{R}^p}{\operatorname{arg\,min}} \left(\frac{1}{n} \| y - X\beta \|_2^2 + 2\lambda \| \beta \|_1^2 \right); \tag{16}$$

$$\check{\sigma}_{\lambda}^{2} = \min_{\beta \in \mathbb{R}^{p}} \left(\frac{1}{n} \|y - X\beta\|_{2}^{2} + 2\lambda \|\beta\|_{1}^{2} \right). \tag{17}$$

Thus, to compute the organic lasso estimator, one simply solves a penalized least squares problem, where the penalty is the square of the ℓ_1 norm. This can be thought of as the exclusive lasso with a single group (Zhou et al. 2010, Campbell et al. 2017). We show in the next section that solving this problem is no harder than solving a standard lasso problem.

One readily sees the connection of the organic lasso to the square-root lasso (11): to get (17), one takes squares of both the loss and the ℓ_1 penalty of (11). However, their origins are actually different in nature: the organic lasso is a maximum of the Gaussian log-likelihood with a scale-equivariant sparsity inducing penalty under parameterization (6), while (11) minimizes the ℓ_1 -penalized Huber concomitant loss function (Antoniadis 2010, Sun & Zhang 2012).

4.2 Algorithm

Coordinate descent is easy to implement and has steadily maintained its place as a start-of-the-art approach for solving lasso-related problems (Friedman et al. 2007). For coordinate descent to work, one typically verifies separability in the non-smooth part of the objective function (Tseng 2001). However, the ℓ_1^2 penalty in (16) is not separable in the coordinates of β . Lorbert et al. (2010) propose a coordinate descent algorithm to solve the Pairwise Elastic Net (PEN) problem, a generalization of (16), and a proof of the convergence of the algorithm is given in Lorbert (2012). In Algorithm 1, we give a coordinate descent algorithm specific to solving (16). The R package natural (Yu 2017) provides a C implementation of Algorithm 1.

Algorithm 1 A coordinate descent algorithm to solve (16)

```
Require: Initial estimate \beta^{(0)} \in \mathbb{R}^p, X \in \mathbb{R}^{n \times p}, y \in \mathbb{R}^n, and \lambda > 0.

Set \beta \leftarrow \beta^{(0)} and r \leftarrow y - X\beta

for j = 1, \dots, p; 1, \dots, p; \dots (until convergence) do
\beta_j^{\text{new}} \leftarrow (2\lambda + n^{-1} \|X_j\|_2^2)^{-1} \mathcal{S}(n^{-1}X_j^T r + n^{-1} \|X_j\|_2^2 \beta_j, 2\lambda \|\beta_{-j}\|_1)
r \leftarrow r + X_j \beta_j - X_j \beta_j^{\text{new}}
\beta_j \leftarrow \beta_j^{\text{new}}
end for
return \beta.
```

Each coordinate update requires O(n) operations. In Algorithm 1, $S(a,b) = \operatorname{sgn}(a)(|a|-b)_+$ is the soft-threshold operator. Empirically Algorithm 1 is as fast as solving a lasso problem. Theorem C.3.9 in Lorbert (2012) shows that, for any initial estimate $\beta^{(0)} \in \mathbb{R}^p$, every limit point of Algorithm 1 is an optimal point of the objective function of (16). This implies that the ℓ_1^2 penalty, although not separable, is well enough behaved that any point that is minimum in every coordinate of the objective function in (16) is indeed a global minimum.

4.3 Theoretical results

A first indication that the organic lasso may succeed where the natural lasso falls short is in terms of scale equivariance. As the design X is usually standardized to be unitless, scale equivariance in this context refers to the effect of scaling y.

Proposition 9. The organic lasso is scale equivariant, i.e., for any t > 0,

$$\check{\beta}_{\lambda}(ty) = t\check{\beta}_{\lambda}(y), \qquad \check{\sigma}_{\lambda}(ty) = t\check{\sigma}_{\lambda}(y).$$

Scale equivariance is a property associated with the ability to prove results in which the tuning parameter λ does not depend on σ . For example, the square-root/scaled lasso (11) is scale equivariant while the lasso, and thus the natural lasso, is not. In particular, $\hat{\beta}_{\lambda}(ty) \neq t\hat{\beta}_{\lambda}(y)$, and $\hat{\sigma}_{\lambda}(ty) \neq t\hat{\sigma}_{\lambda}(y)$ for some t > 0.

In Lemma 2, we saw how expressing an estimator as the optimal value of a convex optimization problem allows us to take full advantage of convex duality in order to derive bounds on the estimator. We therefore start our analysis of (17) by characterizing its dual problem.

Lemma 10. The dual problem of (17) is

$$\max_{u \in \mathbb{R}^n} \left\{ \frac{1}{n} \left(\|y\|_2^2 - \|y - u\|_2^2 \right) - \frac{1}{2\lambda} \left\| \frac{X^T u}{n} \right\|_{\infty}^2 \right\}.$$

Similar arguments as in Lemma 2 give a bound expressing $\check{\sigma}_{\lambda}^2$'s closeness to the oracle estimator of σ^2 .

Lemma 11. If $\lambda \geq n^{-1} \|X^T(\sigma^{-1}\varepsilon)\|_{\infty}$, then

$$-2\lambda\sigma^{2}\left(\frac{\left\|\beta^{*}\right\|_{1}}{\sigma}+\frac{1}{4}\right)\leq\check{\sigma}_{\lambda}^{2}-\frac{1}{n}\left\|\varepsilon\right\|_{2}^{2}\leq2\lambda\left\|\beta^{*}\right\|_{1}^{2}.$$

Comparing with Lemma 2, we see that the condition on λ depends only on a quantity $\sigma^{-1}\varepsilon \sim N(0, I_n)$ that is independent of σ^2 . Indeed, this leads to a mean squared error bound with the desired property of λ not depending on σ .

Theorem 12. Suppose that each column X_j of the matrix $X \in \mathbb{R}^{n \times p}$ has been scaled so that $\|X_j\|_2^2 = n$ for all $j = 1, \ldots, p$, and $\varepsilon \sim N\left(0, \sigma^2 I_n\right)$. Then, for any constant M > 1, the organic lasso estimator (17) with $\lambda = (2Mn^{-1}\log p)^{1/2}$ satisfies the following relative mean squared error bound:

$$E\left\{ \left(\frac{\check{\sigma}_{\lambda}^{2}}{\sigma^{2}} - 1 \right)^{2} \right\} \leq \left\{ \left(8M + 8\frac{p^{1-8M}}{\log p} \right)^{1/2} \max\left(\frac{\|\beta^{*}\|_{1}^{2}}{\sigma^{2}}, \frac{\|\beta^{*}\|_{1}}{\sigma} + \frac{1}{4} \right) \left(\frac{\log p}{n} \right)^{1/2} + \left(\frac{2}{n} \right)^{1/2} \right\}^{2}. \tag{18}$$

Compared with Theorem 3, the organic lasso estimator of σ^2 retains the same rate in terms of n and p but has a slower rate in terms of $\sigma^{-1} \|\beta^*\|_1$. Importantly, though, the value of λ attaining (18) does not depend on σ . This tuning-insensitive property is also enjoyed by the square-root/scaled lasso estimate of σ^2 , as shown in Proposition 7. As in Remark 5, similar high-probability bounds can be obtained for ε with bounded polynomial moments.

Although not central to our main purpose, the organic lasso estimator (16) of β^* is interesting in its own right. The following theorem gives a slow rate bound in prediction error.

Theorem 13. For any L > 0, the solution to (16) with $\lambda = \{2n^{-1}(\log p + L)\}^{1/2}$ has the following bound on the prediction error with probability greater than $1 - e^{-L}$:

$$\frac{1}{n} \| X \check{\beta}_{\lambda} - X \beta^* \|_2^2 \le \left(\sigma^2 + 4 \| \beta^* \|_1^2 \right) \left(\frac{2 \log p + 2L}{n} \right)^{1/2}.$$

In Appendix J, we provide mappings between the path of the natural lasso, $\{\hat{\beta}_{\lambda} : \lambda > 0\}$, and the path of the organic lasso $\{\check{\beta}_{\lambda} : \lambda > 0\}$. We also include a fast-rate prediction error bound of (16) under a compatibility condition in Appendix K.

5 Simulation studies

5.1 Simulation settings

Reid et al. (2016) carry out an extensive simulation study to compare many error variance estimators. We have matched their simulation settings fairly closely, so that the performance comparison with various other methods mentioned in Reid et al. (2016) can be inferred. Specifically, all simulations are run with p = 500 and n = 100. Each row of the design X is generated from a multivariate $N(0, \Sigma)$, with $\Sigma_{ij} = \rho \in (0, 1)$ for $i \neq j$ and $\Sigma_{ii} = 1$. To generate β^* , we randomly select the indices of $\lceil n^{\alpha} \rceil$ nonzero elements out of p variables where $\alpha \in (0, 1)$, and each of the nonzero elements has a value that is randomly drawn from a Laplace distribution with rate 1. The error variance is generated using $\sigma^2 = \tau^{-1}\beta^{*T}\Sigma\beta^*$ for $\tau > 0$. Finally, y is generated following (1).

Each model is indexed by a triplet (ρ, α, τ) , where ρ captures the correlation among features, α determines the sparsity of β^* , and τ characterizes the signal-to-noise ratio. We vary $\rho, \alpha \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$ and $\tau \in \{0.3, 1, 3\}$. We compute a Monte Carlo estimate of both the mean squared error $E\{(\sigma^{-1}\hat{\sigma}-1)^2\}$ and $E(\sigma^{-1}\hat{\sigma})$ as the measure of performance. The methods in comparison include (a) the naive estimator (3) with $\hat{\beta}_{\lambda}$ in (2), (b) the degrees of freedom adjusted estimator $\hat{\sigma}_R^2$ in (5) (Reid et al. 2016), (c) the square-root/scaled lasso (Belloni et al. 2011, Sun & Zhang 2013), (d) the natural lasso (7), and (e) the organic lasso (17). As a benchmark, we also include the oracle $n^{-1}\|\varepsilon\|_2^2$. The simulator R package (Bien 2016) was used for all simulations.

5.2 Methods with regularization parameter selected by cross-validation

We carry out two sets of simulations. In the first set, we compare the performance of the aforementioned methods with regularization parameter selected in a data-adaptive way. In particular, five-fold cross-validation is used to select the tuning parameter for each method.

