



**Problem:** High-dimensional (HD) linear regression,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \in \mathbb{R}^n, \quad \boldsymbol{\beta} \in \mathbb{R}^p, \quad p \gg n$$

## Goals

- Prediction
- Estimation of the coefficient vector
- Variable selection
- Statistical inference: confidence intervals and p-values (time permits)
- Closely related problems (time permits)

## Basic assumptions

- Standardized design:  $\|\mathbf{X}_j\|_2^2 = n$  where  $\mathbf{X}_j$  are columns of  $\mathbf{X}$
- Noise level  $\sigma$ , e.g. sufficiently large or  $\varepsilon \sim N(0, \sigma^2 I)$  or sub-Gaussian
- Universal threshold level:  $\lambda_{univ} = \sigma \sqrt{(2/n) \log p}$
- Condition on noise: For a certain target  $\beta^*$ , e.g.  $\beta$  or oracle estimator,

$$\|\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\beta^*)/n\|_\infty \leq \lambda_{univ}$$

- Approximate sparsity for prediction: For a certain  $\bar{\beta}$

$$\|\mathbf{X}\beta - \mathbf{X}\bar{\beta}\|_2^2 / (n\lambda_{univ}^2) + \|\bar{\beta}_{S^c}\|_1 / \lambda_{univ} + |S| \leq s^*$$

- Capped- $\ell_1$  sparsity for coefficient estimation:

$$\|\beta_{S^c}\|_1 / \lambda_{univ} + |S| \leq s^*$$

- Hard sparsity for variable selection:

$$\|\beta\|_0 = s \leq s^*$$

- Complexity relative to sample size:  $s^* \log p \leq a_0 n$ , with (small) fixed  $a_0$
- Replacement of  $\log p$  by  $\log(p/s^*)$  and adjust  $\|\cdot\|_\infty$  in noise condition

## Minimax rates for the estimation of HD objects

- Universal penalty level:  $\lambda_{univ} = \sigma \sqrt{(2/n) \log p}$
- Prediction loss: When  $\log(p/s^*) \asymp \log p$ ,

$$\|\mathbf{X}\hat{\beta} - \mathbf{X}\beta\|_2^2/n \leq C_{\text{pred}} s^* \lambda_{univ}^2$$

- Estimation loss: When  $\log(p/s^*) \asymp \log p$ ,

$$\|\hat{\beta} - \beta\|_q^q \leq C_{\text{est},q} s^* \lambda_{univ}^q$$

- Variable selection: Signal strength for selection/sign consistency,

$$\beta_{\min} = \min_{\beta_j \neq 0} |\beta_j| \geq C_{\text{select}} \lambda_{univ}$$

$$\Rightarrow \mathbb{P}\left\{\text{sgn}(\hat{\beta}) = \text{sgn}(\beta)\right\} \rightarrow 1$$

- Proper conditions on  $\mathbf{X}$  in terms of  $C_{\text{pred}}$ ,  $C_{\text{est},q}$  and  $C_{\text{select}}$

## Noise inflation factor due to model uncertainty:

$$\sqrt{2 \log(p/s^*)}$$

**Lower bounds.** We assume

- $\varepsilon \sim N(0, \sigma^2 I_{n \times n})$
- $\|\mathbf{X}_j\|_2^2 = n \ \forall j$

**Theorem 1.** Let  $\mathbb{P}_G$  correspond to a prior  $G$  under which  $\beta_j$  are iid random variables and  $Z_j | \beta_j \sim N(\beta_j, \sigma^2/n)$ . Then, for any nonnegative loss function  $L_0$ ,

$$\inf_{\hat{\beta} = \hat{\beta}(\mathbf{X}, \mathbf{Y})} \mathbb{E}_G \sum_{j=1}^p L_0(\hat{\beta}_j, \beta_j) \geq p \inf_{t(\cdot)} \mathbb{E}_G L_0(t(Z_1), \beta_1).$$

**Proof:** The l.h.s. equals to

$$\begin{aligned} & \sum_{j=1}^p \mathbb{E}_G \left\{ \inf_{\hat{\beta} = \hat{\beta}(\mathbf{X}, \mathbf{Y})} \mathbb{E}_G \left[ L_0(\hat{\beta}_j, \beta_j) \middle| \mathbf{X}, \mathbf{Y} \right] \right\} \\ & \geq \sum_{j=1}^p \mathbb{E}_G \left\{ \inf_{\hat{\beta} = \hat{\beta}(\mathbf{X}, \mathbf{Y})} \mathbb{E}_G \left[ L_0(\hat{\beta}_j, \beta_j) \middle| \mathbf{X}, \mathbf{Y}, \beta_k, k \neq j \right] \right\} \\ & = \sum_{j=1}^p \mathbb{E}_G \left\{ \inf_{t(\cdot)} \mathbb{E}_G \left[ L_0(t(Z_j), \beta_j) \middle| Z_j \right] \right\} \end{aligned}$$

**Lemma 1 (cf. Donoho-Johnstone, 94).** Let  $\sigma_n = \sigma/\sqrt{n}$ . Suppose  $\mathbb{P}_G\{\beta_1 \neq 0\} = \mathbb{P}_G\{\beta_1 = \mu\} = \pi_0$  with  $\pi_0 < \epsilon_0 \leq 1 - \pi_0$ ,  $\mu = \mu_0\sigma_n$  and  $\mu_0 = \sqrt{2\log(1/\pi_0)} - \sqrt{2\log(1/\epsilon_0)}$ . Then,

$$\inf_{t(\cdot)} \mathbb{E}_G |t(Z_1) - \beta_1|^q \geq (1 - \epsilon_0/2)(1 - \epsilon_1)\pi_0(\sigma_n\mu_0)^q$$

where  $\epsilon_1 = 1 - 1/[\{\epsilon_0/(1 - \pi_0)\}^{1/(q-1)} + 1]^{q-1}$  for  $q > 1$  and  $\epsilon_1 = 0$  for  $q = 1$ .

**Sketch of proof:** Let  $\xi = Z_1/\sigma_n$  and  $\delta_1 = I\{\beta_1 \neq 0\}$ . Then,  $\xi|\delta_1 \sim N(\mu_0\delta_1, 1)$  and  $\delta_1 \sim \text{Bernoulli}(\pi_0)$  under  $\mathbb{E}_G$ , so that

$$\inf_{t(\cdot)} \frac{\mathbb{E}_G |t(Z_1) - \beta_1|^q}{\sigma_n^q \mu_0^q} = \inf_{t(\cdot)} \mathbb{E}_G |t(\xi) - \delta_1|^q = \int \frac{\pi_0 \varphi(\xi - \mu_0)}{\{\omega_1^{1/(q-1)}(\xi) + 1\}^{q-1}} d\xi$$

with  $\omega_1(\xi) = \mathbb{P}_G\{\delta_1 = 1|\xi\}/\mathbb{P}_G\{\delta_1 = 0|\xi\} = \pi_0\varphi(\xi - \mu_0)/\{(1 - \pi_0)\varphi(\xi)\}$ . Let  $a_0 = \sqrt{2\log(1/\epsilon_0)} = \sqrt{2\log(1/\pi_0)} - \mu_0$ . For  $\xi \leq \mu_0 + a_0$ , we have  $\omega_1(\xi) \leq \pi_0\varphi(a_0)/\{(1 - \pi_0)\varphi(\mu_0 + a_0)\} = \epsilon_0/(1 - \pi_0)$ . Thus,

$$\begin{aligned} \min_{t(\cdot)} \mathbb{E}_G |t(\xi) - \delta_1|^q &\geq (1 - \epsilon_1) \int_{-\infty}^{\mu_0 + a_0} \pi_0 \varphi(\xi - \mu_0) d\xi \\ &\geq (1 - \epsilon_1)\pi_0(1 - \epsilon_0/2). \end{aligned}$$

Let  $p \gg s^* \rightarrow \infty$ ,  $\lambda_{mm} = \sigma \sqrt{(2/n) \log(p/s^*)}$  and

$$\Theta_{n,p,s^*} = \{\mathbf{v} \in \mathbb{R}^p : \|\mathbf{v}\|_0 \leq s^*, \|\mathbf{v}\|_\infty \leq \lambda_{mm}\}.$$

**Theorem 2. Estimation lower bounds:** Let  $\epsilon < 1 \leq q < \infty$ . Then,

$$\inf_{\hat{\beta}} \sup_{\beta \in \Theta_{n,p,s^*}} \mathbb{E}_\beta \|\hat{\beta} - \beta\|_q^q \geq (1 + o(1)) s^* \lambda_{mm}^q,$$

$$\inf_{\hat{\beta}} \sup_{\beta \in \Theta_{n,p,s^*}} \mathbb{P}_\beta \left\{ \|\hat{\beta} - \beta\|_q^q \geq (1 - \epsilon) s^* \lambda_{mm}^q \right\} \geq (\epsilon + o(1)) / 3^q.$$

**Proof:** For the  $\mathbb{P}_G$  in Theorem 1 and Lemma 1 and  $s^* = (1 + \epsilon_0) \pi_0 p$ ,

$$\inf_{\hat{\beta} = \hat{\beta}(\mathbf{X}, \mathbf{Y})} \mathbb{E}_G \|\hat{\beta} - \beta\|_q^q \geq (1 + o(1)) p \pi_0 \mu^q \approx s^* \mu^q \approx s^* \lambda_{mm}^q.$$

Let  $N = \|\beta\|_0$  and  $\delta^* = \arg \min_{\mathbf{b}} \mathbb{E}_G [\|\beta - \mathbf{b}\|_q^q | \mathbf{X}, \mathbf{Y}, N \leq s^*]$ . We have

$$(1 + o(1)) s^* \mu^q \leq \mathbb{E}_G [\|\delta^* - \beta\|_q^q | N \leq s^*] + \mathbb{E}_G \|\delta^* - \beta\|_q^q I\{N > s^*\}.$$

As  $\|\beta\|_q^q \leq s^* \mu^q$  when  $N \leq s^*$ ,  $\|\delta^*\|_q^q \leq s^* \mu^q$  by the convexity of  $\|\cdot\|_q^q$ . Thus,

$$\mathbb{E}_G \|\delta^* - \beta\|_q^q I\{N > s^*\} \leq 2^{q-1} \mu^q \mathbb{E}_G (s^* + N) I\{N > (1 + \epsilon_0) \mathbb{E}_G N\}.$$

As  $N \sim \text{binomial}(p, \pi_0)$ , the r.h.s. is  $o(p \pi_0 \mu^q) = o(s^* \mu^q)$ . Thus,

$$(\text{minimax } \ell_q^q \text{ risk})(\Theta_{n,p,s^*}) \geq \mathbb{E}_G [\|\delta^* - \beta\|_q^q | N \leq s^*] \geq (1 + o(1)) s^* \lambda_{mm}^q.$$

**Proof of Theorem 2 (continuation):** Let  $c = (1 - \epsilon)^{1/q}$ ,

$$\begin{aligned} L(\hat{\beta}, \beta) &= I\{\|\hat{\beta} - \beta\|_q \geq c(s^*)^{1/q} \lambda_{mm}\}, \\ \tilde{\beta} &= \hat{\beta} I\{\|\hat{\beta}\|_q \leq (1 + c)(s^*)^{1/q} \lambda_{mm}\}. \end{aligned}$$

