

# A Comparative Framework for Preconditioned Lasso Algorithms

**Authored by:**

Nebojsa Jojic  
 Fabian L. Wauthier  
 Michael I. Jordan

## Abstract

The Lasso is a cornerstone of modern multivariate data analysis, yet its performance suffers in the common situation in which covariates are correlated. This limitation has led to a growing number of *Preconditioned Lasso* algorithms that pre-multiply  $X$  and  $y$  by matrices  $P_X$ ,  $P_y$  prior to running the standard Lasso. A direct comparison of these and similar Lasso-style algorithms to the original Lasso is difficult because the performance of all of these methods depends critically on an auxiliary penalty parameter  $\lambda$ . In this paper we propose an agnostic, theoretical framework for comparing Preconditioned Lasso algorithms to the Lasso without having to choose  $\lambda$ . We apply our framework to three Preconditioned Lasso instances and highlight when they will outperform the Lasso. Additionally, our theory offers insights into the fragilities of these algorithms to which we provide partial solutions.

## 1 Paper Body

Variable selection is a core inferential problem in a multitude of statistical analyses. Confronted with a large number of (potentially) predictive variables, the goal is to select a small subset of variables that can be used to construct a parsimonious model. Variable selection is especially relevant in linear observation models of the form  $y = X\beta + w$

with  $w \sim N(0, \sigma^2 I_n)$ ,

(1)  
 $\beta$

where  $X$  is an  $n \times p$  matrix of features or predictors,  $\beta$  is an unknown  $p$ -dimensional regression parameter, and  $w$  is a noise vector. In high-dimensional settings where  $n \approx p$ , ordinary least squares is generally inappropriate. Assuming that  $\beta$  is sparse (i.e., the support set  $S(\beta) = \{i : \beta_i \neq 0\}$  has cardinality  $k \ll p$ ), a mainstay algorithm for such settings is the Lasso [10]:  $\hat{\beta} = \arg\min_{\beta} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$ .

(2)

For a particular choice of  $\lambda$ , the variable selection properties of the Lasso can be analyzed by quantifying how well the estimated support  $S(\lambda)$  approximates the true support  $S$ . More carefully, how well the estimated support  $S(\lambda)$  approximates the signed support  $S^s(\lambda)$ ,  $(S^s(\lambda))_i = +1$  if  $x_i \geq 0$  and  $-1$  if  $x_i < 0$ .  $S^s(\lambda)$  is the signed support of  $x$ . Theoretical developments during the last decade have shed light onto the support recovery properties of the Lasso and highlighted practical difficulties when the columns of  $X$  are correlated. These developments have led to various conditions on  $X$  for support recovery, such as the mutual incoherence or the irrepresentable condition [1, 3, 8, 12, 13].

In recent years, several modifications of the standard Lasso have been proposed to improve its support recovery properties [2, 7, 14, 15]. In this paper we focus on a class of Preconditioned Lasso algorithms [5, 6, 9] that pre-multiply  $X$  and  $y$  by suitable matrices  $PX$  and  $Py$  to yield  $\tilde{X} = PX$ ,  $\tilde{y} = Py$ , prior to running Lasso. Thus, the general strategy of these methods is

$\hat{\beta} = \arg\min_{\beta} \|\tilde{y} - \tilde{X}\beta\|_2^2 + \lambda \|\beta\|_1$ . Preconditioned Lasso:  $\hat{\beta} = \arg\min_{\beta} \|\tilde{y} - \tilde{X}\beta\|_2^2 + \lambda \|\beta\|_1$ . Although this class of algorithms often compares favorably to the Lasso in practice, our theoretical understanding of them is at present still fairly poor. Huang and Jojic [5], for example, consider only empirical evaluations, while both Jia and Rohe [6] and Paul et al. [9] consider asymptotic consistency under various assumptions. Important and necessary as they are, consistency results do not provide insight into the relative performance of Preconditioned Lasso variants for finite data sets. In this paper we provide a new theoretical basis for making such comparisons. Although the focus of the paper is on problems of the form of Eq. (4), we note that the core ideas can also be applied to algorithms that right-multiply  $X$  and/or  $y$  with some matrices (e.g., [4, 11]). For particular instances of  $X$ ,  $y$ , we want to discover whether a given Preconditioned Lasso algorithm following Eq. (4) improves or degrades signed support recovery relative to the standard Lasso of Eq. (2). A major roadblock to a one-to-one comparison are the auxiliary penalty parameters,  $\lambda$ ,  $\lambda_0$ . A correct choice of penalty parameter is essential for signed support recovery: If it is too small, the algorithm behaves like ordinary least squares; if it is too large, the estimated support may be empty. Unfortunately, in all but the simplest cases, pre-multiplying data  $X$ ,  $y$  by matrices  $PX$ ,  $Py$  changes the relative geometry of the  $\lambda$  penalty contours to the elliptical objective contours in a nontrivial way. Suppose we wanted to compare the Lasso to the Preconditioned Lasso by choosing  $\lambda$  in Eq. (4). For a fair comparison, the resulting mapping for each  $\lambda$  in Eq. (2) a suitable, matching  $\lambda_0$  would have to capture the change of relative geometry induced by preconditioning of  $X$ ,  $y$ ,  $\lambda_0 = f(\lambda, X, y, PX, Py)$ . It seems difficult to theoretically characterize such a mapping. Furthermore, it seems unlikely that a comparative framework could be built by independently choosing  $\lambda$ . Meinshausen and Bühlmann [8], for example, demonstrate that a ‘‘ideal’’ penalty parameters  $\lambda$ ,  $\lambda_0$ : seemingly reasonable oracle estimator of  $\lambda$  will not lead to consistent support recovery in the Lasso. In the Preconditioned Lasso literature this problem is commonly sidestepped

