

# The Spike-and-Slab LASSO

By V. Ročková and E.I. George

*University of Chicago and University of Pennsylvania*

## Supplemental Material

### A Proofs of Section 6

#### A.1 Proofs of Section 6.2.1

Throughout this section, we denote by  $Q(\beta) = -\frac{1}{2}\|\mathbf{Y} - \mathbf{X}\beta\|^2 + \text{pen}_S(\beta | \theta)$  the log-posterior under the separable *SSL* penalty.

##### A.1.1 Proof of Theorem 6.1

*Proof.* Denote by  $\Theta = \hat{\beta} - \beta_0$ . Because  $0 \geq Q(\beta^0) - Q(\hat{\beta})$ , we can write

$$0 \geq \|\mathbf{X}\Theta\|^2 - 2\epsilon'\mathbf{X}\Theta + 2\log \left[ \frac{\pi(\beta_0 | \theta)}{\pi(\hat{\beta} | \theta)} \right]. \quad (\text{A.1})$$

Using the fact  $p_\theta^*(\hat{\beta}_j) > c_+$  when  $\hat{\beta}_j \neq 0$ , we can write

$$\begin{aligned} \log \left[ \frac{\pi(\beta_0 | \theta)}{\pi(\hat{\beta} | \theta)} \right] &\geq -\lambda_1 |\beta_0 - \hat{\beta}_0| + \sum_{j=1}^p \log \left[ \frac{p_\theta^*(\hat{\beta}_j)}{p_\theta^*(0)} \right] + \sum_{j=1}^p \log \left[ \frac{p_\theta^*(0)}{p_\theta^*(\beta_{0j})} \right] \\ &\geq -\lambda_1 |\beta_0 - \hat{\beta}_0| + \hat{q}b + (\hat{q} - q) \log[1/p_\theta^*(0)], \end{aligned}$$

where  $0 > b = \log c_+ > \log 0.5$  is a constant very close to 0. Because  $\|\mathbf{X}'\epsilon\|_\infty \leq \eta\Delta$  under  $\eta$ -NC condition, we can use the Hölder inequality  $|\alpha'\beta| \leq |\alpha|_\infty |\beta|$  to find that

$$0 \geq \|\mathbf{X}\Theta\|^2 - 2(\eta\Delta + \lambda_1)|\Theta| + 2\hat{q}b + 2(\hat{q} - q) \log[1/p_\theta^*(0)]. \quad (\text{A.2})$$

From Lemma 1 we know that  $\Theta$  lives inside the cone  $C(\eta; \beta_0)$ . Thus, we can use Definition 1 to find that  $\|\mathbf{X}\Theta\|^2 \geq c(\eta; \beta_0)^2 \|\Theta\|^2 \|\mathbf{X}\|^2$ . Denote by  $c = c(\eta; \beta_0)$ . Using the fact  $|\Theta| \leq \|\Theta\| \|\Theta\|_0^{1/2}$ , we have

$$0 \geq c^2 \|\Theta\|^2 \|\mathbf{X}\|^2 - 2(\eta\Delta + \lambda_1) \|\Theta\| \|\Theta\|_0^{1/2} + 2\hat{q}b + 2(\hat{q} - q) \log[1/p_\theta^*(0)],$$

which is equivalent to writing

$$\left[ c \|\Theta\| \|\mathbf{X}\| - \frac{(\eta\Delta + \lambda_1)}{c \|\mathbf{X}\|} \|\Theta\|_0^{1/2} \right]^2 - \frac{(\eta\Delta + \lambda_1)^2}{c^2 \|\mathbf{X}\|^2} \|\Theta\|_0 + 2\hat{q} + 2(\hat{q} - q) \log[1/p_\theta^*(0)] \leq 0.$$

This yields

$$(\hat{q} - q) \log[1/p_\theta^*(0)] + \hat{q}b \leq \frac{(\eta\Delta + \lambda_1)^2}{2c^2\|\mathbf{X}\|^2} \|\boldsymbol{\Theta}\|_0.$$

By noting  $\|\boldsymbol{\Theta}\|_0 \leq \hat{q} + q$  and  $\|\mathbf{X}\|^2 = n$ , we can write

$$\hat{q} \leq q \left( 1 + \frac{2A - b}{B + b - A} \right),$$

where  $A = \frac{(\eta\Delta + \lambda_1)^2}{2c^2\|\mathbf{X}\|^2}$  and  $B = \log[1/p_\theta^*(0)]$ . Assume, for simplicity, that  $(1 - \theta)/\theta = C_1 p^a$ ,  $\lambda_0 = C_2 p^d$  and  $\lambda_1 < 4\sqrt{n \log p} < 4p$  we have  $B = \log \left( 1 + \frac{1-\theta}{\theta} \frac{\lambda_0}{\lambda_1} \right) > \log(C_1 C_2/4) + (a + d - 1) \log p$ . With  $C_1 C_2/4 > 0$  we obtain  $\lambda_1 < 4\sqrt{n \log p} < 4\sqrt{n B/(a + d - 1)}$  and  $\frac{A}{B} < \left( \frac{\eta}{c} + \frac{(\eta+1)2\sqrt{2}}{c\sqrt{a+d-1}} \right)^2 \equiv D$ . We can then write  $\hat{q} \leq q \left( 1 + M \frac{D}{1-D} \right)$ .  $\square$

