

# Coordinate descent for SLOPE

Anonymous Author  
Anonymous Institution

## Abstract

Table 1: Example of the permutation operator  $(i)$  and its inverse  $(i)^{-}$

| $i$ | $\beta_i$ | $(i)$ | $(i)^{-}$ |
|-----|-----------|-------|-----------|
| 1   | 0.5       | 2     | 3         |
| 2   | -5        | 3     | 1         |
| 3   | 4         | 1     | 2         |

## 1 Introduction

In this paper we present a novel numerical algorithm for Sorted L-One Penalized Estimation (SLOPE) [Bog+13; Bog+15; ZF14], defined as

$$\min_{\beta \in \mathbb{R}^p} P(\beta) = \ell(\beta) + J(\beta) \quad (1)$$

where we take  $\ell$  to be smooth and twice differentiable and

$$J(\beta) = \sum_{j=1}^p \lambda_j |\beta_{(j)}| \quad (2)$$

is the *sorted*  $\ell_1$  norm, defined through

$$|\beta_{(1)}| \geq |\beta_{(2)}| \geq \dots \geq |\beta_{(p)}|,$$

with  $\lambda$  being a fixed non-increasing and non-negative sequence.

The sorted  $\ell_1$  norm is a sparsity-enforcing penalty that has become increasingly popular due to several appealing properties, such as its ability to control false discovery rate [Bog+15; KB20], cluster coefficients [FN16; ST20], and recover sparsity and ordering patterns in the solution [Bog+22]. Unlike other competing sparse regularization methods such as MCP [Zha10] and SCAD [FL01], SLOPE is also a convex problem [Bog+15].

In spite of the availability of predictor screening rules [LBW20; EH22], which help speed up SLOPE in the high-dimensional regime, current state-of-the-art algorithms for SLOPE perform poorly in comparison to those of more established penalization methods such as the lasso ( $\ell_1$ -norm regularization) and ridge regression

( $\ell_2$ -norm regularization). As a small illustration of this issue, we compared the speed at which the SLOPE and glmnet packages fit a complete regularization path for the bcTCGA data set. SLOPE takes x seconds to fit the full path, whilst glmnet requires only y seconds.

This lack of cluster performance has hampered the applicability of SLOPE to many real-world applications. A major reason for why algorithms for solving  $\ell_1$ -, MCP-, or SCAD-regularized problems enjoy better performance is that they use coordinate descent [Tse01; FHT10; BH11]. Current SLOPE solvers, on the other hand, rely on proximal gradient descent algorithms such as FISTA [BT09] and the alternating direction method of multipliers method (ADMM) [Boy+10], which have proven to be less efficient than coordinate descent in empirical benchmarks on related problems, such as the lasso [Mor+22]. Applying coordinate descent to SLOPE is not, however, straightforward since convergence guarantees for coordinate descent require the objective to be separable, which is not the case for SLOPE. As a result, naive coordinate descent schemes can get stuck (figure 1).

In this article we address this problem by introducing a new, highly effective algorithm for SLOPE based on a hybrid proximal gradient and coordinate descent scheme. Our method features convergence guarantees and reduces the time required to fit SLOPE by orders of magnitude in our empirical experiments.

**Notation** Let  $(i)^{-}$  be the inverse of  $(i)$  such that  $((i)^{-})^{-} = (i)$ . See table 1 for an example of this operator for a particular  $\beta$ . This means that

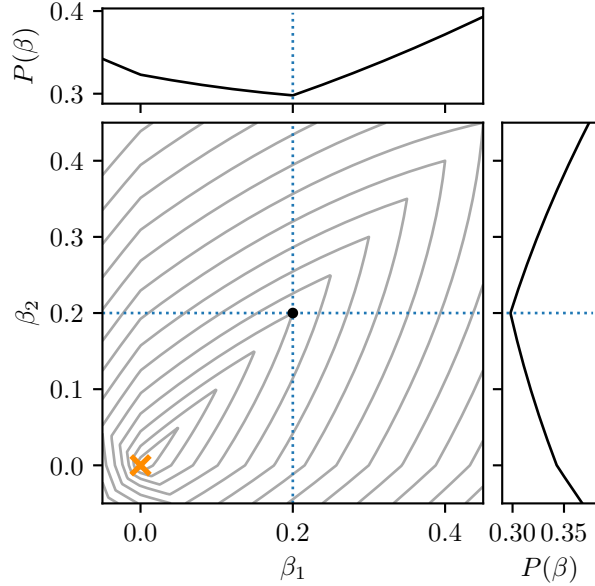


Figure 1: An example of (standard) coordinate descent getting stuck on a two-dimensional SLOPE problem. The main plot shows level curves for the primal objective (1), with the optimum  $\beta^* = [0, 0]^T$  indicated by the orange cross. The marginal plots displays objective values at  $\beta_1 = 0.2$  when optimizing over  $\beta_2$  and vice versa. At  $\beta = [0.2, 0.2]^T$ , a naive coordinate descent algorithm can only move in the directions indicated by the dashed lines—neither of which are descent directions for the objective. As a result, the algorithm is stuck.

$$J(\beta) = \sum_{j=1}^p \lambda_j |\beta_{(j)}| = \sum_{j=1}^p \lambda_{(j)} |\beta_j|.$$

Sorted  $\ell_1$  norm penalization produces vectors with clustered coefficients. To this end, for a fixed  $\beta$  such that  $|\beta_j|$  takes  $m$  values, we introduce  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m$  and  $c_1, c_2, \dots, c_m$  for the indices and coefficients respectively of the  $m$  clusters of  $\beta$ , such that  $\mathcal{C}_i = \{j : |\beta_j| = c_i\}$  and  $c_1 > c_2 > \dots > c_m \geq 0$ . We also let  $\bar{\mathcal{C}}$  denote the complement of  $\mathcal{C}$ . Let  $(e_i)_{i \in [d]}$  denote the canonical basis of  $\mathbb{R}^d$ . Finally, let  $\text{sign}(x) = x/|x|$  (with the convention  $0/0 = 1$ ) be the scalar sign, that acts entrywise on vectors.

## 2 Theory

### 2.1 Coordinate Descent for SLOPE

Proximal coordinate descent cannot be applied to Problem (1) because it is not separable. However, if the clusters  $\mathcal{C}_1^*, \dots, \mathcal{C}_{m^*}^*$  of the solution  $\beta^*$  were known together with their signs, then the values  $c_1^*, \dots, c_{m^*}^*$  taken by  $\beta^*$  on the clusters could be equivalently computed by solving:

$$\min_{z \in \mathbb{R}^{m^*}} \ell \left( \sum_{i=1}^{m^*} \sum_{j \in \mathcal{C}_i} z_i \text{sign}(\beta_j^*) e_j \right) + \sum_{i=1}^{m^*} |z_i| \sum_{\lambda \in \lambda_{\mathcal{C}_i^*}} \lambda, \quad (3)$$

Conditionally on the knowledge of the clusters and the signs of the coefficients, the penalty becomes separable and coordinate descent could be used. Hence, the idea of our algorithm is to intertwine steps that identify the clusters, and large coordinate-wise steps on these clusters.