either by resorting to asymptotic comparisons [6, 9], empirically comparing regularization paths [5], or using modelselection techniques which aim to choose reasonably “good” matching penalty parameters [6]. We deem these approaches to be unsatisfactory: asymptotic and empirical analyses provide limited insight, and model selection strategies add a layer of complexity that may lead to unfair comparisons. It is our view that all of these approaches place unnecessary emphasis on particular choices of penalty parameter. In this paper we propose an alternative strategy that instead compares the Lasso to the Preconditioned Lasso by comparing data-dependent upper and lower penalty parameter bounds. Specifically, we give bounds  $(\lambda_u, \lambda_l)$  on  $\lambda$  so that the Lasso in Eq. (2) is guaranteed to recover the signed support iff  $\lambda_l \leq \lambda \leq \lambda_u$ . Consequently, if  $\lambda \notin [\lambda_l, \lambda_u]$  signed support recovery is not possible.  $\tilde{X} = PX^T X$ ,  $\tilde{y} = Py^T y$  and will thus induce new The Preconditioned Lasso in Eq. (4) uses data  $X$  and  $\tilde{X}$  bounds  $(\lambda_u, \lambda_l)$  on  $\lambda$ . The comparison of Lasso and Preconditioned Lasso on an instance  $X, y$  is then proceeds by suitably comparing the bounds on  $\lambda_u$  and  $\lambda_l$ . The advantage of this approach is that the upper and lower bounds are easy to compute, even though a general mapping between specific penalty parameters cannot be readily derived. To demonstrate the effectiveness of our framework, we use it to analyze three Preconditioned Lasso algorithms [5, 6, 9]. Using our framework we make several contributions: (1) We confirm intuitions about advantages and disadvantages of the algorithms proposed in [5, 9]; (2) We show that for an SVD-based construction of  $n \times p$  matrices  $X$ , the algorithm in [6] changes the bounds deterministically; (3) We show that in the context of our framework, this SVD-based construction can be thought of as a limit point of a Gaussian construction. The paper is organized as follows. In Section 2 we will discuss three recent instances of Eq. (4). We outline our comparative framework in Section 3 and highlight some immediate consequences for [5] and [9] on general matrices  $X$  in Section 4. More detailed comparisons can be made by considering a generative model for  $X$ . In Section 5 we introduce such a model based on a block-wise SVD of  $X$  and then analyze [6] for specific instances of this generative model. Finally, we show that in terms of signed support recovery, this generative model can be thought of as a limit point of a Gaussian construction. Section 6 concludes with some final thoughts. The proofs of all lemmas and theorems are in the supplementary material.

## 2

### Preconditioned Lasso Algorithms

Our interest lies in the class of Preconditioned Lasso algorithms that is summarized by Eq. (4). Extensions to related algorithms, such as [4, 11] will follow readily. In this section we focus on three recent Preconditioned Lasso examples and instantiate the matrices  $PX$ ,  $Py$  appropriately. Detailed derivations can be found in the supplementary material. For later reference, we will denote each algorithm by the author initials. Huang and Jovic [5] (HJ). Huang and Jovic proposed Correlation Sifting [5], which, although not presented as a preconditioning algorithm, can be rewritten as one. Let the SVD of  $X$  be  $X = UDV^T$ . Given an algorithm parameter  $q$ , let  $U_A$  be the set of  $q$  smallest left singular vectors of  $X$ . Then HJ amounts to setting  $\tilde{X} = PX = U_A U_A^T X$ .

