

Supplement to “Bounded support in linear random coefficient models: Identification and variable selection”

Philipp Hermann and Hajo Holzmann*
Department of Mathematics and Computer Science
Philipps-Universität Marburg
{herm, holzmann}@mathematik.uni-marburg.de

February 26, 2024

Abstract

As supplementary material for the paper ‘Bounded support in linear random coefficient models: Identification and variable selection’ we provide the proofs of the results in Section 3.1, a primal-dual witness condition of the adaptive LASSO as well as an analysis for estimating the mean vector with diverging number of parameters.

6 Proofs for Section 3.1

Proof of Proposition 3.1

Proof of Proposition 3.1. From Theorem 2.4, under the assumptions of the proposition the matrix

$$S = \left[v\left((1, \mathbf{W}_1^\top)^\top\right), \dots, v\left((1, \mathbf{W}_{p(p+1)/2}^\top)^\top\right) \right]^\top$$

is of full rank with positive probability. Therefore, the random positive semi-definite matrix

$$\frac{1}{n} (\mathbb{X}_n^\sigma)^\top \mathbb{X}_n^\sigma = \frac{1}{n} \sum_{i=1}^n v\left((1, \mathbf{W}_i^\top)^\top\right) v\left((1, \mathbf{W}_i^\top)^\top\right)^\top$$

for $n \geq p(p+1)/2$ is positive definite with positive probability. Hence its expected value, which equals C^σ , is positive definite. \square

Proof of Theorem 3.2

Turning to the proof of Theorem 3.2, recall the decomposition

$$\varepsilon_n^\sigma = \delta_n + \zeta_n + \xi_n \tag{6.1}$$

*Corresponding author. Prof. Dr. Hajo Holzmann, Department of Mathematics and Computer Science, Philipps-Universität Marburg, Hans-Meerweinstr., 35043 Marburg, Germany

with

$$\begin{aligned}\delta_n &:= \left(v(\mathbf{X}_1)^\top \text{vec}(D_1 - \Sigma^*), \dots, v(\mathbf{X}_n)^\top \text{vec}(D_n - \Sigma^*) \right)^\top, \\ \zeta_n &:= \left(v(\mathbf{X}_1)^\top \text{vec}(E_n), \dots, v(\mathbf{X}_n)^\top \text{vec}(E_n) \right)^\top, \\ \xi_n &:= \left(v(\mathbf{X}_1)^\top \text{vec}(F_{n,1}), \dots, v(\mathbf{X}_n)^\top \text{vec}(F_{n,n}) \right)^\top.\end{aligned}\tag{6.2}$$

of the error term (3.3).

Lemma 6.1. *Under the conditions of Theorem 3.2, we have that*

$$\frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top (\zeta_n + \xi_n) = o_{\mathbb{P}}(1).$$

The proofs of the previous as well as the following lemma are deferred to the end of this section.

Lemma 6.2. *Set $Z_n^{\sigma,1} = \frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \delta_n$, then*

$$\mathbb{E}[Z_n^{\sigma,1} | \mathbb{X}_n^\sigma] = \mathbf{0}_{p(p+1)/2} \quad \text{and} \quad \text{Cov}(Z_n^{\sigma,1} | \mathbb{X}_n^\sigma) = \frac{1}{n} (\mathbb{X}_n^\sigma)^\top \Omega_n^\sigma \mathbb{X}_n^\sigma,$$

where Ω_n^σ is a diagonal matrix with entries $v(\mathbf{X}_1)^\top \Psi^* v(\mathbf{X}_1), \dots, v(\mathbf{X}_n)^\top \Psi^* v(\mathbf{X}_n)$. In particular, $\text{Cov}(Z_n^{\sigma,1}) = \mathbf{B}^\sigma$ and $Z_n^{\sigma,1} = \mathcal{O}_{\mathbb{P}}(1)$.

Proof of Theorem 3.2. We shall use the primal-dual witness characterization of the adaptive LASSO in Lemma 7.1 in this supplement, Section 7, to prove the sign-consistency (3.7), and the Lindeberg-Feller central limit theorem for random vectors, see van der Vaart (1998, Proposition 2.27), to prove the asymptotic normality (3.8). For more details see also the proof of Theorem 3.5. By Lemmas 6.1 and 6.2, setting

$$\mathbf{P}_{\mathbb{X}_{n,S_\sigma}^\sigma} = \mathbf{I}_n - \mathbb{X}_{n,S_\sigma}^\sigma \left((\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbb{X}_{n,S_\sigma}^\sigma \right)^{-1} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top,$$

we have that

$$\frac{1}{\sqrt{n}} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbf{P}_{\mathbb{X}_{n,S_\sigma}^\sigma} \varepsilon_n^\sigma = \mathcal{O}_{\mathbb{P}}(1).$$

In addition, the requirements $\sqrt{n} \lambda_n^\sigma \rightarrow 0$ and $\sqrt{n} (\hat{\sigma}_n^{\text{init}} - \sigma^*) = \mathcal{O}_{\mathbb{P}}(1)$ in Theorem 3.2 lead to

$$0 \leq \frac{\sqrt{n} \lambda_n^\sigma}{|\hat{\sigma}_{n,k}^{\text{init}}|} \leq \frac{\sqrt{n} \lambda_n^\sigma}{\left| |\sigma_k^*| - |\hat{\sigma}_{n,k}^{\text{init}} - \sigma_k^*| \right|} \xrightarrow{\mathbb{P}} 0\tag{6.3}$$

for all $k \in S_\sigma$ since $|\sigma_k^*| > 0$ for these k . This implies

$$\begin{aligned}\sqrt{n} \left[(\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbb{X}_{n,S_\sigma}^\sigma \left((\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbb{X}_{n,S_\sigma}^\sigma \right)^{-1} \left(\lambda_n^\sigma \left(\frac{1}{|\hat{\sigma}_{n,S_\sigma}^{\text{init}}|} \odot \text{sign}(\sigma_{S_\sigma}^*) \right) \right) \right] + \frac{1}{n} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbf{P}_{\mathbb{X}_{n,S_\sigma}^\sigma} \varepsilon_n^\sigma \\ = \mathcal{O}_{\mathbb{P}}(1) o_{\mathbb{P}}(1) + \mathcal{O}_{\mathbb{P}}(1) = \mathcal{O}_{\mathbb{P}}(1).\end{aligned}\tag{6.4}$$