Based on this observation, we derive a coordinate descent algorithm for minimizing the SLOPE problem (1) with respect to the coefficients of a single cluster at a time.

In all the sequel,  $\beta$  is fixed, with  $m$  clusters  $\mathcal{C}_1, \dots, \mathcal{C}_m$  corresponding to values  $c_1, \dots, c_m$ . Let also  $k \in [m]$  be fixed, let  $s = \text{sign} \beta_{\mathcal{C}_k}$ . We are interested in updating  $\beta$  by changing only the value taken on the  $k$ -th cluster. To this end, we let

$$\beta_i(z) = \begin{cases} s_k z, & \text{if } i \in \mathcal{C}_k, \\ \beta_i, & \text{otherwise.} \end{cases} \quad (4)$$

Minimizing the objective in this directions amounts to solving the following one-dimensional problem:

$$\min_{z \in \mathbb{R}} \left( G(z) = P(\beta(z)) = F(\beta(z)) + \phi(z) \right), \quad (5)$$

Mathurin: T  
prop and writ

where

$$\phi(z) = |z| \sum_{j \in \mathcal{C}_k} \lambda_{(j)z}^- + \sum_{j \notin \mathcal{C}_k} |\beta_j| \lambda_{(j)z}^-$$

is the *partial sorted  $\ell_1$  norm* with respect to the  $k$ -th cluster and where we write  $\lambda_{(j)z}^-$  to indicate that the inverse sorting permutation  $(j)_z^-$  is defined with respect to  $\beta(z)$ . The optimality condition for [Problem \(5\)](#) is

$$G'(z; \delta) \geq 0,$$

where  $G'(z; \delta)$  is the directional derivative in the direction  $\delta$ . Since  $F$  is differentiable we have

$$G'(z; \delta) = \delta \nabla F(\beta(z)) + \phi'(z; \delta),$$

where  $\phi'(z; \delta)$  is the directional derivative of  $\phi$ .

Throughout the rest of this section we derive a solution to (5). To do so, we will introduce the directional derivative for the sorted  $\ell_1$  norm with respect to the coefficient of the  $k$ -th cluster. First, let  $C(z)$  be a function that returns the cluster corresponding to  $z$ , that is

$$C(z) = \{j : |\beta(z)_j| = z\}.$$

Next, let  $\varepsilon_c$  be an arbitrary positive value such that

$$\varepsilon_c < |c_i - c_j|, \quad \forall i \neq j \text{ and } \varepsilon_c < c_m \text{ if } c_m \neq 0. \quad (6)$$

Then for  $\delta \in \{-1, 1\}$  we have

$$\lim_{h \downarrow 0} C(z+h\delta) = C(z+\varepsilon_c\delta) \quad \text{and} \quad \lim_{h \downarrow 0} \lambda_{(i)z+h\delta}^- = \lambda_{(i)z+\varepsilon_c\delta}^-. \quad (7)$$

In other words, the order permutation corresponding to  $\beta(z+h\delta)$  depends only on  $\delta$  as  $h$  tends to 0.

*Remark 2.1.* As consequence of the definition of  $\varepsilon_c$ , the order permutations for  $\beta(z)$  and  $\beta(z+\varepsilon_c\delta)$  differ only for a subset of the permutation vectors. The permutation for  $\beta(h\delta)$  corresponds to

$$\begin{cases} \mathcal{C}_1, \dots, \mathcal{C}_{k-1}, \mathcal{C}_{k+1}, \dots, \mathcal{C}_m, C(\varepsilon_c\delta) & \text{if } c_m > 0, \\ \mathcal{C}_1, \dots, \mathcal{C}_{k-1}, \mathcal{C}_{k+1}, \dots, C(\varepsilon_c\delta), \mathcal{C}_m & \text{if } c_m = 0. \end{cases}$$

If  $z \neq 0$ , the permutation for  $\beta(z+\varepsilon_c\delta)$  corresponds to

$$\begin{cases} \mathcal{C}_1, \dots, \mathcal{C}_{k-1}, \mathcal{C}_{k+1}, \dots, C(c_i + \varepsilon_c\delta), \mathcal{C}_i, \dots, \mathcal{C}_m & \text{if } \delta = 1, \\ \mathcal{C}_1, \dots, \mathcal{C}_{k-1}, \mathcal{C}_{k+1}, \dots, \mathcal{C}_i, C(c_i + \varepsilon_c\delta), \dots, \mathcal{C}_m & \text{if } \delta = -1. \end{cases}$$

We are now ready to state the directional derivative of  $\phi$ .

**Theorem 2.2.** Let  $c^{\setminus k}$  be the  $m-1$  length version of  $c$  where the  $k$ -th coordinate has been removed:  $c^{\setminus k} = (c_1, \dots, c_{k-1}, c_{k+1}, \dots, c_m)$ . Let  $\varepsilon_c$  defined as in (6).

The directional derivative of the partial sorted  $\ell_1$  norm with respect to the  $k$ -th cluster  $\phi$ , in the direction  $\delta$  is

$$\phi'(z; \delta) = \begin{cases} \sum_{j \in C(\varepsilon_c)} \lambda_{(j)\varepsilon_c}^- & \text{if } z = 0, \\ \text{sign}(z)\delta \sum_{j \in C(z+\varepsilon_c\delta)} \lambda_{(j)z+\varepsilon_c\delta}^- & \text{if } |z| = c_i^{(k)} > 0, \\ \text{sign}(z)\delta \sum_{j \in C(z)} \lambda_{(j)z}^- & \text{otherwise.} \end{cases}$$

*Proof.* By definition of the directional derivative,

$$\begin{aligned} \phi'(z; \delta) &= \lim_{h \downarrow 0} \frac{\phi(z+h\delta) - \phi(z)}{h} \\ &= \lim_{h \downarrow 0} \frac{1}{h} \left( |z+h\delta| \sum_{j \in C(z+h\delta)} \lambda_{(j)z+h\delta}^- + \sum_{j \notin C(z+h\delta)} |\beta_j| \lambda_{(j)z+h\delta}^- - |z| \sum_{j \in C(z)} \lambda_{(j)z}^- \right) \\ &= \lim_{h \downarrow 0} \frac{1}{h} \left( |z+h\delta| \sum_{j \in C(z+\varepsilon_c\delta)} \lambda_{(j)z+\varepsilon_c\delta}^- + \sum_{j \notin C(z+\varepsilon_c\delta)} |\beta_j| \lambda_{(j)z+\varepsilon_c\delta}^- - |z| \sum_{j \in C(z)} \lambda_{(j)z}^- \right) \end{aligned} \quad (8)$$

We have the following cases to consider:  $|z| \in \{c_i^{(k)}\}_{i=1}^{m-1}$ ,  $|z| = 0$ , and  $|z| \notin \{0, \{c_i^{(k)}\}_{i=1}^{m-1}\}$ .