Due to space constraints, we present a subset of the results in Fig 1. Additional results are presented in Appendix L. The result for the square-root/scaled lasso is averaged over 100 repetitions due to the large computational time. For all other methods, the results are averaged over 1000 repetitions. Overall, the natural lasso does well in adjusting the downward bias of the naive estimator, while other methods tend to produce under-estimates. In each panel, we fix signal-to-noise ratio (τ) and correlations among features (ρ) , and vary model sparsity (α) . All estimates get worse with growing α , except for the natural lasso, which improves as the true β^* gets denser. In particular, both the natural lasso and the organic lasso gain performance advantage over other methods when the underlying models do not satisfy conditions for the support recovery of the lasso solution. From left to right, Fig 1 illustrates the effect of increasing ρ . As observed in Reid et al. (2016), high correlations can be helpful: All curves approach the oracle as ρ increases. Finally, we find that the organic lasso is uniformly better or equivalent to $\hat{\sigma}_R^2$.

Paired t-tests and Wilcoxon signed-rank tests show that the differences in mean squared errors of different methods are significant at the 5% level for almost all points shown in Fig 1.

Results in Appendix L also show the natural lasso estimator doing well when the signal-to-noise ratio is low: the performances of all methods degrade as τ gets large. This is expected from Theorem 3 and Theorem 12, and is also observed in Reid et al. (2016).

5.3 Methods with fixed choice of regularization parameter

Although solving (17) is fast enough for one to use cross-validation with the organic lasso, Theorem 12 implies that $\lambda_0 = (2n^{-1}\log p)^{1/2}$ is a theoretically sound choice of regularization parameter. We also conjecture that a sharper rate may be obtainable at $\lambda_1 \geq n^{-2} \|X^T \epsilon\|_{\infty}^2$, where $\epsilon \sim N(0,1)$. With high probability, $n^{-2} \|X^T \epsilon\|_{\infty}^2 \approx n^{-1} \log(p)$. Thus, we also show the

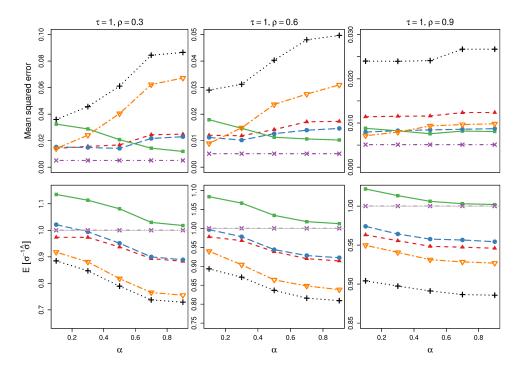


Figure 1: Simulation results of methods using cross-validation. From left to right, columns show Monte Carlo estimates of the mean squared error, in the top panel, and $E(\sigma^{-1}\hat{\sigma})$, in the bottom panel, of various methods in three simulation settings. Line styles and their corresponding methods: + for naive, \wedge for $\hat{\sigma}_R^2$, \nearrow for the square-root/scaled lasso, \neg for the natural lasso, \wedge for the oracle.

performance of the organic lasso with tuning parameter values equal to $\lambda_2 = n^{-1} \log(p)$, and λ_3 , which is a Monte Carlo estimate of $E(n^{-2}||X^T\epsilon||_{\infty}^2)$, where the expectation is with respect to $\epsilon \sim N(0,1)$.

We compare the organic lasso at these three fixed values of tuning parameter to the square-root/scaled lasso estimator (11) of error variance, which is another method whose theoretical choice of λ does not depend on σ . Sun & Zhang (2012) find that λ_0 works very well for (11), which we denote by scaled(1), and Sun & Zhang (2013) propose a refined choice of λ , which is proved to attain a sharper rate, denoted by scaled(2). The results of all the methods are averaged over 1000 repetitions.

Fig 2 shows similar patterns as Fig 1. Specifically, large value of ρ helps all methods, while performance generally degrades for denser β^* . Although not shown here, all methods struggle as τ increases. The theoretically justified tuning parameter λ_0 for the organic lasso appears in practice to overshrink the estimate of β^* and thus to overestimate σ^2 , leading to poor performance; however, the organic lasso with the smaller tuning parameter values λ_2 and λ_3 do quite well, generally outperforming the square-root/scaled lasso based methods.

6 Error estimation for Million Song dataset

We apply our error variance estimators to the Million Song dataset.¹ The data consist of information about 463715 songs, and the primary goal is to model the release year of a song using p=90 of its timbre features. The dataset has a very large sample size so that we can reliably estimate the ground truth of the target of estimation on a very large set of held out data. In particular, we randomly select half of the songs for this purpose and use $\bar{\sigma}^2 =$

 $^{^1}$ The whole data set can be obtained at https://labrosa.ee.columbia.edu/millionsong/. We consider a subset of the whole data, which is available at https://archive.ics.uci.edu/ml/datasets/yearpredictionmsd.

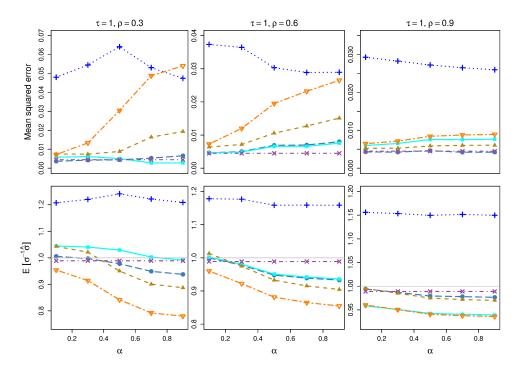


Figure 2: Simulation results of methods using pre-specified regularization parameter values. From left to right, columns show Monte Carlo estimates of the mean squared error, in the top panel, and $E(\sigma^{-1}\hat{\sigma})$, in the bottom panel, of various methods in three simulation settings. Line styles and their corresponding methods: + for organic (λ_0) , - for organic (λ_2) , - for scaled (1), - for scaled (2), - for the oracle.

 $(n-p)^{-1}||y-X\hat{\beta}_{LS}||_2^2$ to form our ground truth, where $\hat{\beta}_{LS}$ is the least-squares estimator of β^* . In practice, model (1) may rarely hold, which alters the interpretation of error variance estimation. Suppose the response vector y has mean μ and covariance matrix Σ . Then $\bar{\sigma}^2$ can be thought of as an estimator of the population quantity

$$\min_{\beta} \frac{1}{n} E\left(\left\| y - X\beta \right\|_2^2 \right) = \frac{1}{n} \operatorname{tr}(\Sigma) + \frac{1}{n} \left\| \left(I - XX^+ \right) \mu \right\|_2^2.$$

In the special case where $\Sigma = \sigma^2 I_n$ and $\mu = X\beta^*$, as in (1), then $\bar{\sigma}^2$ reduces to the linear model noise variance σ^2 .

From the remaining data that was not previously used to yield $\bar{\sigma}^2$, we randomly form training datasets of size n and compare the performance of various error variance estimators. We vary n in $\{20, 40, 60, 80, 100, 120\}$ to gauge the performance of these methods in situations in which n < p and $n \approx p$. For each n, we repeat the data selection and error variance estimation on 1000 disjoint training sets, and report estimates of the mean squared error $E\{(\bar{\sigma}^{-1}\hat{\sigma}-1)^2\}$ in Table 1 and estimates of $E(\bar{\sigma}^{-1}\hat{\sigma})$ in Appendix L.

All methods produce a substantial performance improvement over the naive estimator for a wide range of values of n. The natural and organic lassos with cross validation perform either better or comparably to $\hat{\sigma}_R^2$ and are in some, but not all, cases outperformed by scaled(2). When n gets large, the natural lasso shows some upward bias, which as we noted before is less problematic than downward bias. The organic lasso with the fixed choices λ_2 or λ_3 performs extremely well for all n.

Future research directions include the analysis of the proposed methods with smaller values of λ , and extending the natural parameterization to penalized non-parametric regression. Finally, an R (R Core Team 2017) package, named natural (Yu 2017), is available on the Comprehensive R Archive Network, implementing our estimators.

Table 1: Mean squared error of noise variance estimation for Million Song dataset

n	20	40	60	80	100	120
naive	17.02 (0.68)	8.48 (0.41)	5.28 (0.26)	3.80 (0.17)	3.03 (0.13)	2.43 (0.10)
$\hat{\sigma}_R^2$	$10.74 \ (0.45)$	5.92(0.29)	3.57(0.17)	2.57(0.11)	2.23(0.10)	1.75 (0.08)
natural(cv)	8.82 (0.38)	5.23(0.27)	3.47(0.16)	2.61 (0.12)	2.39(0.11)	2.01 (0.09)
$\operatorname{organic}(\operatorname{cv})$	8.08(0.32)	4.23(0.20)	2.59(0.12)	2.00(0.08)	1.72(0.08)	1.54 (0.07)
scaled(1)	7.43(0.37)	4.92(0.25)	3.84(0.17)	3.08(0.13)	2.94(0.12)	2.75(0.11)
scaled(2)	7.11(0.28)	3.36 (0.15)	2.23(0.10)	2.57(0.83)	1.61 (0.07)	1.46 (0.07)
$\operatorname{organic}(\lambda_2)$	5.87(0.24)	3.17(0.14)	1.93(0.09)	$1.40 \ (0.06)$	$1.20 \ (0.05)$	1.02(0.05)
$\operatorname{organic}(\lambda_3)$	5.72(0.24)	3.15(0.14)	1.99(0.09)	1.45 (0.07)	1.28 (0.05)	1.12 (0.05)

Mean and standard errors, over 1000 replications, of the squared error of various methods. Each entry is multiplied by 100 to convey information more compactly.

Acknowledgement

We thank Irina Gaynanova for a useful conversation that helped us prove Propositions 6 and 7. JB was supported by an NSF CAREER grant, DMS-1653017.

Appendices

A Proof of Lemma 2

From (2) in the paper, it follows that

$$\hat{\sigma}_{\lambda}^{2} \leq \frac{1}{n} \|y - X\beta^{*}\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1} = \frac{1}{n} \|\varepsilon\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1}.$$

By introducing the dual variable $2n^{-1}u \in \mathbb{R}^n$,

$$\begin{split} \hat{\sigma}_{\lambda}^{2} &= \min_{\beta} \left(\frac{1}{n} \left\| y - X\beta \right\|_{2}^{2} + 2\lambda \left\| \beta \right\|_{1} \right) = \min_{\beta, z} \max_{u} \left\{ \frac{1}{n} \left\| y - z \right\|_{2}^{2} + \frac{2}{n} u^{T} \left(z - X\beta \right) + 2\lambda \left\| \beta \right\|_{1} \right\} \\ &\geq \max_{u} \min_{\beta, z} \left\{ \frac{1}{n} \left\| y - z \right\|_{2}^{2} + \frac{2}{n} u^{T} \left(z - X\beta \right) + 2\lambda \left\| \beta \right\|_{1} \right\} \\ &= \max_{u} \left(\frac{1}{n} \left\| y \right\|_{2}^{2} - \frac{1}{n} \left\| y - u \right\|_{2}^{2}, \text{subject to } \left\| X^{T} u \right\|_{\infty} \leq n\lambda \right). \end{split}$$

By assumption, ε is dual feasible, which means that

$$\hat{\sigma}_{\lambda}^{2} \geq \frac{1}{n} \|y\|_{2}^{2} - \frac{1}{n} \|y - \varepsilon\|_{2}^{2} \geq \frac{1}{n} \|\varepsilon\|_{2}^{2} + \frac{2}{n} \varepsilon^{T} X \beta^{*} \geq \frac{1}{n} \|\varepsilon\|_{2}^{2} - 2\lambda \|\beta^{*}\|_{1},$$

where in the last step we applied Hölder's inequality.