### A.1.2 Proof of Theorem 6.2

*Proof.* With  $\boldsymbol{\Theta} = \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0$  and by noting  $\log \left[ \frac{\pi(\boldsymbol{\beta}_0|\theta)}{\pi(\boldsymbol{\beta}|\theta)} \right] > -\lambda_1 |\boldsymbol{\Theta}| + q \log p_\theta^*(0)$  we can write

$$0 \geq \|\mathbf{X}\boldsymbol{\Theta}\|^2 - 2(\eta\Delta + \lambda_1)|\boldsymbol{\Theta}| + 2q \log[p_\theta^*(0)], \quad (\text{A.3})$$

where  $\Delta$  is the selection threshold. From Theorem 6.1 we have  $\|\boldsymbol{\Theta}\|_0 \leq (K + 1)q$  under the  $\eta$ -NC condition. Denote by  $\phi = \phi[C(\eta; \boldsymbol{\beta}_0)]$ . From the definition of the compatibility number  $\phi$ , and using  $4uv \leq u^2 + 4v^2$ , we find that

$$\begin{aligned} 2(\eta\Delta + \lambda_1)|\boldsymbol{\Theta}| &\leq 3(\eta\Delta + \lambda_1) \frac{\|\mathbf{X}\boldsymbol{\Theta}\| \sqrt{(K+1)q}}{\|\mathbf{X}\|\phi} - (\eta\Delta + \lambda_1)|\boldsymbol{\Theta}| \\ &\leq \frac{\|\mathbf{X}\boldsymbol{\Theta}\|^2}{2} + \frac{5(K+1)q(\eta\Delta + \lambda_1)^2}{\|\mathbf{X}\|^2 \phi^2} - (\eta\Delta + \lambda_1)|\boldsymbol{\Theta}|. \end{aligned}$$

Thus it follows from (A.3) that

$$\frac{1}{2} \|\mathbf{X}\boldsymbol{\Theta}\|^2 + (\eta\Delta + \lambda_1)|\boldsymbol{\Theta}| \leq \frac{5(K+1)q(\eta\Delta + \lambda_1)^2}{\|\mathbf{X}\|^2 \phi^2} + 2q \log[1/p_\theta^*(0)]. \quad (\text{A.4})$$

With  $(1 - \theta)/\theta = C_1 p^a$ ,  $\lambda_0 = C_2 p^d$  and  $\sqrt{n}/p < \lambda_1 < 4\sqrt{n \log p}$  we have  $(\eta\Delta + \lambda_1) < C_3 \eta \sqrt{n \log p}$  and  $\log[1/p_\theta^*(0)] < C_4 \log p$ . With  $\|\mathbf{X}\|^2 = n$ , the first two statements of the theorem follow directly from (A.4). Let  $c = c(\eta, \boldsymbol{\beta}_0)$  be the minimal restricted eigenvalue. Then the last statement is obtained from  $\|\mathbf{X}\boldsymbol{\Theta}\| > c\|\mathbf{X}\|\|\boldsymbol{\Theta}\|$ .  $\square$

## A.2 Proofs of Section 6.2.2

The construction of the proof follows Castillo et al. (2015), where suitable modifications are required when using the notion of generalized dimensionality. Before proceeding, we

need to introduce some more notation. Let

$$\Lambda_{n,\beta,\beta_0} = e^{-\frac{1}{2}\|\mathbf{X}(\beta-\beta_0)\|^2 + (\mathbf{y}-\mathbf{X}\beta_0)'\mathbf{X}(\beta-\beta_0)}$$

and

$$\Pi(\beta \mid \theta) = \prod_{i=1}^p [\theta \psi_1(\beta_i) + (1-\theta) \psi_0(\beta_i)].$$

Throughout this section we will denote by  $\bar{\lambda} = 2\sqrt{n \log p}$  the universal threshold. The rates in this section will be expressed in terms of slightly different compatibility and minimal eigenvalue numbers. Following Castillo (2015), for  $S \subset \{1, \dots, p\}$ , we define: the compatibility number  $\tilde{\phi}(S)$  of a model  $S$  by

$$\tilde{\phi}(S) = \inf \left\{ \frac{\|\mathbf{X}\beta\| |S|^{1/2}}{\|\mathbf{X}\| |\beta_S|} : |\beta_{S^c}| \leq 5|\beta_S|, \beta_S \neq 0 \right\}, \quad (\text{A.5})$$

the compatibility in vectors of generalized dimension  $s$  by

$$\bar{\phi}(s) = \inf \left\{ \frac{\|\mathbf{X}\beta\| s^{1/2}}{\|\mathbf{X}\| |\beta|} : 0 < |\gamma(\beta)| \leq s \right\} \quad (\text{A.6})$$

and the minimal eigenvalue restricted to vectors  $\beta$  of generalized dimensionality at most  $s$  by

$$\bar{c}(s) = \inf \left\{ \frac{\|\mathbf{X}\beta\|}{\|\mathbf{X}\| |\beta|} : 0 < |\gamma(\beta)| \leq s \right\} \quad (\text{A.7})$$

For  $S \subset \{1, \dots, p\}$ , let  $\beta_S \in \mathbb{R}^p$  be a subset of  $\beta$  with coordinates in  $S$ . Denote by  $\Pi_S(\beta \mid \theta)$  the marginal prior confined to coordinates in  $S$ . Denote by  $\delta$  the intersection point between  $SSL$  densities, by  $\pi \equiv P(|\beta_1| \leq \delta)$  and by  $\pi(s \mid \theta) = \binom{p}{s} \pi^s (1-\pi)^{p-s} = P[|\gamma(\beta)| = s \mid \theta]$  the prior distribution on the effective dimensionality. By assumption, we have  $\|\mathbf{X}\|^2 = \max_{1 \leq i \leq p} \|\mathbf{X}_i\|^2 = n$ .

We will need the following analogue of Lemma 2 of Castillo et al. (2015) for the separable  $SSL$  prior.