Starting with  $|z| = c_i^{(k)}$ , note that we have, as a result of the definition of  $\varepsilon_c$ , the following identities:

$$\begin{aligned} C(c_i^{(k)} + \varepsilon_c\delta) &\subseteq C(c_i^{(k)}), \\ \tilde{C}_i &= \overline{C(c_i^{(k)} + \varepsilon_c\delta)} \cap C(c_i^{(k)}), \\ C(c_i^{(k)}) &= \tilde{C}_i \cup (C(c_i^{(k)} + \varepsilon_c\delta) \cap C(c_i^{(k)})) = \tilde{C}_i \cup C(c_i^{(k)} + \varepsilon_c\delta), \\ \overline{C(c_i^{(k)} + \varepsilon_c\delta)} &= \overline{C(c_i^{(k)})} \cup \tilde{C}_i. \end{aligned}$$

Using this, we can rewrite (8) as

$$J'_k(z; \delta) = \lim_{h \downarrow 0} \frac{1}{h} \left( |z+h\delta| \sum_{j \in C(z+\varepsilon_c\delta)} \lambda_{(j)z+\varepsilon_c\delta}^- + |c_i^{(k)}| \sum_{j \in \tilde{C}_i} \lambda_{(j)z+\varepsilon_c\delta}^- + \sum_{j \in \overline{C(z)}} |\beta_j| \lambda_{(j)z}^- - |z| \sum_{j \in \tilde{C}_i} \lambda_{(j)z}^- - |c_i^{(k)}| \sum_{j \in C(z+\varepsilon_c\delta)} \lambda_{(j)z+\varepsilon_c\delta}^- - \sum_{j \in \overline{C(z)}} |\beta_j| \lambda_{(j)z}^- \right)$$

Next, observe that  $\lambda_{(j)z+\varepsilon_c\delta}^- = \lambda_{(j)z}^-$  for all  $j \in \overline{C(z)}$  and consequently

$$\sum_{j \in \overline{C(z)}} |\beta_j| \lambda_{(j)z+\varepsilon_c\delta}^- = \sum_{j \in \overline{C(z)}} |\beta_j| \lambda_{(j)z}^-.$$

Moreover, note that, since  $z = \pm c_i$ , there exists a permutation corresponding to  $\lambda_{(j)z}^-$  such that

$$|z| \sum_{j \in \tilde{C}_i} \lambda_{(j)z+\varepsilon_c\delta}^- = |c_i^{(k)}| \sum_{j \in \tilde{C}_i} \lambda_{(j)z}^-$$

and consequently

$$J'_k(z; \delta) = \lim_{h \downarrow 0} \frac{1}{h} \left( |z + h\delta| \sum_{j \in C(z + \varepsilon_c \delta)} \lambda_{(j)}^- - |c_i^{(k)}| \sum_{j \in C(z + \varepsilon_c \delta)} \lambda_{(j)}^- \right). \quad (9)$$

Now, since  $c_i^{(k)} + h\delta > 0$  and  $-c_i^{(k)} + h\delta < 0$  in the limit as  $h$  goes to 0 for  $c_i \neq 0$ , we have

$$\lim_{h \downarrow 0} |-c_i + h\delta| = \lim_{h \downarrow 0} (|c_i^{(k)}| - h\delta) \quad \text{and} \quad \lim_{h \downarrow 0} |c_i^{(k)} + h\delta| = \lim_{h \downarrow 0} (|c_i^{(k)}| + h\delta)$$

which means that

$$\begin{aligned} J'_k(z; \delta) &= \lim_{h \downarrow 0} \frac{1}{h} \left( (|z| + \text{sign}(z)h\delta) \sum_{j \in C(z + \varepsilon_c \delta)} \lambda_{(j)}^- - |c_i| \sum_{j \in C(z + \varepsilon_c \delta)} \lambda_{(j)}^- \right) \\ &= \lim_{h \downarrow 0} \frac{1}{h} \text{sign}(z)h\delta \sum_{j \in C(z + \varepsilon_c \delta)} \lambda_{(j)}^- \\ &= \text{sign}(z)\delta \sum_{j \in C(z + \varepsilon_c \delta)} \lambda_{(j)}^- \end{aligned}$$

For the case when  $z = 0$  there are two sub-cases to consider. First, when  $c_m^{(k)} = 0$  then (9) is simply

$$J'_k(0; \delta) = \lim_{h \downarrow 0} \frac{1}{h} |h\delta| \sum_{j \in C(\varepsilon_c \delta)} \lambda_{(j)}^- = \sum_{j \in C(\varepsilon_c \delta)} \lambda_{(j)}^-$$

since  $|\delta| = 1$  by definition.

The other sub-case is  $c_m^{(k)} \neq 0$ . Here (8) reduces to

$$J'_k(0; \delta) = \lim_{h \downarrow 0} \frac{1}{h} \left( |h\delta| \sum_{j \in C(\varepsilon_c \delta)} \lambda_{(j)}^- + \sum_{j \notin C(\varepsilon_c \delta)} |\beta_j| \lambda_{(j)}^- - \sum_{j \notin C(0)} |\beta_j| \lambda_{(j)}^- \right)$$

In this case, we have  $C(0) = C(\varepsilon_c \delta)$  since  $c_m^{(k)} \neq 0$ , and therefore

$$\sum_{j \notin C(\varepsilon_c \delta)} |\beta_j| \lambda_{(j)}^- = \sum_{j \notin C(0)} |\beta_j| \lambda_{(j)}^-$$

which means that

$$J'_k(0; \delta) = \lim_{h \downarrow 0} \frac{1}{h} |h\delta| \sum_{j \in C(\varepsilon_c \delta)} \lambda_{(j)}^- = |\delta| \sum_{j \in C(\varepsilon_c \delta)} \lambda_{(j)}^- = \sum_{j \in C(0)} \lambda_{(j)}^-.$$

Finally, note that  $J_k$  is differentiable in every other case, i.e.  $|z| \notin \{0, \{c_i^{(k)}\}_{i=1}^{m-1}\}$ . Here since  $C(z + \varepsilon_c \delta) = C(z)$ , has the directional derivative