B Proof of Propositions 1 and 8

We prove in this section that both the natural lasso and the organic lasso estimates of error variance can be simply expressed as the minimizing values of certain convex optimization problems. To do so, we exploit the first order optimality condition of each convex program.

We start with proving that the natural lasso estimate of σ^2 is the minimal value of a lasso problem (2). The following lemma characterizes the conditions for which $(\hat{\theta}_{\lambda}, \hat{\phi}_{\lambda})$ is a solution to (8) with $\Omega(\theta, \phi) = \|\theta\|_1$.

Lemma 14 (Optimality condition of the natural lasso). For any $\lambda > 0$, $(\hat{\theta}_{\lambda}, \hat{\phi}_{\lambda})$ is a solution to (8) with $\Omega(\theta, \phi) = \|\theta\|_1$ if and only if

$$-\frac{1}{\hat{\phi}_{\lambda}} + \frac{1}{n} \|y\|_2^2 - \frac{\left\|X\hat{\theta}_{\lambda}\right\|_2^2}{n\hat{\phi}_{\lambda}^2} = 0, \qquad -X^T y + X^T X \frac{\hat{\theta}_{\lambda}}{\hat{\phi}_{\lambda}} + n\lambda \hat{g} = 0$$

where $\hat{g} \in \partial(\|\hat{\theta}_{\lambda}\|_1)$.

Given $(\hat{\theta}_{\lambda}, \hat{\phi}_{\lambda})$, we reverse the natural parameterization to get $\hat{\beta}_{\lambda} = \hat{\phi}_{\lambda}^{-1} \hat{\theta}_{\lambda}$ and $\hat{\sigma}_{\lambda}^{2} = \hat{\phi}_{\lambda}^{-1}$. From Lemma 14,

$$\hat{\sigma}_{\lambda}^{2} = \frac{1}{n} \left(\|y\|_{2}^{2} - \left\| X \hat{\beta}_{\lambda} \right\|_{2}^{2} \right) \quad \text{and} \quad 0 = -\hat{\beta}_{\lambda}^{T} X^{T} y + \left\| X \hat{\beta}_{\lambda} \right\|_{2}^{2} + n\lambda \left\| \hat{\beta}_{\lambda} \right\|_{1}^{2}.$$

Note that

$$\left\|y-X\hat{\beta}_{\lambda}\right\|_{2}^{2}=\left\|y\right\|_{2}^{2}-\left\|X\hat{\beta}_{\lambda}\right\|_{2}^{2}+2\left(\left\|X\hat{\beta}_{\lambda}\right\|_{2}^{2}-y^{T}X\hat{\beta}_{\lambda}\right)=\left\|y\right\|_{2}^{2}-\left\|X\hat{\beta}_{\lambda}\right\|_{2}^{2}-2n\lambda\left\|\hat{\beta}_{\lambda}\right\|_{1}.$$

We have

$$\hat{\sigma}_{\lambda}^2 = \frac{1}{n} \left(\|y\|_2^2 - \left\| X \hat{\beta}_{\lambda} \right\|_2^2 \right) = \frac{1}{n} \left\| y - X \hat{\beta}_{\lambda} \right\|_2^2 + 2\lambda \left\| \hat{\beta}_{\lambda} \right\|_1.$$

We show that the organic lasso estimate of σ^2 is the minimal value of the ℓ_1^2 -penalized least squares problem. As the natural lasso, we start with studying the following optimality condition:

Lemma 15 (Optimality condition of the organic lasso). For any $\lambda > 0$, $(\check{\theta}_{\lambda}, \check{\phi}_{\lambda})$ is a solution to (15) if and only if

$$-\frac{1}{\check{\phi}_{\lambda}} + \frac{1}{n} \|y\|_2^2 - \frac{\left\|X\check{\theta}_{\lambda}\right\|_2^2}{n\check{\phi}_{\lambda}^2} - 2\lambda \frac{\left\|\check{\theta}_{\lambda}\right\|_1^2}{\check{\phi}_{\lambda}^2} = 0, \qquad -X^T y + X^T X \frac{\check{\theta}_{\lambda}}{\check{\phi}_{\lambda}} + 2n\lambda \frac{\left\|\check{\theta}_{\lambda}\right\|_1}{\check{\phi}_{\lambda}} \check{g} = 0$$

where $\check{g} \in \partial(||\check{\theta}||_1)$.

So following the natural parameterization, we have that $\check{\beta}_{\lambda} = \check{\theta}_{\lambda}^{-1} \check{\rho}_{\lambda}$ and $\check{\sigma}_{\lambda}^{2} = \check{\rho}_{\lambda}^{-1}$, and

$$\check{\sigma}_{\lambda}^{2} = \frac{1}{n} \left(\|y\|_{2}^{2} - \|X\check{\beta}_{\lambda}\|_{2}^{2} - 2n\lambda \|\check{\beta}_{\lambda}\|_{1}^{2} \right)
0 = -\check{\beta}_{\lambda}^{T} X^{T} y + \|X\check{\beta}_{\lambda}\|_{2}^{2} + 2n\lambda \|\check{\beta}_{\lambda}\|_{1}^{2}.$$

Note that

$$\begin{aligned} \|y - X \check{\beta}_{\lambda}\|_{2}^{2} &= \|y\|_{2}^{2} + \|X \check{\beta}_{\lambda}\|_{2}^{2} - 2y^{T} X \check{\beta}_{\lambda} \\ &= \|y\|_{2}^{2} - \|X \check{\beta}_{\lambda}\|_{2}^{2} + 2\left(\|X \check{\beta}_{\lambda}\|_{2}^{2} - y^{T} X \check{\beta}_{\lambda}\right) \\ &= \|y\|_{2}^{2} - \|X \check{\beta}_{\lambda}\|_{2}^{2} - 4n\lambda \|\check{\beta}_{\lambda}\|_{1}^{2}. \end{aligned}$$

We have

$$\check{\sigma}_{\lambda}^2 = \frac{1}{n} \left(\left\| y \right\|_2^2 - \left\| X \check{\beta}_{\lambda} \right\|_2^2 - 2n\lambda \left\| \check{\beta}_{\lambda} \right\|_1^2 \right) = \frac{1}{n} \left\| y - X \check{\beta}_{\lambda} \right\|_2^2 + 2\lambda \left\| \check{\beta}_{\lambda} \right\|_1^2.$$

C Proof of Lemma 10: the dual problem of the ℓ_1^2 -penalized least squares

The primal problem of the ℓ_1^2 -penalized least squares (16) in the paper can be written as an equality constrained minimization problem:

$$\min_{\beta \in \mathbb{R}^p} \left(\frac{1}{n} \|y - z\|_2^2 + 2\lambda \|\beta\|_1^2 \quad \text{s.t.} \quad \frac{2}{n} z = \frac{2}{n} X \beta \right).$$

The Lagrange dual function is

$$\begin{split} g\left(u\right) &= \min_{\beta \in \mathbb{R}^{p}, z \in \mathbb{R}^{n}} \left\{ \frac{1}{n} \left\| y - z \right\|_{2}^{2} + 2\lambda \left\| \beta \right\|_{1}^{2} + \frac{2u^{T}}{n} \left(z - X\beta \right) \right\} \\ &= \min_{z \in \mathbb{R}^{n}} \left(\frac{1}{n} \left\| y - z \right\|_{2}^{2} + \frac{2}{n} u^{T} z \right) + \min_{\beta \in \mathbb{R}^{p}} \left\{ 2\lambda \left\| \beta \right\|_{1}^{2} - 2\left(\frac{X^{T} u}{n} \right)^{T} \beta \right\}. \end{split}$$

The minimization of u is

$$\min_{z \in \mathbb{R}^n} \left(\frac{1}{n} \|y - z\|_2^2 + \frac{2}{n} u^T z \right) = \frac{2}{n} u^T y - \frac{1}{n} \|u\|_2^2 = \frac{1}{n} \left(\|y\|_2^2 - \|y - u\|_2^2 \right),$$

where the minimum is attained at

$$\hat{z} = y - u.$$

The minimization problem of β can be written as

$$\min_{\beta \in \mathbb{R}^p} \left\{ 2\lambda \left\| \beta \right\|_1^2 - 2\left(\frac{X^T u}{n}\right)^T \beta \right\} = -2\lambda \max_{\beta \in \mathbb{R}^p} \left\{ \left(\frac{X^T u}{\lambda n}\right)^T \beta - \left\| \beta \right\|_1^2 \right\}.$$

Observe that the maximum is the Fenchel conjugate function of $\|\cdot\|_1^2$, evaluated at $(\lambda n)^{-1}X^Tu$. By Boyd & Vandenberghe (2004, Example 3.27, pp. 92-93),

$$-2\lambda \max_{\beta \in \mathbb{R}^p} \left\{ \left(\frac{X^T u}{\lambda n}\right)^T \beta - \left\|\beta\right\|_1^2 \right\} = -\frac{2\lambda}{4} \left\|\frac{X^T u}{\lambda n}\right\|_{\infty}^2 = -\frac{1}{2\lambda} \left\|\frac{X^T u}{n}\right\|_{\infty}^2.$$

So

$$g(u) = \frac{1}{n} \left(\|y\|_2^2 - \|y - u\|_2^2 \right) - \frac{1}{2\lambda} \left\| \frac{X^T u}{n} \right\|_{\infty}^2.$$

D Proof of Lemma 11

A direct upper bound is

$$\check{\sigma}_{\lambda}^{2} \leq \frac{1}{n} \|y - X\beta^{*}\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1}^{2} = \frac{1}{n} \|\varepsilon\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1}^{2}.$$

To get a lower bound of $\hat{\sigma}^2$, note that the dual problem in Lemma 10 and the strong duality imply that

$$\begin{split} \check{\sigma}_{\lambda}^2 &= \min_{\beta \in \mathbb{R}^p} \left(\frac{1}{n} \left\| y - X \beta \right\|_2^2 + 2\lambda \left\| \beta \right\|_1^2 \right) = \max_{u \in \mathbb{R}^n} \left(\frac{1}{n} \left\| y \right\|_2^2 - \frac{1}{n} \left\| y - u \right\|_2^2 - \frac{1}{2\lambda} \left\| \frac{X^T u}{n} \right\|_{\infty}^2 \right) \\ &\geq \frac{1}{n} \left\| y \right\|_2^2 - \frac{1}{n} \left\| y - \varepsilon \right\|_2^2 - \frac{1}{2\lambda} \left\| \frac{X^T \varepsilon}{n} \right\|_{\infty}^2 = \frac{1}{n} \left\| \varepsilon \right\|_2^2 + \frac{2}{n} \varepsilon^T X \beta^* - \frac{1}{2\lambda} \left\| \frac{X^T \varepsilon}{n} \right\|_{\infty}^2 \\ &\geq \frac{1}{n} \left\| \varepsilon \right\|_2^2 - 2 \left\| \frac{X^T \varepsilon}{n} \right\|_{\infty} \left\| \beta^* \right\|_1 - \frac{1}{2\lambda} \left\| \frac{X^T \varepsilon}{n} \right\|_{\infty}^2 \geq \frac{1}{n} \left\| \varepsilon \right\|_2^2 - 2\lambda \sigma^2 \left(\frac{\left\| \beta^* \right\|_1}{\sigma} + \frac{1}{4} \right), \end{split}$$

where the last inequality holds for

$$\lambda \geq \frac{\|X^T \varepsilon\|_{\infty}}{n\sigma}.$$

E Proof of Theorem 3 and Theorem 12

We present in this section the proof of Theorem 12. The proof of Theorem 3 follows the same set of arguments. First we use the following lemma to characterize the event that $\lambda \geq n^{-1}\sigma^{-1}\|X^T\varepsilon\|_{\infty}$ is true, so that we can use Lemma 11 to prove a high probability bound.