(5)

Paul et al. [9] (PBHT). An earlier instance of the preconditioning idea was put forward by Paul et al. [9]. For some algorithm parameter  $q$ , let  $A$  be the  $q$  column indices of  $X$  with largest absolute correlation to  $y$ , (i.e., where  $|X_j|_2 / \|y\|_2$  is largest). Define  $U_A$  to be the  $q$  largest left singular vectors of  $X_A$ . With this, PBHT can be expressed as setting  $\hat{y} = U_A U_A^T y$ .

$$P_X = I_n - U_A U_A^T$$

(6)

Jia and Rohe [6] (JR). Jia and Rohe [6] propose a preconditioning method that amounts to whitening the matrix  $X$ . If  $X = U D V^T$  is full rank, then JR defines  $\hat{y} = (X^T X + \lambda I)^{-1} X^T y$ . (7)  $\hat{y} = (X^T X + \lambda I)^{-1} X^T y$  and if  $n \geq p$  then  $X^T X = X P X^T$  where  $P = I_p$ . If  $n < p$  then  $X^T X = X P X^T$  Both HJ and PBHT estimate a basis  $U_A$  for a  $q$ -dimensional subspace onto which they project  $y$  and/or  $X$ . However, since the methods differ substantially in their assumptions, the estimators differ also. Empirical results in [5] and [9] suggest that the respective assumptions are useful in a variety of situations. In contrast, JR reweights the column space directions  $U$  and requires no extra parameter  $q$  to be estimated.

3

### Comparative Framework

In this section we propose a new comparative approach for Preconditioned Lasso algorithms which avoids choosing particular penalty parameters  $\lambda, \gamma$ . We first derive upper and lower bounds for  $\lambda$  and  $\gamma$  respectively so that signed support recovery can be guaranteed iff  $\lambda$  and  $\gamma$  satisfy the bounds. We then compare estimators by comparing the resulting bounds. 3.1

#### Conditions for signed support recovery

Before proceeding, we make some definitions motivated by Wainwright [12]. Suppose that the support set of  $\beta$  is  $S$ ,  $|S| = k$ , with  $|S^c| = n - k$ . To simplify notation, we will assume throughout that  $S = \{1, \dots, k\}$  so that the corresponding off-support set is  $S^c = \{k+1, \dots, n\}$ , with  $|S^c| = n - k$ . Denote by  $X_j$  column  $j$  of  $X$  and by  $X_A$  the submatrix of  $X$  consisting of columns indexed by set  $A$ . Define the following variables: For all  $j \in S^c$  and  $i \in S$ , let  $w_{ij} = X_j^T X_i / (X_j^T X_j)^{1/2} (X_i^T X_i)^{1/2}$ . (8)  $w_{ij} = X_j^T X_i / (X_j^T X_j)^{1/2} (X_i^T X_i)^{1/2}$

1

$$\|w\|_1 = \sum_{i \in S} \sum_{j \in S^c} |w_{ij}| = \sum_{i \in S} \sum_{j \in S^c} |X_j^T X_i| / (X_j^T X_j)^{1/2} (X_i^T X_i)^{1/2}. \quad (9)$$

The choice of smallest singular vectors is considered for matrices  $X$  with sharply decaying spectrum. We note that Jia and Rohe [6] let  $D$  be square, so that it can be directly inverted. If  $X$  is not full rank, the pseudo-inverse of  $D$  can be used. 2

3

$$\|w\|_1 = \sum_{i \in S} \sum_{j \in S^c} |w_{ij}| = \sum_{i \in S} \sum_{j \in S^c} |X_j^T X_i| / (X_j^T X_j)^{1/2} (X_i^T X_i)^{1/2}$$

$$\|w\|_1 = \sum_{i \in S} \sum_{j \in S^c} |w_{ij}| = \sum_{i \in S} \sum_{j \in S^c} |X_j^T X_i| / (X_j^T X_j)^{1/2} (X_i^T X_i)^{1/2}$$

$$1 \ 0.8 \ 0.6 \ 0.4 \ 0.2 \ 0 \ 0.5$$

$$1 \ f$$

$$0.8 \ 0.6 \ 0.4 \ 0.2 \ 0 \ 0.5$$

1.5

(a) Signed support recovery around  $\lambda_l$ .

1 f

1.5

(b) Signed support recovery around  $\lambda_u$ .

Figure 1: Empirical evaluation of the penalty parameter bounds of Lemma

1. For each of 500 synthetic Lasso problems ( $n = 300$ ,  $p = 1000$ ,  $k = 10$ ) we computed  $\lambda_l$ ,  $\lambda_u$  as per Lemma 1. Then we ran Lasso using penalty parameters  $f \lambda_l$  in Figure (a) and  $f \lambda_u$  in Figure (b), where the factor  $f = 0.5, \dots, 1.5$ . The figures show the empirical probability of signed support recovery as a function of the factor  $f$  for both  $\lambda_l$  and  $\lambda_u$ . As expected, the probabilities change sharply at  $f = 1$ . For the traditional Lasso of Eq. (2), results in (for example) Wainwright [12] connect settings of  $\lambda$  with instances of  $X, \beta, w$  to certify whether or not Lasso will recover the signed support. We invert these results and, for particular instances of  $X, \beta, w$ , derive bounds on  $\lambda$  so that signed support recovery is guaranteed if and only if the bounds are satisfied. Specifically, we prove the following Lemma in the supplementary material. Lemma 1. Suppose that  $XS \succeq XS$  is invertible,  $\|x_j\|_1 \leq 1$ ,  $\|x_j\|_2 \leq S$ , and  $\text{sgn}(\beta_i) \beta_i \leq 0$ ,  $\beta_i \in S$ . Then  $\lambda = S^2 (\|x_j\|_1 + \|x_j\|_2)$  if and the Lasso has a unique solution  $\hat{\beta}$  which recovers the signed support (i.e.,  $S^2 (\lambda_l \leq \lambda \leq \lambda_u)$  only if  $\lambda_l \leq \lambda \leq \lambda_u$ , where  $\lambda_l = \|x_j\|_1 + \|x_j\|_2$ ,  $\lambda_u = \max_j (2\|x_j\|_2^2 / \|x_j\|_1) + \|x_j\|_2^2$ ).

