

# Supplement to "SLOPE – Adaptive Variable Selection via Convex Optimization"

Małgorzata Bogdan, Ewout van den Berg, Chiara Sabatti, Weijie Su and Emmanuel J. Candès

June 9, 2015

## A Formulation of SLOPE

We here provide proofs of some properties of the SLOPE optimization problem; we show that the SLOPE regularizing function  $J_\lambda(b)$  is a norm and that the SLOPE proximal operator reduces to the solution of the optimization problem equation (2.3).

### A.1 Proof of Proposition 1.2

**Proof** It is clear that  $J_\lambda(b) = 0$  if and only if  $b = 0$ , and that for any scalar  $t \in \mathbb{R}$ ,  $J_\lambda(tb) = |t|J_\lambda(b)$ . Thus it remains to prove that  $J_\lambda(b)$  is convex. For this, observe that by the Hardy-Littlewood-Pólya inequality

$$J_\lambda(b) = \max_{\pi} \sum_{i=1}^p \lambda_{\pi(i)} |b_i|,$$

where the maximum is over all permutations of  $p$  objects. Thus  $J_\lambda(b)$  is convex since it is a maximum over a collection of convex functions. ■

### A.2 Proof of Proposition 2.2

**Proof** It is enough to show that under Assumption 2.1, the solution  $x$  to equation (2.2) satisfies

$$x_1 \geq x_2 \geq \dots \geq x_p \geq 0. \tag{A.1}$$

Suppose that (A.1) does not hold so that there exists a pair of indices  $i < j$  such that  $x_i < x_j$  (and  $y_i > y_j$ ). Form a copy  $x'$  of  $x$  with entries  $i$  and  $j$  exchanged. Letting  $f$  be the objective functional in equation (2.2), we have

$$f(x) - f(x') = \frac{1}{2}(y_i - x_i)^2 + \frac{1}{2}(y_j - x_j)^2 - \frac{1}{2}(y_i - x_j)^2 - \frac{1}{2}(y_j - x_i)^2.$$

This follows from the fact that the sorted  $\ell_1$  norm takes on the same value at  $x$  and  $x'$  and that all the quadratic terms cancel but those for  $i$  and  $j$ . This gives

$$f(x) - f(x') = x_j y_i - x_i y_i + x_i y_j - x_j y_j = (x_j - x_i)(y_i - y_j) > 0,$$

which shows that the objective  $x'$  is strictly smaller, thereby contradicting the optimality of  $x$ . ■

## B FDR Control Under Orthogonal Designs

In this section, we prove FDR control in the orthogonal design, namely, Theorem 1.1. As we have seen in Section 1, the SLOPE solution reduces to

$$\min_{b \in \mathbb{R}^p} \frac{1}{2} \|\tilde{y} - b\|_{\ell_2}^2 + \sum_{i=1}^p \lambda_i |b|_{(i)},$$

where  $\tilde{y} = X'y \sim \mathcal{N}(\beta, I_p)$ . From this, it is clear that it suffices to consider the setting in which  $n = p$  and  $y \sim \mathcal{N}(\beta, I_p)$ , which we assume from now on.

We are thus testing the  $p$  hypotheses  $H_i : \beta_i = 0$ ,  $i = 1, \dots, p$  and set things up so that the first  $p_0$  hypotheses are null, i.e.  $\beta_i = 0$  for  $i \leq p_0$ . The SLOPE solution is

$$\hat{\beta} = \arg \min \frac{1}{2} \|y - b\|_{\ell_2}^2 + \sum_{i=1}^p \lambda_i |b|_{(i)} \quad (\text{B.1})$$

with  $\lambda_i = \Phi^{-1}(1 - iq/2p)$ . We reject  $H_i$  if and only if  $\hat{\beta}_i \neq 0$ . Letting  $V$  (resp.  $R$ ) be the number of false rejections (resp. the number of rejections) or, equivalently, the number of indices in  $\{1, \dots, p_0\}$  (resp. in  $\{1, \dots, p\}$ ) for which  $\hat{\beta}_i \neq 0$ , we have

$$\text{FDR} = \mathbb{E} \left[ \frac{V}{R \vee 1} \right] = \sum_{r=1}^p \mathbb{E} \left[ \frac{V}{r} \mathbb{1}_{\{R=r\}} \right] = \sum_{r=1}^p \frac{1}{r} \mathbb{E} \left[ \sum_{i=1}^{p_0} \mathbb{1}_{\{H_i \text{ is rejected}\}} \mathbb{1}_{\{R=r\}} \right]. \quad (\text{B.2})$$

The proof of Theorem 1.1 now follows from the two key lemmas below.