Lemma 16 (Corollary 4.3, Giraud (2014)). Assume that each column X_j of the design matrix $X \in \mathbb{R}^{n \times p}$ satisfies $||X_j||_2^2 = n$ for all $j = 1, \ldots, p$, and $\varepsilon \sim N\left(0, \sigma^2 I_n\right)$. Then for any L > 0,

$$\operatorname{pr}\left\{\frac{\left\|X^T\varepsilon\right\|_{\infty}}{n\sigma} > \left(\frac{2\log p + 2L}{n}\right)^{1/2}\right\} \leq e^{-L}.$$

Lemma 16 implies that a good choice of the value of λ would be $\{n^{-1}(2 \log p + 2L)\}^{1/2}$, which does not depend on any parameter of the underlying model. The following corollary shows that with this value of λ , the organic lasso estimate of σ^2 is close to the oracle estimator with high probability.

Corollary 17. Assume that each column X_j of the design matrix $X \in \mathbb{R}^{n \times p}$ satisfies $||X_j||_2^2 = n$ for all $j = 1, \ldots, p$, and $\varepsilon \sim N\left(0, \sigma^2 I_n\right)$. Then for any L > 0, the organic lasso with

$$\lambda = \left(\frac{2\log p + 2L}{n}\right)^{1/2}$$

has the following bound

$$\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \|\varepsilon\|_{2}^{2}\right)^{2} \leq 8 \max\left\{\left\|\beta^{*}\right\|_{1}^{2}, \sigma^{2}\left(\frac{\|\beta^{*}\|_{1}}{\sigma} + \frac{1}{4}\right)\right\}^{2} \frac{\log p + L}{n}$$

with probability greater than $1 - e^{-L}$.

In general, a high probability bound does not necessarily imply an expectation bound. However, when the probability bound holds with an exponential tail, it implies an expectation bound with essentially the same rate.

Theorem 18. Assume that each column X_j of the design matrix $X \in \mathbb{R}^{n \times p}$ satisfies $||X_j||_2^2 = n$ for all $j = 1, \ldots, p$, and $\varepsilon \sim N\left(0, \sigma^2 I_n\right)$. Then, for any constant M > 1, the organic lasso estimate with

$$\lambda = \left(\frac{2M\log p}{n}\right)^{1/2}$$

satisfies the following bound in expectation:

$$\mathbb{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n}\left\|\varepsilon\right\|_{2}^{2}\right)^{2}\right\} \leq 8\left(M + \frac{p^{1-M}}{\log p}\right) \max\left\{\left\|\beta^{*}\right\|_{1}^{2}, \sigma^{2}\left(\frac{\left\|\beta^{*}\right\|_{1}}{\sigma} + \frac{1}{4}\right)\right\}^{2} \frac{\log p}{n}.$$

Proof. For any M > 1, take $L = (M - 1) \log p$ in Corollary 17. Denote $X_n = (\check{\sigma}_{\lambda}^2 - n^{-1} \|\varepsilon\|^2)^2$, and $r_n = 8 \max(\|\beta^*\|_1^2, \sigma \|\beta^*\|_1 + 4^{-1}\sigma^2)^2 n^{-1} \log p$. Then we have

$$\operatorname{pr}(X_n > Mr_n) \le e^{-(M-1)\log p}.$$

So

$$E\left(\frac{X_n}{r_n}\right) = \int_0^\infty \operatorname{pr}\left(\frac{X_n}{r_n} > t\right) dt = \int_0^M \operatorname{pr}\left(\frac{X_n}{r_n} > t\right) dt + \int_M^\infty \operatorname{pr}\left(\frac{X_n}{r_n} > t\right) dt$$

$$\leq M + \int_M^\infty e^{-(t-1)\log p} dt = M + \frac{p^{1-M}}{\log p},$$

and the expectation bound follows.

Now we are ready to present the proof of Theorem 12. Since $\sigma^{-2} \|\varepsilon\|_2^2 \sim \chi^2(n)$, we have

$$\mathrm{E}\left(\frac{1}{n}\left\|\varepsilon\right\|_{2}^{2}\right) = \sigma^{2}, \qquad \mathrm{var}\left(\frac{1}{n}\left\|\varepsilon\right\|_{2}^{2}\right) = \frac{2\sigma^{4}}{n},$$

Therefore,

$$\begin{split} & \operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \sigma^{2}\right)^{2}\right\} = \operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2} + \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2} - \sigma^{2}\right)^{2}\right\} \\ & = \operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)^{2}\right\} + \operatorname{E}\left\{\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2} - \sigma^{2}\right)^{2}\right\} + 2\operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2} - \sigma^{2}\right)\right\} \\ & \leq \operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)^{2}\right\} + \operatorname{var}\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right) + 2\left\{\operatorname{var}\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)\operatorname{var}\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)\right\}^{1/2} \\ & \leq \operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)^{2}\right\} + \operatorname{var}\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right) + 2\left[\operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)^{2}\right\}\operatorname{var}\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)\right]^{1/2} \\ & = \left[\left[\operatorname{E}\left\{\left(\check{\sigma}_{\lambda}^{2} - \frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)^{2}\right\}\right]^{1/2} + \left\{\operatorname{var}\left(\frac{1}{n} \left\|\varepsilon\right\|_{2}^{2}\right)\right\}^{1/2}\right]^{2} \\ & \leq \left[\left\{8\left(M + \frac{p^{1-M}}{\log p}\right)\right\}^{1/2} \operatorname{max}\left\{\left\|\beta^{*}\right\|_{1}^{2}, \sigma^{2}\left(\frac{\left\|\beta^{*}\right\|_{1}}{\sigma} + \frac{1}{4}\right)\right\}\left(\frac{\log p}{n}\right)^{1/2} + \sigma^{2}\left(\frac{2}{n}\right)^{1/2}\right]^{2}, \end{split}$$

where the last inequality holds from Theorem 18.

F Proof of Remark 5

For the independent zero-mean noise ε_i with variance σ^2 and bounded m-th order moment $(m=3,4,\ldots)$

$$E |\varepsilon_i|^m \le \frac{m!}{2} K^{m-2}$$

for some constant K>0, a Bernstein's type inequality (Bühlmann & Van De Geer 2011, Lemma 14.13) implies that

$$\operatorname{pr}\left[\max_{1\leq j\leq p}\frac{1}{n\sigma}\left\|X_{j}^{T}\varepsilon\right\|_{\infty}\geq \frac{2K\log p}{n}+2\left\{\frac{\log(2p)}{n}\right\}^{1/2}\right]\leq \frac{1}{p}.$$

Then the proof of Corollary 17 goes through.

G Proof of Proposition 6 and Proposition 7

The following lemma gives a general result on the estimation error of $\hat{\sigma}^2$ of the form (3) in the paper based on $\hat{\beta}$:

Lemma 19.

$$\left| \hat{\sigma}^2 - \frac{1}{n} \|\varepsilon\|_2^2 \right| \le \frac{1}{n} \left\| X \hat{\beta} - X \beta^* \right\|_2^2 + \frac{2}{n} \left\| X^T \varepsilon \right\|_{\infty} \left(\|\beta^*\|_1 + \|\hat{\beta}\|_1 \right)$$

Proof. First by definition

$$\hat{\sigma}^2 = \frac{1}{n} \left\| y - X \hat{\beta} \right\|_2^2 = \frac{1}{n} \left\| \varepsilon + X \beta^* - X \hat{\beta} \right\|_2^2 = \frac{1}{n} \| \varepsilon \|_2^2 + \frac{1}{n} \left\| X \hat{\beta} - X \beta^* \right\|_2^2 + \frac{2}{n} \varepsilon^T X \left(\hat{\beta} - \beta^* \right).$$

Note that

$$\left| \varepsilon^T X \left(\hat{\beta} - \beta^* \right) \right| \le \| X^T \varepsilon \|_{\infty} \| \hat{\beta} - \beta^* \|_1,$$

and the result follows.

G.1 Slow rate bound for the naive estimator of σ^2

We now give the proof of Proposition 6. From the basic inequality

$$\frac{1}{n} \left\| y - X \hat{\beta}_{\lambda} \right\|_{2}^{2} + 2\lambda \|\hat{\beta}_{\lambda}\|_{1} \le \frac{1}{n} \|y - X \beta^{*}\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1},$$

which implies that

$$\frac{1}{n} \left\| X \hat{\beta}_{\lambda} - X \beta^* \right\|_{2}^{2} + 2\lambda \|\hat{\beta}_{\lambda}\|_{1} \leq \frac{2}{n} \left| \varepsilon^{T} X \left(\hat{\beta}_{\lambda} - \beta^* \right) \right| + 2\lambda \|\beta^*\|_{1} \\
\leq \frac{2}{n} \left\| X^{T} \varepsilon \right\|_{\infty} \left\| \hat{\beta}_{\lambda} - \beta^* \right\|_{1} + 2\lambda \|\beta^*\|_{1}.$$

We thank Irina Gaynanova (Gaynanova 2018) for showing us the technique of taking λ to be twice its usual size. For $\lambda \geq 2n^{-1} ||X^T \varepsilon||_{\infty}$, we have that

$$\frac{1}{n} \left\| X \hat{\beta}_{\lambda} - X \beta^* \right\|_2^2 + 2\lambda \|\hat{\beta}_{\lambda}\|_1 \le \lambda \|\hat{\beta}_{\lambda} - \beta^*\|_1 + 2\lambda \|\beta^*\|_1 \le \lambda \|\hat{\beta}_{\lambda}\|_1 + 3\lambda \|\beta^*\|_1,$$

so $n^{-1} \| X \hat{\beta}_{\lambda} - X \beta^* \|_2^2 + \lambda \| \hat{\beta}_{\lambda} \|_1 \le 3\lambda \| \beta^* \|_1$. So by Lemma 19 we have

$$\left| \hat{\sigma}_{\text{naive}}^{2} - \frac{1}{n} \|\varepsilon\|_{2}^{2} \right| \leq \frac{1}{n} \left\| X \hat{\beta}_{\lambda} - X \beta^{*} \right\|_{2}^{2} + \frac{2}{n} \left\| X^{T} \varepsilon \right\|_{\infty} \left(\|\beta^{*}\|_{1} + \|\hat{\beta}_{\lambda}\|_{1} \right)$$

$$\leq \frac{1}{n} \left\| X \hat{\beta}_{\lambda} - X \beta^{*} \right\|_{2}^{2} + \lambda \|\beta^{*}\|_{1} + \lambda \|\hat{\beta}_{\lambda}\|_{1} \leq 4\lambda \|\beta^{*}\|_{1}.$$

Finally, taking $\lambda = 2\sigma \{n^{-1}(2\log p + 2L)\}^{1/2}$ with $L = \log p$, the result follows from Lemma 16.