$\lambda_l = \max_j (\|x_j\|_1 + \|x_j\|_2)$ ,  $\lambda_u = \max_j (2\|x_j\|_2^2 / \|x_j\|_1) + \|x_j\|_2^2$ .  $J_K$  denotes the indicator function and  $\max(0, \lambda) = \max(0, \lambda)$  denotes the hinge function. On the other hand, if  $XS \succeq XS$  is not invertible, then the signed support cannot in general be recovered. Lemma 1 recapitulates well-worn intuitions about when the Lasso has difficulty recovering the signed support. For instance, assuming that  $w$  has symmetric distribution with mean 0, if  $\|x_j\|_1 - \|x_j\|_2$  is small (i.e., the irrepresentable condition almost fails to hold), then  $\lambda_l$  will tend to be large. In extreme cases we might have  $\lambda_l \geq \lambda_u$  so that signed support recovery is impossible. Figure 1 empirically validates the bounds of Lemma 1 by estimating probabilities of signed support recovery for a range of penalty parameters on synthetic Lasso problems.

### Comparisons

In this paper we propose to compare a preconditioning algorithm to the traditional Lasso by comparing the penalty parameter bounds produced by Lemma 1. As highlighted in Eq. 4, the preconditioned Lasso  $\tilde{X} \tilde{\beta} = y$  can be written as  $\tilde{X} \tilde{\beta} = y$ . For the purpose of applying the framework runs Lasso on modified variables  $\tilde{X}$  Lemma 1, these transformations induce a new noise vector  $\tilde{w} = Py (X\beta + w) - \tilde{P}X X\beta - w = y - \tilde{X}\tilde{\beta}$ . Note that if  $PX = Py$  then  $w = Py w$ . Provided the conditions of Lemma 1 hold for  $X, \beta, w$  we can define updated variables  $\tilde{x}_j, \tilde{\beta}_i, \tilde{w}_i$  from which the bounds  $\tilde{\lambda}_u, \tilde{\lambda}_l$  on the penalty parameter  $\tilde{\lambda}$  be derived. In order for our comparison to be scale-invariant, we will compare algorithms by ratios of resulting penalty parameter bounds. That is, we deem a Preconditioned Lasso algorithm to be  $\tilde{\lambda}_u / \tilde{\lambda}_l \leq \lambda_u / \lambda_l$ . Intuitively, the upper bound  $\tilde{\lambda}_u$  is then more effective than the traditional Lasso if  $\tilde{\lambda}_u$  is disproportionately larger than  $\lambda_l$  relative to  $\lambda_u$  and  $\lambda_l$ , which in principle

allows easier tuning of  $\gamma$ . We will later encounter the special case  $u = 0$ ,  $v = 0$  in which case we define  $\gamma(u, v)$  to be  $\frac{1}{\|u\|_1}$ ; then signed support recovery indicates that the preconditioned problem is very easy. If  $u \neq 0$ ,  $v \neq 0$ , if  $u = v$  or  $u = -v$  it is in general impossible. Finally, to match this intuition, we define  $\gamma(u, v) = \frac{\|u\|_1}{\|u\|_1 + \|v\|_1}$  could also be considered. However, we find the ratio to be a particularly intuitive measure.

44

## General Comparisons

We begin our comparisons with some immediate consequences of Lemma 1 for HJ and PBHT. In order to highlight the utility of the proposed framework, we focus in this section on special cases of  $\mathbf{P}\mathbf{X}$ ,  $\mathbf{P}\mathbf{y}$ . The framework can of course also be applied to general matrices  $\mathbf{P}\mathbf{X}$ ,  $\mathbf{P}\mathbf{y}$ . As we will see, both HJ and PBHT have the potential to improve signed support recovery relative to the traditional Lasso, provided the matrices  $\mathbf{P}\mathbf{X}$ ,  $\mathbf{P}\mathbf{y}$  are suitably estimated. The following notation will be used during our comparisons: We will write  $\mathbf{A} \succcurlyeq \mathbf{A}'$  to indicate that random variable  $\mathbf{A}$  stochastically  $\succcurlyeq$  that is,  $\Pr(\mathbf{A} \leq t) \geq \Pr(\mathbf{A}' \leq t)$ . We also let  $\mathbf{U}\mathbf{S}$  be a minimal basis for the column dominates  $\mathbf{A}$ .