Moreover, $\sqrt{n}(\hat{\sigma}_n^{\text{init}} - \sigma^*) = \mathcal{O}_{\mathbb{P}}(1)$ implies also $\sqrt{n}\hat{\sigma}_{n,k}^{\text{init}} = \mathcal{O}_{\mathbb{P}}(1)$ for all $k \in S_\sigma^c$ since $\sigma_k^* = 0$ for these k . Thus, by the second requirement $n\lambda_n^\sigma \rightarrow \infty$ on the regularization parameter it follows that

$$\frac{\sqrt{n}\lambda_n^\sigma}{|\hat{\sigma}_{n,k}^{\text{init}}|} = \frac{n\lambda_n^\sigma}{\sqrt{n}|\hat{\sigma}_{n,k}^{\text{init}}|} \xrightarrow{\mathbb{P}} \infty$$

for all $k \in S_\sigma^c$. Together with (6.4) this implies the first condition (7.1) of Lemma 7.1 with high probability for a sufficient large number n of observations. Furthermore, let

$$\tilde{\sigma}_{n,S_\sigma} = \sigma_{S_\sigma}^* + \left(\frac{1}{n} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbb{X}_{n,S_\sigma}^\sigma \right)^{-1} \left(\frac{1}{n} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \varepsilon_n^\sigma - \lambda_n^\sigma \left(\frac{1}{|\hat{\sigma}_{n,S_\sigma}^{\text{init}}|} \odot \text{sign}(\sigma_{S_\sigma}^*) \right) \right).$$

Then we obtain

$$\sqrt{n}(\tilde{\sigma}_{n,S_\sigma} - \sigma_{S_\sigma}^*) = \left(\frac{1}{n} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbb{X}_{n,S_\sigma}^\sigma \right)^{-1} \frac{1}{\sqrt{n}} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \varepsilon_n^\sigma + \mathcal{O}_{\mathbb{P}}(1)$$

by (6.3). Moreover, with Lemmas 6.1 and 6.2 it follows that

$$\begin{aligned} \sqrt{n}(\tilde{\sigma}_{n,S_\sigma} - \sigma_{S_\sigma}^*) &= \left(\frac{1}{n} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \mathbb{X}_{n,S_\sigma}^\sigma \right)^{-1} \frac{1}{\sqrt{n}} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \delta_n + \mathcal{O}_{\mathbb{P}}(1) \\ &= \mathcal{O}_{\mathbb{P}}(1) + \mathcal{O}_{\mathbb{P}}(1) = \mathcal{O}_{\mathbb{P}}(1), \end{aligned} \quad (6.5)$$

which leads to $\tilde{\sigma}_{n,S_\sigma} - \sigma_{S_\sigma}^* = \mathcal{O}_{\mathbb{P}}(1)$. Therefore the second condition, $\text{sign}(\tilde{\sigma}_{n,S_\sigma}) = \text{sign}(\sigma_{S_\sigma}^*)$, of Lemma 7.1 is also satisfied with high probability for large n . Sign-consistency of the adaptive LASSO and $\hat{\sigma}_{n,S_\sigma}^{\text{AL}} = \tilde{\sigma}_{n,S_\sigma}$ is the consequence.

Note that for the asymptotic normality (3.8) of the rescaled estimation error only the first term in (6.5) is crucial. Hence we consider the random vectors

$$Z_n^{\sigma,1} = \frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \delta_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n (e_i^\top \delta_n) \mathbf{v}(\mathbf{X}_i) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\mathbf{v}(\mathbf{X}_i)^\top \text{vec}(D_i - \Sigma^*) \right) \mathbf{v}(\mathbf{X}_i),$$

where $D_i = (\mathbf{A}_i - \mu^*)(\mathbf{A}_i - \mu^*)^\top$ and δ_n is defined in (6.2). Now we want to apply the Lindeberg-Feller central limit theorem for the array

$$Q_{n,i} = \frac{1}{\sqrt{n}} \left(\mathbf{v}(\mathbf{X}_i)^\top \text{vec}(D_i - \Sigma^*) \right) \mathbf{v}(\mathbf{X}_i), \quad i = 1, \dots, n,$$

of random vectors. These are independent and identically distributed in each row (for fixed n) since $(\mathbf{X}_1^\top, \mathbf{A}_1^\top)^\top, \dots, (\mathbf{X}_n^\top, \mathbf{A}_n^\top)^\top$ are independent and identically distributed. Furthermore, they are centered,

$$\mathbb{E}[Q_{n,i}] = \frac{1}{\sqrt{n}} \mathbb{E} \left[\mathbb{E} \left[\mathbf{v}(\mathbf{X}_i)^\top \text{vec}(D_i - \Sigma^*) \mid \mathbb{X}_n^\sigma \right] \mathbf{v}(\mathbf{X}_i) \right] = \frac{1}{\sqrt{n}} \mathbb{E}[0 \cdot \mathbf{v}(\mathbf{X}_i)] = \mathbf{0}_{p(p+1)/2},$$