G.2 Slow rate bound for the square-root/scaled lasso estimator of σ^2

As shown in Lederer et al. (2016) (proof of Lemma A.3), we note that with probability 1, $||y - X\tilde{\beta}_{SQRT}||_2 > 0$ for $\lambda > 0$. So the first order optimality condition of the square-root/scaled lasso is

$$\frac{1}{n^{1/2}} \frac{-X^T \left(y - X \tilde{\beta}_{SQRT} \right)}{\left\| y - X \tilde{\beta}_{SQRT} \right\|_2} + \lambda \hat{g} = 0$$

for some $\hat{g} \in \partial \|\tilde{\beta}_{SQRT}\|_1$. Taking an inner product with $\tilde{\beta}_{SQRT} - \beta^*$ on both sides, we have

$$-\frac{1}{n^{1/2}} \frac{\left(\tilde{\beta}_{SQRT} - \beta^*\right)^T X^T \left(y - X\tilde{\beta}_{SQRT}\right)}{\left\|y - X\tilde{\beta}_{SQRT}\right\|_2} + \lambda \hat{g}^T \left(\tilde{\beta}_{SQRT} - \beta^*\right) = 0,$$

which implies that

$$\frac{\left\|X\left(\beta^* - \tilde{\beta}_{SQRT}\right)\right\|_{2}^{2}}{n^{1/2} \left\|y - X\tilde{\beta}_{SQRT}\right\|_{2}^{2}} - \frac{\left(\tilde{\beta}_{SQRT} - \beta^{*}\right)^{T} X^{T} \varepsilon}{n^{1/2} \left\|y - X\tilde{\beta}_{SQRT}\right\|_{2}^{2}} \leq \lambda \hat{g}^{T} \left(\beta^{*} - \tilde{\beta}_{SQRT}\right) \leq \lambda \|\beta^{*}\|_{1} - \lambda \|\tilde{\beta}_{SQRT}\|_{1},$$

and thus

$$\begin{split} &\frac{1}{n} \left\| X \left(\beta^* - \tilde{\beta}_{\mathrm{SQRT}} \right) \right\|_{2}^{2} \leq \frac{1}{n} \left| \varepsilon^{T} X \left(\tilde{\beta}_{\mathrm{SQRT}} - \beta^* \right) \right| + \frac{\lambda}{n^{1/2}} \left\| y - X \tilde{\beta}_{\mathrm{SQRT}} \right\|_{2} \left(\| \beta^* \|_{1} - \| \tilde{\beta}_{\mathrm{SQRT}} \|_{1} \right) \\ &\leq \frac{1}{n} \left\| X^{T} \varepsilon \right\|_{\infty} \left\| \tilde{\beta}_{\mathrm{SQRT}} - \beta^* \right\|_{1} + \frac{\lambda}{n^{1/2}} \left\| y - X \tilde{\beta}_{\mathrm{SQRT}} \right\|_{2} \left(\| \beta^* \|_{1} - \| \tilde{\beta}_{\mathrm{SQRT}} \|_{1} \right) \\ &\leq \frac{1}{n} \left\| X^{T} \varepsilon \right\|_{\infty} \left(\| \tilde{\beta}_{\mathrm{SQRT}} \|_{1} + \| \beta^* \|_{1} \right) + \frac{\lambda}{n^{1/2}} \left\| y - X \tilde{\beta}_{\mathrm{SQRT}} \right\|_{2} \left(\| \beta^* \|_{1} - \| \tilde{\beta}_{\mathrm{SQRT}} \|_{1} \right) \\ &\leq \left(\frac{1}{n} \left\| X^{T} \varepsilon \right\|_{\infty} + \frac{\lambda}{n^{1/2}} \left\| y - X \tilde{\beta}_{\mathrm{SQRT}} \right\|_{2} \right) \| \beta^* \|_{1} + \left(\frac{1}{n} \left\| X^{T} \varepsilon \right\|_{\infty} - \frac{\lambda}{n^{1/2}} \left\| y - X \tilde{\beta}_{\mathrm{SQRT}} \right\|_{2} \right) \| \tilde{\beta}_{\mathrm{SQRT}} \|_{1}. \end{split}$$

Taking $\lambda = 3n^{-1/2} ||y - X\tilde{\beta}_{SQRT}||_2^{-1} ||X^T \varepsilon||_{\infty}$, which is 3 times what is suggested in Lederer et al. (2016), we have

$$\frac{1}{n} \left\| X \left(\beta^* - \tilde{\beta}_{SQRT} \right) \right\|_2^2 \le \frac{4 \left\| X^T \varepsilon \right\|_{\infty}}{n} \|\beta^*\|_1 - \frac{2 \|X^T \varepsilon\|_{\infty}}{n} \|\tilde{\beta}_{SQRT}\|_1.$$

By Lemma 19

$$\left| \tilde{\sigma}_{\text{SQRT}}^2 - \frac{1}{n} \|\varepsilon\|_2^2 \right| \leq \frac{1}{n} \left\| X \tilde{\beta}_{\text{SQRT}} - X \beta^* \right\|_2^2 + \frac{2}{n} \left\| X^T \varepsilon \right\|_{\infty} \left(\|\beta^*\|_1 + \|\tilde{\beta}_{\text{SQRT}}\|_1 \right)$$

$$\leq \frac{6}{n} \left\| X^T \varepsilon \right\|_{\infty} \|\beta^*\|_1.$$

The result then follows from Lemma 16 by taking $L = \log p$.

H Proof of Proposition 9: scale-equivariance of the organic lasso

Proof. Suppose $\check{\beta}_{\lambda}(y)$ is a solution to the organic lasso, where we write out explicitly the dependence of the solution on the response y. Then using notation from previous section,

$$L\left(t\check{\beta}_{\lambda}\left(y\right)|ty,\lambda\right) = \frac{1}{n}\left\|ty - tX\check{\beta}_{\lambda}\left(y\right)\right\|_{2}^{2} + 2\lambda\left\|t\check{\beta}_{\lambda}\left(y\right)\right\|_{1}^{2}$$
$$= t^{2}L\left(\check{\beta}_{\lambda}\left(y\right)|y,\lambda\right).$$

This implies that $t\check{\beta}_{\lambda}(y)$ is a solution to the problem with response ty, i.e., $\check{\beta}_{\lambda}(ty) = t\check{\beta}(y)$. Consequently,

$$\begin{split} \check{\sigma}_{\lambda}^{2}\left(ty\right) &= \min_{\beta} L\left(\beta_{\lambda}|ty,\lambda\right) \\ &= L\left(t\check{\beta}_{\lambda}\left(y,\lambda\right)|ty,\lambda\right) = t^{2}L\left(\check{\beta}_{\lambda}\left(y,\lambda\right)|y,\lambda\right) = t^{2}\check{\sigma}_{\lambda}^{2}\left(y,\lambda\right), \end{split}$$

which establishes the theorem.

I Proof of Theorem 13

Proof. We start from the basic inequality

$$\frac{1}{n} \|y - X \check{\beta}_{\lambda}\|_{2}^{2} + 2\lambda \|\check{\beta}_{\lambda}\|_{1}^{2} \leq \frac{1}{n} \|y - X \beta^{*}\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1}^{2},$$

which leads to

$$\frac{1}{n} \|X\check{\beta}_{\lambda} - X\beta^*\|_{2}^{2} \leq 2 \left(\frac{X^{T}\varepsilon}{n}\right)^{T} \left(\check{\beta}_{\lambda} - \beta^*\right) + 2\lambda \left(\|\beta^*\|_{1}^{2} - \|\check{\beta}_{\lambda}\|_{1}^{2}\right)
\leq 2 \left\|\frac{X^{T}\varepsilon}{n}\right\|_{\infty} \|\check{\beta}_{\lambda} - \beta^*\|_{1} + 2\lambda \left(\|\beta^*\|_{1}^{2} - \|\check{\beta}_{\lambda}\|_{1}^{2}\right).$$

If

$$\left\| \frac{X^T \varepsilon}{n} \right\|_{\infty} \le \sigma \lambda,$$

then

$$\frac{1}{n} \| X \check{\beta}_{\lambda} - X \beta^* \|_{2}^{2} \leq 2\sigma\lambda \| \check{\beta}_{\lambda} - \beta^* \|_{1} + 2\lambda \left(\| \beta^* \|_{1}^{2} - \| \check{\beta}_{\lambda} \|_{1}^{2} \right)
\leq \sigma^{2}\lambda + \lambda \| \check{\beta}_{\lambda} - \beta^* \|_{1}^{2} + 2\lambda \left(\| \beta^* \|_{1}^{2} - \| \check{\beta}_{\lambda} \|_{1}^{2} \right)
\leq \sigma^{2}\lambda + \lambda \left(\| \check{\beta}_{\lambda} \|_{1} + \| \beta^* \|_{1} \right)^{2} + 2\lambda \left(\| \beta^* \|_{1}^{2} - \| \check{\beta}_{\lambda} \|_{1}^{2} \right)
\leq \sigma^{2}\lambda + 2\lambda \left(\| \check{\beta}_{\lambda} \|_{1}^{2} + \| \beta^* \|_{1}^{2} \right) + 2\lambda \left(\| \beta^* \|_{1}^{2} - \| \check{\beta}_{\lambda} \|_{1}^{2} \right)
= \sigma^{2}\lambda + 4\lambda \| \beta^* \|_{1}^{2}.$$

The result then holds from Lemma 16.