We have  $\mu \leq \lambda_{mm}$  and  $\|\beta\|_q^q \leq s^* \mu^q$  for  $N \leq s^*$ . It follows that

$$\|\tilde{\beta} - \beta\|_q^q \leq c^q s^* \lambda_{mm}^q \{1 - L(\hat{\beta}, \beta)\} + (2 + c)^q s^* \lambda_{mm}^q L(\hat{\beta}, \beta).$$

As  $\|\beta\|_q^q = N\mu^q \leq N\lambda_{mm}^q$ ,

$$\begin{aligned} \mathbb{E}_G \|\tilde{\beta} - \beta\|_q^q &\leq c^q s^* \lambda_{mm}^q + \{(2 + c)^q - c^q\} s^* \lambda_{mm}^q \max_{\beta \in \Theta} \mathbb{E}_\beta L(\hat{\beta}, \beta) \\ &\quad + 2^{q-1} s^* \lambda_{mm}^q \mathbb{E}_G \left( N/s^* + (1 + c)^q \right) I\{N > (1 + \epsilon_0)k\}. \end{aligned}$$

As  $\mathbb{E}_G (N/s^* + (1 + c)^q) I\{N > (1 + \epsilon_0)k\} \rightarrow 0$ , Lemma 1 yields

$$\begin{aligned} \{(2 + c)^q - c^q\} \max_{\beta \in \Theta} \mathbb{E}_\beta L(\hat{\beta}, \beta) &\geq \mathbb{E}_G \|\tilde{\beta} - \beta\|_q^q / (s^* \lambda_{mm}^q) - c^q + o(1) \\ &\geq 1 - c^q + o(1). \end{aligned}$$

This gives  $\max_{\beta \in \Theta} \mathbb{E}_\beta L(\hat{\beta}, \beta) \geq (1 - c^q + o(1))/3^q$ .

See Ye-Z (2010, JMLR).



Let  $\bar{\Sigma} = \mathbf{X}^T \mathbf{X} / n$  and  $\phi_{\pm}(m; \Sigma)$  be the upper and lower sparse eigenvalues,

$$\begin{aligned}\phi_+(m; \Sigma) &= \max_{|A|=m, \|\mathbf{u}\|_2=1} \mathbf{u}^T \Sigma_{A,A} \mathbf{u}, \\ \phi_-(m; \Sigma) &= \min_{|A|=m, \|\mathbf{u}\|_2=1} \mathbf{u}^T \Sigma_{A,A} \mathbf{u}.\end{aligned}$$

**Theorem 3. Prediction lower bounds:** *Under the assumptions of Theorem 2,*

$$\inf_{\hat{\beta}} \sup_{\beta \in \Theta_{n,p,s^*}} \mathbb{E}_{\beta} \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta\|_2^2 / n \geq (1 + o(1)) s^* \lambda_{mm}^2 \phi_-(2s^*; \bar{\Sigma}) / 4,$$

and

$$\inf_{\hat{\beta}} \sup_{\beta \in \Theta_{n,p,s^*}} \mathbb{P}_{\beta} \left\{ \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta\|_2^2 / n \geq \frac{(1 - \epsilon)^2 s^* \lambda_{mm}^2}{4 / \phi_-(2s^*; \bar{\Sigma})} \right\} \geq \epsilon / 4 + o(1).$$

**Remark:** When  $\mathbf{X}$  has iid  $N(0, \Sigma)$  rows,

$$\frac{\phi_-(m; \bar{\Sigma})}{\phi_-(m; \Sigma)} \geq 1 - \epsilon(a), \quad \frac{\phi_+(m; \bar{\Sigma})}{\phi_+(m; \Sigma)} \leq 1 + \epsilon(a),$$

with  $\epsilon(0+) = 0$  not dependent on  $\Sigma$ , where  $a = m(\log p)/n$ . Thus, conditions on  $\phi_{\pm}(m, \bar{\Sigma})$ , e.g. the restricted isometry property (RIP) and sparse Reisz condition (SRC), are of  $\ell_2$ -type.

**Theorem 4. Selection lower bound:** Let  $\mathcal{S}(s, \beta_*)$  be the set of all coefficient vectors  $\beta$  satisfying  $\|\beta\|_0 = s$  and  $\min_{\beta_j \neq 0} |\beta_j| \geq \beta_*$ , with  $s < p - 1$ . Then,

$$\inf_{\hat{\beta}} \sup_{\beta \in \mathcal{S}(s, \beta_*)} \mathbb{P}\{\text{supp}(\hat{\beta}) \neq \text{supp}(\beta)\} \geq 1 - \frac{2\beta_*^2 + (\sigma^2/n) \log 2}{(\sigma^2/n) \log(p - s)}.$$

**Proof:** Let  $m = p - s + 1$ . For  $j = 1, \dots, m$  let  $\mathbb{P}_j$  be the probability under which  $\beta_j = \beta_{p-s+2} = \dots = \beta_p = \beta_*$  and  $\beta_k = 0$  for all  $k \neq j$  and  $k \leq p - s + 1$ . Clearly  $\beta \in \mathcal{S}(s, \beta_*)$  under  $\mathbb{P}_j$ . Since the Kullback-Leibler information between  $N(\mathbf{u}, \sigma^2 \mathbf{I}_{n \times n})$  and  $N(\mathbf{v}, \sigma^2 \mathbf{I}_{n \times n})$  is  $\|\mathbf{u} - \mathbf{v}\|_2^2 / (2\sigma^2)$ ,

$$K(\mathbb{P}_j, \mathbb{P}_k) = \mathbb{E}\|\mathbf{X}_j \beta_* - \mathbf{X}_k \beta_*\|_2^2 / (2\sigma^2) \leq 2n\beta_*^2 / \sigma^2.$$

Let  $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$  and  $\delta(\mathbf{Z}) = \min\{j : \hat{\beta}_j \neq 0\}$ . It follows from Fano's lemma that

$$\min_j \mathbb{P}_j\{\text{supp}(\hat{\beta}) = \text{supp}(\beta)\} \leq \frac{2\beta_*^2 / (\sigma^2/n) + \log 2}{\log(m - 1)}.$$

See Wainwright (2007, Tech Rep) and Z (2007, Tech Rep; 2010, AOS)

**Lemma 2 (Fano 61).** *Let  $\mathbb{P}_1, \dots, \mathbb{P}_m$  be probability measures for data  $\mathbf{Z}$ ,  $\delta(\mathbf{z})$  a  $\{1, \dots, m\}$ -valued statistic, and  $K(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{\mathbb{P}} \log(d\mathbb{P}/d\mathbb{Q})$  the Kullback-Leibler information. Then,*

$$\frac{1}{m} \sum_{j=1}^m \mathbb{P}_j \left\{ \delta(\mathbf{Z}) = j \right\} \leq \frac{1}{m^2} \sum_{1 \leq j, k \leq m} \frac{K(\mathbb{P}_j, \mathbb{P}_k) + \log 2}{\log(m-1)}.$$

**Penalized LSE:** Global or suitable local minimizers of

$$\mathcal{L}_\lambda(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/(2n) + \|p_\lambda(\beta)\|_1$$

where  $\|p_\lambda(\beta)\|_1 = \sum_{j=1}^p p_\lambda(\beta_j)$

**Basic properties of penalty functions**

- Zero baseline penalty:  $p_\lambda(0) = 0$
- Symmetry:  $p_\lambda(-t) = p_\lambda(t)$
- Monotonicity:  $p_\lambda(x) \leq p_\lambda(y)$  for all  $0 \leq x < y < \infty$
- Sparsity:  $p_\lambda(0+) + p'_\lambda(0+) > 0$
- Subadditivity:  $p_\lambda(x+y) \leq p_\lambda(x) + p_\lambda(y)$  for all  $0 \leq x \leq y < \infty$
- One-sided differentiability:  $p'_\lambda(t\pm)$  are defined for all real  $t$

Convention:  $p'_\lambda(t) = x$  means

$$\min \{p'_\lambda(t+), p'_\lambda(t-)\} \leq x \leq \max \{p'_\lambda(t+), p'_\lambda(t-)\}$$

## Important quantities

- Maximum concavity:  $\kappa(p_\lambda) = \sup_{t, \epsilon \neq 0} \{p'_\lambda(t - \epsilon) - p'_\lambda(t)\} / \epsilon$
- Bias threshold:  $a_\lambda = \inf \{t > 0 : p_\lambda(x) = p_\lambda(t) \ \forall |x| \geq t\}$ 
  - For  $p = 1$ , PLSE = LSE when  $|\hat{\beta}^{(lse)}| \geq a_\lambda$
- Global threshold level:  $\lambda_* = \lambda_*(p_\lambda) = \inf_{t > 0} \{t/2 + p_\lambda(t)/t\}$ 
  - For  $p = 1$ ,  $\hat{\beta} = 0$  attains global minimum iff  $|\hat{\beta}^{(lse)}| \leq \lambda_*$

## Standardization of penalty functions (scaling key quantities in $\lambda$ )

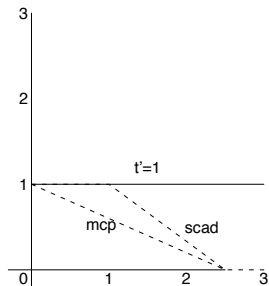
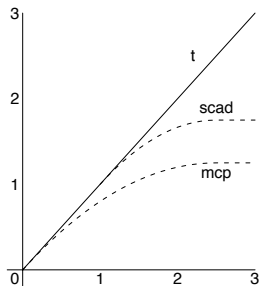
- Standardized local threshold level:  $p'_\lambda(0+) = \lambda$  if  $\max |p'_\lambda(0)| < \infty$
- Standardized global threshold level:  $\lambda_*(p_\lambda) = \lambda$  if  $\max |p'_\lambda(0)| = \infty$
- Standardized bias threshold:  $a = a_\lambda / \lambda$
- Standardized local threshold level:  $a^* = p'_\lambda(0+) / \lambda$
- Standardized global threshold level:  $a_* = \lambda_* / \lambda$
- Standardized maximum penalty:  $\gamma^* = \|p_\lambda(t)\|_\infty / \lambda^2$

The above quantities are indeed constants when  $p_\lambda(t) = \lambda^2 p_1(t/\lambda)$ , as in the following examples. See Fan-Li (2001, JASA), Z (2010, AOS), and Z-Zhang (2012, StatSci) for general discussion of penalty functions

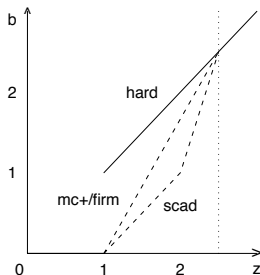
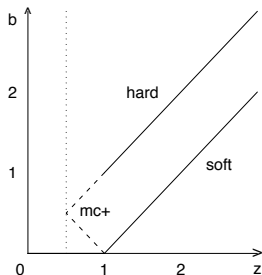
Penalty	$p_1(t)$	Bias $a$	Maximum Concavity $\kappa(p_\lambda)$	Thresh Levels $a^*/a_*$	Max Penalty $\gamma^*$
$\ell_0$	$2^{-1}I\{t \neq 0\}$	0	$\infty$	$\infty/1$	1/2
Bridge/ $\ell_\alpha$	$C_\alpha t ^\alpha$	$\infty$	$\infty$	$\infty/1$	$\infty$
Lasso/ $\ell_1$	$ t $	$\infty$	0	1/1	$\infty$
SCAD	$\int_0^{ t } \left(1 - \frac{(x-1)_+}{a-1}\right) dx$	$a \geq 2$	$1/(a-1)$	1/1	$a/2 + 1/2$
MCP	$\int_0^{ t } \left(1 - x/a\right)_+ dx$	$a \geq 1$	$1/a$	1/1	$a/2$
Capped- $\ell_1$	$\min( t , a)$	$a \geq 1/2$	$\infty$	1/1	$a$