and for the sum of the covariance matrices

$$\sum_{i=1}^n \text{Cov}(Q_{n,i}) = \text{Cov} \left(\sum_{i=1}^n Q_{n,i} \right) = \text{Cov}(Z_n^{\sigma,1})$$

we get by Lemma 6.2

$$\sum_{i=1}^n \text{Cov}(Q_{n,i}) = \mathbf{B}^\sigma.$$

Moreover, we obtain for arbitrary $\delta > 0$ the equation

$$\begin{aligned} \sum_{i=1}^n \mathbb{E} \left[\|Q_{n,i}\|_2^2 \mathbb{1} \{ \|Q_{n,i}\|_2 > \delta \} \right] &= \mathbb{E} \left[\mathbf{v}(\mathbf{X})^\top \text{vec}(D - \Sigma^*) \mathbf{v}(\mathbf{X})^\top \text{vec}(D - \Sigma^*) \mathbf{v}(\mathbf{X})^\top \mathbf{v}(\mathbf{X}) \right. \\ &\quad \left. \cdot \mathbb{1} \{ \mathbf{v}(\mathbf{X})^\top \text{vec}(D - \Sigma^*) \mathbf{v}(\mathbf{X})^\top \text{vec}(D - \Sigma^*) \mathbf{v}(\mathbf{X})^\top \mathbf{v}(\mathbf{X}) > \delta^2 n \} \right]. \end{aligned}$$

The expected mean $\mathbb{E}[\mathbf{v}(\mathbf{X})^\top \text{vec}(D - \Sigma^*) \mathbf{v}(\mathbf{X})^\top \text{vec}(D - \Sigma^*) \mathbf{v}(\mathbf{X})^\top \mathbf{v}(\mathbf{X})]$ exists because of Assumption 1 and the Cauchy Schwarz inequality. Thus we get

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{E} \left[\|Q_{n,i}\|_2^2 \mathbb{1} \{ \|Q_{n,i}\|_2 > \delta \} \right] = 0$$

by Lebesgue's dominated convergence theorem, which coincides with Lindeberg's condition, see van der Vaart (1998, Proposition 2.27). Hence the mentioned proposition implies the weak convergence

$$Z_n^{\sigma,1} = \frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \delta_n = \sum_{i=1}^n Q_{n,i} \xrightarrow{d} Q \sim \mathcal{N}_{p(p+1)/2}(\mathbf{0}_{p(p+1)/2}, \mathbf{B}^\sigma),$$

respectively

$$\frac{1}{\sqrt{n}} (\mathbb{X}_{n,S_\sigma}^\sigma)^\top \delta_n \xrightarrow{d} Q_{S_\sigma} \sim \mathcal{N}_{s_\sigma}(\mathbf{0}_{s_\sigma}, \mathbf{B}_{S_\sigma S_\sigma}^\sigma).$$

So all in all a multivariate version of Slutsky's theorem, see for example van der Vaart (1998, Theorem 2.7, Lemma 2.8), together with equation (6.5) leads to

$$\sqrt{n} (\hat{\sigma}_{n,S_\sigma}^{\text{AL}} - \sigma_{S_\sigma}^*) \xrightarrow{d} (\mathbf{C}_{S_\sigma S_\sigma}^\sigma)^{-1} Q_{S_\sigma}.$$

In addition, it follows that

$$(\mathbf{C}_{S_\sigma S_\sigma}^\sigma)^{-1} Q_{S_\sigma} \sim \mathcal{N}_{s_\sigma}(\mathbf{0}_{s_\sigma}, (\mathbf{C}_{S_\sigma S_\sigma}^\sigma)^{-1} \mathbf{B}_{S_\sigma S_\sigma}^\sigma (\mathbf{C}_{S_\sigma S_\sigma}^\sigma)^{-1})$$

by the symmetry of $\mathbf{C}_{S_\sigma S_\sigma}^\sigma$ and the properties of the multivariate normal distribution, and hence the asymptotic normality (3.8). \square

Proof of Lemma 6.1. We prove Lemma 6.1 in two steps. First we show that

$$\frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \zeta_n = o_{\mathbb{P}}(1). \quad (6.6)$$

We obtain

$$\begin{aligned}
\frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \zeta_n &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (e_i^\top \zeta_n) \mathbf{v}(\mathbf{X}_i) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\mathbf{v}(\mathbf{X}_i)^\top \text{vec}(E_n) \right) \mathbf{v}(\mathbf{X}_i) \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\sum_{q=1}^{p(p+1)/2} \mathbf{v}(\mathbf{X}_i)_q \text{vec}(E_n)_q \right) \mathbf{v}(\mathbf{X}_i) \\
&= \sum_{q=1}^{p(p+1)/2} \sqrt{n} \text{vec}(E_n)_q \left(\frac{1}{n} \sum_{i=1}^n \mathbf{v}(\mathbf{X}_i)_q \mathbf{v}(\mathbf{X}_i) \right), \tag{6.7}
\end{aligned}$$

where

$$E_n = (\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}_n)(\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}_n)^\top.$$