**Lemma A.1.** *Assume  $\beta_0 \in \mathbb{R}^p$  has a support  $S_0 \subset \{1, \dots, p\}$ , where  $|S_0| = q$ . Assume  $\lambda_0 = (1-\theta)/\theta = Cp^a$ , where  $a \geq 2$  and  $C > 0$ , and  $\sqrt{n}/p < \lambda_1 \leq 4\sqrt{n \log p}$ . Assume  $p > n$ . Then*

$$\int \Lambda_{n,\beta,\beta_0} \Pi(\beta \mid \theta) d\beta \geq \frac{\pi(q \mid \theta)}{p^{2q}} e^{-1-D-\lambda_1|\beta|_1},$$

where  $D > 0$ .

*Proof.* Denote by  $g(\beta) = e^{-\|\mathbf{X}\beta\|^2 + (\mathbf{y} - \mathbf{X}\beta_0)' \mathbf{X}\beta}$ . Using the fact  $\|\mathbf{X}\beta\|^2 \leq 2\|\mathbf{X}\beta_{S_0}\|^2 + 2\|\mathbf{X}\beta_{S_0^c}\|^2$ , we can write

$$\Lambda_{n,\beta,\beta_0} > g(\beta_{S_0^c})g(\beta_{S_0} - \beta_{0S_0}).$$

By the Jensen's inequality we have

$$\int g(\beta)\Pi(\beta|\theta)d\beta \geq \int e^{-\|\mathbf{X}\beta\|^2}\Pi(\beta|\theta)d\beta.$$

Conditionally on  $\theta$ , the SSL prior is separable, implying  $\Pi(\beta|\theta) = \Pi_{S_0}(\beta|\theta)\Pi_{S_0^c}(\beta|\theta)$ . Changing variables  $\mathbf{b} \rightarrow (\beta - \beta_0)$  and noting  $\Pi_{S_0}(\beta|\theta) > \theta^q \left(\frac{\lambda_1}{2}\right)^q e^{-\lambda_1|\beta_{S_0}|}$ , we can write

$$\int \Lambda_{n,\beta,\beta_0}\Pi(\beta|\theta)d\beta > \int e^{-\|\mathbf{X}\beta_{S_0^c}\|^2}\Pi_{S_0^c}(\beta|\theta)d\beta_{S_0^c} \quad (\text{A.8})$$

$$\times \theta^q e^{-\lambda_1|\beta_0|} \int e^{-\|\mathbf{X}\mathbf{b}_{S_0}\|^2} \left(\frac{\lambda_1}{2}\right)^q e^{-\lambda_1|\mathbf{b}_{S_0}|} d\mathbf{b}_{S_0}. \quad (\text{A.9})$$

To simplify the integral in (A.9) we use arguments of Castillo (2015) in the proof of Lemma

2. Under the assumption  $\|\mathbf{X}\|/p < \lambda_1 < 4\|\mathbf{X}\|\sqrt{\log p}$ , we obtain

$$\int e^{-\|\mathbf{X}\mathbf{b}_{S_0}\|^2} \left(\frac{\lambda_1}{2}\right)^q e^{-\lambda_1|\mathbf{b}_{S_0}|} d\mathbf{b}_{S_0} > e^{-1} \left(\frac{\lambda_1}{\|\mathbf{X}\|}\right)^q \frac{e^{-\lambda_1/\|\mathbf{X}\|}}{q!} > \frac{e^{-1}}{p^q q!} \quad (\text{A.10})$$

To simplify the integral in (A.8), we use  $\|\mathbf{X}\beta\| \leq \|\mathbf{X}\|\|\beta\|$  to find that

$$\int e^{-\|\mathbf{X}\beta_{S_0^c}\|^2}\Pi_{S_0^c}(\beta|\theta)d\beta_{S_0^c} > \int_{|\beta_i| \leq \delta; i \notin S_0} e^{-\|\mathbf{X}\|^2\|\beta_{S_0^c}\|^2}\Pi_{S_0^c}(\beta|\theta)d\beta_{S_0^c} \quad (\text{A.11})$$

$$\geq e^{-(p-q)^2\|\mathbf{X}\|^2\delta^2} \mathbf{P}(|\beta_1| \leq \delta)^{p-q}. \quad (\text{A.12})$$

Combining (A.10) and (A.12) and noting  $\theta > \pi \equiv \mathbf{P}(|\beta_1| > \delta|\theta)$ , we obtain

$$\int \Lambda_{n,\beta,\beta_0}\Pi(\beta|\theta)d\beta > e^{-(p-q)^2\|\mathbf{X}\|^2\delta^2} e^{-1-\lambda_1|\beta_0|} \pi^q (1-\pi)^{p-q} \frac{1}{p^q q!}.$$

Recall that  $\pi(q|\theta) = \binom{p}{q} \pi^q (1-\pi)^{p-q}$  is the prior probability of the effective dimensionality  $q$ . Since  $\binom{p}{q} q! \leq p^q$ , we can write

$$\int \Lambda_{n,\beta,\beta_0}\Pi(\beta|\theta)d\beta > e^{-(p-q)^2\|\mathbf{X}\|^2\delta^2} e^{-1-\lambda_1|\beta_0|} \frac{\pi(q|\theta)}{p^{2q}}.$$

Using the fact  $\delta = \frac{1}{\lambda_0 - \lambda_1} \log[1/p^*(0) - 1]$ , we obtain

$$(p-q)^2\|\mathbf{X}\|^2\delta^2 = \frac{(p-q)^2\|\mathbf{X}\|^2}{(\lambda_0 - \lambda_1)^2} \log^2[1/p^*(0) - 1] \quad (\text{A.13})$$