$$J'_k(z; \delta) = \delta \text{sign}(z) \sum_{j \in C(z)} \lambda_{(j)}^-.$$

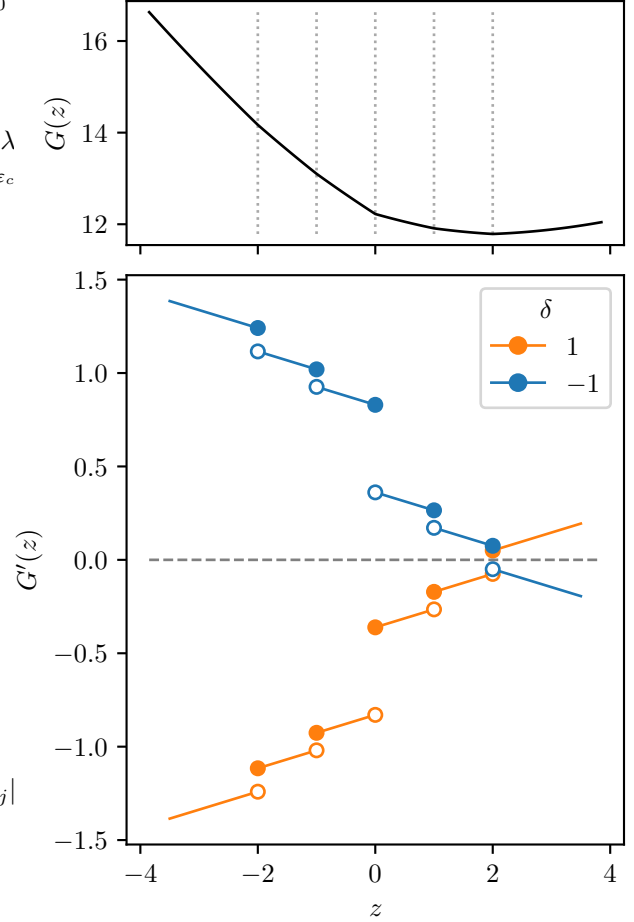


Figure 2: The function  $P_k(z)$  and its directional derivative  $P'_k(z; \delta)$  for an example with  $\beta = [-3.0, 1.0, 3.0, 2.0]^T$ ,  $k = 1$ , and consequently  $c^{(k)} = [2, 1]^T$ . The solution corresponds to the value of  $z$  for which  $P'_k(z; \delta) \geq 0$  for  $\delta \in \{-1, 1\}$ , which holds only at  $z = 2$ , which must therefore be the solution.

□

In figure 2, we show an example of the directional derivative and the objective function.

We are now ready to present the SLOPE thresholding operator.

**Theorem 2.3** (The SLOPE Thresholding Operator).

$$T_k(\gamma; \omega, c, \lambda) = \begin{cases} 0 & \text{if } |\gamma| \leq \sum_{j \in C(\varepsilon_c)} \lambda_{(j)\varepsilon_c}^- \\ \text{sign}(\gamma)c_i & \text{if } \omega c_i + \sum_{j \in C(c_i - \varepsilon_c)} \lambda_{(j)c_i - \varepsilon_c}^- \leq |\gamma| \leq \omega c_i + \sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- \\ \frac{\text{sign}(\gamma)}{\omega} \left( |\gamma| - \sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- \right) & \text{if } \omega c_i + \sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- < |\gamma| < \omega c_{i-1} + \sum_{j \in C(c_{i-1} - \varepsilon_c)} \lambda_{(j)c_{i-1} - \varepsilon_c}^- \\ \frac{\text{sign}(\gamma)}{\omega} \left( |\gamma| - \sum_{j \in C(c_{i-1} - \varepsilon_c)} \lambda_{(j)c_{i-1} - \varepsilon_c}^- \right) & \text{if } |\gamma| \geq \omega c_{i-1} + \sum_{j \in C(c_{i-1} - \varepsilon_c)} \lambda_{(j)c_{i-1} - \varepsilon_c}^- \end{cases}$$

Furthermore, since  $P_k$  is differentiable in  $(c_i^{(k)}, c_{i-1}^{(k)})$ , we have  $\frac{\partial}{\partial z} P_k(z) \Big|_{z=|z^*|} = \omega z^* + \sum_{j \in C(z^*)} \text{sign}(z^*) \lambda_{(j)z^*}^- = 0$ , and therefore (11) must be the solution. The solution for the last case follows using reasoning analogous to that of the third case.  $\square$

with  $\varepsilon_c$  defined as in (6) and let  $\gamma = \tilde{r}^T x$ ,  $\omega = \tilde{x}^T \tilde{x}$ . Then  $T(\gamma; \omega, c^{(k)}, \lambda) \in \arg \min_{z \in \mathbb{R}} P_k(z)$ .

*Proof.* Recall that  $P_k(z) : \mathbb{R} \mapsto \mathbb{R}$  is a convex, continuous piecewise-differentiable function with kinks whenever  $|z| = c_i^{(k)}$  or  $z = 0$ . Let  $\gamma = \tilde{r}^T \tilde{x}$  and  $\omega = \tilde{x}^T \tilde{x}$  and note that the optimality criterion for (5) is

$$\delta(\omega z - \gamma) + J'_k(z; \delta) \geq 0, \quad \forall \delta \in \{-1, 1\},$$

which is equivalent to

$$\omega z - J'_k(z; -1) \leq \gamma \leq \omega z + J'_k(z; 1). \quad (10)$$

We now proceed to show that there is a solution  $z^* \in \arg \min_{z \in \mathbb{R}} J_k(z)$  for every interval over  $\gamma \in \mathbb{R}$ .

First, assume that the first case in the definition of  $T_k$  holds and note that this is equivalent to (10) with  $z = 0$  since  $C(\varepsilon_c) = C(-\varepsilon_c)$  and  $\lambda_{(j)\varepsilon_c}^- = \lambda_{(j)-\varepsilon_c}^-$ . This is sufficient for  $z^* = 0$ .

Next, assume that the second case holds and observe that this is equivalent to (10) with  $z = c_i^{(k)}$ , since  $C(c_i + \varepsilon_c) = C(-c_i - \varepsilon_c)$  and  $C(-c_i + \varepsilon_c) = C(c_i - \varepsilon_c)$ . Thus  $z^* = \text{sign}(\gamma)c_i^{(k)}$ .

For the third case, we have

$$\sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- = \sum_{j \in C(c_{i-1} - \varepsilon_c)} \lambda_{(j)c_{i-1} - \varepsilon_c}^-$$

and therefore (10) is equivalent to

$$c_i < \frac{1}{\omega} \left( |\gamma| - \sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- \right) < c_{i-1}.$$

Now let

$$z^* = \frac{\text{sign}(\gamma)}{\omega} \left( |\gamma| - \sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- \right) \quad (11)$$

and note that  $|z^*| \in (c_i^{(k)}, c_{i-1}^{(k)})$  and hence

$$\frac{1}{\omega} \left( |\gamma| - \sum_{j \in C(c_i + \varepsilon_c)} \lambda_{(j)c_i + \varepsilon_c}^- \right) = \frac{1}{\omega} \left( |\gamma| - \sum_{j \in C(z^*)} \lambda_{(j)z^*}^- \right).$$

Furthermore, since  $P_k$  is differentiable in  $(c_i^{(k)}, c_{i-1}^{(k)})$ ,

we have  $\frac{\partial}{\partial z} P_k(z) \Big|_{z=|z^*|} = \omega z^* + \sum_{j \in C(z^*)} \text{sign}(z^*) \lambda_{(j)z^*}^- = 0$ , and therefore (11) must be the solution. The solution for the last case follows using reasoning analogous to that of the third case.  $\square$

In figure 3, we visualize the SLOPE thresholding operator.