By the assumption on $\widehat{\boldsymbol{\mu}}_n$ we get $e_k^\top E_n e_l = (\widehat{\mu}_{n,k} - \mu_k^*)(\widehat{\mu}_{n,l} - \mu_l^*) = \mathcal{O}_{\mathbb{P}}(1/n)$ for $k, l \in \{1, \dots, p\}$, and hence also

$$\sqrt{n} \text{vec}(E_n)_q = \mathcal{O}_{\mathbb{P}}\left(\frac{1}{\sqrt{n}}\right) \tag{6.8}$$

for all $q \in \{1, \dots, p(p+1)/2\}$. Furthermore, the random vectors $Q_i^q = \mathbf{v}(\mathbf{X}_i)_q \mathbf{v}(\mathbf{X}_i)$ are independent and identically distributed with

$$\mathbb{E}[\|Q_i^q\|_2] \leq \mathbb{E}[\|Q_i^q\|_1] = \mathbb{E}\left[\left\|\mathbf{v}(\mathbf{X}_i)_q \mathbf{v}(\mathbf{X}_i)\right\|_1\right] = \sum_{r=1}^{p(p+1)/2} \mathbb{E}\left[\left|\mathbf{v}(\mathbf{X}_i)_r \mathbf{v}(\mathbf{X}_i)_q\right|\right] < \infty,$$

so that by the law of large numbers

$$\frac{1}{n} \sum_{i=1}^n \mathbf{v}(\mathbf{X}_i)_q \mathbf{v}(\mathbf{X}_i) = \mathcal{O}_{\mathbb{P}}(1) \tag{6.9}$$

for all $q \in \{1, \dots, p(p+1)/2\}$ follows. In summary, (6.7), (6.8) and (6.9) lead to (6.6).

In the second step, consider

$$\frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \xi_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\sum_{q=1}^{p(p+1)/2} \mathbf{v}(\mathbf{X}_i)_q \text{vec}(F_{n,i})_q \right) \mathbf{v}(\mathbf{X}_i),$$

where $F_{n,i} = (\mathbf{A}_i - \boldsymbol{\mu}^*)(\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}_n)^\top + (\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}_n)(\mathbf{A}_i - \boldsymbol{\mu}^*)^\top$. Then we obtain analogously

$$\begin{aligned}
\frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \xi_n &= \sum_{k,l=1}^p \sqrt{n} (\widehat{\mu}_{n,k} - \mu_k^*) \left(-\frac{2}{n} \sum_{i=1}^n X_{i,k} X_{i,l} (A_{i,l} - \mu_l^*) \mathbf{v}(\mathbf{X}_i) \right) \\
&= \mathcal{O}_{\mathbb{P}}(1) \circ_{\mathbb{P}}(1) = \circ_{\mathbb{P}}(1),
\end{aligned}$$

since

$$\mathbb{E}\left[X_{1,k} X_{1,l} (A_{1,l} - \mu_l^*) \mathbf{v}(\mathbf{X}_1)\right] = 0$$

by the independence of \mathbf{X}_1 and \mathbf{A}_1 . □

Proof of Lemma 6.2. We obtain by simple calculation $\mathbb{E}[\delta_n | \mathbb{X}_n^\sigma] = \mathbf{0}_n$ and $\text{Cov}(\delta_n | \mathbb{X}_n^\sigma) = \Omega_n^\sigma$, hence

$$\mathbb{E}[Z_n^{\sigma,1} | \mathbb{X}_n^\sigma] = \frac{1}{\sqrt{n}} (\mathbb{X}_n^\sigma)^\top \mathbb{E}[\delta_n | \mathbb{X}_n^\sigma] = \mathbf{0}_{p(p+1)/2}$$

and

$$\text{Cov}(Z_n^{\sigma,1} | \mathbb{X}_n^\sigma) = \frac{1}{n} (\mathbb{X}_n^\sigma)^\top \text{Cov}(\delta_n | \mathbb{X}_n^\sigma) \mathbb{X}_n^\sigma = \frac{1}{n} (\mathbb{X}_n^\sigma)^\top \Omega_n^\sigma \mathbb{X}_n^\sigma.$$

For random variables Q_1, Q_2 and Q_3 the law of total covariance implies the decomposition

$$\text{Cov}(Q_1, Q_2) = \mathbb{E}[\text{Cov}(Q_1, Q_2 | Q_3)] + \text{Cov}(\mathbb{E}[Q_1 | Q_3], \mathbb{E}[Q_2 | Q_3]).$$

This can be extended to random vectors and covariance matrices and hence we obtain

$$\begin{aligned} \text{Cov}(Z_n^{\sigma,1}) &= \mathbb{E}[\text{Cov}(Z_n^{\sigma,1} | \mathbb{X}_n^\sigma)] + \text{Cov}(\mathbb{E}[Z_n^{\sigma,1} | \mathbb{X}_n^\sigma]) \\ &= \mathbb{E}\left[\frac{1}{n} (\mathbb{X}_n^\sigma)^\top \Omega_n^\sigma \mathbb{X}_n^\sigma\right] = \mathbf{B}^\sigma. \end{aligned}$$

Boundedness in probability follows since by the law of large numbers,

$$\frac{1}{n} (\mathbb{X}_n^\sigma)^\top \Omega_n^\sigma \mathbb{X}_n^\sigma = \frac{1}{n} \sum_{i=1}^n \left(\mathbf{v}(\mathbf{X}_i)^\top \Psi^* \mathbf{v}(\mathbf{X}_i) \right) \mathbf{v}(\mathbf{X}_i) \mathbf{v}(\mathbf{X}_i)^\top \rightarrow \mathbf{B}^\sigma.$$