Since  $\|\mathbf{X}\| = \sqrt{n} < \sqrt{p}$  we have  $\lambda_1 \leq 4\|\mathbf{X}\|\sqrt{\log p} < 4p$ . Because  $\lambda_0 \geq p^d$  with  $d \geq 2$ , we have  $\frac{(p-q)^2}{\lambda_0 - \lambda_1} \preceq \frac{1}{p^{d-2}}$ . Because  $(1-\theta)/\theta \sim p^a$ , we have

$$\log^2[1/p^*(0) - 1] \preceq \log^2 p + \log^2 \lambda_0.$$

Therefore, with  $p > n$  and  $\lambda_0 \succeq p^d$  with  $d \geq 2$  we obtain

$$(p - q)^2 \|\mathbf{X}\|^2 \delta^2 \preceq \frac{n (\log^2 p + \log^2 \lambda_0)}{p^{d-2} \lambda_0} \rightarrow 0. \quad \square$$

### A.2.1 Proof of Theorem 6.3

*Proof.* Denote by  $\mathcal{B} = \{\boldsymbol{\beta} : |\gamma(\boldsymbol{\beta})| > R\}$ . Then  $\mathbb{E}_{\boldsymbol{\beta}_0} \mathbf{P}(\mathcal{B} | \mathbf{Y}, \theta) \leq \mathbb{E}_{\boldsymbol{\beta}_0} \mathbf{P}(\mathcal{B} | \mathbf{Y}, \theta) \mathbb{I}_{\tau_0} + \frac{2}{p}$ , where  $\tau_0 = \{ \|\mathbf{X}'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}_0)\|_\infty \leq \bar{\lambda} \}$  and  $\bar{\lambda} = 2\sqrt{n \log p}$ . Then

$$\mathbf{P}(\mathcal{B} | \mathbf{Y}, \theta) = \frac{\int_{\mathcal{B}} \Lambda_{n,\boldsymbol{\beta},\boldsymbol{\beta}_0} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}}{\int \Lambda_{n,\boldsymbol{\beta},\boldsymbol{\beta}_0} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}} \leq A e^{\lambda_1 |\boldsymbol{\beta}_0|} \int_{\mathcal{B}} e^{-\frac{1}{2} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 + (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}_0)' \mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}, \quad (\text{A.14})$$

where  $A = \frac{p^{2q}}{\pi(q|\theta)} e^{1+D}$ . Similarly as in the proof of Theorem 12 of Castillo et al. (2015), we use Hölder's inequality to obtain on  $\tau_0$

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}_0)' \mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0) \leq \bar{\lambda} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|. \quad (\text{A.15})$$

Therefore, the expectation under  $\boldsymbol{\beta}_0$  of the integrand satisfies

$$e^{-\frac{1}{2} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2} \mathbb{E}_{\boldsymbol{\beta}_0} \left[ e^{\left(1 - \frac{\lambda_1}{2\bar{\lambda}}\right) (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}_0)' \mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)} \mathbb{I}_{\tau_0} \right] e^{\frac{\lambda_1}{2} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|} \quad (\text{A.16})$$

$$\leq e^{-\frac{1}{2} \left[ 1 - \left(1 - \frac{\lambda_1}{2\bar{\lambda}}\right)^2 \right] \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2} e^{\frac{\lambda_1}{2} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|} \quad (\text{A.17})$$

$$\leq e^{-\frac{\lambda_1}{4\bar{\lambda}} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2} e^{\frac{\lambda_1}{2} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|}, \quad (\text{A.18})$$

where we used  $\lambda_1 \leq 2\bar{\lambda}$  and invoked the expectation of a log-normally distributed r.v. Thus,

$$\mathbb{E}_{\boldsymbol{\beta}_0} \mathbf{P}(\mathcal{B} | \mathbf{Y}, \theta) \mathbb{I}_{\tau_0} \leq A e^{\lambda_1 |\boldsymbol{\beta}_0|} \int_{\mathcal{B}} e^{-\frac{\lambda_1}{4\bar{\lambda}} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2} e^{\frac{\lambda_1}{2} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|} d\Pi(\boldsymbol{\beta} | \theta). \quad (\text{A.19})$$

Now, when  $5|\boldsymbol{\beta}_{S_0} - \boldsymbol{\beta}_0| \leq |\boldsymbol{\beta}_{S_0^c}|$ , then

$$|\boldsymbol{\beta}_0| + \frac{1}{2} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| \leq |\boldsymbol{\beta}_{S_0}| + \frac{5}{4} |\boldsymbol{\beta}_S - \boldsymbol{\beta}_0| + \frac{3}{4} |\boldsymbol{\beta}_{S_0^c}| - \frac{1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| < -\frac{1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| + |\boldsymbol{\beta}| \quad (\text{A.20})$$

$$< -\frac{1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| + |\boldsymbol{\beta}| + \frac{1}{4\bar{\lambda}} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 + 2 \frac{\bar{\lambda} |S_0|}{\|\mathbf{X}\|^2 \tilde{\phi}(S_0)^2}. \quad (\text{A.21})$$

When  $5|\boldsymbol{\beta}_{S_0} - \boldsymbol{\beta}_0| > |\boldsymbol{\beta}_{S_0^c}|$ , we use the definition of the compatibility number to find that

$$|\boldsymbol{\beta}_{S_0}| + \frac{5}{4} |\boldsymbol{\beta}_{S_0} - \boldsymbol{\beta}_0| + \frac{3}{4} |\boldsymbol{\beta}_{S_0^c}| - \frac{1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| < -\frac{1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| + |\boldsymbol{\beta}| + \frac{5 \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\| |S_0|^{1/2}}{4 \|\mathbf{X}\| \tilde{\phi}(S_0)}.$$

Invoking the inequality  $4uv \leq u^2 + 4v^2$ , we can bound the last display from above by

$$-\frac{1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0| + |\boldsymbol{\beta}| + \frac{1}{4\bar{\lambda}} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 + 2 \frac{\bar{\lambda} q}{\|\mathbf{X}\|^2 \tilde{\phi}(S_0)^2}.$$