### 2.1.1 Naive Updates

As in Friedman, Hastie and Tibshirani [FHT10], we can improve the efficiency of updates by observing that

$$\begin{aligned} \tilde{r}_k &= y - \tilde{y}_k \\ &= y - X_{\bar{c}_k} \beta_{\bar{c}_k} - \tilde{x} c_k + \tilde{x} c_k \\ &= r + \tilde{x} c_k \end{aligned}$$

and therefore that

$$\tilde{x}_k^T (y - \tilde{y}_k) = \tilde{x}_k^T r + \tilde{x}_k^T \tilde{x}_k c_k. \quad (12)$$

### 2.1.2 Caching Reductions

Observe that  $\tilde{x}_k$  only changes between subsequent coordinate updates provided that the members of the cluster  $k$  change, for instance if two clusters are merged, a predictor leaves a cluster, or the signs flip (through an update of  $\alpha_k$ ). As a result, it is possible to obtain computational gains by caching  $\tilde{x}_k$  and  $\tilde{x}_k^T \tilde{x}_k$  for each cluster (except the zero cluster, which we do not consider in our coordinate descent step). When there is no change in the clusters, there is no need to recompute these quantities. And even when there are changes, we can still reduce the costs involved since  $\tilde{x}_k$  can be updated in place. If a large cluster is joined by few new predictors, then the cost of updating may be much lower than recomputing the quantities for the entire cluster. Also note that, for single-member clusters we only need to store  $\tilde{x}_k^T \tilde{x}_k$  since  $\tilde{x}_k$  is just a column in  $X$  times the corresponding sign.

Letting  $\tilde{x}_k^{\text{old}}$  correspond to the value of  $\tilde{x}_k$  before the update, we note that  $\tilde{x}_k \leftarrow \tilde{x}_k^{\text{old}} + x_j \text{sign}(\beta_j)$  for each  $j \in C_k^{\text{new}} \setminus C_k^{\text{old}}$  and  $\tilde{x}_k \leftarrow \tilde{x}_k^{\text{old}} - x_j \text{sign}(\beta_j)$  for each  $j \in C_k^{\text{old}} \setminus C_k^{\text{new}}$ . If only the signs flip, we simply have to also flip the signs in  $\tilde{x}_k$ .

**JL:** Consider showing that generalizes the operator.

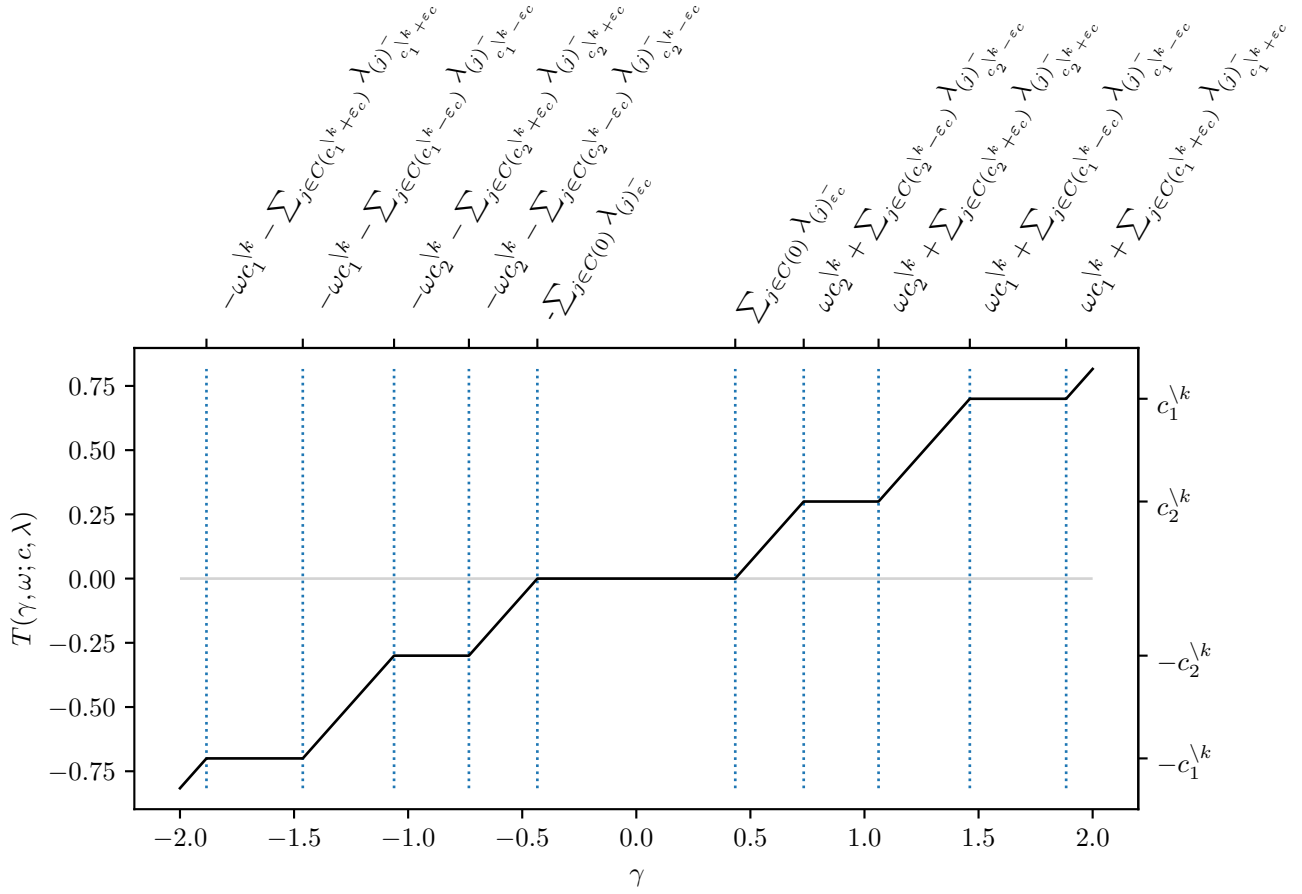


Figure 3: An example of the SLOPE thresholding operator. The result corresponds to an example for  $\beta = [0.5, -0.5, 0.3, 0.7]^T$ ,  $c = [0.7, 0.5, 0.3]^T$  with an update for the second cluster ( $k = 2$ ), such that  $c^{\setminus k} = [0.5, 0.3]^T$ . Across regions where the function is constant, the operator sets the result to be either exactly 0 or to the value of one of the elements of  $c^{\setminus k}$ .