□

7 The adaptive LASSO

We look for a fixed number $n \in \mathbb{N}$ of observations at the ordinary linear regression model

$$\mathbb{Y}_n = \mathbb{X}_n \beta^* + \varepsilon_n,$$

where $\mathbb{Y}_n \in \mathbb{R}^n$ is the vector of the response variables, $\mathbb{X}_n \in \mathbb{R}^{n \times p}$ the deterministic design matrix, $\beta^* \in \mathbb{R}^p$ the unknown coefficient vector and $\varepsilon_n \in \mathbb{R}^n$ represents additive noise. Moreover, we allow the coefficients β^* to be sparse, in other words it is $s \leq p$ for

$$S = \text{supp}(\beta^*) = \left\{ k \in \{1, \dots, p\} \mid \beta_k^* \neq 0 \right\}, \quad s = |S|.$$

In addition, let $S^c = \{1, \dots, p\} \setminus S$ be the relative complement of S . Because of the sparsity of the coefficients the linear regression model can also be expressed by

$$\mathbb{Y}_n = \mathbb{X}_{n,S} \beta_S^* + \varepsilon_n.$$

Consider the adaptive LASSO estimator with regularization parameter $\lambda_n > 0$, given by

$$\hat{\beta}_n^{\text{AL}} \in \rho_{n,\lambda_n}^{\text{AL}} := \arg \min_{\beta \in \mathbb{R}^p} \left(\frac{1}{n} \|\mathbb{Y}_n - \mathbb{X}_n \beta\|_2^2 + 2\lambda_n \sum_{k=1}^p \frac{|\beta_k|}{|\hat{\beta}_{n,k}^{\text{init}}|} \right),$$

where $\hat{\beta}_n^{\text{init}} \in \mathbb{R}^p$ is an initial estimator of β^* . If $\hat{\beta}_{n,k}^{\text{init}} = 0$, we require $\beta_k = 0$ in the above definition.

Lemma 7.1 (Primal-dual witness characterization of the adaptive LASSO). *Assume $s \leq n$ and $\text{rank}(\mathbb{X}_{n,S}) = s$. If*

$$\left| \mathbb{X}_{n,S^c}^\top \mathbb{X}_{n,S} \left(\mathbb{X}_{n,S}^\top \mathbb{X}_{n,S} \right)^{-1} \lambda_n \left(\frac{1}{|\widehat{\beta}_{n,S}^{\text{init}}|} \odot \text{sign}(\beta_S^*) \right) + \frac{1}{n} \mathbb{X}_{n,S^c}^\top \mathbb{P}_{\mathbb{X}_{n,S}^\perp} \varepsilon_n \right| < \frac{\lambda_n}{|\widehat{\beta}_{n,S^c}^{\text{init}}|} \quad (7.1)$$

with

$$\mathbb{P}_{\mathbb{X}_{n,S}^\perp} := \mathbf{I}_n - \mathbb{X}_{n,S} \left(\mathbb{X}_{n,S}^\top \mathbb{X}_{n,S} \right)^{-1} \mathbb{X}_{n,S}^\top$$

holds, and

$$\widetilde{\beta}_{n,S} = \beta_S^* + \left(\frac{1}{n} \mathbb{X}_{n,S}^\top \mathbb{X}_{n,S} \right)^{-1} \left(\frac{1}{n} \mathbb{X}_{n,S}^\top \varepsilon_n - \lambda_n \left(\frac{1}{|\widehat{\beta}_{n,S}^{\text{init}}|} \odot \text{sign}(\beta_S^*) \right) \right)$$

satisfies $\text{sign}(\widetilde{\beta}_{n,S}) = \text{sign}(\beta_S^*)$, then the unique adaptive LASSO solution $\rho_{n,\lambda_n}^{\text{AL}} = \{\widehat{\beta}_n^{\text{AL}}\}$ satisfies

$$\text{sign}(\widehat{\beta}_n^{\text{AL}}) = \text{sign}(\beta^*), \quad \widehat{\beta}_{n,S}^{\text{AL}} = \widetilde{\beta}_{n,S} \text{ and } \widehat{\beta}_{n,S^c}^{\text{AL}} = \mathbf{0}_{|S^c|}.$$

Proof. Cf. Lemma 12.1 in Zhou et al. (2009) with $\vec{w} = (1/|\widehat{\beta}_{n,1}^{\text{init}}|, \dots, 1/|\widehat{\beta}_{n,p}^{\text{init}}|)^\top \in \mathbb{R}^p$. \square

8 Estimating the means with diverging number p of parameters

The model is given in vector-matrix form by

$$\mathbb{Y}_n^\mu = \mathbb{X}_n^\mu \mu^* + \varepsilon_n^\mu,$$

where

$$\mathbb{Y}_n^\mu := (Y_1, \dots, Y_n)^\top, \quad \mathbb{X}_n^\mu := [\mathbf{X}_1, \dots, \mathbf{X}_n]^\top, \quad \varepsilon_n^\mu := \left(\mathbf{X}_1^\top (\mathbf{A}_1 - \mu^*), \dots, \mathbf{X}_n^\top (\mathbf{A}_n - \mu^*) \right)^\top.$$

Then the adaptive LASSO estimator with regularization parameter $\lambda_n^\mu > 0$ is given by

$$\widehat{\mu}_n^{\text{AL}} \in \rho_{\mu,n,\lambda_n^\mu}^{\text{AL}} := \arg \min_{\beta \in \mathbb{R}^p} \left(\frac{1}{n} \|\mathbb{Y}_n^\mu - \mathbb{X}_n^\mu \beta\|_2^2 + 2\lambda_n^\mu \sum_{k=1}^p \frac{|\beta_k|}{|\widehat{\mu}_{n,k}^{\text{init}}|} \right), \quad (8.1)$$

where $\widehat{\mu}_n^{\text{init}} \in \mathbb{R}^p$ is an initial estimator of μ^* . Note that if $\widehat{\mu}_{n,k}^{\text{init}} = 0$, we require again $\beta_k = 0$. Further, let

$$\mathbf{C}^\mu := \mathbb{E}[\mathbf{X} \mathbf{X}^\top], \quad \mathbf{B}^\mu := \mathbb{E}[(\mathbf{X}^\top \Sigma^* \mathbf{X}) \mathbf{X} \mathbf{X}^\top],$$

and we denote by

$$S_\mu := \text{supp}(\mu^*) = \left\{ k \in \{1, \dots, p\} \mid \mu_k^* \neq 0 \right\}, \quad s_\mu := |S_\mu|,$$

the support of the mean vector μ^* . $S_\mu^c := \{1, \dots, p\} \setminus S_\mu$ is again the corresponding relative complement.