space of the submatrix  $\mathbf{X}\mathbf{S}$ , and define  $\text{span}(\mathbf{US}) = \mathbf{x} \in \mathbb{R}^n$  s.t.  $\mathbf{x} = \mathbf{US}\mathbf{c}$  for  $\mathbf{c} \in \mathbb{R}^n$ . Finally, we let  $\mathbf{US}^\perp$  be a minimal basis for the orthogonal complement of  $\text{span}(\mathbf{US})$ .  $\square$  Consequences for HJ. Recall from Section 2 that HJ uses  $\mathbf{P}\mathbf{X} = \mathbf{P}\mathbf{y} = \mathbf{U}\mathbf{A}\mathbf{U}^\top$ , where  $\mathbf{U}\mathbf{A}$  is a column basis estimated from  $\mathbf{X}$ . We have the following theorem: Theorem 1. Suppose that the conditions of Lemma 1 are met for a fixed instance of  $\mathbf{X}$ ,  $\mathbf{y}$ . If  $\text{span}(\mathbf{US}) \subseteq \text{span}(\mathbf{U}\mathbf{A})$ , then after preconditioning using HJ the conditions continue to hold, and  $\|\mathbf{u}\|_2 \leq \sqrt{\kappa(\mathbf{X})} \|\mathbf{y}\|_2$ , (12)  $\square$

where the stochasticity on both sides is due to independent noise vectors  $w$ . On the other hand, if  $\mathbf{X}^\top \mathbf{X}$  is not invertible, then HJ cannot in general recover the signed support.  $\mathbf{X}^\top \mathbf{X}$  We briefly sketch the proof of Theorem 1. If  $\text{span}(\mathbf{U} \mathbf{S}) \neq \text{span}(\mathbf{U} \mathbf{A})$  then plugging in the definition of  $\mathbf{X}$  into  $\mathbf{y} = \mathbf{X}^\top \mathbf{z}$ ,  $\mathbf{y}_i = \mathbf{X}_i^\top \mathbf{z}$ , one can derive the following  $\mathbf{y} = \mathbf{X}^\top \mathbf{z}$   $\mathbf{y}_i = \mathbf{X}_i^\top \mathbf{z}$  (13)

Let  $\tilde{X}_j = X_j - \frac{1}{n} \sum_{i=1}^n X_i$ . If  $\text{span}(U) = \text{span}(US)$ , then it is easy to see that  $\tilde{X}_j = 0$ . Notice that because  $\tilde{X}_j$  and  $\tilde{X}_i$  are unchanged, if the conditions of Lemma 1 hold for the original Lasso problem (i.e.,  $X^T X$  is invertible,  $\|X_j\|_1 \leq \frac{1}{2} \|X_i\|_1$  and  $\text{sgn}(\tilde{X}_j) = \text{sgn}(\tilde{X}_i)$ ), they will continue to hold for the preconditioned problem. Suppose then that the conditions set forth in Lemma 1 are met. With some additional work one can show that

$$u_i + v_j - l = \max_{u,v} \{u_i + v_j - l\} = u_i + v_j - l. \quad (15)$$

c i?S j?S (2J? ??i ?j ) 0K ? 1) ? ? ?j + ? 1 , ?l are both independent of ? ? u = ?u . Note that if The result then follows by showing that ? ? ? ? span(UA ) = span(US ), then ?l = 0 and so ?u /?l , ?. In the more common case when span(US ) 6? span(UA ) the performance of the Lasso depends on how misaligned UA and US are. In extreme cases, XS XS PX PX XS is singular and so signed support recovery is not in general possible. Consequences for PBHT.

Recall from Section 2 that PBHT uses  $PX = I_n \gamma_n$ ,  $Py = UA UA^T$ , where  $UA$  is a column basis estimated from  $X$ . We have the following theorem. Theorem 2. Suppose that the conditions of Lemma 1 are met for a fixed instance of  $X$ ,  $\gamma$ . If  $\text{span}(US) \subset \text{span}(UA)$ , after preconditioning using PBHT the conditions continue to hold, and  $\gamma_u \gamma_u^T = (16) \gamma$ ,  $\gamma_l \gamma_l^T$

where the stochasticity on both sides is due to independent noise vectors  $w$ . On the other hand, if  $\text{span}(US^c) = \text{span}(UA)$ , then PBHT cannot recover the signed support. As before, we sketch the proof to build some intuition. Because PBHT does not set  $PX = Py$  as HJ  $\gamma$  does, there is no danger of  $XS \gamma$   $PX$   $PX$   $XS$  becoming singular. On the other hand, this complicates  $\gamma$  the form of the induced noise vector  $w$ . Plugging  $PX$  and  $Py$  into Eq. (11), we find  $w = (UA UA^T \gamma \gamma^T I_n)X\gamma + UA UA^T w$ . However, even though the noise has a more complicated form, derivations in the supplementary material show that if  $\text{span}(US) \subset \text{span}(UA)$ , then  $\gamma_j = \gamma_j$   $\gamma_i = \gamma_i$  (17)

$$\gamma \gamma^T w \gamma_j = X_j I_n \gamma \gamma^T US US^T UA UA^T \gamma_i = \gamma_i. \quad (18) \quad n = 5$$