### 2.1.3 Covariance Updates

Notice that we can rewrite the first term in (12) as

$$\begin{aligned}
 \tilde{x}_k^T r &= \tilde{x}_k^T y - \sum_{j:\beta_j \neq 0} \tilde{x}_k^T x_j \beta_j \\
 &= \tilde{x}_k^T y - \sum_{j:c_j \neq 0} \tilde{x}_k^T \tilde{x}_j c_j \\
 &= s_{\tilde{C}_k}^T X_{\tilde{C}_k}^T y - \sum_{j:\beta_j \neq 0} s_{\tilde{C}_k}^T X_{\tilde{C}_k}^T x_j \beta_j \\
 &= s_{\tilde{C}_k}^T (X^T y)_{\tilde{C}_k} - \sum_{j:\beta_j \neq 0} s_{\tilde{C}_k}^T X_{\tilde{C}_k}^T x_j \beta_j \\
 &= \sum_{j \in \tilde{C}_k} \left( s_j x_j^T y - \sum_{t:\beta_t \neq 0} s_j x_j^T x_t \beta_t \right)
 \end{aligned} \tag{13}$$

As in Friedman, Hastie and Tibshirani [FHT10], this formulation can be used to achieve so-called *covariance updates*. We compute  $X^T y$  once at the start. Then, each time a new predictor becomes non-zero, we compute its inner product with all other predictors, caching these products.

## 2.2 Hybrid proximal coordinate descent strategy

We propose an iterative solver that alternates between proximal gradient step and proximal coordinate descent. Since the regularization term for SLOPE is not separable, applying PCD does not guarantee convergence. However, [DT22] showed that once the clusters are known, the subdifferential of  $J$  can be written as the cartesian product of the subdifferential of  $J$  restricted to the clusters. Hence, if one knew the clusters, PCD updates could be applied on each cluster.

The notion of clusters for SLOPE extends the notion of sparsity coming from the LASSO. Identification of the sparsity pattern throughout the iterative algorithm have largely been studied. Talk about Partly Smooth functions, related Manifold and that support transpose to cluster for SLOPE regularization. DO the maths.

Then identification of this underlying structure occurs when applying PGD. Hence the idea to alternate, PGD and PCD step to take advantage of the speed of PCD and ensure convergence via the identification of the right structure with PGD steps.

In figure 4, we show how algorithm 1 works in practice on a two-dimensional SLOPE problem.

**Lemma 2.4.** *Let  $\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(k)}$  be a sequence of iterates generated by algorithm 1,  $1/t$  the frequency of proximal gradient descent iterates in algorithm 1, and*

**Algorithm 1** Hybrid coordinate descent and proximal gradient descent algorithm for SLOPE

---

**input:**  $X \in \mathbb{R}^{n \times p}, y \in \mathbb{R}^n, \beta \in \mathbb{R}^p, \lambda \in \{\mathbb{R}^p : \lambda_1 \geq \lambda_2 \geq \dots > 0\}, m \in \mathbb{N}$

**init** :  $t \leftarrow 0, \beta \leftarrow 0, L \leftarrow \|X\|_2^2$

- 1 **repeat**
- 2      $t \leftarrow t + 1$
- 3     **if**  $t \bmod m = 0$  **then**
- 4          $\beta \leftarrow \text{prox}_{J/L}(\beta - \frac{1}{L} \nabla f(\beta))$
- 5     **else**
- 6          $k \leftarrow 0$
- 7         **while**  $k \leq |C|$  **do**
- 8              $k \leftarrow k + 1$
- 9              $s \leftarrow \text{sign}(\beta_{\tilde{C}_k})$
- 10            **if**  $s \neq 0$  **then**
- 11                 $\tilde{x}_k \leftarrow X_{\tilde{C}_k} s$
- 12                 $\tilde{y}_k \leftarrow X_{\tilde{C}_k}^T \beta_{\tilde{C}_k}$
- 13                 $\tilde{r}_k \leftarrow y - \tilde{y}_k$
- 14                 $\beta_{\tilde{C}_k} \leftarrow sT\left(\frac{\tilde{r}_k^T \tilde{x}_k}{\tilde{x}_k^T \tilde{x}_k}, \frac{\lambda}{\tilde{x}_k^T \tilde{x}_k}, k, \tilde{\beta}, \tilde{C}\right) \triangleright \tilde{C} \text{ is updated at this step.}$
- 15 **until** convergence;
- 16 **return**  $\beta$

---

$L$  the Lipschitz constant of  $\nabla \ell$ . Then

$$P(\beta^{(k)}) - P(\beta^*) \leq \frac{L \|\beta^{(0)} - \beta^*\|_2^2}{2 \lfloor k/t \rfloor}.$$

Here  $\beta^* \in \arg \min_{\beta \in \mathbb{R}^p} P$ .

*Proof.* Assume that we are at iteration  $k + 1$  of algorithm 1 and observe that problem (5) can be written as a constrained form of problem (1):

$$\begin{aligned}
 &\min_{\beta \in \mathbb{R}^p} P(\beta), \\
 &\text{subject to } |\beta_i| = |\beta_j| \quad i, j \in \tilde{C}_k, i \neq j \\
 &\quad \text{sign}(\beta_i) \text{sign}(\beta_i^{(k)}) = \text{sign}(\beta_j) \text{sign}(\beta_j^{(k)}) \quad i, j \in \tilde{C}_k, \\
 &\quad \beta_i = \beta_i^{(k)} \quad i \notin \tilde{C}_k.
 \end{aligned}$$

And note that these constraints hold for any  $\beta^{(k)}$  used as input to this problem. In theorem 2.3, we showed that the thresholding operator  $T$  minimizes this problem for any  $\beta^{(k)}$ , which means that  $P(\beta^{(k)}) \leq P(\beta^{(k+1)})$  for any coordinate descent iteration  $k + 1$  in algorithm 1.

The convergence properties of proximal gradient descent are well known [BT09; DDD04], and hence we know that for a sequence of iterates of the proximal gradient step (line 4),  $\beta^{(t)}, \beta^{(2t)}, \dots, \beta^{(\lfloor k/t \rfloor)}$ , it holds that

$$P(\beta^{(\lfloor k/t \rfloor)}) - P(\beta^*) \leq \frac{L \|\beta^{(0)} - \beta^*\|_2^2}{2 \lfloor k/t \rfloor}.$$

JW: We only

JW: I don't t  
this level of d  
that line 9-14  
 $\beta^{t+k/|C|} = T$   
and by theore  
that  $P(\beta^{t+k/|C|})$   
 $P(\beta^{t+(k-1)/|C|})$   
approximately