Thus (A.19) can be bounded by

$$A e^{\frac{2\lambda_1 \bar{\lambda} q}{\|\mathbf{X}\|^2 \phi(S_0)^2}} \int_{\mathcal{B}} e^{\lambda_1 |\boldsymbol{\beta}| - \frac{\lambda_1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}.$$

Note that  $\pi(\boldsymbol{\beta} | \theta) < \theta \lambda_1 e^{-\lambda_1 |\boldsymbol{\beta}|}$  when  $|\boldsymbol{\beta}| > \delta$ . For  $\mathcal{B} = \{\boldsymbol{\beta} : |\boldsymbol{\gamma}(\boldsymbol{\beta})| > R\}$ , we can write

$$\int_{\mathcal{B}} e^{\lambda_1 |\boldsymbol{\beta}| - \frac{\lambda_1}{4} |\boldsymbol{\beta} - \boldsymbol{\beta}_0|} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta} \leq \sum_{S: |S| > R} \theta^{|S|} \lambda_1^{|S|} \int_{|\boldsymbol{\beta}_i| > \delta; i \in S} e^{-\frac{\lambda_1}{4} |\boldsymbol{\beta}_S - \boldsymbol{\beta}_{0S}|} d\boldsymbol{\beta}_S \quad (\text{A.22})$$

$$\times \int_{|\boldsymbol{\beta}_i| \leq \delta; i \in S^c} e^{\lambda_1 |\boldsymbol{\beta}_{S^c}| - \frac{\lambda_1}{4} |\boldsymbol{\beta}_{S^c} - \boldsymbol{\beta}_{0S^c}|} \Pi_{S^c}(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}_{S^c} \quad (\text{A.23})$$

$$< \sum_{k=R+1}^p \binom{p}{k} (8\theta)^k e^{\lambda_1 \delta(p-k)} (1-\pi)^{p-k}. \quad (\text{A.24})$$

Recall that  $\pi \equiv P(|\beta_1| > \delta | \theta) = \theta e^{-\delta \lambda_1} \left(1 + \frac{\lambda_1}{\lambda_0}\right) < \theta$ . Because  $\theta < \pi e^{\delta \lambda_1}$ , we can bound the last display by

$$e^{\lambda_1 p \delta} \sum_{k=R+1}^p 8^k \binom{p}{k} \pi^k (1-\pi)^{p-k} = e^{\lambda_1 p \delta} \sum_{k=R+1}^p 8^k \pi(k | \theta).$$

Because  $\pi/(1-\pi) < \theta/(1-\theta) \leq 1/p^a$  for  $a \geq 2$ , we have

$$\pi(k | \theta) \leq \left(\frac{1}{p}\right)^{a-1} \pi(k-1 | \theta) \quad \text{for } k \geq 1.$$

Thereby, we can write for  $R > q$

$$e^{\lambda_1 p \delta} \sum_{k=R+1}^p 8^k \pi(k | \theta) < e^{\lambda_1 p \delta} 8^q \pi(q | \theta) \left(\frac{8}{p^{a-1}}\right)^{R+1-q} \sum_{k=0}^{\infty} \left(\frac{8}{p^{a-1}}\right)^k.$$

With  $(1-\theta)/\theta \sim p^a$  and  $\lambda_0 \geq p^d$  and  $\|\mathbf{X}\|^2 = n$ , we have  $\lambda_1 \delta p \leq 1$ . Altogether

$$\begin{aligned} P(\mathcal{B} | \mathbf{Y}, \theta) &\leq e^{2q \log p + \lambda_1 p \delta + \frac{2\lambda_1 \bar{\lambda} q}{n \phi(S_0)^2}} \left(\frac{8}{p^{a-1}}\right)^{R+1-q} + \frac{2}{p} \\ &\leq e^{(R+1-q) \log 8 + 2q \log p [1 + 4\lambda_1 / (\bar{\lambda} \phi(S_0)^2)] - (R+1-q)(a-1) \log p} + \frac{2}{p}. \end{aligned}$$

The right side of the above display goes to zero when  $R > q \left[1 + \frac{M}{a-1} \left(1 + \frac{4\lambda_1}{\bar{\lambda} \phi(S_0)^2}\right)\right]$  for some  $M > 2$ .  $\square$

## A.2.2 Proof of Theorem 6.4

*Proof.* By Theorem 6.3, the posterior distribution is asymptotically supported on the event  $E = \{\boldsymbol{\beta} : |\boldsymbol{\gamma}(\boldsymbol{\beta})| \leq q(1+K)\}$ , where  $K = \frac{M}{a-1} \left(1 + \frac{4\lambda_1}{\bar{\lambda} \phi(S_0)^2}\right)$ . Thus, we confine attention to

$E^* = E \cap \tau_0$ , where  $\tau_0$  was defined in the proof of Theorem 6.3. From (A.14) and (A.15), we can see that

$$\Pi(\mathcal{B} | \mathbf{Y}, \theta) \mathbb{I}_{\tau_0} \leq \frac{p^{2q} e^{1+D}}{\pi(q | \theta)} \int_{\mathcal{B}} e^{-\frac{1}{2} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 + 3\bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| + \lambda_1 |\boldsymbol{\beta}|} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}. \quad (\text{A.25})$$

We now use the definition of the compatibility number in vectors of generalized dimensionality (A.6). With the inequality  $4uv \leq u^2 + 4v^2$ , we can then write

$$\begin{aligned} (4-1)\bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| &\leq \frac{4\bar{\lambda} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\| |\gamma(\boldsymbol{\beta} - \boldsymbol{\beta}_0)|^{1/2}}{\sqrt{n} \bar{\phi}(|\gamma(\boldsymbol{\beta} - \boldsymbol{\beta}_0)|)} - \bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| \\ &\leq \frac{1}{4} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 + \frac{16\bar{\lambda}^2 |\gamma(\boldsymbol{\beta} - \boldsymbol{\beta}_0)|}{n \bar{\phi}(|\gamma(\boldsymbol{\beta} - \boldsymbol{\beta}_0)|)^2} - \bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| \end{aligned}$$