**Table:** Specific penalty functions:  $C_\alpha = \{2(1-\alpha)\}^{1-\alpha}/(2-\alpha)^{2-\alpha}, 0 < \alpha < 1$

The Lasso, SCAD and MCP:  $p_1(t)$  and  $p'_1(t)$  for  $t \geq 0$



The Lasso, SCAD and MC+ estimators for  $p = 1$   
(vertical line indicates the bias threshold  $a$ )





**The Lasso:** With the  $\ell_1$  penalty  $p_\lambda(t) = \lambda|t|$ ,

$$\hat{\beta} = \hat{\beta}(\lambda) = \hat{\beta}^{(\ell_1)} = \arg \min_{\beta} \left\{ \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 / (2n) + \lambda \|\beta\|_1 \right\}$$

## Motivation

- A convex minimization problem
  - Easier to compute than the traditional subset selection ( $\ell_0$  penalty)
  - Do not have to deal with (potentially bad) local minima
- Rate optimality in prediction, coefficient estimation and variable selection in the rate minimax sense (adaptive in complexity  $s^*$ )

## Some drawbacks

- Biasedness
  - Does not take advantage of signal strength well
  - Requires more restrictive  $\ell_\infty$ -type conditions for selection consistency
- Performance unclear for individual  $\hat{\beta}_j$

## Convex minimization

- With convex loss  $\mathcal{L}_0$  and convex regularizer  $R_\lambda$ , let

$$\mathcal{L}_\lambda(\beta) = \mathcal{L}_0(\beta) + R_\lambda(\beta)$$

- The KKT condition:**  $\hat{\beta}$  is a minimizer of  $\mathcal{L}_\lambda(\beta)$  iff

$$\partial\{\mathcal{L}_0(\hat{\beta}) + R_\lambda(\hat{\beta})\} \ni \mathbf{0}$$

- Sub-differential: For any convex function  $f$ ,

$$\partial f(x) = \left\{ v : f(x+h) - f(x) - \langle v, h \rangle \geq 0 \quad \forall h \right\}$$

- Bregman divergence: For convex  $f$  and  $\nabla f(y) \in \partial f(y)$

$$D_f(x, y) = f(x) - f(y) - \langle \nabla f(y), x - y \rangle \geq 0$$

- Symmetric Bregman divergence:

$$D_f^s(x, y) = D_f(x, y) + D_f(y, x) = \langle \nabla f(x) - \nabla f(y), x - y \rangle \geq 0$$

**Proposition 1. KKT for the Lasso:** A vector  $\hat{\beta} \in \mathbb{R}^p$  is a Lasso solution iff

$$g_j = \mathbf{X}_j^T (\mathbf{Y} - \mathbf{X}\hat{\beta})/n = p'_\lambda(\hat{\beta}_j) \begin{cases} = \lambda \operatorname{sgn}(\hat{\beta}_j) & \hat{\beta}_j \neq 0 \\ \in \lambda [-1, 1] & \forall j. \end{cases}$$

Moreover,  $\hat{\beta}$  is the unique solution of the Lasso iff  $\|\mathbf{X}\mathbf{u}\|_2^2 > 0$  for all  $\mathbf{u} \neq 0$  satisfying the conditions  $\operatorname{supp}(\mathbf{u}) \subseteq \{j : |g_j| = \lambda\}$  and  $g_j u_j \geq 0$  for all  $j$  with  $\hat{\beta}_j = 0$ , e.g.  $\operatorname{rank}(\mathbf{X}_{\mathcal{S}_1}) = |\mathcal{S}_1|$  with  $\mathcal{S}_1 = \{j : |g_j| = \lambda\}$ .

**Proof:**

- Recall that  $\mathcal{L}_\lambda(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/(2n) + \lambda\|\beta\|_1$
- Joint differentiability of  $\mathcal{L}_\lambda(\hat{\beta} + \mathbf{u})$  in  $\mathbf{u}$  for small  $\mathbf{u}$  with fixed  $\operatorname{sgn}(\mathbf{u})$
- $\partial p_\lambda(t) = \lambda \operatorname{sgn}(t)$  for  $t \neq 0$
- $\partial p_\lambda(t) = \lambda[-1, 1]$  for  $t = 0$
- Uniqueness: Check  $(\partial/\partial t)^2 \mathcal{L}_\lambda(\hat{\beta} + t\mathbf{u}) = \|\mathbf{X}\mathbf{u}\|_2^2/n$  for small  $t > 0$

## Selection consistency of the Lasso: Define

- Gram matrix:  $\bar{\Sigma} = \mathbf{X}^T \mathbf{X} / n$
- Support:  $\mathcal{S} = \text{supp}(\beta)$
- Orthogonal projection to the column space of  $\mathbf{X}_A$ :  $\mathbf{P}_A$
- Oracle estimator  $\hat{\beta}^o$ :  $\hat{\beta}_{\mathcal{S}}^o = \mathbf{X}_{\mathcal{S}}^{\dagger} \mathbf{Y}$ ,  $\hat{\beta}_{\mathcal{S}^c}^o = \mathbf{0}$
- Neighborhood stability/exact recovery/irrepresentability coefficient (Meinshausen-Bühlmann, 06; Tropp, 06; Zhao-Yu, 06):

$$\kappa_{\text{select}} = \|\bar{\Sigma}_{\mathcal{S}^c, \mathcal{S}} \bar{\Sigma}_{\mathcal{S}, \mathcal{S}}^{-1} \text{sgn}(\beta_{\mathcal{S}})\|_{\infty}$$

**Theorem 5 (i) Sufficient conditions:** Suppose  $\text{rank}(\mathbf{X}_{\mathcal{S}}) = |\mathcal{S}|$  and

$$\begin{aligned} \max_{j \notin \mathcal{S}} |\mathbf{X}_j^T \mathbf{P}_{\mathcal{S}}^{\perp} (\mathbf{Y} - \mathbf{X}\beta) / n| &\leq (1 - \kappa_{\text{select}}) \lambda, \\ \text{sgn}(\beta_{\mathcal{S}}) &= \text{sgn}(\hat{\beta}_{\mathcal{S}}^o - \lambda \bar{\Sigma}_{\mathcal{S}, \mathcal{S}}^{-1} \text{sgn}(\beta_{\mathcal{S}})). \end{aligned}$$

Then, a Lasso solution is sign consistent,

$$\text{sgn}(\hat{\beta}) = \text{sgn}(\beta).$$

Moreover, the solution is unique if the first condition holds strictly.

**Uniform signal strength/ $\beta_{\min}$  condition:**

$$\min_{j \in \mathcal{S}} \left( |\beta_j| - \text{sgn}(\beta_j) \lambda (\bar{\Sigma}_{\mathcal{S}, \mathcal{S}}^{-1} \text{sgn}(\beta_{\mathcal{S}}))_j - \lambda' (\bar{\Sigma}_{\mathcal{S}, \mathcal{S}}^{-1})_{j,j}^{1/2} \right) > 0,$$

**Theorem 5 (ii) Necessary conditions:** Suppose  $\text{rank}(\mathbf{X}_{\mathcal{S}}) = |\mathcal{S}|$  and  $\mathbf{Y} - \mathbf{X}\beta$  is a symmetric continuous random vector. If either  $\kappa_{\text{select}} \geq 1$  or the  $\beta_{\min}$  condition fails to hold for  $\lambda' = 0$ , then

$$\mathbb{P}\{\text{sgn}(\hat{\beta}) = \text{sgn}(\beta)\} \leq 1/2.$$

If  $\kappa_{\text{select}} < 1$  and the  $\beta_{\min}$  condition fails to hold for a certain  $\lambda'$ , then

$$\begin{aligned} & \mathbb{P}\{\text{supp}(\hat{\beta}) = \text{supp}(\beta)\} \\ & \leq \max_{j \in \mathcal{S}} \mathbb{P}\left\{\text{sgn}(\beta_j)(\hat{\beta}_j^{\circ} - \beta_j)/(\bar{\Sigma}_{\mathcal{S}, \mathcal{S}}^{-1})_{j,j}^{1/2} > -\lambda'\right\}. \end{aligned}$$

**Corollary:** Suppose  $\text{rank}(\mathbf{X}_{\mathcal{S}}) = |\mathcal{S}|$ ,  $\varepsilon = \mathbf{Y} - \mathbf{X}\beta$  is sub-Gaussian,  $\max_j \|\mathbf{P}_{\mathcal{S}}^{\perp} \mathbf{X}_j\|_2^2 \leq n$ ,  $\lambda \geq (1 - \kappa_{\text{select}})_+^{-1} \sigma \sqrt{(2/n) \log((p - |\mathcal{S}|)/\epsilon)}$ , and the  $\beta_{\min}$  condition holds with  $\lambda' = \sigma \sqrt{(2/n) \log(|\mathcal{S}|/\epsilon)}$ . Then,

$$\mathbb{P}\{\text{sgn}(\hat{\beta}) = \text{sgn}(\beta)\} \geq 1 - 3\epsilon.$$

## Sufficient and nearly necessary conditions for Lasso selection consistency:

- Strong irrepresentable condition:  $\kappa_{\text{select}} = \|\bar{\Sigma}_{S^c, S} \bar{\Sigma}_{S, S}^{-1} \text{sgn}(\beta_S)\|_\infty < 1$
- Penalty level:  $\lambda \geq \lambda_{\text{univ}} / (1 - \kappa_{\text{select}})$
- Signal strength:  $\beta_{\min} \geq \lambda \|\bar{\Sigma}_{S, S}^{-1} \text{sgn}(\beta_S)\|_\infty + \lambda' \|\text{diag}(\bar{\Sigma}_{S, S}^{-1})\|_\infty^{1/2}$

The signal strength requirement is of optimal rate, but the  $\ell_\infty$ -type conditions do not scale well in  $\|\beta\|_0 = \|\text{sgn}(\beta_S)\|_2^2$

## Proof of Theorem 5:

- By KKT,  $\hat{\beta}$  is a Lasso solution with  $\text{sgn}(\hat{\beta}) = \text{sgn}(\beta)$  iff

$$\begin{aligned} \mathbf{X}_j^T (\mathbf{Y} - \mathbf{X}_S \hat{\beta}_S) / n &= \lambda \text{sgn}(\beta_j), \quad j \in S, \\ |\mathbf{X}_j^T (\mathbf{Y} - \mathbf{X}_S \hat{\beta}_S) / n| &\leq \lambda \text{ and } \hat{\beta}_j = 0, \quad j \notin S. \end{aligned}$$

- The first part yields  $\hat{\beta}_S = \hat{\beta}_S^o - \lambda \bar{\Sigma}_{S, S}^{-1} \text{sgn}(\beta_S)$
- The KKT (i.e. nasc) for  $\text{sgn}(\hat{\beta}) = \text{sgn}(\beta)$  can be written as

$$\begin{aligned} \text{sgn}(\hat{\beta}_S^o - \lambda \bar{\Sigma}_{S, S}^{-1} \text{sgn}(\beta_S)) &= \text{sgn}(\beta_S), \\ \|\mathbf{X}_{S^c}^T \mathbf{P}_S^\perp \epsilon / n + \lambda \bar{\Sigma}_{S^c, S} \bar{\Sigma}_{S, S}^{-1} \text{sgn}(\beta_S)\|_\infty &\leq \lambda. \end{aligned}$$