JW: I suggest  
formal. Cite o  
article where i



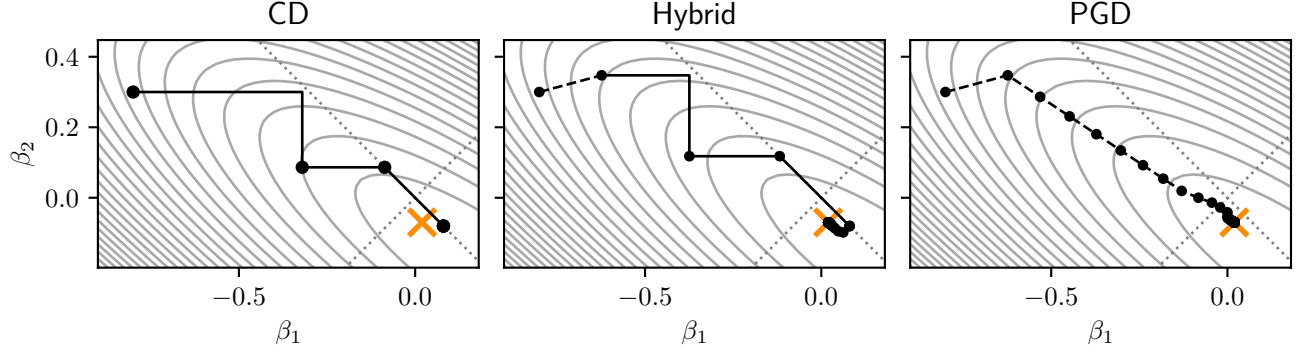


Figure 4: Illustration of the proposed solver. The figures show progress until convergence for the coordinate descent (CD) solver that we use as part of the hybrid method, our hybrid method, and proximal gradient descent (PGD). The orange cross marks the optimum. Dotted lines indicate where the coefficients are equal in absolute value. The dashed lines indicate PGD steps and solid lines CD steps. Each dot marks a complete epoch, which may correspond to only a single coefficient update for the CD and hybrid solvers if the coefficients flip order. Each solver was run until the duality gap was smaller than  $10^{-10}$ . Note that the CD algorithm cannot split clusters and is therefore stuck after the third epoch. The hybrid and PGD algorithms, meanwhile, reach convergence after 67 and 156 epochs respectively.

Combining this and the result from the previous paragraph yields the desired result.  $\square$

### 3 Experiments

The proposed algorithm is part of a python package relying on numpy and numba [Har+20; LPS15]. It will be made open-source upon publication. To compare its efficiency, we used several public datasets described in section 3. We performed an extensive benchmark with the following competitors:

- Alternating Direction Method of Multipliers (admm) [Boy+10]
- Anderson acceleration for proximal gradient descent (anderson) [ZOB20]
- Proximal Gradient Descent (pgd) [CW05]
- Fast Iterative Shrinkage-Thresholding Algorithm (fista) [BT09]
- Semismooth Newton-Based Augmented Lagrangian (newt-alt) [Luo+19]
- The Oracle solver (oracle) uses the clusters obtained via another solver to compute coordinate descent updates from the known solver.
- The Hybrid (ours) (hybrid) solver (see algorithm 1) combines proximal gradient descent and coordinate descent to overcome the non-separability of the SLOPE problem.

We used the benchopt [Mor+22] tool to obtain the convergence curves for the different solvers. Benchopt launches each solver several times increasing the number of iterations and store the objective value, dual gap and time to reach it. The repository to reproduce

| Datasets    | #samples n | #features p | density |
|-------------|------------|-------------|---------|
| Simulated 1 | 200        | 20 000      | 1       |
| Simulated 2 | 20 000     | 200         | 1       |
| Simulated 3 | 200        | 2 000 000   | 0.001   |
| Rhee2006    | 842        | 361         | ?       |
| bcTCGA      | 536        | 17 322      | 1       |
| Scheetz2006 | 120        | 18 975      | ?       |

the benchmark is available at XXX.

#### 3.1 Simulated data

The design matrix  $X$  is simulated with correlation between features  $j$  and  $j'$  equal to  $\rho^{|j-j'|}$  where  $\rho$  is a parameter that can be chosen in  $[0, 1[$ . The true regression vector  $\beta^*$  contains a certain number of non-zero entries that are obtained by simulating a gaussian distribution with zero mean and unit variance. The observations are equal to  $y = X\beta^* + \varepsilon$  where  $\varepsilon$  is a centered gaussian distribution with variance such that  $\|X\beta^*\|/\|\varepsilon\| = 3$ .

#### 3.2 Real data

### 4 Discussion

### References

- [BH11] Patrick Breheny and Jian Huang. ?Coordinate Descent Algorithms for Nonconvex Penalized Regression, with Applications to Biological Feature Selection? in *The Annals of*

Klopfe: Do v  
three different  
that changes  
lambdas?



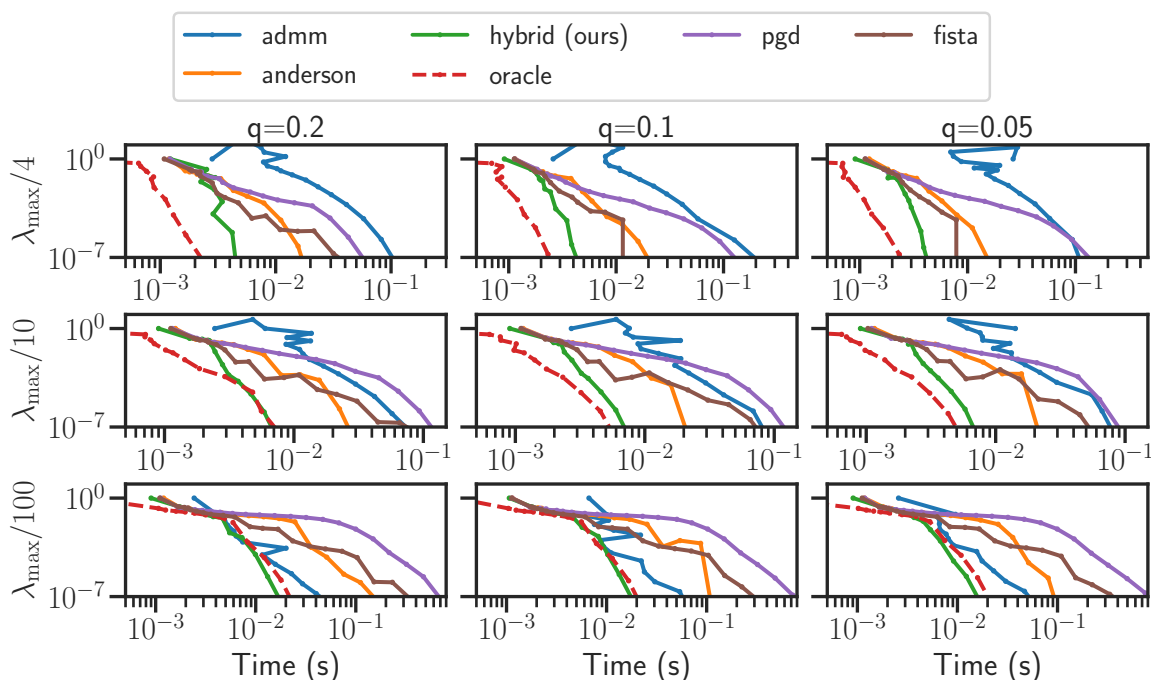


Figure 5: Benchmark on the Rhee2006 dataset.

*Applied Statistics*: 5.1 (march 2011), pages 232–253. ISSN: 1932-6157, 1941-7330. DOI: 10 / dxzfz. URL: <https://projecteuclid.org/euclid.aoas/1300715189> (urlseen 12/03/2018).

[Bog+13] Małgorzata Bogdan and others. ?Statistical Estimation and Testing via the Sorted L1 Norm? 29 october 2013. arXiv: 1310.1969 [math, stat]. URL: <http://arxiv.org/abs/1310.1969> (urlseen 16/04/2020).