Thus,

$$\Pi(\mathcal{B} | \mathbf{Y}, \theta) \mathbb{I}_{\tau_0} \leq \frac{p^{2q} e^{1+D}}{\pi(q | \theta)} e^{\frac{16\bar{\lambda}q(2+K)}{n[\bar{\phi}(2q+Kq)]^2}} \int_{\mathcal{B}} e^{-\frac{1}{4} \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 - \bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| + \lambda_1 |\boldsymbol{\beta}|} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta}.$$

Denote now  $\mathcal{B} = \{\boldsymbol{\beta} \in E^* : \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\| > R\}$ . Then

$$\begin{aligned} \Pi(\mathcal{B} | \mathbf{Y}, \theta) \mathbb{I}_{\tau_0} &\leq \frac{p^{2q} e^{1+D}}{\pi(q | \theta)} e^{\frac{16\bar{\lambda}q(2+K)}{n[\bar{\phi}(2q+Kq)]^2}} e^{-\frac{R^2}{4}} \int_{\mathcal{B}} e^{-\bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| + \lambda_1 |\boldsymbol{\beta}|} \Pi(\boldsymbol{\beta} | \theta) d\boldsymbol{\beta} \\ &\leq \frac{p^{2q} e^{1+D}}{\pi(q | \theta)} e^{\frac{16\bar{\lambda}q(1+K)}{n[\bar{\phi}(2q+Kq)]^2}} e^{-\frac{R^2}{4}} e^{\lambda_1 \delta p} \sum_{k=0}^p 8^k \pi(k | \theta) \end{aligned}$$

Now, because  $\pi > \theta e^{-\delta\lambda_1}$  and  $\theta \sim 1/p^a$  we can write

$$\pi(q | \theta) \succeq \pi(q-1 | \theta) \frac{e^{-\delta\lambda_1}}{p^a} \succeq \pi(q-1 | \theta) \frac{C_2}{p^a},$$

where we used the fact  $e^{-\delta\lambda_1} > C_2$ . Thus,  $\pi(q | \theta) \succeq C_2^q / p^{a^q} \pi(0 | \theta)$ . Thereby,

$$\Pi(\mathcal{B} | \mathbf{Y}, \theta) \mathbb{I}_{\tau_0} \leq C_2^{-q} p^{q(2+a)} e^{1+D} e^{\frac{16\bar{\lambda}^2 q(2+K)}{n[\bar{\phi}(2q+Kq)]^2}} e^{-\frac{R^2}{4}} e^{\lambda_1 \delta p} \sum_{k=0}^p \left( \frac{8}{p^{a-1}} \right)^k.$$

This quantity will tend to zero for

$$R^2 \succeq 4q(2+a) \log p + \frac{16\bar{\lambda}^2 q(2+K)}{n[\bar{\phi}(2q+Kq)]^2} \succeq \frac{q(2+K) \log p}{[\bar{\phi}(2q+Kq)]^2}.$$

This proves the first assertion. The second assertion follows from

$$\bar{\lambda}|\boldsymbol{\beta} - \boldsymbol{\beta}_0| \leq \|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\|^2 + \frac{\bar{\lambda}^2 q(2+K)}{2n\bar{\phi}(2q+Kq)^2}$$

and the last one from the definition of a minimal eigenvalue restricted to  $\boldsymbol{\beta}$  of generalized dimensionality at most  $s$  in (A.7), which yields  $\|\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\| > \bar{c}(2q+Kq)\|\mathbf{X}\| \|\boldsymbol{\beta} - \boldsymbol{\beta}_0\|$ .  $\square$

## B Proof of Lemma 6.2

*Proof.* Denote by  $\mathbf{e}_j$  the  $j^{th}$  canonical vector. The global optimality of  $\widehat{\boldsymbol{\beta}}$  yields

$$-||\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}||/2 + \text{pen}_{NS}(\widehat{\boldsymbol{\beta}}) \geq -||\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}} - t\mathbf{X}_j||/2 + \text{pen}_{NS}(\widehat{\boldsymbol{\beta}} + t\mathbf{e}_j)$$

for any  $t \in \mathbb{R}$ . Because  $||\mathbf{X}_j||^2 = n$ , this is equivalent to

$$t\mathbf{X}_j'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) \leq t^2n/2 + \log \left[ \frac{\pi(\widehat{\boldsymbol{\beta}})}{\pi(\widehat{\boldsymbol{\beta}} + t\mathbf{e}_j)} \right] < nt^2/2 - \tilde{\rho}(t; \widehat{\boldsymbol{\beta}}_{\setminus j}),$$

where we used the definition of the conditional singleton (3.10) together with (3.8) and (3.9). The statement of the lemma follows from the definition  $\Delta_j$ .  $\square$

## C Simulation Study: Initialization

In order to assess the sensitivity of the *Spike-and-Slab LASSO* to initialization, we repeated the simulation study from Section 5 (with the correlated block design), considering different starting vectors  $\boldsymbol{\beta}^0$  as well as different spike penalty sequences  $\{\lambda_0^1, \dots, \lambda_0^L\}$ . In particular, we generated 100 datasets and initialized *SSL* at (a) the zero vector, (b) a random vector from  $\mathcal{N}_p(0, I)$  and (c) a random vector from  $\mathcal{N}_p(0, 10 \times I)$ . Different random starting vectors were used for each dataset. The results are summarized in Table 1. With a fine grid of  $\lambda_0$  values (e.g. choosing  $\lambda_0 \in \{1, 2, \dots, 100\}$ ), we do not observe any sensitivity to initialization! With a crude ladder of  $\lambda_0$  values (e.g. choosing  $\lambda_0 \in \{10, 20, 30, \dots, 100\}$ ), there is indeed sensitivity to initialization. However, we observed that the zero starting vector performed the best regardless of the choice of  $\lambda_0$  values.