2.5

0.6 0.4 0.2 0

Orthogonal data Gaussian data

2 Lasso 55 35 25 15 10 0

1000

2000

3000

4000

$\gamma_u / \gamma_l \gamma \gamma_u / \gamma_l$

$P(\gamma_u / \gamma_l \gamma_t)$

1 0.8

1.5 1 0.5 0 0.2

5000

t

0.4

0.6

0.8

1

1.2

1.4

f

(a) Empirical validation of Theorems 1 and 2.

(b) Evaluation of JR on Gaussian ensembles.

Figure 2: Experimental evaluations. Figure (a) shows empirical c.d.f.s of penalty parameter bounds ratios estimated from 1000 variable selection problems. Each problem consists of Gaussians  $X$  and  $w$ , and  $\gamma$ , with  $n = 100$ ,  $p = 300$ ,  $k = 5$ . The blue curve shows the c.d.f. for  $\gamma_u / \gamma_l$  estimated on  $\gamma$  the original data (Lasso). Then we projected the data using  $PX = Py = UA UA^T$ , where  $\text{span}(US) \subset \text{span}(UA)$  but  $\dim(UA) = \dim(\text{span}(UA))$  is variable (see legend), and estimated the resulting  $\gamma_u / \gamma_l$ . As predicted by Theorems 1 and 2,  $\gamma_u / \gamma_l \gamma \gamma_u / \gamma_l$ . In c.d.f. for the updated bounds ratio  $\gamma$  2 Figure

(b) the blue curve shows the scale factor  $(p - k)/(n + p - k)$  predicted by Theorem 3 for problems constructed from Eq. (19) for  $f = 1 - (n/p)$ . The red curve plots the corresponding factor estimated from the Gaussian construction in Eq. (25) ( $n = 100$ ,  $m = 2000$ ,  $p = 200$ ,  $k = 5$ ) using the same  $S$ ,  $S_c$  as in Theorem 3, averaged over 50 problem instances and with error bars for one standard deviation. As in Theorem 3, the factor is approximately 1 if  $f = 1$ .

As with HJ, if  $\text{span}(U_A) = \text{span}(U_S)$ , then  $\beta_j = 0$ . Because  $\beta_j$  and  $\beta_i$  are again unchanged, the conditions of Lemma 1 will continue to hold for the preconditioned problem if they hold for the original Lasso problem. With the previous equalities established, the remainder of the proof is identical to that of Theorem 1. The fact that the above  $\beta_j$ ,  $\beta_j$ ,  $\beta_i$ ,  $\beta_i$  are identical to those of HJ depends crucially on the fact that  $\text{span}(U_S) = \text{span}(U_A)$ . In general the values will differ because PBHT sets  $PX = I_n^n$ , but HJ does not. On the other hand, if  $\text{span}(U_S) \neq \text{span}(U_A)$  then the distribution of  $\beta_i$  depends on how misaligned  $U_A$  and  $U_S$  are. In the extreme case when  $\text{span}(U_S) = \text{span}(U_A)$ , one can show that  $\beta_i = \beta_i$ ,  $\beta_u = 0$ ,  $\beta_l = 1$ . Because  $P(\beta_l = 0) = 1$ , signed support recovery is not possible. which results in ? Remarks. Our theoretical analyses show that both HJ and PBHT can indeed lead to improved signed support recovery relative to the Lasso on finite datasets. To underline our findings, we empirically validate Theorems 1 and 2 in Figure 2(a), where we plot estimated c.d.f.s for penalty parameter bounds ratios of Lasso and Preconditioned Lasso for various subspaces  $U_A$ . Our theorems focussed on specific settings of  $PX$ ,  $P_y$  and ignored others. In general, the gains of HJ and PBHT over Lasso depend on how much the decoy signals in  $X_S$  are suppressed and how much of the true signal due to  $X_S$  is preserved. Further comparison of HJ and PBHT must thus analyze how the subspaces  $\text{span}(U_A)$  are estimated in the context of the assumptions made in [5] and [9]. A final note concerns the dimension of the subspace  $\text{span}(U_A)$ . Both HJ and PBHT were proposed with the implicit goal of finding a basis  $U_A$  that has the same span as  $U_S$ . This of course requires estimating  $|S| = k$  by  $q$ , which adds another layer of complexity to these algorithms. Theorems 1 and 2 suggest that underestimating  $k$  can be more detrimental to signed support recovery than overestimating it. By overestimating  $q \geq k$ , we can trade off milder improvement when  $\text{span}(U_S) \neq \text{span}(U_A)$  against poor behavior should we have  $\text{span}(U_S) = \text{span}(U_A)$ .