## Prediction and coefficient estimation

- Noise condition:  $\|\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta^*)/n\|_\infty \leq \eta\lambda$  with  $\eta \leq 1$ , e.g.  $\lambda \geq \lambda_{univ}/\eta$
- Let  $\mathbf{w}_S = \mathbf{X}_S^T(\mathbf{Y} - \mathbf{X}\beta^*)/(\lambda n) - \text{sgn}(\bar{\beta}_S)$ , with  $\|\mathbf{w}_S\|_\infty \leq 1 + \eta$ , and

$$\psi(\mathbf{u}) = \mathbf{w}_S^T \mathbf{u}_S - (1 - \eta)\|\mathbf{u}_{S^c}\|_1.$$

**Lemma 3. A basic inequality:** For any  $S$  and  $\bar{\beta}$ ,

$$\begin{aligned} & \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2/n + \|\mathbf{X}\hat{\beta} - \mathbf{X}\bar{\beta}\|_2^2/n \\ & \leq \|\mathbf{X}\bar{\beta} - \mathbf{X}\beta^*\|_2^2/n + 4\lambda\|\bar{\beta}_{S^c}\|_1 + 2\lambda\psi(\hat{\beta} - \bar{\beta}). \end{aligned}$$

**Proof:** Let  $\mathbf{h}^* = \hat{\beta} - \beta^*$  and  $\bar{\mathbf{h}} = \hat{\beta} - \bar{\beta}$ . We have

$$\begin{aligned} \text{KKT: } \lambda \partial\|\hat{\beta}\|_1 &= -\partial\mathcal{L}_0(\hat{\beta}) = \mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\beta})/n, \\ \bar{\mathbf{h}}^T \partial\|\hat{\beta}\|_1 &= \bar{\mathbf{h}}^T \mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta^*)/(\lambda n) - (\mathbf{X}\bar{\mathbf{h}})^T(\mathbf{X}\mathbf{h}^*)/(\lambda n), \\ \bar{\mathbf{h}}^T \partial\|\hat{\beta}\|_1 &\geq \|\bar{\mathbf{h}}_{S^c}\|_1 - 2\|\bar{\beta}_{S^c}\|_1 + \bar{\mathbf{h}}_S^T \partial\|\hat{\beta}_S\|_1, \\ \bar{\mathbf{h}}_S^T \partial\|\hat{\beta}_S\|_1 &\geq \bar{\mathbf{h}}_S^T \text{sgn}(\bar{\beta}_S) = \bar{\mathbf{h}}_S^T \mathbf{X}_S^T(\mathbf{Y} - \mathbf{X}\beta^*)/(\lambda n) - \mathbf{w}_S^T \bar{\mathbf{h}}_S, \\ \bar{\mathbf{h}}^T \partial\|\hat{\beta}\|_1 &\geq (1 - \eta)\|\bar{\mathbf{h}}_{S^c}\|_1 - 2\|\bar{\beta}_{S^c}\|_1 + \bar{\mathbf{h}}^T \mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta^*)/(\lambda n) - \mathbf{w}_S^T \bar{\mathbf{h}}_S, \\ (\mathbf{X}\bar{\mathbf{h}})^T(\mathbf{X}\mathbf{h}^*)/n &\leq -2\lambda\|\bar{\beta}_{S^c}\|_1 + \lambda\psi(\bar{\mathbf{h}}), \\ 2ab &= a^2 + b^2 - (a - b)^2 \end{aligned}$$

## Prediction error bounds

- Let  $\psi(\mathbf{u}) = \mathbf{w}_S^T \mathbf{u}_S - (1 - \eta) \|\mathbf{u}_{S^c}\|_1$  and define

$$\bar{C}_{\text{pred}}^{(\ell_1)}(\mathcal{S}, \bar{\Delta}) = \sup \left\{ \frac{[\psi(\mathbf{u}) + \{\psi^2(\mathbf{u}) + 2\bar{\Delta}\}^{1/2}]^2}{4(1 \vee |\mathcal{S}|)} : \mathbf{u}^T \bar{\Sigma} \mathbf{u} = 1 \right\}$$

- Basic inequality: With  $\bar{\Delta} = \|\mathbf{X}\bar{\beta} - \mathbf{X}\beta^*\|_2^2/(\lambda^2 n) + 4\|\bar{\beta}_{S^c}\|_1/\lambda$ ,

$$\|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2/n \leq \lambda^2 \bar{\Delta} + 2\lambda \psi(\hat{\beta} - \bar{\beta}) - \|\mathbf{X}\hat{\beta} - \mathbf{X}\bar{\beta}\|_2^2/n$$

**Theorem 6.** For any  $\bar{\beta}$  and  $\mathcal{S}$ ,

$$\begin{cases} \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2/n + \|\mathbf{X}\hat{\beta} - \mathbf{X}\bar{\beta}\|_2^2/n \leq \lambda^2 \bar{\Delta}, & \mathcal{S} = \emptyset, \\ \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2/n \leq \lambda^2 \bar{\Delta} + \lambda^2 |\mathcal{S}| \bar{C}_{\text{pred}}^{(\ell_1)}(\mathcal{S}, 0), & \mathcal{S} \neq \emptyset, \\ \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2/n \leq \lambda^2 (1 \vee |\mathcal{S}|) \bar{C}_{\text{pred}}^{(\ell_1)}(\mathcal{S}, \bar{\Delta}), & \bar{\beta} = \beta^*. \end{cases}$$

**Proof:** With  $\bar{\mathbf{h}}/\lambda = t\mathbf{u}$  at the largest possible  $t$  and  $c = 0$  or  $c = \bar{\Delta}$ ,  
 $(\lambda^2 c + 2\lambda \psi(\bar{\mathbf{h}}) - \bar{\mathbf{h}}^T \bar{\Sigma} \bar{\mathbf{h}})/\lambda^2 = c + 2t\psi(\mathbf{u}) - t^2 \leq (1 \vee |\mathcal{S}|) \bar{C}_{\text{pred}}^{(\ell_1)}(\mathcal{S}, c).$



## Remarks

- Approximate sparsity: For certain  $\beta^*$  satisfying the noise condition and  $\bar{\beta}$

$$\|\mathbf{X}\beta^* - \mathbf{X}\bar{\beta}\|_2^2 / (n\lambda_{univ}^2) + \|\bar{\beta}_{Sc}\|_1 / \lambda_{univ} + |S| \leq s^*$$

- Prediction error bound of order  $\lambda^2 s^*$ :

$$\|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2 / n \leq \|\mathbf{X}\bar{\beta} - \mathbf{X}\beta^*\|_2^2 / n + 4\lambda\|\bar{\beta}_{Sc}\|_1 + \lambda^2|S|\bar{C}_{pred}^{(\ell_1)}(\mathcal{S}, 0)$$

- Regularity condition on  $\mathbf{X}$ : With  $\mathbf{w}_S = \mathbf{X}_S^T(\mathbf{Y} - \mathbf{X}\beta^*) / (\lambda n) - \text{sgn}(\bar{\beta}_S)$ ,

$$\bar{C}_{pred}^{(\ell_1)}(\mathcal{S}, 0) = \sup_{\mathbf{u} \neq 0} \left\{ \frac{[\mathbf{w}_S^T \mathbf{u}_S - (1 - \eta)\|\mathbf{u}_{Sc}\|_1]_+^2}{(1 \vee |S|)\mathbf{u}^T \bar{\Sigma} \mathbf{u}} \right\} = O(1)$$

- Compatibility condition: As  $\|\mathbf{w}_S\|_\infty \leq 1 + \eta$ , a sufficient condition is

$$\sup \left\{ \frac{(1 + \eta)^2 \|\mathbf{u}_S\|_1^2}{(1 \vee |S|)\mathbf{u}^T \bar{\Sigma} \mathbf{u}} : (1 - \eta)\|\mathbf{u}_{Sc}\|_1 < (1 + \eta)\|\mathbf{u}_S\|_1 \right\} = O(1)$$

- Restricted eigenvalue (RE) condition: Another sufficient condition is

$$\sup \left\{ \frac{(1 + \eta)^2 \|\mathbf{u}_S\|_2^2}{\mathbf{u}^T \bar{\Sigma} \mathbf{u}} : (1 - \eta)\|\mathbf{u}_{Sc}\|_1 < (1 + \eta)\|\mathbf{u}_S\|_1 \right\} = O(1)$$

- All the above conditions are of  $\ell_2$ -type

## Error bounds for coefficient estimation

- Let  $\psi(\mathbf{u}) = \mathbf{w}_S^T \mathbf{u}_S - (1 - \eta) \|\mathbf{u}_{S^c}\|_1$  and define

$$\overline{C}_{\text{est},q}^{(\ell_1)}(\mathcal{S}, \overline{\Delta}) = \sup \left\{ \frac{\|\mathbf{u}\|_q [\psi(\mathbf{u}) + \{\psi^2(\mathbf{u}) + 2\overline{\Delta}\}^{1/2}]}{2(1 \vee |\mathcal{S}|)^{1/(q \wedge 2)}} : \mathbf{u}^T \overline{\Sigma} \mathbf{u} = 1 \right\}$$

- Recall the basic inequality: With  $\overline{\Delta} = \|\mathbf{X}\overline{\beta} - \mathbf{X}\beta^*\|_2^2/(\lambda^2 n) + 4\|\overline{\beta}_{S^c}\|_1/\lambda$ ,

$$\|\mathbf{X}\widehat{\beta} - \mathbf{X}\beta^*\|_2^2/n \leq \lambda^2 \overline{\Delta} + 2\lambda \psi(\widehat{\beta} - \overline{\beta}) - \|\mathbf{X}\widehat{\beta} - \mathbf{X}\overline{\beta}\|_2^2/n$$

**Theorem 7.** For any  $\overline{\beta}$  and  $\mathcal{S}$ ,

$$\begin{cases} \|\widehat{\beta} - \overline{\beta}\|_q \leq 2\lambda |\mathcal{S}|^{1/(q \wedge 2)} \overline{C}_{\text{est},q}^{(\ell_1)}(\mathcal{S}, \overline{\Delta}/2), & \forall \overline{\beta} \\ \|\widehat{\beta} - \beta^*\|_q \leq \lambda |\mathcal{S}|^{1/(q \wedge 2)} \overline{C}_{\text{est},q}^{(\ell_1)}(\mathcal{S}, \overline{\Delta}), & \overline{\beta} = \beta^*. \end{cases}$$

**Proof:** With  $\overline{\mathbf{h}}/\lambda = t\mathbf{u}$  at the largest  $t$ ,  $\|\overline{\mathbf{h}}\|_q = \lambda t \|\mathbf{u}\|_q$ .