J Mapping between the paths of the natural and organic lasso

In this section, we draw a connection between the natural lasso and the organic lasso estimates

Theorem 20. Letting $\hat{\beta}_s$ and $\check{\beta}_t$ denote the lasso and organic lasso estimates of β^* with tuning parameters s and t,

$$\hat{\beta}_{\lambda} = \check{\beta}_{(2\|\hat{\beta}_{\lambda}\|_{1})^{-1}\lambda}, \qquad \check{\beta}_{\nu} = \hat{\beta}_{2\nu\|\check{\beta}_{\nu}\|_{1}}. \tag{19}$$

This result implies that one can start with a lasso solution $\hat{\beta}_{\lambda}$ with tuning parameter λ , and then report a solution to the organic lasso with tuning parameter $(2\|\hat{\beta}_{\lambda}\|_1)^{-1}\lambda$. Likewise, an organic lasso solution $\dot{\beta}_{\nu}$ is equivalent to a standard lasso solution with tuning parameter $2\nu \|\dot{\beta}_{\nu}\|_{1}$. This equivalence is also observed in Lorbert et al. (2010) that considers a more general penalty.

Although the methods' paths are the same, this does not imply that the cross-validated methods will be the same. In K-fold cross-validation, the natural lasso estimator is evaluated on K differing datasets for a fixed value of λ . A fixed tuning parameter λ for the natural lasso over multiple datasets corresponds to running the organic lasso with a different λ on each fold. Thus, the two methods in fact have different cross-validation performance.

Proof. Let $\hat{\beta}_{\lambda}$ be a solution to (2) with tuning parameter λ , and $\tilde{\beta}_{\nu}$ be a solution to (16) with tuning parameter ν , then they satisfy optimality conditions

$$-\frac{1}{n}X^{T}\left(y - X\hat{\beta}_{\lambda}\right) + \lambda\hat{g} = 0 \quad \text{where} \quad \hat{g} \in \partial\left(\left\|\hat{\beta}_{\lambda}\right\|_{1}\right), \tag{20}$$

$$-\frac{1}{n}X^{T}\left(y - X\hat{\beta}_{\lambda}\right) + \lambda\hat{g} = 0 \quad \text{where} \quad \hat{g} \in \partial\left(\left\|\hat{\beta}_{\lambda}\right\|_{1}\right), \tag{20}$$
$$-\frac{1}{n}X^{T}\left(y - X\tilde{\beta}_{\nu}\right) + 2\nu\left\|\tilde{\beta}_{\nu}\right\|_{1}\tilde{g} = 0 \quad \text{where} \quad \tilde{g} \in \partial\left(\left\|\tilde{\beta}_{\nu}\right\|_{1}\right). \tag{21}$$

If $\hat{\beta}_{\lambda} = \tilde{\beta}_{\nu}$, then simply comparing (20) and (21) we have that $\lambda = 2\nu \|\tilde{\beta}_{\nu}\|_{1}$, and $\nu = 2\nu \|\tilde{\beta}_{\nu}\|_{1}$

Now for $\hat{\beta}_{\lambda}$ that satisfies (20), by plugging $\lambda = 2\nu \|\hat{\beta}_{\lambda}\|_1$, we have that $\hat{\beta}_{\lambda}$ satisfies (21), i.e., $\tilde{\beta}_{\nu} = \hat{\beta}_{\lambda}$ where $\lambda = 2\nu \|\hat{\beta}_{\lambda}\|_1$. Following the same argument, for $\tilde{\beta}_{\nu}$ that satisfies (21), we take $\nu = (2\|\tilde{\beta}_{\nu}\|_1)^{-1}\lambda$, and find that $\tilde{\beta}_{\nu}$ satisfies (20). This implies that $\hat{\beta}_{\lambda} = \tilde{\beta}_{\nu}$, where $\nu = (2\|\tilde{\beta}_{\nu}\|_1)^{-1}\lambda.$

K Fast rate in prediction error of the squared lasso

Recall the squared lasso estimate of β^* :

$$\check{\beta} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^p} \frac{1}{n} \|y - X\beta\|_2^2 + 2\lambda \|\beta\|_1^2. \tag{22}$$

It is well known that the fast rate is built on the compatibility condition of the lasso problem. Let $S = \text{supp}(\beta^*)$, i.e., the support of the true regression coefficient β^* , the compatibility condition of the squared lasso problem requires that for all $\mu \in \mathbb{R}^p$ such that $\|\mu_{S^c}\|_1 - \sigma \leq 3\|\mu_S\|_1$,

$$\|\mu_{\mathcal{S}}\|_{1} + \frac{1}{4}\sigma \le |\mathcal{S}|^{1/2} \frac{\|X\mu\|_{2}}{n^{1/2}\phi_{0}}.$$
 (23)

The following theorem establishes that the fast rate prediction error and an estimation error rate of $\check{\beta}$ in (22) can be attained with a value of λ that does not depend on any unknown parameters.

Theorem 21. Suppose that each column X_j of the matrix $X \in \mathbb{R}^{n \times p}$ has been scaled so that $||X_j||_2^2 = n$ for all $j = 1, \ldots, p$, and $\varepsilon \sim N\left(0, \sigma^2 I_n\right)$. If compatibility condition (23) holds, then for any L > 0, the solution $\check{\beta}$ in (22) with tuning parameter

$$\lambda = \left(\frac{2\log p + 2L}{n}\right)^{1/2} \tag{24}$$

attains the following estimation error rate and fast rate bound in prediction with probability greater than $1 - e^{-L}$:

$$\frac{1}{2n} \| X \check{\beta} - X \beta^* \|_2^2 \le \frac{64 \max(\|\beta^*\|_1, \sigma)^2 |\mathcal{S}| (\log p + L)}{\phi_0^2 n};$$
$$\| \beta^* - \check{\beta} \|_1 \le \frac{16 \max(\|\beta^*\|_1, \sigma) |\mathcal{S}|}{\phi_0^2} \left(\frac{2 \log p + 2L}{n} \right)^{1/2}.$$

Proof. First by the optimality of $\check{\beta}$, we have

$$\frac{1}{n} \|y - X\check{\beta}\|_{2}^{2} + 2\lambda \|\check{\beta}\|_{1}^{2} \le \frac{1}{n} \|y - X\beta^{*}\|_{2}^{2} + 2\lambda \|\beta^{*}\|_{1}^{2},$$

which implies that

$$\frac{1}{n} \left\| X \check{\beta} - X \beta^* \right\|_2^2 \le \frac{2}{n} \left(\check{\beta} - \beta^* \right)^T X^T \varepsilon + 2\lambda \left\| \beta^* \right\|_1^2 - 2\lambda \left\| \check{\beta} \right\|_1^2. \tag{25}$$

The following proof is considered in two cases:

(1). When $\|\beta^*\|_1 \geq \sigma$: Note that $\|\cdot\|_1^2$ is convex and by chain rule, for any $g \in \partial(\|\beta^*\|_1)$,

$$\|\check{\beta}\|_{1}^{2} - \|\beta^{*}\|_{1}^{2} \ge 2 \|\beta^{*}\|_{1} g^{T} (\check{\beta} - \beta^{*}).$$

For $j \in \mathcal{S}$, we have that $g_j = \text{sign}(\beta_j^*)$. For any $j \in \mathcal{S}^C$, we let

$$g_j = \operatorname{sign} \left(\check{\beta}_j - \beta_j^* \right) = \operatorname{sign} \left(\check{\beta}_j \right).$$

Then g is still a valid sub-differential of $\|\beta^*\|_1$. Moreover, conditional on the event

$$\mathcal{T} = \left\{ \frac{1}{n} \| X^T \varepsilon \|_{\infty} \le \lambda \sigma \right\},\,$$

from (25) we have

$$\frac{1}{n} \| X \check{\beta} - X \beta^* \|_{2}^{2} \leq \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 4\lambda \| \beta^* \|_{1} g^{T} \left(\beta^* - \check{\beta} \right)
= \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 4\lambda \| \beta^* \|_{1} g^{T}_{\mathcal{S}} \left(\beta^*_{\mathcal{S}} - \check{\beta}_{\mathcal{S}} \right) + 4\lambda \| \beta^* \|_{1} g^{T}_{\mathcal{S}^{C}} \left(\beta^*_{\mathcal{S}^{C}} - \check{\beta}_{\mathcal{S}^{C}} \right)
= \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 4\lambda \| \beta^* \|_{1} g^{T}_{\mathcal{S}} \left(\beta^*_{\mathcal{S}} - \check{\beta}_{\mathcal{S}} \right) - 4\lambda \| \beta^* \|_{1} \| \beta^*_{\mathcal{S}^{C}} - \check{\beta}_{\mathcal{S}^{C}} \|_{1}
\leq \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 4\lambda \| \beta^* \|_{1} \| \beta^*_{\mathcal{S}} - \check{\beta}_{\mathcal{S}} \|_{1} - 4\lambda \| \beta^* \|_{1} \| \beta^*_{\mathcal{S}^{C}} - \check{\beta}_{\mathcal{S}^{C}} \|_{1}.$$

Since $\sigma \leq \|\beta^*\|_1$ and \mathcal{T} holds, we have that $n^{-1} \|X^T \varepsilon\|_{\infty} \leq \lambda \sigma \leq \lambda \|\beta^*\|_1$, and thus

$$\frac{1}{n} \| X \check{\beta} - X \beta^* \|_{2}^{2} \leq 2\lambda \| \beta^* \|_{1} \| \beta^* - \check{\beta} \|_{1} + 4\lambda \| \beta^* \|_{1} \| \beta^*_{S} - \check{\beta}_{S} \|_{1} - 4\lambda \| \beta^* \|_{1} \| \beta^*_{S^{C}} - \check{\beta}_{S^{C}} \|_{1}
= 2\lambda \| \beta^* \|_{1} \left(3 \| \beta^*_{S} - \check{\beta}_{S} \|_{1} - \| \beta^*_{S^{C}} - \check{\beta}_{S^{C}} \|_{1} \right).$$

This first implies that $3\|\beta_{\mathcal{S}}^* - \check{\beta}_{\mathcal{S}}\|_1 \ge \|\beta_{\mathcal{S}^C}^* - \check{\beta}_{\mathcal{S}^C}\|_1$, and that

$$\frac{1}{n} \| X \check{\beta} - X \beta^* \|_2^2 + 2\lambda \| \beta^* \|_1 \| \beta_{\mathcal{S}^C}^* - \check{\beta}_{\mathcal{S}^C} \|_1 \le 6\lambda \| \beta^* \|_1 \| \beta_{\mathcal{S}}^* - \check{\beta}_{\mathcal{S}} \|_1.$$