Assumption 3 (Growing p). We assume that $(\mathbf{X}_i^\top, \mathbf{A}_i^\top)^\top$, $i = 1, \dots, n$, are identically distributed, and that

(A9) the random coefficients \mathbf{A} have finite second moments,

(A10) the covariate vector \mathbf{X} is sub-Gaussian,

(A11) $c_{C^\mu, l} \leq \lambda_{\min}(C^\mu) \leq \lambda_{\max}(C^\mu) \leq c_{C^\mu, u}$ for some positive constants $0 < c_{C^\mu, l} \leq c_{C^\mu, u} < \infty$, where $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ denote the minimal and maximal eigenvalues of a symmetric matrix A ,

(A12) $\lambda_{\max}(B^\mu) \leq c_{B^\mu, u}$ for some positive constant $c_{B^\mu, u} > 0$,

(A13) $\lim_{n \rightarrow \infty} p/n = 0$.

Theorem 8.1 (Variable selection for growing p). *Let Assumption 3 be satisfied, and assume that for the initial estimator $\hat{\mu}_n^{\text{init}}$ in the adaptive LASSO $\hat{\mu}_n^{\text{AL}}$ in (8.1) we have $\sqrt{n/p} \|\hat{\mu}_n^{\text{init}} - \mu^*\|_2 = \mathcal{O}_{\mathbb{P}}(1)$. Moreover, if the regularization parameter is chosen as $\lambda_n^\mu \rightarrow 0$,*

$$\sqrt{s_\mu n} \lambda_n^\mu / (\mu_{\min}^* \sqrt{p}) \rightarrow 0, \quad \sqrt{p} / (\mu_{\min}^* \sqrt{n}) \rightarrow 0, \quad n \lambda_n^\mu / p \rightarrow \infty$$

with $\mu_{\min}^* := \min_{k \in S_\mu} |\mu_k^*|$, then it follows that $\hat{\mu}_n^{\text{AL}}$ is sign-consistent,

$$\mathbb{P}\left(\text{sign}(\hat{\mu}_n^{\text{AL}}) = \text{sign}(\mu^*)\right) \rightarrow 1. \quad (8.2)$$

For the proof of Theorem 8.1 we need the following auxiliary lemma.

Lemma 8.2. *Set $Z_n^\mu = \frac{1}{n} (\mathbb{X}_n^\mu)^\top \varepsilon_n^\mu$, then $\|Z_n^\mu\|_2 = \mathcal{O}_{\mathbb{P}}(\sqrt{p/n})$.*

Proof of Lemma 8.2. It is

$$\begin{aligned} \mathbb{E}\left[\|Z_n^\mu\|_2^2\right] &= \frac{1}{n^2} \mathbb{E}\left[(\varepsilon_n^\mu)^\top \mathbb{X}_n^\mu (\mathbb{X}_n^\mu)^\top \varepsilon_n^\mu\right] = \frac{1}{n^2} \mathbb{E}\left[\text{trace}((\mathbb{X}_n^\mu)^\top \varepsilon_n^\mu (\varepsilon_n^\mu)^\top \mathbb{X}_n^\mu)\right] \\ &= \frac{1}{n^2} \mathbb{E}\left[\text{trace}((\mathbb{X}_n^\mu)^\top \mathbb{E}[\varepsilon_n^\mu (\varepsilon_n^\mu)^\top | \mathbb{X}_n^\mu] \mathbb{X}_n^\mu)\right] \\ &= \frac{1}{n} \text{trace}\left(\mathbb{E}\left[\frac{1}{n} (\mathbb{X}_n^\mu)^\top \Omega_n^\mu \mathbb{X}_n^\mu\right]\right), \end{aligned}$$

where $\Omega_n^\mu = \text{Cov}(\varepsilon_n^\mu | \mathbb{X}_n^\mu)$ is a diagonal matrix with entries $\mathbf{X}_1^\top \Sigma^* \mathbf{X}_1^\top, \dots, \mathbf{X}_n^\top \Sigma^* \mathbf{X}_n^\top$. It is obvious that

$$\mathbb{E}\left[\frac{1}{n} (\mathbb{X}_n^\mu)^\top \Omega_n^\mu \mathbb{X}_n^\mu\right] = B^\mu,$$

and hence we obtain by Assumption (A12) the estimate

$$\mathbb{E}\left[\|Z_n^\mu\|_2^2\right] = \frac{\text{trace}(B^\mu)}{n} \leq \frac{\lambda_{\max}(B^\mu) p}{n} \leq \frac{c_{B^\mu, u} p}{n}.$$