## Remarks

- Approximate sparsity: For certain  $\beta^*$  satisfying the noise condition and  $\bar{\beta}$

$$\|\mathbf{X}\beta^* - \mathbf{X}\bar{\beta}\|_2^2 / (n\lambda_{univ}^2) + \|\bar{\beta}_{Sc}\|_1 / \lambda_{univ} + |\mathcal{S}| \leq s^*$$

- The Lasso has prediction error  $\|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2 / n \lesssim s^* \lambda^2$
- The Lasso actually estimates  $\bar{\beta}$  with  $\|\hat{\beta} - \bar{\beta}\|_q^q \lesssim s^* \lambda^q$ ,  $1 \leq q \leq 2$
- $\ell_q$  error bounds for  $2 < q \leq \infty$  based on the (sign-restricted) CIF
- The error bounds match the rate of the minimax lower bounds
- The error bounds in Theorems 6 and 7 are the sharpest possible based on the basic inequality
- Many error bounds for prediction and coefficient estimation have been derived in the literature, mostly based on cruder versions of the basic inequality and under the stronger hard sparsity condition with  $\bar{\Delta} = 0$
- Similar prediction and estimation error bounds can be established for the Dantzig selector

$$\hat{\beta}^{(DS)} = \arg \min_b \left\{ \|b\|_1 : \|\mathbf{X}^T(\mathbf{Y} - \mathbf{X}b)/n\|_\infty \leq \lambda \right\}, \lambda \geq \lambda_{univ}$$

## The false positive and sparsity of the Lasso

- Sparse condition number of the Gram matrix  $\bar{\Sigma} = \mathbf{X}^T \mathbf{X} / n$ ,

$$\phi_{\text{cond}}(m; \mathcal{S}, \bar{\Sigma}) = \max_{|B \setminus \mathcal{S}| \leq (1 \vee m)} \{ \phi_{\max}(\bar{\Sigma}_{B,B}) / \phi_{\min}(\bar{\Sigma}_{B,B}) \}$$

- Constant factor for controlling the false positive

$$C_{\text{FP}}^{(\ell_1)}(\mathcal{S}, \bar{\Sigma}) = \inf \left\{ t : \frac{\phi_{\text{cond}}(t|\mathcal{S}|; \mathcal{S}, \bar{\Sigma}) - 1}{2(1 - \eta)^2 / (1 + \eta)^2} + \frac{1}{|\mathcal{S}|} \leq t \right\}.$$

**Theorem 8.** Let  $\mathcal{S} \supseteq \text{supp}(\beta) \cup \text{supp}(\beta^*)$  and  $\hat{\mathcal{S}} = \text{supp}(\hat{\beta})$ . Then,

$$|\hat{\mathcal{S}} \setminus \mathcal{S}| \leq |\mathcal{S}| C_{\text{FP}}^{(\ell_1)}(\mathcal{S}, \bar{\Sigma}).$$

**Heuristics:** The KKT implies that for any  $B \subseteq \hat{\mathcal{S}} \setminus \mathcal{S}$ ,

$$\begin{aligned} \lambda^2 |B| &= \|\mathbf{X}_B(\mathbf{Y} - \mathbf{X}\hat{\beta})/n\|_2^2 \\ &\leq \phi_{\max}(\bar{\Sigma}_{B,B}) \|\mathbf{P}_B(\mathbf{Y} - \mathbf{X}\hat{\beta})\|_2^2 / n. \end{aligned}$$

Done if the r.h.s. is smaller than  $\lambda^2 |\mathcal{S}| C_{\text{FP}}^{(\ell_1)}(\mathcal{S}, \bar{\Sigma})$  when  $|B| \leq |\mathcal{S}| C_{\text{FP}}^{(\ell_1)}(\mathcal{S}, \bar{\Sigma})$ .  
See Z-Huang (2008), Z (2010) and Z-Zhang (2012).

## LSE after Lasso selection

- LSE in a given model  $\mathcal{M} \subset \{1, \dots, p\}$ :

$$\hat{\beta}^{(\mathcal{M})} = \arg \min_b \left\{ \|\mathbf{Y} - \mathbf{X}b\|_2^2 : \text{supp}(\mathbf{b}) \subset \mathcal{M} \right\}$$

- Relaxed lower sparse eigenvalue: With  $\bar{\Sigma} = \mathbf{X}^T \mathbf{X} / n$ ,

$$\phi_*(m; \mathcal{S}, \bar{\Sigma}) = \min_{|B \setminus \mathcal{S}| \leq (1 \vee m)} \phi_{\min}(\bar{\Sigma}_{B,B})$$

**Theorem 9.** Let  $\mathcal{S} = \text{supp}(\beta)$  and  $\hat{\mathcal{S}} = \text{supp}(\hat{\beta})$  and  $m_1 = \lfloor C_{\text{pred}}^{(\ell_1)}(\mathcal{S}) |\mathcal{S}| \rfloor$ . Suppose that  $\mathbf{Y} - \mathbf{X}\beta$  is sub-Gaussian and  $\lambda = \eta^{-1} \sigma \sqrt{(2/n) \log(p/\epsilon)}$ . Then,

$$\begin{aligned} \frac{\|\hat{\beta}^{(\hat{\mathcal{S}})} - \beta\|_2^2}{\phi_*(m_1; \mathcal{S}, \bar{\Sigma})} &\leq \|\mathbf{X}\hat{\beta}^{(\hat{\mathcal{S}})} - \mathbf{X}\beta\|_2^2 / n \\ &\leq |\mathcal{S}| \left\{ C_{\text{pred}}^{(\ell_1)}(\mathcal{S}) (\lambda^2 + 7.5\sigma^2/n) + 2.5 C_{\text{FP}}^{(\ell_1)}(\mathcal{S}, \bar{\Sigma}) (\lambda_1^2 + 3\sigma^2/n) \right\} \end{aligned}$$

with at least probability  $1 - 2\epsilon - \epsilon^{m_1}$ , where  $\lambda_1 = \sigma \sqrt{(2/n) \log(ep/(\epsilon m_1))}$

**Proof:**  $|\hat{\mathcal{S}} \setminus \mathcal{S}| \leq m_1$  and  $\|\mathbf{X}\hat{\beta}^{(\hat{\mathcal{S}})} - \mathbf{X}\beta^*\|_2^2 / n \leq \|\mathbf{P}_{\hat{\mathcal{S}}}^\perp \mathbf{X}\beta^*\|_2^2 / n + \|\mathbf{P}_{\hat{\mathcal{S}}} \epsilon\|_2^2 / n$ .

## Some references:

- The Lasso: Tibshirani (96), Chen-Donoho-Saunders (01)
- Prediction/persistence: Greenshtein-Ritov (04)
- Dantzig selector/restricted isometry condition: Candes-Tao (05,07)
- Lasso/sparse Riesz condition: Z-Huang (08), Zhang (09)
- RE condition: Bickel-Ritov-Tsybakov (09), Koltchinskii (09)
- Compatibility condition: van de Geer (07), van de Geer-Bühlmann (09)
- Cone invertibility factors (CIF): Ye-Z (10)
- Capped- $\ell_1$  sparsity: Z-Huang (08), Z-Zhang (12)
- Approximate sparsity: Sun-Z (12)
- False positive/Post Lasso LSE: Sun-Z (12)

**Restricted eigenvalue:** Let  $\xi = (1 + \eta)/(1 - \eta)$ ,  $\bar{\mathbf{\Sigma}} = \mathbf{X}^T \mathbf{X}/n$  and

$$\text{RE}^2(\mathcal{S}, \xi; \bar{\mathbf{\Sigma}}) = \inf \left\{ \mathbf{u}^T \bar{\mathbf{\Sigma}} \mathbf{u} : \|\mathbf{u}_{\mathcal{S}}\|_2 = 1, \|\mathbf{u}_{\mathcal{S}^c}\|_1 < \xi \|\mathbf{u}_{\mathcal{S}}\|_1 \right\}.$$

**Lemma 4.** For integer  $m > 0$  and  $\mathbf{u} \in \mathbb{R}^p$ , there exists  $\mathbf{v} \in \mathbb{R}^p$  such that

$$\mathbf{v}_{\mathcal{S}} = \mathbf{u}_{\mathcal{S}}, \text{supp}(\mathbf{v}) \subseteq \text{supp}(\mathbf{u}), \|\mathbf{v}_{\mathcal{S}^c}\|_0 \leq m, \|\mathbf{v}_{\mathcal{S}^c}\|_1 = \|\mathbf{u}_{\mathcal{S}^c}\|_1,$$

and

$$\mathbf{v}^T \bar{\mathbf{\Sigma}} \mathbf{v} - \mathbf{u}^T \bar{\mathbf{\Sigma}} \mathbf{u} \leq \|\mathbf{u}_{\mathcal{S}^c}\|_1^2/m \leq |\mathcal{S}| \xi^2 \|\mathbf{u}_{\mathcal{S}}\|_2^2/m.$$

**Proof:** Assume  $\|\mathbf{u}_{\mathcal{S}^c}\|_1 > 0$ . Let  $\pi_j = |u_j|/\|\mathbf{u}_{\mathcal{S}^c}\|_1$  and  $\mathbf{Z}^i$  be a iid vectors with

$$\mathbb{P} \left\{ \mathbf{Z}^i = \mathbf{u}_{\mathcal{S}} + \|\mathbf{u}_{\mathcal{S}^c}\|_1 \text{sgn}(u_j) \mathbf{e}_j \right\} = \pi_j, \quad j \notin \mathcal{S}.$$

Let  $\mathbf{V} = \sum_{i=1}^m \mathbf{Z}^i/m$ . We have  $\mathbb{E} \mathbf{Z} = \mathbf{u}$  and

$$\mathbb{E} \mathbf{V}^T \bar{\mathbf{\Sigma}} \mathbf{V} - \mathbf{u}^T \bar{\mathbf{\Sigma}} \mathbf{u} = \frac{1}{m} \mathbb{E} (\mathbf{Z} - \mathbf{u})_{\mathcal{S}^c}^T \bar{\mathbf{\Sigma}}_{\mathcal{S}^c, \mathcal{S}^c} (\mathbf{Z} - \mathbf{u})_{\mathcal{S}^c} \leq \frac{\|\mathbf{u}_{\mathcal{S}^c}\|_1^2}{m} \|\bar{\mathbf{\Sigma}}_{\mathcal{S}^c, \mathcal{S}^c}\|_{\max}.$$

Thus,  $\mathbf{v}$  can be one of the realizations of  $\mathbf{V}$ .