**Lemma B.1.** *Let  $H_i$  be a null hypothesis and let  $r \geq 1$ . Then*

$$\{y: H_i \text{ is rejected and } R = r\} = \{y: |y_i| > \lambda_r \text{ and } R = r\}.$$

**Lemma B.2.** *Consider applying the SLOPE procedure to  $\tilde{y} = (y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_p)$  with weights  $\tilde{\lambda} = (\lambda_2, \dots, \lambda_p)$  and let  $\tilde{R}$  be the number of rejections this procedure makes. Then with  $r \geq 1$ ,*

$$\{y: |y_i| > \lambda_r \text{ and } R = r\} \subset \{y: |y_i| > \lambda_r \text{ and } \tilde{R} = r - 1\}.$$

To see why these intermediate results give Theorem 1.1, observe that if  $H_i$  is a null, then

$$\begin{aligned} \mathbb{P}(H_i \text{ rejected and } R = r) &\leq \mathbb{P}(|y_i| \geq \lambda_r \text{ and } \tilde{R} = r - 1) \\ &= \mathbb{P}(|y_i| \geq \lambda_r) \mathbb{P}(\tilde{R} = r - 1) \\ &= \frac{q^r}{p} \mathbb{P}(\tilde{R} = r - 1), \end{aligned}$$

where the inequality is a consequence of the lemmas above and the first equality follows from the independence between  $y_i$  and  $\tilde{y}$ . Plugging this inequality into equation (B.2) gives

$$\text{FDR} = \sum_{r=1}^p \frac{1}{r} \sum_{i=1}^{p_0} \mathbb{P}(H_i \text{ rejected and } R = r) \leq \sum_{r \geq 1} \frac{qp_0}{p} \mathbb{P}(\tilde{R} = r - 1) = \frac{qp_0}{p},$$

which finishes the proof.

## B.1 Proof of Lemma B.1

We begin with a lemma we shall use more than once, and which characterizes the solution to equation (2.3).

**Lemma B.3.** *Consider a pair of nonincreasing and nonnegative sequences  $y_1 \geq y_2 \geq \dots \geq y_p \geq 0$ ,  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ , and let  $\hat{b}$  be the solution to*

$$\begin{aligned} & \text{minimize} && f(b) = \frac{1}{2} \|y - b\|_{\ell_2}^2 + \sum_{i=1}^p \lambda_i b_i \\ & \text{subject to} && b_1 \geq b_2 \geq \dots \geq b_p \geq 0. \end{aligned}$$

If  $\hat{b}_r > 0$  and  $\hat{b}_{r+1} = 0$ , then for every  $j \leq r$ , it holds that

$$\sum_{i=j}^r (y_i - \lambda_i) > 0 \tag{B.3}$$

and for every  $j \geq r + 1$ ,

$$\sum_{i=r+1}^j (y_i - \lambda_i) \leq 0. \tag{B.4}$$

**Proof** To prove equation (B.3), consider a new feasible—i.e. nonnegative and nonincreasing—sequence  $b^*$ , which differs from  $\hat{b}$  only by subtracting a small positive scalar  $h < \hat{b}_r$  from  $\hat{b}_j, \dots, \hat{b}_r$ . Now

$$f(b^*) - f(\hat{b}) = h \sum_{i=j}^r (y_i - \lambda_i - \hat{b}_i) + h^2 \sum_{i=j}^r \frac{1}{2}.$$