Then by compatibility condition,

$$\frac{1}{n} \| X \check{\beta} - X \beta^* \|_2^2 + 2\lambda \| \beta^* \|_1 \| \beta^* - \check{\beta} \|_1
= \frac{1}{n} \| X \check{\beta} - X \beta^* \|_2^2 + 2\lambda \| \beta^* \|_1 \| \beta_{\mathcal{S}}^* - \check{\beta}_{\mathcal{S}} \|_1 + 2\lambda \| \beta^* \|_1 \| \beta_{\mathcal{S}^C}^* - \check{\beta}_{\mathcal{S}^C} \|_1
\leq 8\lambda \| \beta^* \|_1 \| \beta_{\mathcal{S}}^* - \check{\beta}_{\mathcal{S}} \|_1 \leq \frac{8\lambda \| \beta^* \|_1 |\mathcal{S}|^{1/2} \| X \beta^* - X \check{\beta} \|_2}{n^{1/2} \phi_0}
\leq \frac{1}{2n} \| X \check{\beta} - X \beta^* \|_2^2 + \frac{32 \| \beta^* \|_1^2 \lambda^2 |\mathcal{S}|}{\phi_0^2}.$$
(26)

(2). When $\|\beta^*\|_1 < \sigma$: We define $\gamma^* \in \mathbb{R}^p$ as

$$\gamma_{j}^{*} = \begin{cases} \beta_{j}^{*} + \frac{\sigma - \|\beta^{*}\|_{1}}{|\mathcal{S}|} & \text{if } \beta_{j}^{*} > 0\\ \beta_{j}^{*} - \frac{\sigma - \|\beta^{*}\|_{1}}{|\mathcal{S}|} & \text{if } \beta_{j}^{*} < 0\\ 0 & \text{if } \beta_{j}^{*} = 0. \end{cases}$$

It is easy to check that $\|\gamma^*\|_1 = \sigma$. Also (25) implies that

$$\frac{1}{n} \| X \check{\beta} - X \beta^* \|_2^2 \le \frac{2}{n} \left(\check{\beta} - \beta^* \right)^T X^T \varepsilon + 2\lambda \left(\| \beta^* \|_1^2 - \| \gamma^* \|_1^2 + \| \gamma^* \|_1^2 - \| \check{\beta} \|_1^2 \right).$$

Then we have that

$$\|\check{\beta}\|_{1}^{2} - \|\gamma^{*}\|_{1}^{2} \ge 2 \|\gamma^{*}\|_{1} g^{T} (\check{\beta} - \gamma^{*})$$

holds for all $g \in \partial(\|\gamma^*\|_1)$, and it further implies that

$$\left\|\gamma^*\right\|_1^2 - \left\|\check{\beta}\right\|_1^2 \leq 2\left\|\gamma^*\right\|_1 g^T\left(\gamma^* - \check{\beta}\right) = 2\sigma g^T\left(\beta^* - \check{\beta}\right) + 2\sigma g^T\left(\gamma^* - \beta^*\right).$$

Note that any $g \in \partial(\|\gamma^*\|_1)$ is also a valid sub-differential of $\|\beta^*\|_1$, and

$$g^{T} \left(\gamma^* - \beta^* \right) = \sigma - \|\beta^*\|_1.$$

Thus we have

$$\frac{1}{n} \left\| X \check{\beta} - X \beta^* \right\|_{2}^{2} \leq \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 2\lambda \left(\left\| \beta^* \right\|_{1}^{2} - \sigma^{2} + 2\sigma g^{T} \left(\beta^* - \check{\beta} \right) + 2\sigma^{2} - 2\sigma \left\| \beta^* \right\|_{1} \right) \\
= \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 4\lambda \sigma g^{T} \left(\beta^* - \check{\beta} \right) + 2\lambda \left(\sigma - \left\| \beta^* \right\|_{1} \right)^{2} \\
\leq \frac{2}{n} \left(\check{\beta} - \beta^* \right)^{T} X^{T} \varepsilon + 4\lambda \sigma g^{T} \left(\beta^* - \check{\beta} \right) + 2\lambda \sigma^{2}$$

Since γ^* and β^* have the same support, we can again choose $g_j = \text{sign}(\check{\beta}_j)$ for $j \in S^c$. Conditional on the event \mathcal{T} , it follows that

$$\frac{1}{n} \left\| X \check{\beta} - X \beta^* \right\|_{2}^{2} \leq 2\lambda \sigma \left\| \check{\beta} - \beta^* \right\|_{1} + 4\lambda \sigma \left\| \check{\beta}_{\mathcal{S}} - \beta^*_{\mathcal{S}} \right\|_{1} - 4\lambda \sigma \left\| \check{\beta}_{\mathcal{S}^{c}} - \beta^*_{\mathcal{S}^{c}} \right\|_{1} + 2\lambda \sigma^{2}$$

$$= 6\lambda \sigma \left\| \check{\beta}_{\mathcal{S}} - \beta^*_{\mathcal{S}} \right\|_{1} - 2\lambda \sigma \left\| \check{\beta}_{\mathcal{S}^{c}} - \beta^*_{\mathcal{S}^{c}} \right\|_{1} + 2\lambda \sigma^{2}.$$

This implies that $3\|\check{\beta}_{\mathcal{S}} - \beta_{\mathcal{S}}^*\| + \sigma \ge \|\check{\beta}_{\mathcal{S}^c} - \beta_{\mathcal{S}^c}^*\|$. And then by the compatibility condition (23),

$$\frac{1}{n} \| X \check{\beta} - X \beta^* \|_{2}^{2} + 2\lambda \sigma \| \check{\beta} - \beta^* \|_{1} = \frac{1}{n} \| X \check{\beta} - X \beta^* \|_{2}^{2} + 2\lambda \sigma \| \check{\beta}_{\mathcal{S}} - \beta^*_{\mathcal{S}} \|_{1} + 2\lambda \sigma \| \check{\beta}_{\mathcal{S}^{c}} - \beta^*_{\mathcal{S}^{c}} \|_{1} \\
\leq 8\lambda \sigma \| \check{\beta}_{\mathcal{S}} - \beta^*_{\mathcal{S}} \|_{1} + 2\lambda \sigma^{2} \\
\leq 8\lambda \sigma |\mathcal{S}|^{1/2} \frac{\| X \left(\check{\beta} - \beta^* \right) \|_{2}}{n^{1/2} \phi_{0}} \\
\leq \frac{1}{2n} \| X \check{\beta} - X \beta^* \|_{2}^{2} + \frac{32\lambda^{2} \sigma^{2} |\mathcal{S}|}{\phi_{0}^{2}}. \tag{27}$$

By the proof of Corollary 4.3 in Giraud (2014), we have

$$\operatorname{pr}\left\{\frac{1}{n}\left\|X^T\varepsilon\right\|_{\infty}>\sigma\left(\frac{2\log p+2L}{n}\right)^{1/2}\right\}\leq e^{-L}.$$

Thus taking λ in (24), we have that

$$\operatorname{pr}(\mathcal{T}^c) = \operatorname{pr}\left(\frac{1}{n} \|X^T \varepsilon\|_{\infty} > \lambda \sigma\right) \le e^{-L}.$$

And the results follow from (26) and (27).

L Additional results in numerical studies

We include in this section some additional results in the numerical studies in Section 5 and Section 6. In particular, Fig 3 and Fig 4 present the complementary results (in different simulation regimes) to Fig 1 and Fig 2 in the paper respectively, and Table 2 shows the p-values of the paired t-tests and the Wilcoxon signed-rank tests of the difference of various methods outputs in Fig 1 in the paper. Finally, Table 3 presents the mean and standard errors of $E(\hat{\sigma}/\sigma)$ of various estimators in the real data example.

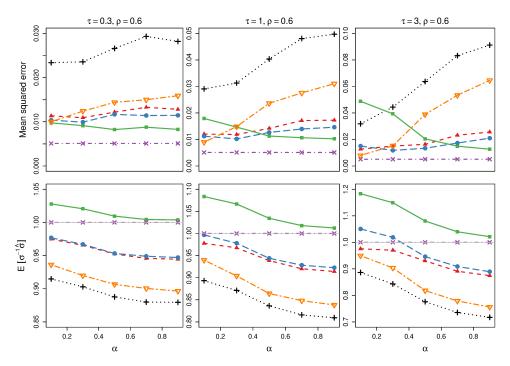


Figure 3: Simulation results of various methods with regularization parameter selected using cross-validation. From left to right, column show the average (over 1000 repetitions) of the mean squared error (top panel) and $E(\sigma^{-1}\hat{\sigma})$ (bottom panel) of various methods in three simulation settings. In each setting, we fix model sparsity (α) and correlations among features (ρ) , and let signal-to-noise ratio(as expressed in τ) change. Line styles and their corresponding methods: + for naive, - for $\hat{\sigma}_R^2$, - for the square-root/scaled lasso, - for the natural lasso, - for the oracle.

Table 2: p-values for testing the difference of various methods outputs

	natural vs. organic	$\hat{\sigma}_R^2$ vs. organic	$\hat{\sigma}_R^2$ vs. natural
$\alpha = 0.1, \rho = 0.3, \tau = 1$	0.00(0.00)	0.07 (0.00)	0.00 (0.00)
$\alpha=0.3, \rho=0.3, \tau=1$	0.00(0.00)	0.19 (0.25)	0.00(0.00)
$\alpha = 0.5, \rho = 0.3, \tau = 1$	0.00(0.00)	0.00(0.00)	0.00(0.00)
$\alpha = 0.7, \rho = 0.3, \tau = 1$	0.00(0.00)	0.00(0.00)	0.00(0.00)
$\alpha = 0.9, \rho = 0.3, \tau = 1$	0.00(0.00)	0.00(0.00)	0.00(0.00)
$\alpha = 0.1, \rho = 0.6, \tau = 1$	0.00(0.00)	0.08(0.01)	0.00(0.00)
$\alpha = 0.3, \rho = 0.6, \tau = 1$	0.00(0.00)	0.00(0.14)	0.00(0.00)
$\alpha = 0.5, \rho = 0.6, \tau = 1$	0.05 (0.10)	0.01 (0.00)	0.00(0.00)
$\alpha = 0.7, \rho = 0.6, \tau = 1$	0.00(0.00)	0.00(0.00)	0.00(0.00)
$\alpha = 0.9, \rho = 0.6, \tau = 1$	0.00(0.00)	0.00(0.00)	0.00(0.00)
$\alpha = 0.1, \rho = 0.9, \tau = 1$	0.06 (0.32)	0.00(0.03)	0.00(0.12)
$\alpha = 0.3, \rho = 0.9, \tau = 1$	0.96 (0.02)	0.00(0.07)	0.00(0.00)
$\alpha = 0.5, \rho = 0.9, \tau = 1$	0.03(0.00)	0.00(0.00)	0.00(0.00)
$\alpha = 0.7, \rho = 0.9, \tau = 1$	0.44(0.00)	0.00(0.00)	0.00(0.00)
$\alpha=0.9, \rho=0.9, \tau=1$	$0.20 \ (0.00)$	0.00(0.01)	0.00(0.00)

In each simulation setting, as characterized by a (α, ρ, τ) triplet, we report p-values of the (two-sided) paired t-tests and the Wilcoxon signed-rank tests (shown in parentheses) for testing the null hypothesis that the output of each pair of methods are the same.