## Restricted eigenvalue:

$$\text{RE}^2(\mathcal{S}, \xi; \bar{\Sigma}) = \inf \left\{ \mathbf{u}^T \bar{\Sigma} \mathbf{u} : \|\mathbf{u}_{\mathcal{S}}\|_2 = 1, \|\mathbf{u}_{\mathcal{S}^c}\|_1 < \xi \|\mathbf{u}_{\mathcal{S}}\|_1 \right\}.$$

A lower bound of the restricted eigenvalue is

$$\text{RE}^2(\mathcal{S}, \xi; \bar{\Sigma}) \geq \inf \left\{ \mathbf{v}^T \bar{\Sigma} \mathbf{v} : \|\mathbf{v}_{\mathcal{S}}\|_2 = 1, \|\mathbf{v}_{\mathcal{S}^c}\|_1 < \xi \|\mathbf{v}_{\mathcal{S}}\|_1, \|\mathbf{v}_{\mathcal{S}^c}\|_0 \leq m \right\} - \xi^2 |\mathcal{S}|/m.$$

If  $\mathbf{X}$  has iid  $N(0, \Sigma)$  (or sub-Gaussian) rows, then

$$\text{RE}^2(\mathcal{S}, \xi; \bar{\Sigma}) \geq \left\{ 1 - \frac{C_0}{n} \binom{p}{|\mathcal{S}| + m} \right\} \text{RE}^2(\mathcal{S}, \xi; \Sigma) - \frac{\xi^2 |\mathcal{S}|}{m}.$$

This strengthens the results in Rudelson-Zhou (2013)



**Concave penalties:** Global or suitable local minimizers of

$$\mathcal{L}_\lambda(\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2/(2n) + \|\mathbf{p}_\lambda(\boldsymbol{\beta})\|_1$$

where  $\|\mathbf{p}_\lambda(\boldsymbol{\beta})\|_1 = \sum_{j=1}^p p_\lambda(\beta_j)$

- Critical point:

$$\liminf_{t \rightarrow 0_+} \inf_{0 < \|\mathbf{u}\|_2 \leq 1} t^{-1} \left\{ \mathcal{L}_\lambda(\hat{\boldsymbol{\beta}} + t\mathbf{u}) - \mathcal{L}_\lambda(\hat{\boldsymbol{\beta}}) \right\} \geq 0$$

- Local minimizer: There exists  $t_0 > 0$  such that

$$\mathcal{L}_\lambda(\mathbf{b}) \geq \mathcal{L}_\lambda(\hat{\boldsymbol{\beta}}), \quad \forall 0 < \|\mathbf{b} - \hat{\boldsymbol{\beta}}\|_2 \leq t_0$$

- Strict local minimizer: There exists  $t_0 > 0$  such that

$$\mathcal{L}_\lambda(\mathbf{b}) > \mathcal{L}_\lambda(\hat{\boldsymbol{\beta}}), \quad \forall 0 < \|\mathbf{b} - \hat{\boldsymbol{\beta}}\|_2 \leq t_0$$

- Global minimizer:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \mathcal{L}_\lambda(\mathbf{b})$$

**Theorem 10. KKT-type conditions for local solutions:** Suppose that  $p_\lambda(t)$  is twice left- and right-differentiable for all real  $t$ .

(i) A vector  $\hat{\beta}$  is a critical point iff

$$\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\hat{\beta})/n = p'_\lambda(\hat{\beta}_j), \quad \forall j = 1, \dots, p.$$

In particular,  $|\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\hat{\beta})/n| \leq p'_\lambda(0+)$  for all  $j$  with  $\hat{\beta}_j = 0$ .

(ii) A critical point  $\hat{\beta}$  is a strict local minimizer if

$$\|\mathbf{X}\mathbf{u}\|_2^2/n + \sum_{u_j > 0} p''_\lambda(\hat{\beta}_j+)u_j^2 + \sum_{u_j < 0} p''_\lambda(\hat{\beta}_j-)u_j^2 > 0$$

for all  $\mathbf{u}$  with  $\{j : u_j > 0\} \subseteq \mathcal{S}_+$  and  $\{j : u_j < 0\} \subseteq \mathcal{S}_-$ , where  $\mathcal{S}_\pm = \{j : \mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\hat{\beta})/n = p'_\lambda(\hat{\beta}_j\pm)\}$ .

(iii) Suppose  $p_\lambda(t)$  is a quadratic spline in  $[0, \infty)$ . Then, a critical point  $\hat{\beta}$  is a strict local minimizer iff the condition in (ii) holds, and it is a local minimizer iff

$$\|\mathbf{X}\mathbf{u}\|_2^2/n + \sum_{u_j > 0} p''_\lambda(\hat{\beta}_j+)u_j^2 + \sum_{u_j < 0} p''_\lambda(\hat{\beta}_j-)u_j^2 \geq 0$$

for the same set of  $\mathbf{u}$ .

**Oracle local solutions:** Let  $\mathcal{S} = \text{supp}(\beta)$  and

$$\theta_{\text{select}}(p_\lambda, \beta) = \inf \left\{ \theta : \frac{\|\mathbf{v}_\mathcal{S} - \beta_\mathcal{S}\|_\infty}{\theta\lambda + \lambda'} \leq 1 \Rightarrow \|\bar{\Sigma}_{\mathcal{S},\mathcal{S}}^{-1} p'_\lambda(\mathbf{v}_\mathcal{S})\|_\infty \leq \theta\lambda \right\},$$

$$\kappa_{\text{select}}(p_\lambda, \beta) = \sup \left\{ \left\| \bar{\Sigma}_{\mathcal{S}^c,\mathcal{S}} \bar{\Sigma}_{\mathcal{S},\mathcal{S}}^{-1} \left( \frac{p'_\lambda(\mathbf{v}_\mathcal{S})}{\lambda} \right) \right\|_\infty : \frac{\|\mathbf{v}_\mathcal{S} - \beta_\mathcal{S}\|_\infty}{\theta_{\text{select}}(p_\lambda, \beta)\lambda + \lambda'} \leq 1 \right\}.$$

**Theorem 11.** Let  $\hat{\beta}^\circ$  be the oracle LSE. Suppose  $p'_\lambda(t)$  is continuous in  $t > 0$ ,  $\text{rank}(\bar{\Sigma}_{\mathcal{S},\mathcal{S}}) = |\mathcal{S}|$ ,  $\min_{j \in \mathcal{S}} |\beta_j| \geq \theta_{\text{select}}(p_\lambda, \beta)\lambda + \lambda'$ , and

$$\|\hat{\beta}^\circ - \beta\|_\infty \leq \lambda', \quad \|\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\beta}^\circ)/n\|_\infty \leq \lambda(1 - \kappa_{\text{select}}(p_\lambda, \beta))_+.$$

Then, there exists a local solution  $\hat{\beta}$  such that

$$\text{sgn}(\hat{\beta}) = \text{sgn}(\beta), \quad \|\hat{\beta} - \beta\|_\infty \leq \theta_{\text{select}}(p_\lambda, \beta)\lambda + \lambda'.$$

Moreover, if in addition  $\theta_{\text{select}}(p_\lambda, \beta) = 0$ , then  $\kappa_{\text{select}}(p_\lambda, \beta) = 0$  and

$\hat{\beta}^\circ$  is a local minimizer.

## Remarks:

- Folded concave penalties provide oracle local solutions under much weaker condition on  $\mathbf{X}$  than the Lasso, e.g. when  $\hat{\beta}^o$  is an oracle solution under the signal strength condition  $\theta_{\text{select}}(p_\lambda, \beta) = 0$
- Due to the multiplicity of solutions, the challenges of concave PLSE are
  - to find such oracle solutions if exists
  - to cause little harm or even some gain in performance for prediction and coefficient estimation, e.g. in view of the rate minimaxity of the Lasso, when such oracle solution does not exist
  - to compute such solutions with reasonable cost

**Outline of the proof of Theorem 11:** Apply Brouwer's fixed-point theorem to

$$\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\hat{\beta})/n = p'_\lambda(\hat{\beta}), \quad \text{supp}(\hat{\beta}) = \mathcal{S},$$

on  $\mathcal{S}$  to find  $\hat{\beta}_{\mathcal{S}}$  and then check the condition on  $\mathcal{S}^c$  using  $\kappa_{\text{select}}(p_\lambda, \beta)$ . See Z-Zhang (2012).

### Proof of Theorem 11:

- Recall that  $\|\hat{\beta}^\circ - \beta\|_\infty \leq \lambda'$ .
- Let  $B = \{\mathbf{h}_S : \|\mathbf{h}_S\|_\infty \leq \theta_{\text{select}}(\mathbf{p}_\lambda, \beta)\lambda\}$ .
- For  $\mathbf{h}_S \in B$ ,  $\mathbf{v}_S = \hat{\beta}_S^\circ + \mathbf{h}_S$  satisfies  $\|\mathbf{v}_S - \beta_S\|_\infty \leq \theta_{\text{select}}(\mathbf{p}_\lambda, \beta)\lambda + \lambda'$ .
- It follows that  $\|\bar{\Sigma}_{S,S}^{-1} \mathbf{p}'_\lambda(\mathbf{v}_S)\|_\infty \leq \theta_{\text{select}}(\mathbf{p}_\lambda, \beta)\lambda$ .
- Thus,  $\mathbf{h}_S = -\bar{\Sigma}_{S,S}^{-1} \mathbf{p}'_\lambda(\hat{\beta}_S^\circ + \mathbf{h}_S)$  has a solution (Brouwer).
- Let  $\hat{\beta}_S = \hat{\beta}_S^\circ + \mathbf{h}_S$  and  $\hat{\beta}_{S^c} = \mathbf{0}$ . We have

$$\mathbf{X}_S^T(\mathbf{Y} - \mathbf{X}\hat{\beta})/n = \mathbf{X}_S^T(\mathbf{X}\hat{\beta}^\circ - \mathbf{X}\hat{\beta})/n = -\bar{\Sigma}_{S,S}\mathbf{h}_S = \mathbf{p}'_\lambda(\hat{\beta}_S).$$

- Moreover, by the definition of  $\kappa_{\text{select}}(\mathbf{p}_\lambda, \beta)$ ,

$$\begin{aligned} & \|\mathbf{X}_{S^c}^T(\mathbf{Y} - \mathbf{X}_S\hat{\beta}_S)/n\|_\infty \\ &= \|\mathbf{X}_{S^c}^T(\mathbf{Y} - \mathbf{X}\hat{\beta}^\circ)/n + \bar{\Sigma}_{S^c,S}\bar{\Sigma}_{S,S}^{-1}\mathbf{p}'_\lambda(\hat{\beta}_S)\|_\infty \\ &\leq \|\mathbf{X}_{S^c}^T(\mathbf{Y} - \mathbf{X}\hat{\beta}^\circ)/n\|_\infty + \kappa_{\text{select}}(\mathbf{p}_\lambda, \beta)\lambda \\ &\leq \lambda. \end{aligned}$$

**Local solution path:** Let  $\kappa(p_\lambda) = \sup_{t_1 < t_2} \{p'_\lambda(t_1) - p'_\lambda(t_2)\} / (t_2 - t_1)$ ,

- $\mathcal{P}(\lambda_0, \kappa_0) = \left\{ p_\lambda(\cdot) : p'_\lambda(0+) = \lambda \geq \lambda_0, \|p'_\lambda(\cdot)\|_\infty \leq \lambda, \kappa(p_\lambda) \leq \kappa_0 \right\}$
- $\mathcal{B}(\lambda_0, \kappa_0)$  = the set of all local solutions for some  $p_\lambda \in \mathcal{P}(\lambda_0, \kappa_0)$
- $\mathcal{B}_0(\lambda_0, \kappa_0)$  = the set of all vectors connected to  $\mathbf{0}$  in  $\mathcal{B}(\lambda_0, \kappa_0)$
- $\beta^*$  be a target vector satisfying  $\|\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta^*)/n\|_\infty < \eta\lambda_0$  with  $\eta < 1$
- $\mathcal{S} = \text{supp}(\beta^*)$
- $\text{RE} = \inf_{\mathbf{u} \neq \mathbf{0}} \left\{ \|\mathbf{X}\mathbf{u}\|_2 / (n^{1/2} \|\mathbf{u}\|_2) : \|\mathbf{u}_{\mathcal{S}^c}\|_1 \leq \|\mathbf{u}_{\mathcal{S}}\|_1 + \eta \|\mathbf{u}\|_1 \right\}$