Taking the limit as  $h$  goes to zero, the optimality of  $\hat{b}$  implies that  $\sum_{i=j}^r (y_i - \lambda_i - \hat{b}_i) \geq 0$ , which gives

$$\sum_{i=j}^r (y_i - \lambda_i) \geq \sum_{i=j}^r \hat{b}_i > 0.$$

For the second claim equation (B.4), consider a feasible sequence  $b^*$ , which differs from  $\hat{b}$  by replacing  $\hat{b}_{r+1}, \dots, \hat{b}_j$  with a positive scalar  $0 < h < \hat{b}_r$ . Now observe that

$$f(b^*) - f(\hat{b}) = -h \sum_{i=r+1}^j (y_i - \lambda_i) + h^2 \sum_{i=r+1}^j \frac{1}{2}.$$

The claim follows from the optimality of  $\hat{b}$ . ■

It is now straightforward so see how these simple relationships give Lemma B.1. Observe that when  $R = r$ , we must have  $|y|_{(r)} > \lambda_r$  and  $|y|_{(r+1)} \leq \lambda_{r+1}$ . Hence, if  $H_1$  is rejected, it must hold that  $|y_1| \geq |y|_{(r)} > \lambda_r$ . This shows that  $\{H_1 \text{ is rejected and } R = r\} \subset \{|y_1| > \lambda_r \text{ and } R = r\}$ . Conversely, assume that  $|y_1| > \lambda_r$  and  $R = r$ . Then  $H_1$  must be rejected since  $|y_1| > |y|_{(r+1)}$ . This shows that  $\{H_1 \text{ is rejected and } R = r\} \supset \{|y_1| > \lambda_r \text{ and } R = r\}$ .

## B.2 Proof of Lemma B.2

We assume without loss of generality that  $y \geq 0$  (see Section 2.2). By assumption the solution to equation (B.1) with  $\lambda_i = \Phi^{-1}(1 - iq/2p)$  has exactly  $r$  strictly positive entries, and we need to show that when  $y_1$  is rejected, the solution to

$$\min J(\tilde{b}) := \sum_{i=1}^{p-1} \frac{1}{2}(\tilde{y}_i - \tilde{b}_i)^2 + \sum_{i=1}^{p-1} \tilde{\lambda}_i |\tilde{b}|_{(i)} \quad (\text{B.5})$$

in which  $\tilde{\lambda}_i = \lambda_{i+1}$  has exactly  $r - 1$  nonzero entries. We prove this in two steps:

- (i) The optimal solution  $\hat{b}$  to equation (B.5) has at least  $r - 1$  nonzero entries.
- (ii) The optimal solution  $\hat{b}$  to equation (B.5) has at most  $r - 1$  nonzero entries.

### B.2.1 Proof of (i)

Suppose by contradiction that  $\hat{b}$  has fewer than  $r - 1$  entries; i.e.,  $\hat{b}$  has  $j - 1$  nonzero entries with  $j < r$ . Letting  $I$  be those indices for which the rank of  $\tilde{y}_i$  is between  $j$  and  $r - 1$ , consider a feasible point  $b$  as in the proof of Lemma B.3 defined as

$$b_i = \begin{cases} h & i \in I, \\ \hat{b}_i & \text{otherwise;} \end{cases}$$

here, the positive scalar  $h$  obeys  $0 < h < b_{(j-1)}$ . By definition,

$$J(b) - J(\hat{b}) = -h \sum_{i=j}^{r-1} (\tilde{y}_{(i)} - \tilde{\lambda}_i) + h^2 \sum_{i=j}^{r-1} \frac{1}{2}.$$

Now

$$\sum_{j \leq i \leq r-1} \tilde{y}_{(i)} - \tilde{\lambda}_i = \sum_{j+1 \leq i \leq r} \tilde{y}_{(i-1)} - \lambda_i \geq \sum_{j+1 \leq i \leq r} y_{(i)} - \lambda_i > 0.$$

The first equality follows from  $\tilde{\lambda}_i = \lambda_{i+1}$ , the first inequality from  $y_{(i)} \leq \tilde{y}_{(i-1)}$  and the last from equation (B.3). By selecting  $h$  small enough, this gives  $J(b) < J(\hat{b})$ , which contradicts the optimality of  $\hat{b}$ .