5

#### Model-Based Comparisons

In the previous section we used Lemma 1 in conjunction with assumptions on  $U_A$  to make statements about HJ and PBHT. Of course, the quality of the estimated  $U_A$  depends on the specific instances  $X$ ,  $\beta$ ,  $w$ , which hinders a general analysis. Similarly, a direct application of Lemma 1 to JR yields bounds that exhibit strong  $X$  dependence. It is possible to crystallize prototypical examples by specializing  $X$  and  $w$  to come from a generative model. In this section we briefly present this model and will show the resulting penalty parameter bounds for JR. 6

##### 5.1

##### Generative model for $X$



As discussed in Section 2, many preconditioning algorithms can be phrased as truncating or reweighting column subspaces associated with  $X$  [5, 6, 9]. This suggests that a natural generative model for  $X$  can be formulated in terms of the SVD of submatrices of  $X$ . Assume  $p \leq k \leq n$  and let  $\Sigma$ ,  $\Sigma_c$  be fixed-spectrum matrices of dimension  $n \times k$  and  $n \times p \times k$  respectively. We will assume throughout this paper that the top left “diagonal” entries of  $\Sigma$ ,  $\Sigma_c$  are positive and the remainder is zero. Furthermore, we let  $U$ ,  $V_S$ ,  $V_{S_c}$  be orthonormal bases of dimension  $n \times n$ ,  $k \times k$  and  $p \times k \times p \times k$  respectively. We assume that these bases are chosen uniformly at random from the corresponding Stiefel manifold. As before and without loss of generality, suppose  $S = \{1, \dots, k\}$ . Then we let the Lasso problem be

$y = X\beta + w$  with  $X = U\Sigma V_S^T, \Sigma_c V_{S_c}^T w \sim N(0, \sigma^2 I_{n \times n})$ , (19)  
To ensure that the column norms of  $X$  are controlled, we compute the spectra  $\Sigma$ ,  $\Sigma_c$  by normalizing  $S$  and  $S_c$  with arbitrary positive elements on the diagonal. Specifically, we let  $\Sigma =$

$$\begin{aligned} \Sigma &= \frac{1}{\sqrt{k}} \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \\ \Sigma_c &= \frac{1}{\sqrt{p \times k}} \begin{bmatrix} \Sigma_{c11} & \Sigma_{c12} \\ \Sigma_{c21} & \Sigma_{c22} \end{bmatrix} \end{aligned} \quad (20)$$

We verify in the supplementary material that with these assumptions the squared column norms of  $X$  are in expectation  $n$  (provided the orthonormal bases are chosen uniformly at random). Intuition.

Note that any matrix  $X$  can be decomposed using a block-wise SVD as

$$X = [X_S, X_{S_c}] = U \begin{bmatrix} \Sigma & \Sigma_c \end{bmatrix} \begin{bmatrix} V_S^T \\ V_{S_c}^T \end{bmatrix}, \quad (21)$$

with orthonormal bases  $U$ ,  $T$ ,  $V_S$ ,  $V_{S_c}$ . Our model in Eq. (19) is only a minor restriction of this model, where we set  $T = I_{n \times n}$ . To develop more intuition, let us temporarily set  $V_S = I_{k \times k}$ ,  $V_{S_c} = I_{p \times k \times p \times k}$ . Then  $X = [X_S, X_{S_c}] = U \begin{bmatrix} \Sigma & \Sigma_c \end{bmatrix}$  and we see that up to scaling  $X_S$  equals the first  $k$  columns of  $X_{S_c}$ . The difficulty for Lasso thus lies in correctly selecting the columns in  $X_S$ , which are highly correlated with the first few columns in  $X_{S_c}$ .