[Bog+15] Małgorzata Bogdan and others. ?SLOPE – Adaptive Variable Selection via Convex Optimization? in *The annals of applied statistics*: 9.3 (september 2015), pages 1103–1140. ISSN: 1932-6157. DOI: 10.1214/15-AOS842. pmid: 26709357. URL: <https://projecteuclid.org/euclid.aoas/1446488733> (urlseen 17/12/2018).

[Bog+22] Małgorzata Bogdan and others. ?Pattern Recovery by SLOPE? 17 may 2022. DOI: 10.48550/arXiv.2203.12086. arXiv: 2203.12086 [math, stat]. URL: <http://arxiv.org/abs/2203.12086> (urlseen 03/06/2022).

[Boy+10] Stephen Boyd and others. ?Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers? in *Foundations and Trends® in*

*Machine Learning*: 3.1 (2010), pages 1–122. ISSN: 1935-8237, 1935-8245. DOI: 10.1561/22000000016. URL: <http://www.nowpublishers.com/article/Details/MAL-016> (urlseen 07/02/2020).

[BT09] A. Beck and M. Teboulle. ?A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems? in *SIAM Journal on Imaging Sciences*: 2.1 (1 january 2009), pages 183–202. DOI: 10.1137/080716542. URL: <https://epubs.siam.org/doi/abs/10.1137/080716542> (urlseen 10/02/2019).

[CW05] Patrick L Combettes and Valérie R Wajs. ?Signal recovery by proximal forward-backward splitting? in *Multiscale modeling & simulation*: 4.4 (2005), pages 1168–1200.

[DDD04] I. Daubechies, M. Defrise and C. De Mol. ?An Iterative Thresholding Algorithm for Linear Inverse Problems with a Sparsity Constraint? in *Communications on Pure and Applied Mathematics*: 57.11 (26 august 2004), pages 1413–1457. ISSN: 1097-0312. DOI: 10.1002/cpa.20042. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpa.20042> (urlseen 27/05/2022).

[DT22] Xavier Dupuis and Patrick J.C. Tardivel. ?Proximal operator for the sorted l1 norm: Application to testing procedures based on

- SLOPE? in *Journal of Statistical Planning and Inference*: 221 (2022), **pages** 1–8. ISSN: 0378-3758. DOI: <https://doi.org/10.1016/j.jspi.2022.02.005>. URL: <https://www.sciencedirect.com/science/article/pii/S0378375822000179>.
- [EH22] Clément Elvira and Cédric Herzet. ?Safe Rules for the Identification of Zeros in the Solutions of the SLOPE Problem? 18 **april** 2022. DOI: 10.48550/arXiv.2110.11784. arXiv: 2110.11784 [cs]. URL: <http://arxiv.org/abs/2110.11784> (urlseen 03/06/2022).
- [FHT10] Jerome Friedman, Trevor Hastie and Robert Tibshirani. ?Regularization Paths for Generalized Linear Models via Coordinate Descent? in *Journal of Statistical Software*: 33.1 (january 2010), **pages** 1–22. DOI: 10.18637/jss.v033.i01. URL: <http://www.jstatsoft.org/v33/i01/>.
- [FL01] Jianqing Fan and Runze Li. ?Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties? in *Journal of the American Statistical Association*: 96.456 (1 december 2001), **pages** 1348–1360. ISSN: 0162-1459. DOI: 10/fd7bfs. URL: <https://doi.org/10.1198/016214501753382273> (urlseen 14/03/2018).
- [FN16] Mario Figueiredo and Robert Nowak. ?Ordered Weighted L1 Regularized Regression with Strongly Correlated Covariates: Theoretical Aspects? in *Artificial Intelligence and Statistics: Artificial Intelligence and Statistics*. 2 may 2016, **pages** 930–938. URL: <http://proceedings.mlr.press/v51/figueiredo16.html> (urlseen 05/11/2019).
- [Har+20] Charles R Harris and others. ?Array programming with NumPy? in *Nature*: 585.7825 (2020), **pages** 357–362.
- [KB20] Michał Kos and Małgorzata Bogdan. ?On the Asymptotic Properties of SLOPE? in *Sankhya A*: 82.2 (11 august 2020), **pages** 499–532. ISSN: 0976-8378. DOI: 10.1007/s13171-020-00212-5. URL: <https://doi.org/10.1007/s13171-020-00212-5> (urlseen 03/06/2022).
- [LBW20] Johan Larsson, Małgorzata Bogdan and Jonas Wallin. ?The Strong Screening Rule for SLOPE? in *Advances in Neural Information Processing Systems* 33: Neural Information Processing Systems 2020. by editor H. Larochelle and others. volume 33. Virtual: Curran Associates, Inc., 6 december 2020–12, **pages** 14592–14603. URL: <https://proceedings.neurips.cc/paper/2020/file/a7d8ae4569120b5bec12e7b6e9648b86-Paper.pdf>.
- [LPS15] Siu Kwan Lam, Antoine Pitrou and Stanley Seibert. ?Numba: A llvm-based python jit compiler? in *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*: 2015, **pages** 1–6.
- [Luo+19] Ziyang Luo and others. ?Solving the OSCAR and SLOPE Models Using a Semismooth Newton-Based Augmented Lagrangian Method? in *Journal of Machine Learning Research*: 20.106 (2019), **pages** 1–25. URL: <http://jmlr.org/papers/v20/18-172.html>.
- [Mor+22] Thomas Moreau and others. ?Benchopt: Reproducible, efficient and collaborative optimization benchmarks? in *arXiv preprint arXiv:2206.13424*: (2022).
- [ST20] Ulrike Schneider and Patrick Tardivel. ?The Geometry of Uniqueness, Sparsity and Clustering in Penalized Estimation? 18 **august** 2020. DOI: 10.48550/arXiv.2004.09106. arXiv: 2004.09106 [math, stat]. URL: <http://arxiv.org/abs/2004.09106> (urlseen 03/06/2022).
- [Tse01] Paul Tseng. ?Convergence of a block coordinate descent method for nondifferentiable minimization? in *Journal of optimization theory and applications*: 109.3 (2001), **pages** 475–494.
- [ZF14] X. Zeng and M. Figueiredo. ?The Ordered Weighted  $\ell_1$  Norm: Atomic Formulation, Projections, and Algorithms? in *arXiv preprint arXiv:1409.4271*: (2014).
- [Zha10] Cun-Hui Zhang. ?Nearly Unbiased Variable Selection under Minimax Concave Penalty? in *The Annals of Statistics*: 38.2 (april 2010), **pages** 894–942. ISSN: 0090-5364, 2168-8966. DOI: 10/bp22zz. URL: <https://projecteuclid.org/euclid.aos/1266586618> (urlseen 14/03/2018).
- [ZOB20] Junzi Zhang, Brendan O’Donoghue and Stephen Boyd. ?Globally convergent type-I Anderson acceleration for nonsmooth fixed-point iterations? in *SIAM Journal on Optimization*: 30.4 (2020), **pages** 3170–3197.

## A Proofs