**Theorem 12. (Feng-Z, 2015).** Let  $\hat{\beta}$  be a local minimizer in  $\mathcal{B}_0(\lambda_0, \kappa_0)$  with a penalty  $p_\lambda \in \mathcal{P}(\lambda_0, \kappa_0)$ . Suppose  $\text{RE}^2 \geq \kappa_0$ . Then, for all seminorms  $\|\cdot\|$ ,

$$\|\hat{\beta} - \beta^*\| \leq \lambda \sup \left\{ \frac{(1+\eta)\|\mathbf{u}\|}{\|\bar{\Sigma}\mathbf{u}\|_\infty} : (1-\eta)\|\mathbf{u}_{\mathcal{S}^c}\|_1 \leq \mathbf{w}_{\mathcal{S}}^T \mathbf{u}_{\mathcal{S}} \right\}$$

where  $\mathbf{w}_{\mathcal{S}}^T = \{\mathbf{X}_{\mathcal{S}}^T(\mathbf{Y} - \mathbf{X}\beta^*)/n - p'_\lambda(\beta_{\mathcal{S}}^*)\}/\lambda$ . In particular,

$$\|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^*\|_2^2/n \leq \frac{4\lambda^2|S|(1+\eta)^2}{(1-\eta)^2\text{RE}^2}, \quad \|\hat{\beta} - \beta^*\|_2 \leq \frac{2\lambda|S|(1+\eta)}{(1-\eta)\text{RE}^2},$$

and when  $\phi_{\min}(\bar{\Sigma}_{\mathcal{S},\mathcal{S}}) > \kappa_0$  and  $\beta^*$  is an oracle solution for  $p_\lambda$ ,

$$\hat{\beta} = \beta^*.$$

**Proof:** Let  $\mathbf{h} = \hat{\beta} - \beta^*$ . For some  $\tilde{\eta} < \eta$ , we have

$$\begin{aligned}
 h_j \mathbf{X}_j^T (\mathbf{X} \hat{\beta} - \mathbf{Y}) / n &= -h_j p'_\lambda(\hat{\beta}_j), && \text{"KKT" for } \hat{\beta} \\
 -h_j p'_\lambda(\hat{\beta}_j) &\leq \kappa_0 h_j^2 - h_j p'_\lambda(\beta_j^*), && \text{bound on concavity} \\
 \|\mathbf{X} \mathbf{h}\|_2^2 / n &\leq \mathbf{h} \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \beta^*) / n + \kappa_0 \|\mathbf{h}\|_2^2 - \mathbf{h}^T p'_\lambda(\beta^*), && \text{summing over } j \\
 -\mathbf{h}^T p'_\lambda(\beta^*) &= -\lambda \|\mathbf{h}_{S^c}\|_1 - \mathbf{h}_S^T p'_\lambda(\beta_S^*), && \text{favorable choice of } p'_\lambda \\
 \mathbf{h} \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \beta^*) / n &\leq \tilde{\eta} \lambda \|\mathbf{h}_{S^c}\|_1 + \mathbf{h}_S \mathbf{X}_S^T (\mathbf{Y} - \mathbf{X} \beta^*) / n, && \text{condition on noise}
 \end{aligned}$$

Combining the above three inequalities, we have a basic inequality

$$\|\mathbf{X} \mathbf{h}\|_2^2 / n + (1 - \tilde{\eta}) \lambda \|\mathbf{h}_{S^c}\|_1 \leq \kappa_0 \|\mathbf{h}\|_2^2 + \lambda \mathbf{w}_S^T \mathbf{h}_S.$$

Key argument: As  $\|\mathbf{w}_S\|_\infty \leq 1 + \tilde{\eta}$ , the above and the RE condition yield

$$\begin{aligned}
 \kappa_0 \|\mathbf{h}\|_2^2 &\leq (\eta - \tilde{\eta}) \|\mathbf{h}\|_1 \\
 \Rightarrow \|\mathbf{h}_{S^c}\|_1 &\leq \|\mathbf{h}_S\|_1 + \eta \|\mathbf{h}\|_1 \\
 \Rightarrow (1 - \tilde{\eta}) \lambda \|\mathbf{h}_{S^c}\|_1 &\leq \lambda \mathbf{w}_S^T \mathbf{h}_S \\
 \Rightarrow \|\mathbf{h}_{S^c}\|_1 &\leq \|\mathbf{h}_S\|_1 + \tilde{\eta} \|\mathbf{h}\|_1.
 \end{aligned}$$

Thus,  $\mathcal{B}_0(\lambda_0, \kappa_0) - \beta^* \subseteq \mathcal{C}(\mathcal{S}, \eta) = \{\mathbf{u} : (1 - \eta) \|\mathbf{u}_{S^c}\|_1 \leq \mathbf{w}_S^T \mathbf{u}_S, \mathbf{w}_S = \mathbf{w}_S^{(p_\lambda)}\}$ .

As  $\|\mathbf{h}\| \leq \|\bar{\Sigma} \mathbf{h}\|_\infty \sup_{\mathbf{u} \in \mathcal{C}(\mathcal{S}, \eta)} \|\mathbf{u}\| / \|\bar{\Sigma} \mathbf{u}\|_\infty$ , the error bound follows from

$$\|\bar{\Sigma} \mathbf{h}\|_\infty \leq \|\mathbf{X}^T (\mathbf{Y} - \mathbf{X} \hat{\beta}) / n\|_\infty + \|\mathbf{X}^T (\mathbf{Y} - \mathbf{X} \beta^*) / n\|_\infty \leq \lambda + \eta \lambda.$$

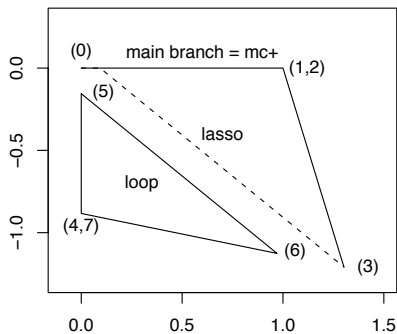
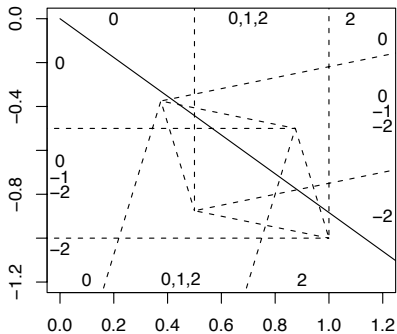
## Remarks:

- $\mathcal{B}_0(\lambda_0, \kappa_0)$  is the set of all local solutions computable by path following algorithms starting from the origin, with the constraints  $\lambda \geq \lambda_0$  and  $\kappa(p_\lambda) \leq \kappa_0$  on the penalty and concavity levels
- The RE condition, which is of the  $\ell_2$  type, is arguably nearly the weakest proven one on  $\mathbf{X}$  for the rate optimal prediction and coefficient estimation error bounds for the Lasso, e.g. in models where the sharper bounds in Theorems 6 and 7 are not much sharper
- Under the RE condition, any path following solution achieves the same or sharper prediction and coefficient estimation error rates as the Lasso
- When a large majority of the components of  $|\hat{\beta}_S^\circ|$  are large,  $\mathbf{w}_S$  can be small for  $\hat{\beta}^\circ$ . In this case, concave penalty outperforms the Lasso in prediction and coefficient estimation
- Under the same RE condition, the path following solution achieves variable selection consistency,  $\hat{\beta} = \beta^*$ , completely removing the need for the restrictive  $\ell_\infty$ -type conditions such as the irrepresentable condition required for the selection consistency of the Lasso



## A path following algorithm:

- $p_\lambda(t) = \lambda^2 p_1(|t|/\lambda)$  with a quadratic spline  $p_1(t)$  in  $[0, \infty)$
- $\hat{\beta}(\lambda)$ : minimizing  $\|\mathbf{Y} - \mathbf{X}\beta\|_2^2/(2n) + \lambda^2 \|p_1(\beta/\lambda)\|_1$
- Rescale:  $\tau = 1/\lambda$ ,  $\mathbf{b} = \mathbf{b}(\tau) = \hat{\beta}(\lambda)/\lambda$ ,  $\mathbf{Z} = \mathbf{X}^T \mathbf{y}/n$
- $\mathbf{b}$ : minimizing  $\mathbf{b}^T \bar{\mathbf{\Sigma}} \mathbf{b}/2 - \tau \mathbf{Z}^T \mathbf{b} + \|p_1(\mathbf{b})\|_1$
- “KKT”:  $\bar{\mathbf{\Sigma}} \mathbf{b}(\tau) - \tau \mathbf{Z} + p'_1(\mathbf{b}(\tau)) = 0$ , linear spline path
- Direction of path:  $\bar{\mathbf{\Sigma}} \mathbf{b}'(\tau) - \mathbf{Z} + p''_1(\mathbf{b}(\tau)) \circ \mathbf{b}'(\tau) = 0$
- **Initialization:**  $\tau^{(0)} = 0$ ,  $\hat{S}^{(0)} = \emptyset$ ,  $\mathbf{b}(\tau) = 0$ ,  $k = 1$
- **Iteration:**
  - Find  $\tau^{(k)}$  as the beginning point for new  $\hat{S}^{(k)}$  or new  $p''_\lambda(\mathbf{b})$
  - Find the new  $\hat{S}^{(k)}$  &  $p''_\lambda(\mathbf{b})$
  - Compute  $\partial \mathbf{b}_{\hat{S}}^{(k)} = (\bar{\mathbf{\Sigma}} + \text{diag}(p''_1(\mathbf{b})))_{\hat{S}, \hat{S}}^{-1} \mathbf{Z}_{\hat{S}}$  with new  $\hat{S}^{(k)}$  &  $p''_\lambda(\mathbf{b})$
  - $\mathbf{b}(\tau) = \mathbf{b}(\tau^{(k)}) + (\tau - \tau^{(k)}) \partial \mathbf{b}^{(k)}$  from  $\tau = \tau^{(k)}$  to  $\tau = \tau^{(k+1)}$
  - $k = k + 1$
- References
  - Osborne et al (2000), Efron et al (2004), LARS for the Lasso
  - Z (2010), PLUS for the Lasso, SCAD, MCP and ...