Markov's inequality implies the assertion. \square

Proof of Theorem 8.1. We shall use the primal-dual witness characterization of the adaptive LASSO in Lemma 7.1 in Section 7 to prove the sign-consistency (8.2). We obtain by Assumption (A10) and Wainwright (2019, Theorem 6.5) that

$$\left\| \frac{1}{n} (\mathbb{X}_n^\mu)^\top \mathbb{X}_n^\mu - C^\mu \right\|_{\text{M},2} = \mathcal{O}_{\mathbb{P}} \left(\sqrt{p/n} \right),$$

which implies together with the Assumptions (A11) and (A13) the invertibility of the Gram matrix for large n , and hence by Loh and Wainwright (2017, Lemma 11) we get also

$$\left\| \left(\frac{1}{n} (\mathbb{X}_n^\mu)^\top \mathbb{X}_n^\mu \right)^{-1} - (C^\mu)^{-1} \right\|_{\text{M},2} = \mathcal{O}_{\mathbb{P}} \left(\sqrt{p/n} \right).$$

Furthermore, basic properties of the ℓ_2 operator norm lead to

$$\left\| (\mathbb{X}_{n,S_\mu^c}^\mu)^\top \mathbb{X}_{n,S_\mu}^\mu \left((\mathbb{X}_{n,S_\mu}^\mu)^\top \mathbb{X}_{n,S_\mu}^\mu \right)^{-1} - C_{S_\mu^c S_\mu}^\mu (C_{S_\mu S_\mu}^\mu)^{-1} \right\|_{\text{M},2} = \mathcal{O}_{\mathbb{P}} \left(\sqrt{p/n} \right).$$

In particular, this implies

$$\left\| \left(\frac{1}{n} (\mathbb{X}_n^\mu)^\top \mathbb{X}_n^\mu \right)^{-1} \right\|_{\text{M},2} = \mathcal{O}_{\mathbb{P}}(1), \quad \left\| (\mathbb{X}_{n,S_\mu^c}^\mu)^\top \mathbb{X}_{n,S_\mu}^\mu \left((\mathbb{X}_{n,S_\mu}^\mu)^\top \mathbb{X}_{n,S_\mu}^\mu \right)^{-1} \right\|_{\text{M},2} = \mathcal{O}_{\mathbb{P}}(1). \quad (8.3)$$

Moreover, let $\hat{\mu}_{n,\min}^{\text{init}} := \min_{k \in S_\mu} |\hat{\mu}_{n,k}^{\text{init}}|$, then

$$\left| \frac{\hat{\mu}_{n,\min}^{\text{init}} - \mu_{\min}^*}{\mu_{\min}^*} \right| \leq \frac{1}{\mu_{\min}^*} \|\hat{\mu}_n^{\text{init}} - \mu^*\|_2 = \mathcal{O}_{\mathbb{P}} \left(\frac{\sqrt{p}}{\mu_{\min}^* \sqrt{n}} \right) = \mathfrak{o}_{\mathbb{P}}(1)$$

since $\sqrt{n/p} \|\hat{\mu}_n^{\text{init}} - \mu^*\|_2 = \mathcal{O}_{\mathbb{P}}(1)$ and $\sqrt{p}/(\mu_{\min}^* \sqrt{n}) \rightarrow 0$. This implies

$$\left(1 + \frac{\hat{\mu}_{n,\min}^{\text{init}} - \mu_{\min}^*}{\mu_{\min}^*} \right)^{-1} = \mathcal{O}_{\mathbb{P}}(1),$$

and hence we obtain

$$\begin{aligned} \sqrt{\frac{n}{p}} \left\| \lambda_n^\mu \left(\frac{1}{|\hat{\mu}_{n,S_\mu}^{\text{init}}|} \odot \text{sign}(\mu_{S_\mu}^*) \right) \right\|_2 &\leq \frac{\sqrt{n} \lambda_n^\mu}{\sqrt{p}} \left\| \frac{1}{|\hat{\mu}_{n,S_\mu}^{\text{init}}|} \right\|_2 \leq \frac{\sqrt{s_\mu n} \lambda_n^\mu}{\sqrt{p}} \left\| \frac{1}{|\hat{\mu}_{n,S_\mu}^{\text{init}}|} \right\|_\infty \\ &= \frac{\sqrt{s_\mu n} \lambda_n^\mu}{\sqrt{p}} (\hat{\mu}_{n,\min}^{\text{init}})^{-1} \\ &= \frac{\sqrt{s_\mu n} \lambda_n^\mu}{\sqrt{p}} (\mu_{\min}^*)^{-1} \left(1 + \frac{\hat{\mu}_{n,\min}^{\text{init}} - \mu_{\min}^*}{\mu_{\min}^*} \right)^{-1} \\ &= \frac{\sqrt{s_\mu n} \lambda_n^\mu}{\mu_{\min}^* \sqrt{p}} \mathcal{O}_{\mathbb{P}}(1) \\ &= \mathfrak{o}_{\mathbb{P}}(1) \end{aligned} \quad (8.4)$$

since $\sqrt{s_\mu n} \lambda_n^\mu / (\mu_{\min}^* \sqrt{p}) \rightarrow 0$ by assumption. It follows that