### B.2.2 Proof of (ii)

The proof is similar to that of (i). Suppose by contradiction that  $\hat{b}$  has more than  $r - 1$  entries; i.e.  $\hat{b}$  has  $j$  nonzero entries with  $j \geq r$ . Letting  $I$  be those indices for which the rank of  $\tilde{y}_i$  is between  $r$  and  $j$ , consider a feasible point  $b$  as in the proof of Lemma B.3 defined as

$$b_i = \begin{cases} \hat{b}_i - h & i \in I \\ \hat{b}_i & \text{otherwise;} \end{cases}$$

here, the positive scalar  $h$  obeys  $0 < h < b_{(j)}$ . By definition,

$$J(b) - J(\hat{b}) = h \sum_{i=r}^j (\tilde{y}_{(i)} - \tilde{\lambda}_i - \hat{b}_{(i)}) + h^2 \sum_{i=r}^j \frac{1}{2}.$$

Now

$$\sum_{r \leq i \leq j} (\tilde{y}_{(i)} - \tilde{\lambda}_i) = \sum_{r+1 \leq i \leq j+1} (y_{(i)} - \lambda_i) \leq 0.$$

The equality follows from the definition and the inequality from equation (B.4). By selecting  $h$  small enough, this gives  $J(b) < J(\hat{b})$ , which contradicts the optimality of  $\hat{b}$ .

## C Algorithmic Issues

### C.1 Duality-based stopping criteria

To derive the dual of equation (1.10) we first rewrite it as

$$\underset{b,r}{\text{minimize}} \quad \frac{1}{2}r'r + J_\lambda(b) \quad \text{subject to} \quad Xb + r = y.$$

The dual is then given by

$$\underset{w}{\text{maximize}} \quad \mathcal{L}(w),$$

where

$$\begin{aligned} \mathcal{L}(w) &:= \inf_{b,r} \left\{ \frac{1}{2}r'r + J_\lambda(b) - w'(Xb + r - y) \right\} \\ &= w'y - \sup_r \left\{ w'r - \frac{1}{2}r'r \right\} - \sup_b \left\{ (X'w)'b - J_\lambda(b) \right\}. \end{aligned}$$

The first supremum term evaluates to  $\frac{1}{2}w'w$  by choosing  $r = w$ . The second term is the conjugate function  $J^*$  of  $J$  evaluated at  $v = X'w$ , which can be shown to reduce to

$$J_\lambda^*(v) := \sup_b \{v'b - J_\lambda(b)\} = \begin{cases} 0 & v \in C_\lambda, \\ +\infty & \text{otherwise,} \end{cases}$$

where the set  $C_\lambda$  is the unit ball of the dual norm to  $J_\lambda(\cdot)$ . In details,

$$w \in C_\lambda \iff \sum_{j \leq i} |w|_{(j)} \leq \sum_{j \leq i} \lambda_j \text{ for all } i = 1, \dots, p.$$

The dual problem is thus given by

$$\underset{w}{\text{maximize}} \quad w'y - \frac{1}{2}w'w \quad \text{subject to} \quad X'w \in C_\lambda.$$

The dual formulation can be used to derive appropriate stopping criteria. At the solution we have  $w = r$ , which motivates estimating a dual point by setting  $\hat{w} = r =: y - Xb$ . At this point the primal-dual gap at  $b$  is the difference between the primal and dual objective:

$$\delta(b) = (Xb)'(Xb - y) + J_\lambda(b).$$

However,  $\hat{w}$  is not guaranteed to be feasible, i.e., we may not have  $\hat{w} \in C_\lambda$ . Therefore we also need to compute a level of infeasibility of  $\hat{w}$ , for example

$$\text{infeasi}(\hat{w}) = \max \left\{ 0, \max_i \sum_{j \leq i} (|\hat{w}|_{(j)} - \lambda_j) \right\}.$$

The algorithm used in the numerical experiments terminates whenever both the infeasibility and primal-dual gap are sufficiently small. In addition, it imposes a limit on the total number of iterations to ensure termination.