## D Further Implementations

### D.1 Implementation via EMVS

The coordinate-wise optimization based on the univariate soft-thresholding operator (as explained in Section 4)

$$S(z, \lambda) = \frac{1}{n}(|z| - \lambda)_+ \text{sign}(z).$$

resembles the LLA algorithm (??), which iterates over joint updates

$$\boldsymbol{\beta}^{(k+1)} = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ -\frac{1}{2}||\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}||^2 - \sum_{j=1}^p \lambda_{\theta}^*(\beta_j^{(k)})|\beta_j| \right\}. \quad (\text{D.1})$$



$\lambda_0$	$\lambda_1$	$\theta$	MSE	FDR	FNR	$\hat{q}$	TRUE	TIME	HAM
$\beta^0 = \mathbf{0}_p$									
$\{1, 2, \dots, 100\}$	1	$\mathcal{B}(1, p)$	3.64	0.288	0.288	6	13	0.63	3.46
$\{1, 12, 23 \dots, 100\}$	1	$\mathcal{B}(1, p)$	4.65	0.358	0.357	6.01	5	0.11	4.29
$\{10, 20, 30 \dots, 100\}$	1	$\mathcal{B}(1, p)$	1.96	0.193	0.192	6.01	22	0.02	2.31
$\beta^0 \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{I}_p)$									
$\{1, 2, \dots, 100\}$	1	$\mathcal{B}(1, p)$	3.64	0.288	0.288	6	13	0.74	3.46
$\{1, 12, 23 \dots, 100\}$	1	$\mathcal{B}(1, p)$	5.78	0.427	0.42	6.07	2	0.14	5.11
$\{10, 20, 30 \dots, 100\}$	1	$\mathcal{B}(1, p)$	5.22	0.38	0.377	6.03	8	0.09	4.55
$\beta^0 \sim \mathcal{N}_p(\mathbf{0}_p, 10 \times \mathbf{I}_p)$									
$\{1, 2, \dots, 100\}$	1	$\mathcal{B}(1, p)$	3.64	0.288	0.288	6	13	0.78	3.46
$\{1, 12, 23 \dots, 100\}$	1	$\mathcal{B}(1, p)$	7.74	0.524	0.517	6.09	2	0.18	6.29
$\{10, 20, 30 \dots, 100\}$	1	$\mathcal{B}(1, p)$	14.38	0.88	0.842	7.93	0	0.16	12.03

Table 1: Simulation study using 100 repetitions; MSE (average mean squared error), FDR (false discovery rate), FNR (false non-discovery rate), DIM (average size of the model), TRUE (# true model detected), TIME (average execution time in seconds), HAM (average Hamming distance)

From another point of view, (D.1) coincides with the M-step of a Bayesian EM algorithm for posterior mode detection under continuous spike-and-slab priors, which treats  $\gamma$  as missing data and keeps  $\theta$  fixed. This connection is made apparent by the fact  $\lambda_\theta^*(\beta_j^{(k)}) = \lambda_1 p_\theta^*(\beta_j^{(k)}) + \lambda_0 [1 - p_\theta^*(\beta_j^{(k)})]$ , and by noting  $p_\theta^*(\beta_j) = \mathbb{E}(\gamma_j | \beta_j^{(k)}, \theta)$  (the E-step calculation). A similar strategy was implemented for a mixture of two Gaussian distributions in the EMVS procedure by ?. Whereas their approach was based on iteratively solving adaptive ridge regressions, here it entails solving weighted LASSO regressions. The advantages of using EMVS with the LASSO updates are (a) automatic variable selection through thresholding, (b) faster speed of convergence (follows from considerations of ?).

Extending LLA to the case of a non-separable penalty is achieved naturally within the Bayesian EM framework by treating  $\theta$  as an additional model parameter. Instead of carrying forward the same fixed value, one now simply updates  $\theta$  throughout the algorithm. The non-separable variant of the M-step thus uses  $\theta = \theta^{(k)}$  to obtain  $\beta^{(k+1)}$  from (D.1). This step is followed by a new update  $\theta^{(k+1)}$  according to

$$\theta^{(k+1)} = \frac{\sum_{j=1}^p p_{\theta^{(k)}}^*(\beta_j^{(k+1)}) + a - 1}{a + b + p - 2}. \quad (\text{D.2})$$

The calculation (D.2) follows directly from equation (3.12) of ?. A variant of this strategy was implemented for sparse factor analysis by ?, where more details on this algorithm can be found.

---



---

<b>EMVS Algorithm:</b> <i>The Spike-and-Slab LASSO</i>
Input a grid of increasing $\lambda_0$ values $I = \{\lambda_0^1, \dots, \lambda_0^L\}$
For each value $l \in \{1, \dots, L\}$
Set $k = 0$
(a) Initialize: $\beta_l^{(k)} = \beta^*, \theta^{(0)} = \theta^*$
(b) While $\text{diff} > \varepsilon$
(i) Increment $k$
(ii) Update $\beta_l^{(k)}$ according to (D.1) with $\theta = \theta^{(k)}$
(iii) Update $\theta^{(k)}$ according to (D.2)
(iv) $\text{diff} = \ \beta^{(k)} - \beta^{(k-1)}\ _2$
(c) Return $\beta_l^{(k)}$
(d) Assign $\beta^* = \beta_l^{(k)}$

---



---

Table 2: The EMVS implementation of the Spike-and-Slab LASSO procedure

## D.2 Posterior Simulation

The *SSL* and *NSSL* priors are also amenable to posterior simulation. Direct Gibbs sampling is available through the exponential scale mixture representation of the Laplace distribution (?), applying the SSVS strategy (?). Alternatively, one could deploy a variant of an orthant sampler developed for the Bayesian LASSO by ?. Whereas simulating from the full-dimensional posterior  $\pi(\beta | \mathbf{Y})$  will only be practical when  $p$  is not overwhelmingly big, initiating the sampler at a posterior mode can save burn-in time and provide a quick insight into uncertainty surrounding the mode. Alternatively, one could confine the simulation to a lower-dimensional subspace, sampling only from active coordinates identified by the mode hunting strategies.

