

# The de-sparsified or de-biased Lasso

Recap: if  $p < n$  and  $\text{rank}(X) = p$ , then:

$$\hat{\beta}_{\text{OLS},j} = Y^T Z^{(j)} / (X^{(j)})^T Z^{(j)}$$

$$Z^{(j)} = X^{(j)} - X^{(-j)} \hat{\gamma}^{(j)}$$

= OLS residuals from  $X^{(j)}$  vs.  $X^{(-j)} = \{X^{(k)}; k \neq j\}$

$$\hat{\gamma}^{(j)} = \text{argmin}_{\gamma} \|X^{(j)} - X^{(-j)} \gamma\|_2^2$$

idea for high-dimensional setting:  
use the Lasso for the residuals  $Z^{(j)}$

## The de-sparsified estimator

consider

$$\begin{aligned} Z^{(j)} &= X^{(j)} - X^{(-j)} \hat{\gamma}^{(j)} \\ &= \text{Lasso residuals from } X^{(j)} \text{ vs. } X^{(-j)} = \{X^{(k)}; k \neq j\} \\ \hat{\gamma}^{(j)} &= \operatorname{argmin}_{\gamma} \|X^{(j)} - X^{(-j)} \gamma\|_2^2 + \lambda_j \|\gamma\|_1 \end{aligned}$$

build projection of  $Y$  onto  $Z^{(j)}$ :

$$\frac{Y^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} \underbrace{=}_{Y=X\beta^0+\varepsilon} \beta_j^0 + \underbrace{\sum_{k \neq j} \frac{(X^{(k)})^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} \beta_k^0}_{\text{bias}} + \frac{\varepsilon^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}}$$

estimate bias and subtract it:

$$\widehat{\text{bias}} = \sum_{k \neq j} \frac{(X^{(k)})^T X^{(j)}}{(X^{(j)})^T Z^{(j)}} \hat{\beta}_k$$

→ de-sparsified estimator

$$\hat{b}_j = \frac{Y^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} - \sum_{k \neq j} \frac{(X^{(k)})^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} \hat{\beta}_k \quad (j = 1, \dots, p)$$

**not sparse!** Never equal to zero for all  $j = 1, \dots, p$

can also be represented as

$$\hat{b}_j = \hat{\beta}_j + \frac{(Y - X\hat{\beta})^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} \quad \text{“de-biased estimator”}$$

using that

$$\frac{Y^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} = \beta_j^0 + \sum_{k \neq j} \frac{(X^{(k)})^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} \beta_k^0 + \frac{\varepsilon^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}}$$

we obtain

$$\sqrt{n}(\hat{b}_j - \beta_j^0) = \underbrace{\sqrt{n} \sum_{k \neq j} \frac{(X^{(k)})^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} (\beta_k^0 - \hat{\beta}_k)}_{\text{bias term}} + \underbrace{\sqrt{n} \frac{\varepsilon^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}}}_{\text{fluctuation term}}$$

so far, this holds for any  $Z^{(j)}$

assume fixed design  $X$ , e.g. condition on  $X$   
Gaussian error  $\varepsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 I)$

fluctuation term:

$$\sqrt{n} \frac{\varepsilon^T \mathbf{Z}^{(j)}}{(\mathbf{X}^{(j)})^T \mathbf{Z}^{(j)}} = \frac{n^{-1/2} \varepsilon^T \mathbf{Z}^{(j)}}{(\mathbf{X}^{(j)})^T \mathbf{Z}^{(j)} / n} \sim \mathcal{N}\left(\mathbf{0}, \frac{\sigma^2 \|\mathbf{Z}^{(j)}\|_2^2 / n}{|(\mathbf{X}^{(j)})^T \mathbf{Z}^{(j)} / n|^2}\right)$$

bias term: we exploit two things

- ▶  $\|\hat{\beta} - \beta^0\|_1 = O_P(s_0 \sqrt{\log(p)/n})$
- ▶ KKT condition for Lasso (on  $X^{(j)}$  versus  $X^{(-j)}$ ):  
 $|(X^{(k)})^T Z^{(j)}/n| \leq \lambda_j/2$

therefore:

$$\begin{aligned} & \sqrt{n} \sum_{k \neq j} \frac{(X^{(k)})^T Z^{(j)}}{(X^{(j)})^T Z^{(j)}} (\beta_k^0 - \hat{\beta}_k) \\ &= \sqrt{n} \sum_{k \neq j} \frac{(X^{(k)})^T Z^{(j)}/n}{(X^{(j)})^T Z^{(j)}/n} (\beta_k^0 - \hat{\beta}_k) \\ &\leq \sqrt{n} \max_{k \neq j} \left| \frac{(X^{(k)})^T Z^{(j)}/n}{(X^{(j)})^T Z^{(j)}/n} \right| \|\hat{\beta} - \beta^0\|_1 \\ &\leq \sqrt{n} \frac{\lambda_j/2}{(X^{(j)})^T Z^{(j)}/n} O_P(s_0 \sqrt{\log(p)/n}) \\ &= O_P(s_0 \log(p)/\sqrt{n}) = o_P(1) \text{ if } s_0 \ll \frac{\sqrt{n}}{\log(p)} \end{aligned}$$

if  $\lambda_j \asymp \sqrt{\log(p)/n}$  and  $(X^{(j)})^T Z^{(j)}/n \asymp O(1)$

summarizing  $\rightsquigarrow$

*Theorem 10.1 in the notes*

assume:

- ▶  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$
- ▶  $\lambda_j = C_j \sqrt{\log(p)/n}$  and  $\|Z^{(j)}\|_2^2/n \geq L > 0$
- ▶  $s_0 = o(\sqrt{n}/\log(p))$
- ▶  $\|\hat{\beta} - \beta^0\|_1 = O_P(s_0 \sqrt{\log(p)/n})$   
(i.e., compatibility constant  $\phi_0^2$  bounded away from zero)

Then:

$$\sigma^{-1} \sqrt{n} \frac{(X^{(j)})^T Z^{(j)}/n}{\|Z^{(j)}\|_2/\sqrt{n}} (\hat{b}_j - \beta_j^0) \implies \mathcal{N}(0, 1) \quad (j = 1, \dots, p)$$

more precisely:

$$\sigma^{-1} \sqrt{n} \frac{(X^{(j)})^T Z^{(j)} / n}{\|Z^{(j)}\|_2 / \sqrt{n}} (\hat{b}_j - \beta_j^0) = W_j + \Delta_j$$
$$(W_1, \dots, W_p)^T \sim \mathcal{N}_p(\mathbf{0}, \sigma^2 \Omega), \quad \max_{j=1, \dots, p} |\Delta_j| = o_P(1)$$

confidence intervals for  $\beta_j^0$ :

$$\hat{b}_j \pm \hat{\sigma} n^{-1/2} \frac{\|Z^{(j)}\|_2 / \sqrt{n}}{|(X^{(j)})^T Z^{(j)} / n|} \Phi^{-1}(1 - \alpha/2)$$

$$\hat{\sigma}^2 = \|Y - X\hat{\beta}\|_2^2 / n \text{ or } \hat{\sigma}^2 = \|Y - X\hat{\beta}\|_2^2 / (n - \|\hat{\beta}\|_0)$$



can also test

$$H_{0,j} : \beta_j^0 = 0 \text{ versus } H_{A,j} : \beta_j^0 \neq 0$$

can also test group hypothesis: for  $G \subseteq \{1, \dots, p\}$

$$H_{0,G} : \beta_j^0 \equiv 0 \forall j \in G$$

$$H_{A,G} : \exists j \in G \text{ such that } \beta_j^0 \neq 0$$

under  $H_{0,G}$ :

$$\max_{j \in G} \sigma^{-1} \sqrt{n} \frac{|(X^{(j)})^T Z^{(j)} / n|}{\|Z^{(j)}\|_2 / \sqrt{n}} |\hat{b}_j| = \max_{j \in G} |W_j + \Delta_j| \asymp \underbrace{\max_{j \in G} |W_j|}_{\text{distr. simulated}}$$

and plug-in  $\hat{\sigma}$  for  $\sigma$

## Choice of tuning parameters

as usual:  $\hat{\beta} = \hat{\beta}(\hat{\lambda}_{CV})$ ; what is the role of  $\lambda_j$ ?

$$\text{variance} = \sigma^2 n^{-1} \frac{\|Z^{(j)}\|_2^2/n}{|(X^{(j)})^T Z^{(j)}/n|^2} \asymp \sigma^2 / \|Z^{(j)}\|_2^2$$

if  $\lambda_j \searrow$  then  $\|Z^{(j)}\|_2^2 \searrow$ , i.e. large variance

error due to bias estimation is bounded by:

$$|\dots| \leq \sqrt{n} \frac{\lambda_j/2}{|(X^{(j)})^T Z^{(j)}/n|} \|\hat{\beta} - \beta^0\|_1 \propto \lambda_j$$

assuming  $\lambda_j$  is not too small

if  $\lambda_j \searrow$  (but not too small) then bias estimation error  $\searrow$

$\leadsto$  inflate the variance a bit to have low error due to bias estimation: control type I error at the price of slightly decreasing power

## How good is the de-biased Lasso?

asymptotic efficiency:

for the de-biased Lasso to “work” we require

- ▶ sparsity:  $s_0 = o(\sqrt{n}/\log(p))$   
this cannot be beaten in a minimax sense
- ▶ compatibility condition for  $X$

for optimality in terms of the lowest possible asymptotic variance achieving the “Cramer-Rao” lower bound:

- ▶ require **in addition** that  $X^{(j)}$  versus  $X^{(-j)}$  is sparse:  
 $s_j \ll n/\log(p)$

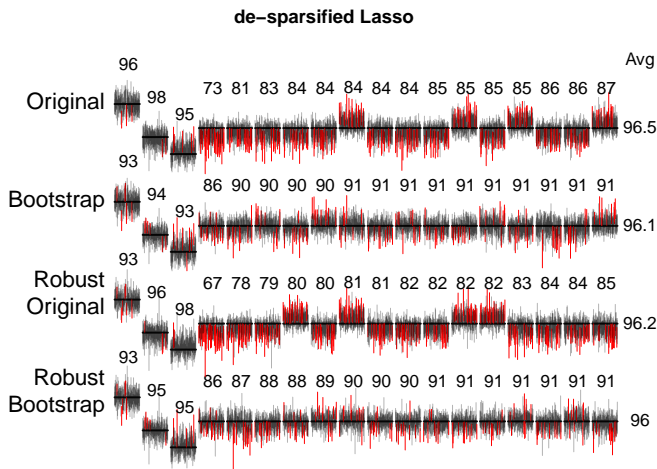
then... skipping details, the de-biased Lasso achieves (see Theorem 10.2):

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j^0) \implies \mathcal{N}\left(0, \underbrace{\sigma^2 \Theta_{jj}}_{\text{Cramer-Rao lower bound}}\right)$$

$\Theta = \Sigma_X^{-1} = \text{Cov}(X)^{-1} \rightsquigarrow$  as for OLS in low dimensions!

# Empirical results

R-software hdi



black: confidence interval covered the true coefficient  
red: confidence interval failed to cover