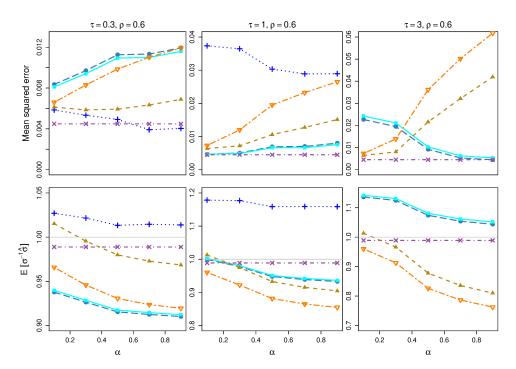


Figure 4: Simulation results of various methods with pre-specified regularization parameter values. From left to right, column show the average (over 1000 repetitions) of the mean squared error (top panel) and $E(\sigma^{-1}\hat{\sigma})$ (bottom panel) of various methods in three simulation settings. In each setting, we fix model sparsity (α) and correlations among features (ρ) , and let signal-to-noise ratio(as expressed in τ) change. Line styles and their corresponding methods: + for organic (λ_0) , - for organic (λ_2) , - for organic (λ_3) , - for scaled (1), - for scaled (2), - for the oracle.

Table 3: $E(\sigma^{-1}\hat{\sigma})$ in MSD dataset

n	20	40	60	80	100	120
naive	80.1 (1.1)	94.2(0.9)	95.8(0.7)	96.4 (0.6)	97.9(0.5)	96.7 (0.5)
$\hat{\sigma}_R^2$	90.0(1.0)	$100.4\ (0.8)$	101.7 (0.6)	$102.3 \ (0.5)$	$103.3 \ (0.5)$	102.4 (0.4)
natural	94.0 (0.9)	$103.3 \ (0.7)$	105.5 (0.6)	$106.0\ (0.5)$	107.0 (0.4)	106.6 (0.4)
organic	86.8 (0.8)	97.6 (0.6)	99.9(0.5)	100.9(0.4)	101.7 (0.4)	101.8 (0.4)
scaled(1)	$106.1\ (0.8)$	109.3 (0.6)	111.2 (0.5)	111.2 (0.4)	111.7(0.4)	111.8 (0.4)
scaled(2)	88.5 (0.8)	$99.0\ (0.6)$	102.9 (0.5)	104.4 (0.5)	105.1 (0.4)	105.5 (0.3)
$\operatorname{organic}(\lambda_2)$	89.7(0.7)	94.7 (0.5)	97.6 (0.4)	98.3 (0.4)	99.2 (0.3)	99.7(0.3)
organic(λ_3)	92.0(0.7)	97.3 (0.6)	$100.1\ (0.4)$	100.7 (0.4)	101.6 (0.4)	102.0 (0.3)

Mean and standard errors (over 1000 replications) of $E(\sigma^{-1}\hat{\sigma})$ of various methods we considered in Section 5. Each entry of the method output is multiplied by 100 to convey information more compactly.

References

- Antoniadis, A. (2010), 'Comments on: ℓ_1 -penalization for mixture regression models', *Test* **19**(2), 257–258.
- Banerjee, O., El Ghaoui, L. & d'Aspremont, A. (2008), 'Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data', *The Journal of Machine Learning Research* 9, 485–516.
- Bayati, M., Erdogdu, M. A. & Montanari, A. (2013), Estimating lasso risk and noise level, *in* 'Advances in Neural Information Processing Systems', pp. 944–952.
- Belloni, A., Chernozhukov, V. & Wang, L. (2011), 'Square-root lasso: pivotal recovery of sparse signals via conic programming', *Biometrika* **98**(4), 791–806.
- Bien, J. (2016), 'The Simulator: An Engine to Streamline Simulations', ArXiv e-prints.
- Bien, J. & Tibshirani, R. J. (2011), 'Sparse estimation of a covariance matrix', *Biometrika* **98**(4), 807–820.
- Boyd, S. & Vandenberghe, L. (2004), Convex optimization, Cambridge university press.
- Bühlmann, P. (2013), 'Statistical significance in high-dimensional linear models', *Bernoulli* 19(4), 1212–1242.
- Bühlmann, P. & Van De Geer, S. (2011), Statistics for high-dimensional data: methods, theory and applications, Springer Science & Business Media.
- Campbell, F., Allen, G. I. et al. (2017), 'Within group variable selection through the exclusive lasso', *Electronic Journal of Statistics* **11**(2), 4220–4257.
- Chatterjee, S. & Jafarov, J. (2015), 'Prediction error of cross-validated lasso', arXiv preprint arXiv:1502.06291.
- Dalalyan, A. & Chen, Y. (2012), Fused sparsity and robust estimation for linear models with unknown variance, in 'Advances in Neural Information Processing Systems', pp. 1259–1267.
- Dalalyan, A. S., Hebiri, M. & Lederer, J. (2017), 'On the prediction performance of the lasso', Bernoulli 23(1), 552–581.
- Dicker, L. H. (2014), 'Variance estimation in high-dimensional linear models', *Biometrika* **101**(2), 269.
- Fan, J., Guo, S. & Hao, N. (2012), 'Variance estimation using refitted cross-validation in ultrahigh dimensional regression', *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **74**(1), 37–65.
- Friedman, J., Hastie, T., Höfling, H., Tibshirani, R. et al. (2007), 'Pathwise coordinate optimization', *The Annals of Applied Statistics* 1(2), 302–332.
- Friedman, J., Hastie, T. & Tibshirani, R. (2008), 'Sparse inverse covariance estimation with the graphical lasso', *Biostatistics* **9**(3), 432–441.
- Gaynanova, I. (2018), 'Prediction and estimation consistency of sparse multi-class penalized optimal scoring', arXiv preprint arXiv:1809.04669.
- Giraud, C. (2014), Introduction to high-dimensional statistics, Vol. 138, CRC Press.

- Hastie, T., Tibshirani, R. & Wainwright, M. (2015), Statistical learning with sparsity: the lasso and generalizations, CRC press.
- Hebiri, M. & Lederer, J. (2013), 'How correlations influence lasso prediction', *IEEE Transactions on Information Theory* **59**(3), 1846–1854.
- Javanmard, A. & Montanari, A. (2014), 'Confidence intervals and hypothesis testing for high-dimensional regression.', *Journal of Machine Learning Research* **15**(1), 2869–2909.
- Lederer, J., Yu, L. & Gaynanova, I. (2016), 'Oracle inequalities for high-dimensional prediction', arXiv preprint arXiv:1608.00624.
- Lee, J. D., Sun, D. L., Sun, Y., Taylor, J. E. et al. (2016), 'Exact post-selection inference, with application to the lasso', *The Annals of Statistics* 44(3), 907–927.
- Lockhart, R., Taylor, J., Tibshirani, R. J. & Tibshirani, R. (2014), 'A significance test for the lasso', *Annals of statistics* **42**(2), 413.
- Lorbert, A. (2012), Alignment and supervised learning with functional neuroimaging data, PhD thesis.
 - **URL:** http://arks.princeton.edu/ark:/88435/dsp01707957683
- Lorbert, A., Eis, D. J., Kostina, V., Blei, D. M. & Ramadge, P. J. (2010), Exploiting covariate similarity in sparse regression via the pairwise elastic net., in 'AISTATS', Vol. 9, pp. 477–484.
- Ning, Y. & Liu, H. (2017), 'A general theory of hypothesis tests and confidence regions for sparse high dimensional models', *Ann. Statist.* **45**(1), 158–195. URL: https://doi.org/10.1214/16-AOS1448
- R Core Team (2017), R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria.

 URL: https://www.R-project.org/
- Reid, S., Tibshirani, R. & Friedman, J. (2016), 'A study of error variance estimation in lasso regression', *Statistica Sinica* pp. 35–67.
- Rigollet, P. & Tsybakov, A. (2011), 'Exponential screening and optimal rates of sparse estimation', *The Annals of Statistics* pp. 731–771.
- Rockafellar, R. T. (2015), Convex analysis, Princeton university press.
- Städler, N., Bühlmann, P. & van de Geer, S. (2010), ' ℓ_1 -penalization for mixture regression models (with discussion)', Test 19, 209–285.
- Sun, T. & Zhang, C.-H. (2010), 'Comments on: ℓ_1 -penalization for mixture regression models', Test 19(2), 270–275.
- Sun, T. & Zhang, C.-H. (2012), 'Scaled sparse linear regression', Biometrika 99(4), 879–898.
- Sun, T. & Zhang, C.-H. (2013), 'Sparse matrix inversion with scaled lasso', *The Journal of Machine Learning Research* **14**(1), 3385–3418.
- Taylor, J. & Tibshirani, R. (2017), 'Post-selection inference for ℓ_1 -penalized likelihood models', Canadian Journal of Statistics .
- Tibshirani, R. (1996), 'Regression shrinkage and selection via the lasso', Journal of the Royal Statistical Society. Series B (Methodological) 58(1), 267–288.

- Tibshirani, R. J., Rinaldo, A., Tibshirani, R., Wasserman, L. et al. (2018), 'Uniform asymptotic inference and the bootstrap after model selection', *The Annals of Statistics* **46**(3), 1255–1287.
- Tibshirani, R. J., Taylor, J., Lockhart, R. & Tibshirani, R. (2016), 'Exact post-selection inference for sequential regression procedures', *Journal of the American Statistical Association* **111**(514), 600–620.
- Tseng, P. (2001), 'Convergence of a block coordinate descent method for nondifferentiable minimization', *Journal of Optimization Theory and Applications* **109**(3), 475–494.
- Van de Geer, S. A. & Bühlmann, P. (2009), 'On the conditions used to prove oracle results for the lasso', *Electronic Journal of Statistics* 3, 1360–1392.
- Van de Geer, S., Bühlmann, P., Ritov, Y., Dezeure, R. et al. (2014), 'On asymptotically optimal confidence regions and tests for high-dimensional models', *The Annals of Statistics* **42**(3), 1166–1202.
- Yu, G. (2017), natural: Estimating the Error Variance in a High-Dimensional Linear Model. R package version 0.9.0.
 - **URL:** https://CRAN.R-project.org/package=natural
- Yuan, M. & Lin, Y. (2007), 'Model selection and estimation in the gaussian graphical model', *Biometrika* **94**(1), 19–35.
- Zhang, C.-H. & Zhang, S. S. (2014), 'Confidence intervals for low dimensional parameters in high dimensional linear models', *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76**(1), 217–242.
- Zhou, Y., Jin, R. & Hoi, S. (2010), Exclusive lasso for multi-task feature selection, in 'International conference on artificial intelligence and statistics', pp. 988–995.