## E Identifiability Considerations

As a followup on Lemma 6.1, it is instructive to compare the geometry of the region  $C(\eta; \beta_0)$  under different penalty functions. This region captures the discrepancy  $\Theta = \hat{\beta} - \beta_0$  between the true vector and its shrunken estimate. The shape of this region thus provides useful insights about the nature of the shrinkage.

For the sake of illustration, we assume  $\beta_0 = (0, 3)'$  and  $\eta = 0.45$ . Under the LASSO penalty (where  $\lambda_1 = \lambda_0$ ), the set  $C(\eta; \beta_0)$  has a diamond shape (Figure 1(a)), embedded within a cone  $\left\{ \Theta \in \mathbb{R}^p : |\Theta_{S^c}| \leq \frac{1+\eta}{1-\eta} |\Theta_S| \right\}$ . On the other hand, for the other limiting case  $\lambda_0 = \infty$  and  $\lambda_1 = 0$ , the *SSL* penalty corresponds to the  $\ell_0$  penalty. The set  $C(\eta; \beta_0)$  then consists of those values  $\Theta = \hat{\beta} - \beta_0$  for which  $\|\hat{\beta}\|_0 < 2$ , as marked by the two solid lines corresponding to  $\Theta_1 = 0$  and  $\Theta_2 = -3$ . The *SSL* penalty yields a compromise between

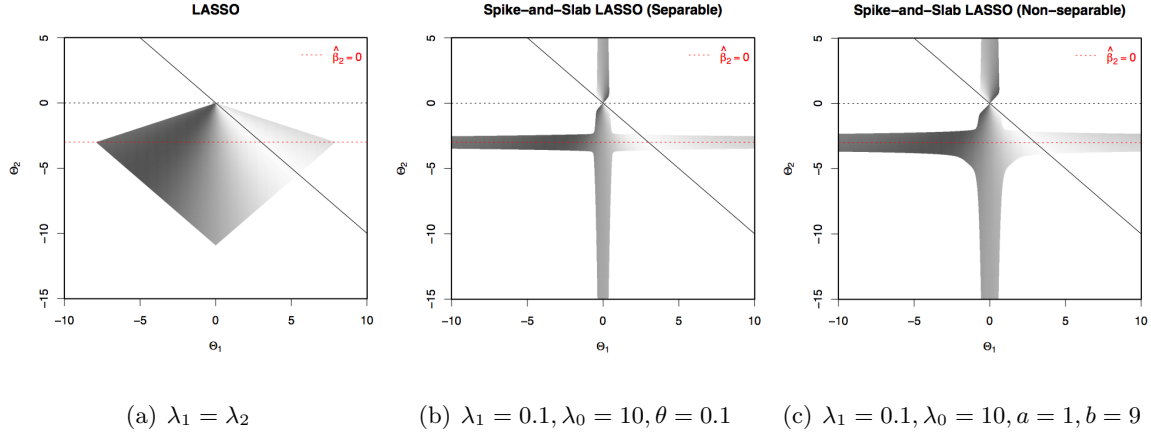


Figure 1: Plots of the feasible regions for  $\Theta$  under the LASSO penalty and the *SSL* penalty.

these two extremes. With  $1 < \lambda_0 < \infty$  and  $0 < \lambda_1 < 1$ ,  $C(\eta; \beta_0)$  is a star-shaped wrap around the set  $\{\Theta : \|\hat{\beta}\|_0 < 2\}$ . With larger  $\lambda_1$ , the set begins to resemble a diamond. The non-separable penalty ties the coordinates together, making the set larger in the center.

Figures 1(a), 1(b) and 1(c) are actual heat-maps of the restricted eigenvalues  $\frac{\|\mathbf{X}\Theta\|}{\|\mathbf{X}\| \|\Theta\|}$  inside  $C(\eta; \beta_0)$ ; the darker the shade of grey, the larger the value. Here,  $\mathbf{X}$  contains  $n = 100$  observations on 2 highly collinear variables (correlation  $\rho = 0.96$ ). The diamond in Figure 1(a) is seen as a continuum of rays of equal eigenvalues, the minimum attained on the ray  $\Theta_2 = -\Theta_1$  (marked by a solid line). This ray dissects all three sets in Figure 1. Under the  $\ell_0$  penalty, the intersection occurs at  $\hat{\beta} = (3, 0)'$ , the “opposite” of  $\beta_0 = (0, 3)$  where the correct variable was mistaken for its “knockoff”. This unfavorable case is assigned a very small  $c(\eta; \beta_0)$  value, an indication that the true variable is not easily distinguishable. Interestingly, compared to the LASSO set, the *SSL* feasible regions indicate that at least one of the coordinates in  $\hat{\beta}$  must be negligible, acknowledging the sparsity of the true  $\beta_0$ . However, despite their different geometries, the minimal eigenvalue taken over these different sets is the same. Thus, regardless of the penalty, the same identifiability condition has to be imposed in this example.