## C.2 Proof of Lemma 2.3

It is useful to think of the prox as the solution to the quadratic program equation (2.3) and we begin by recording the Karush-Kuhn-Tucker (KKT) optimality conditions for this QP.

*Primal feasibility:*  $x_1 \geq x_2 \geq \dots \geq x_n \geq 0$ .

*Dual feasibility:*  $\mu = (\mu_1, \dots, \mu_n)$  obeys  $\mu \geq 0$ .

*Complementary slackness:*  $\mu_i(x_i - x_{i+1}) = 0$  for all  $i = 1, \dots, n$  (with the convention  $x_{n+1} = 0$ ).

*Stationarity of the Lagrangian:* with the convention that  $\mu_0 = 0$ ,

$$x_i - y_i + \lambda_i - (\mu_i - \mu_{i-1}) = 0.$$

We now turn to the proof of the second claim of the lemma. Set  $x = (y - \lambda)_+$ , which by assumption is primal feasible, and let  $i_0$  be the last index such that  $y_i - \lambda_i > 0$ . Set  $\mu_1 = \mu_2 = \dots = \mu_{i_0} = 0$  and for  $j > i_0$ , recursively define

$$\mu_j = \mu_{j-1} - (y_j - \lambda_j) \geq 0.$$

Then it is straightforward to check that the pair  $(x, \mu)$  obeys the KKT optimality conditions. Hence  $x$  is solution.

Consider now the first claim. We first argue that the prox has to be constant over any monotone segment of the form

$$y_i - \lambda_i \leq y_{i+1} - \lambda_{i+1} \leq \dots \leq y_j - \lambda_j.$$

To see why this is true, set  $x = (y; \lambda)$  and suppose the contrary: then over a segment as above, there is  $k \in \{i, i+1, \dots, j-1\}$  such that  $x_k > x_{k+1}$  (we cannot have a strict inequality in the other direction since  $x$  has to be primal feasible). By complementary slackness,  $\mu_k = 0$ . This gives

$$\begin{aligned} x_k &= y_k - \lambda_k - \mu_{k-1} \\ x_{k+1} &= y_{k+1} - \lambda_{k+1} + \mu_{k+1}. \end{aligned}$$

Since  $y_{k+1} - \lambda_{k+1} \geq y_k - \lambda_k$  and  $\mu \geq 0$ , we have  $x_k \leq x_{k+1}$ , which is a contradiction.

Now an update replaces an increasing segment as in equation (2.4) with a constant segment and we have just seen that both proxies must be constant over such segments. Now consider the cost function associated with the prox with parameter  $\lambda$  and input  $y$  over an increasing segment as in equation (2.4),

$$\sum_{i \leq k \leq j} \left\{ \frac{1}{2}(y_k - x_k)^2 + \lambda_k x_k \right\}. \tag{C.1}$$

Since all the variables  $x_k$  must be equal to some value  $z$  over this block, this cost is equal to

$$\begin{aligned} \sum_{i \leq k \leq j} \left\{ \frac{1}{2}(y_k - z)^2 + \lambda_k z \right\} &= \sum_k \frac{1}{2}(y_k - \bar{y})^2 + \sum_{i \leq k \leq j} \left\{ \frac{1}{2}(\bar{y} - z)^2 + \bar{\lambda} z \right\} \\ &= \sum_k \frac{1}{2}(y_k - \bar{y})^2 + \sum_{i \leq k \leq j} \left\{ \frac{1}{2}(y_k^+ - z)^2 + \bar{\lambda}_k^+ z \right\}, \end{aligned}$$

where  $\bar{y}$  and  $\bar{\lambda}$  are block averages. The second term in the right-hand side is the cost function associated with the prox with parameter  $\lambda^+$  and input  $y^+$  over the same segment since all the variables over this segment must also take on the same value. Therefore, it follows that replacing each appearance of block sums as in equation (C.1) in the cost function yields the same minimizer. This proves the claim.