$$\begin{aligned}
& \sqrt{\frac{n}{p}} \left\| \left(\mathbb{X}_{n, S_\mu^c}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \left(\left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \right)^{-1} \left(\lambda_n^\mu \left(\frac{1}{|\widehat{\mu}_{n, S_\mu}^{\text{init}}|} \odot \text{sign}(\mu_{S_\mu}^*) \right) \right) + \frac{1}{n} \left(\mathbb{X}_{n, S_\mu^c}^\mu \right)^\top \mathbb{P}_{\mathbb{X}_{n, S_\mu}^\mu} \varepsilon_n^\mu \right\|_2 \\
& \leq \left\| \left(\mathbb{X}_{n, S_\mu^c}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \left(\left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \right)^{-1} \right\|_{\text{M},2} \sqrt{\frac{n}{p}} \left\| \lambda_n^\mu \left(\frac{1}{|\widehat{\mu}_{n, S_\mu}^{\text{init}}|} \odot \text{sign}(\mu_{S_\mu}^*) \right) \right\|_2 \\
& \quad + \sqrt{\frac{n}{p}} \left\| \frac{1}{n} \left(\mathbb{X}_{n, S_\mu^c}^\mu \right)^\top \varepsilon_n^\mu \right\|_2 + \left\| \left(\mathbb{X}_{n, S_\mu^c}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \left(\left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \right)^{-1} \right\|_{\text{M},2} \sqrt{\frac{n}{p}} \left\| \frac{1}{n} \left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \varepsilon_n^\mu \right\|_2 \\
& = \mathcal{O}_{\mathbb{P}}(1) \circ_{\mathbb{P}}(1) + \mathcal{O}_{\mathbb{P}}(1) + \mathcal{O}_{\mathbb{P}}(1) \\
& = \mathcal{O}_{\mathbb{P}}(1)
\end{aligned} \tag{8.5}$$

by Lemma 8.2 and (8.3), where

$$\mathbb{P}_{\mathbb{X}_{n, S_\mu}^\mu} = \mathbf{I}_n - \mathbb{X}_{n, S_\mu}^\mu \left(\left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \right)^{-1} \left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top.$$

Furthermore, it is

$$\frac{|\widehat{\mu}_{n, k}^{\text{init}}|}{\lambda_n^\mu} \leq \frac{\|\widehat{\mu}_{n, S_\mu^c}^{\text{init}}\|_2}{\lambda_n^\mu} = \frac{\|\widehat{\mu}_{n, S_\mu^c}^{\text{init}} - \mu_{S_\mu^c}^*\|_2}{\lambda_n^\mu} \leq \frac{\|\widehat{\mu}_n^{\text{init}} - \mu^*\|_2}{\lambda_n^\mu} = \frac{\sqrt{n/p} \|\widehat{\mu}_n^{\text{init}} - \mu^*\|_2}{\sqrt{n/p} \lambda_n^\mu}$$

for all $k \in S^c$. The condition $\sqrt{n/p} \|\widehat{\mu}_n^{\text{init}} - \mu^*\|_2 = \mathcal{O}_{\mathbb{P}}(1)$ together with $n \lambda_n^\mu / p \rightarrow \infty$ implies the convergence

$$\frac{|\widehat{\mu}_{n, k}^{\text{init}}|}{\sqrt{n/p} \lambda_n^\mu} = \frac{1}{n \lambda_n^\mu / p} \mathcal{O}_{\mathbb{P}}(1) = \circ_{\mathbb{P}}(1).$$

Hence it follows by (8.5) that the first condition (7.1) of Lemma 7.1 is satisfied with high probability for a sufficient large sample size n . Furthermore, let

$$\widetilde{\mu}_{n, S_\mu} = \mu_{S_\mu}^* + \left(\frac{1}{n} \left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \right)^{-1} \left(\frac{1}{n} \left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \varepsilon_n^\mu - \lambda_n^\mu \left(\frac{1}{|\widehat{\mu}_{n, S_\mu}^{\text{init}}|} \odot \text{sign}(\mu_{S_\mu}^*) \right) \right).$$

Then we obtain

$$\begin{aligned}
\sqrt{\frac{n}{p}} \left\| \widetilde{\mu}_{n, S_\mu} - \mu_{S_\mu}^* \right\|_2 & \leq \left\| \left(\frac{1}{n} \left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \mathbb{X}_{n, S_\mu}^\mu \right)^{-1} \right\|_{\text{M},2} \left(\sqrt{\frac{n}{p}} \left\| \frac{1}{n} \left(\mathbb{X}_{n, S_\mu}^\mu \right)^\top \varepsilon_n^\mu \right\|_2 \right. \\
& \quad \left. + \sqrt{\frac{n}{p}} \left\| \lambda_n^\mu \left(\frac{1}{|\widehat{\mu}_{n, S_\mu}^{\text{init}}|} \odot \text{sign}(\mu_{S_\mu}^*) \right) \right\|_2 \right) \\
& = \mathcal{O}_{\mathbb{P}}(1) (\mathcal{O}_{\mathbb{P}}(1) + \circ_{\mathbb{P}}(1)) = \mathcal{O}_{\mathbb{P}}(1)
\end{aligned}$$

by (8.3), (8.4) and Lemma 8.2. In particular, this implies

$$\left\| \widetilde{\mu}_{n, S_\mu} - \mu_{S_\mu}^* \right\|_2 = \mathcal{O}_{\mathbb{P}}(\sqrt{p/n}) = \circ_{\mathbb{P}}(1)$$

by Assumption (A13), and hence the second condition, $\text{sign}(\widetilde{\mu}_{n, S_\mu}) = \text{sign}(\mu_{S_\mu}^*)$, of Lemma 7.1 is also satisfied with high probability for large sample sizes n . Sign-consistency of the adaptive LASSO and $\widehat{\mu}_{n, S_\mu}^{\text{AL}} = \widetilde{\mu}_{n, S_\mu}$ is the consequence. \square

References

- Loh, P.-L. and M. J. Wainwright (2017). Support recovery without incoherence: A case for nonconvex regularization. *The Annals of Statistics* 45(6), 2455–2482.
- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, Volume 48. Cambridge University Press.
- Zhou, S., S. van de Geer, and P. Bühlmann (2009). Adaptive lasso for high dimensional regression and Gaussian graphical modeling. *arXiv preprint arXiv:0903.2515*.