## 5.2 Piecewise constant spectra

For notational clarity we will now focus on a special case of the above model. To begin, we develop some notation. In previous sections we used  $U_S$  to denote a basis for the column space of  $X_S$ . We will continue to use this notation, and let  $U_S$  contain the first  $k$  columns of  $U$ . Accordingly, we let  $\Sigma$ ,  $\Sigma_c$ ,  $\Sigma_{c,c}$  denote the last  $n \times k$  columns of  $U$  by  $U_{S_c}$ . We let the diagonal elements of  $\Sigma$ ,  $\Sigma_c$  be identified by their column indices. That is, the diagonal entries  $\Sigma_{c,c}$  of  $\Sigma$  and  $\Sigma_{c,c}$  of  $\Sigma_c$  are indexed by  $c \in \{1, \dots, k\}$ ; the diagonal entries  $\Sigma_{c,c}$  of  $\Sigma_c$  and  $\Sigma_{c,c}$  of  $\Sigma_c$  are indexed by  $c \in \{1, \dots, n\}$ . Each of the diagonal entries in  $\Sigma$ ,  $\Sigma_c$  is associated with a column of  $U$ . The set of diagonal entries of  $\Sigma$  and  $\Sigma_c$  associated with  $U_S$  is  $\Sigma(S) = \{1, \dots, k\}$  and the set of diagonal entries in  $\Sigma_c$  associated with  $U_{S_c}$  is  $\Sigma(S_c) = \{1, \dots, n\} \setminus \Sigma(S)$ . We will construct spectrum matrices  $\Sigma$ ,  $\Sigma_c$  that are piecewise constant on their diagonals. For some  $\epsilon > 0$ , we let  $\Sigma_{c,i} = 1$ ,  $\Sigma_{c,i}$



$\frac{\|W\|_F}{\|W_m\|_F} \leq \frac{\|W\|_F}{\|W_m\|_F} \leq \frac{\|W\|_F}{\|W_m\|_F}$ , (26) where the stochasticity on the left is due to  $W_m$ ,  $w_m$  and on the right is due to  $w$ . Thus, with respect to the bounds ratio  $\frac{\|W\|_F}{\|W_m\|_F}$ , the construction of Eq. (19) can be thought of as the limiting construction of Gaussian Lasso problems in Eq. (25) for large  $m$ . As such, we believe that Eq. (19) is a useful proxy for less restrictive generative models. Indeed, as the experiment in Figure 2(b) shows, Theorem 3 can be used to predict the scaling factor for penalty parameter bounds  $\frac{\lambda}{\lambda_0} \approx \frac{\|W\|_F}{\|W_m\|_F}$  with good accuracy even for Gaussian ensembles. ratios (i.e.,  $\frac{\lambda}{\lambda_0}$ )

6

## Conclusions

This paper proposes a new framework for comparing Preconditioned Lasso algorithms to the standard Lasso which skirts the difficulty of choosing penalty parameters. By eliminating this parameter from consideration, finite data comparisons can be greatly simplified, avoiding the use of model selection strategies. To demonstrate the framework's usefulness, we applied it to a number of Preconditioned Lasso algorithms and in the process confirmed intuitions and revealed fragilities and mitigation strategies. Additionally, we presented an SVD-based generative model for Lasso problems that can be thought of as the limit point of a less restrictive Gaussian model. We believe this work to be a first step towards a comprehensive theory for evaluating and comparing Lasso-style algorithms and believe that the strategy can be extended to comparing other penalized likelihood methods on finite datasets. 8

## 2 References

- [1] D.L. Donoho, M. Elad, and V.N. Temlyakov. Stable recovery of sparse over-complete representations in the presence of noise. *Information Theory, IEEE Transactions on*, 52(1):6?18, 2006.
- [2] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96:1348?1360, 2001.
- [3] J.J. Fuchs. Recovery of exact sparse representations in the presence of bounded noise. *Information Theory, IEEE Transactions on*, 51(10):3601?3608, 2005.
- [4] H.-C. Huang, N.-J. Hsu, D.M. Theobald, and F.J. Breidt. Spatial Lasso with applications to GIS model selection. *Journal of Computational and Graphical Statistics*, 19(4):963?983, 2010.
- [5] J.C. Huang and N. Jojic. Variable selection through Correlation Sifting. In V. Bafna and S.C. Sahinalp, editors, *RECOMB*, volume 6577 of *Lecture Notes in Computer Science*, pages 106?123. Springer, 2011.
- [6] J. Jia and K. Rohe. ?Preconditioning? to comply with the irrerepresentable condition. 2012.
- [7] N. Meinshausen. Lasso with relaxation. Technical Report 129, Eidgen?ossische Technische Hochschule, Z?urich, 2005.
- [8] N. Meinshausen and P. B?uhlmann. High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34(3):1436?1462, 2006.
- [9] D. Paul, E. Bair, T. Hastie, and R. Tibshirani. ?Preconditioning? for feature selection and regression in high-dimensional problems. *Annals of Statistics*, 36(4):1595?1618, 2008.
- [10] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal*

of the Royal Statistical Society, Series B, 58(1):267-288, 1994. [11] R.J. Tibshirani. The solution path of the Generalized Lasso. Stanford University, 2011. [12] M.J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (Lasso). IEEE Transactions on Information Theory, 55(5):2183-2202, 2009. [13] P. Zhao and B. Yu. On model selection consistency of Lasso. Journal of Machine Learning Research, 7:2541-2563, 2006. [14] H. Zou. The Adaptive Lasso and its oracle properties. Journal of the American Statistical Association, 101:1418-1429, 2006. [15] H. Zou and T. Hastie. Regularization and variable selection via the Elastic Net. Journal of the Royal Statistical Society, Series B, 67:301-320, 2005.