## Majorize-Minimization (MM) algorithms:

- $\mathcal{L}_\lambda(\beta) = \mathcal{L}_0(\beta) + R_\lambda(\beta)$ ,  $\hat{\beta}^{(old)}$
- $f(\beta) \geq \mathcal{L}_\lambda(\beta)$ ,  $f(\hat{\beta}^{(old)}) = \mathcal{L}_\lambda(\hat{\beta}^{(old)})$
- $\hat{\beta}^{(new)} = \arg \min_{\beta} f(\beta)$
- $\mathcal{L}_\lambda(\hat{\beta}^{(new)}) \leq f(\hat{\beta}^{(new)}) \leq f(\hat{\beta}^{(old)}) = \mathcal{L}_\lambda(\hat{\beta}^{(old)})$

- Majorizing the penalty

- Local quadratic approximation (LQA; Fan-Li, 2001):

$$p_\lambda(|\beta_j|) \leq p_\lambda(|\hat{\beta}_j^{(old)}|) + p'_\lambda(|\hat{\beta}_j^{(old)}|)(\beta_j^2 - |\hat{\beta}_j^{(old)}|^2)/(2|\hat{\beta}_j^{(old)}|)$$

- Local linear approximation (LLA; Zou-Li, 2008):

$$p_\lambda(|\beta_j|) \leq p_\lambda(|\hat{\beta}_j^{(old)}|) + p'_\lambda(|\hat{\beta}_j^{(old)}|)(|\beta_j| - |\hat{\beta}_j^{(old)}|) = p'_\lambda(|\hat{\beta}_j^{(old)}|)|\beta_j| + C_j$$

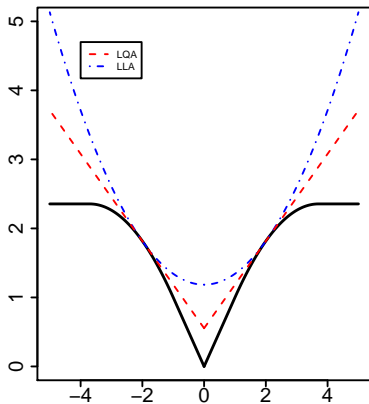
- Majorizing the loss

- (Fast) iterative shrinkage-thresholding algorithm (FISTA; Nesterov, 83; Daubechies et al, 2004; Beck-Teboulle, 2009)

$$\mathcal{L}_0(\beta) \leq \mathcal{L}_0(\hat{\beta}^{(old)}) + \left\langle \nabla \mathcal{L}_0(\hat{\beta}^{(old)}), \beta - \hat{\beta}^{(old)} \right\rangle + \|\beta - \hat{\beta}^{(old)}\|_2^2/(2s)$$

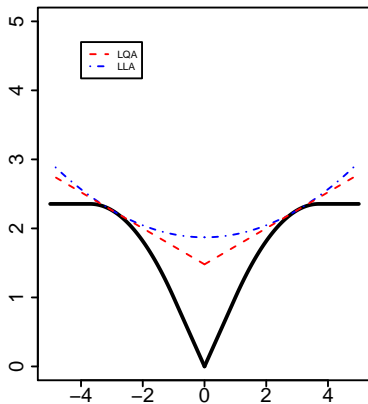
$$\text{with } \sup_{\beta} \|\nabla^{\otimes 2} \mathcal{L}_0(\beta)\|_{\text{spectrum}} \leq 1/s$$

## LQA and LLA



(a)

## LQA and LLA



(a)

See Zou-Li 08; Thanks to J. Fan

**Properties of LLA:** Let  $\tilde{\beta} = \hat{\beta}^{(old)}$  and  $\hat{\beta}$  be an estimator satisfying

$$|X_j^T(X\hat{\beta} - Y)/n + p'_\lambda(\tilde{\beta}_j)\partial|\tilde{\beta}_j|| \leq \lambda\nu_j$$

**Basic inequality:** Suppose  $|p'_\lambda(t_1) - p'_\lambda(t_2)| \leq \kappa(p_\lambda)|t_1 - t_2| \forall 0 < t_1 < t_2$ . Similar to the basic inequality for the path following algorithm, we have

$$\mathbf{h}^T \bar{\Sigma} \mathbf{h} + (1 - \eta)\lambda \|\mathbf{h}_{S^c}\|_1 \leq \lambda \mathbf{w}_S^T \mathbf{h}_S + \|\mathbf{h}\|_2 (\kappa(p_\lambda) \|\tilde{\mathbf{h}}\|_2 + \lambda \|\boldsymbol{\nu}\|_2)$$

with  $\mathbf{h} = \hat{\beta} - \beta^*$ ,  $\tilde{\mathbf{h}} = \tilde{\beta} - \beta^*$ , and  $\mathbf{w}_S = \{X_S^T(Y - X\beta^*)/n - p'_\lambda(\beta_S^*)\}/\lambda$ .

**Theorem 13:** (i) Suppose  $p'_\lambda$  satisfies the Lipschitz condition. Then,

$$\|\mathbf{h}\| \leq \lambda \sup \left\{ \|\mathbf{u}\| \psi(\mathbf{u}) : \mathbf{u}^T \bar{\Sigma} \mathbf{u} = 1 \right\}$$

with  $\psi(\mathbf{u}) = \mathbf{w}_S^T \mathbf{u}_S - (1 - \eta)\|\mathbf{u}_{S^c}\|_1 + \|\mathbf{u}\|_2 \{ \kappa(p_\lambda) \|\tilde{\mathbf{h}}\|_2 / \lambda + \|\boldsymbol{\nu}\|_2 \}$ .

(ii) Let  $\rho_0 = \sup \{ \|\mathbf{u}\|_2 : \psi(\mathbf{u}) > 0, \mathbf{u}^T \bar{\Sigma} \mathbf{u} = 1 \}$ . Then,

$$\|\mathbf{h}\|_2 \leq \rho_0 \kappa(p_\lambda) \|\tilde{\mathbf{h}}\|_2 + \lambda \sup_{\mathbf{u}^T \bar{\Sigma} \mathbf{u} = 1} \left\{ \|\mathbf{u}\|_2 (\mathbf{w}_S^T \mathbf{u}_S - (1 - \eta)\|\mathbf{u}_{S^c}\|_1 + \|\mathbf{u}\|_2 \|\boldsymbol{\nu}\|_2) \right\}$$

- Multistage LLA/approximate solution (Zhang, 2010, 2013; Z-Zhang, 2012)
- RSC (Negahban et al, 2012)

## Estimation of HD objects

- Universal penalty level:  $\lambda_{univ} = \sigma \sqrt{(2/n) \log p}$
- Sparsity condition:  $\sum_{j=1}^p \min(1, |\beta_j|/\lambda_{univ}) \leq s^*$
- Prediction: With at least probability  $1 - p^{-a_0}$ ,

$$\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\boldsymbol{\beta}\|_2^2/n \leq C_{\text{pred}} s^* \lambda_{univ}^2$$

- Estimation: With at least probability  $1 - p^{-a_0}$ ,

$$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_q^q \leq C_{\text{est},q} s^* \lambda_{univ}^q$$

- Variable selection:

$$\begin{aligned} \beta_{\min} &= \min_{\beta_j \neq 0} |\beta_j| \geq C_{\text{select}} \lambda_{univ} \\ \Rightarrow \mathbb{P}\left\{\text{sgn}(\hat{\boldsymbol{\beta}}) = \text{sgn}(\boldsymbol{\beta})\right\} &\geq 1 - p^{-a_0} \end{aligned}$$

- Proper conditions on  $\mathbf{X}$  yield bounded  $C_{\text{pred}}$ ,  $C_{\text{est},q}$  and  $C_{\text{select}}$

**How about confidence interval and p-value for  $\beta_j$ ?**

## Statistical inference after model selection

- A model  $M$  is selected based on screening data
- Fresh data  $(\mathbf{X}, \mathbf{Y})$  are observed independent of  $M$ ,  $n \gg |M|$
- Confidence interval for  $\theta = \mathbf{a}^T \boldsymbol{\beta}$  after model selection:

$$\mathbf{a}_M^T \hat{\boldsymbol{\beta}}_M^{(lse)} \pm 1.96\sigma \left( \mathbf{a}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{a}_M \right)^{1/2}$$

## Statistical inference after model selection

- A model  $M$  is selected based on screening data
- Fresh data  $(\mathbf{X}, \mathbf{Y})$  are observed independent of  $M$ ,  $n \gg |M|$
- Confidence interval for  $\theta = \mathbf{a}^T \boldsymbol{\beta}$  after model selection:

$$\mathbf{a}_M^T \hat{\boldsymbol{\beta}}_M^{(lse)} \pm 1.96\sigma \left( \mathbf{a}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{a}_M \right)^{1/2}$$

A superefficiency paradox

- If an oracle model  $A$  is available such that

$$A \supseteq M \cup \{j : \beta_j \neq 0\}, \quad |M| \vee \|\boldsymbol{\beta}\|_0 < |A| \ll n,$$

- do we use

$$\mathbf{a}_A^T \hat{\boldsymbol{\beta}}_A^{(lse)} \pm 1.96\sigma \left( \mathbf{a}_A (\mathbf{X}_A^T \mathbf{X}_A)^{-1} \mathbf{a}_A \right)^{1/2}?$$

- If we do, then the CI with the oracle information is wider:

$$\left( \mathbf{a}_A (\mathbf{X}_A^T \mathbf{X}_A)^{-1} \mathbf{a}_A \right)^{1/2} > \left( \mathbf{a}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{a}_M \right)^{1/2}.$$

Conservative approach: Berk et al (2013)



## Estimation of the noise level

- Theoretical results for the Lasso suggests

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \{ \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 / (2n) + \lambda \|\beta\|_1 \}$$

with  $\lambda = \sigma \lambda_0$ ,  $\lambda_0 = \eta^{-1} \sqrt{(2/n) \log(p/k)}$ ,  $1 \leq k \ll n / \log p$

- A “naive estimate” of  $\sigma$  (Sun-Z, 2010): A recursive solution of

$$\hat{\sigma}^2 = \|\mathbf{Y} - \mathbf{X}\hat{\beta}(\hat{\sigma}\lambda_0)\|_2^2 / n$$

- Convex minimization formulation for the same  $\hat{\sigma}$  (Antoniadis, 2010):

$$\{\hat{\beta}, \hat{\sigma}\} = \arg \min_{\{\beta, \sigma\}} \{ \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 / (2\sigma n) + \sigma/2 + \lambda_0 \|\beta\|_1 \}$$

- Theory for this scaled Lasso (Sun-Z, 2012, 2013):

$$\begin{aligned} |\hat{\sigma}/\sigma^* - 1| &\leq O(1) C_{\text{est},1} (s^* + k) \lambda_0^2 \\ (s^* + k) \log p &\ll n^{1/2} \Rightarrow \sqrt{n}(\hat{\sigma}/\sigma - 1) \rightarrow N(0, 1/2) \end{aligned}$$

where  $(\sigma^*)^2 = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 / n = \sigma^2 \chi_n^2$

- An equivalent “square-root Lasso” formulation of  $\hat{\beta}$  (Belloni et al, 2011)
  - An earlier scaled Lasso proposal (Städler et al, 2010) yields
- $$\|\hat{\sigma}/\sigma^* - 1\| \lesssim s^* \lambda_0^2 + \lambda_0 \|\beta/\sigma\|_1$$

## Bias correction/Statistical inference after regularized estimation

- Linear bias correction from an initial estimator, e.g.  $\hat{\beta}^{(init)} = \hat{\beta}^{(lasso)}$ :

$$\hat{\beta}_j = \hat{\beta}_j^{(init)} + \mathbf{w}_j^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}^{(init)})$$

- Error decomposition:

$$\hat{\beta}_j - \beta_j = \mathbf{w}_j^\top \boldsymbol{\varepsilon} - (\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j)^\top (\hat{\beta}^{(init)} - \boldsymbol{\beta})$$

- The LS property:  $\mathbf{X}^\top \mathbf{w}_j = \mathbf{e}_j$ , i.e.  $\mathbf{w}_j \propto \mathbf{X}_j^\perp \perp \mathbf{X}_k \forall k \neq j$
- A relaxed LS property is sufficient
  - Noise factor:  $\tau_j = \|\mathbf{w}_j\|_2$ , hopefully  $\tau_j \asymp n^{-1/2}$
  - Bias factor:  $\eta_j = \|\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j\|_\infty / \|\mathbf{w}_j\|_2$ , hopefully  $\eta_j \leq \sqrt{2 \log p}$
  - Asymptotic theory based on  $\ell_\infty$ - $\ell_1$  split:

$$\eta_j \|\hat{\beta}^{(init)} - \boldsymbol{\beta}\|_1 \ll 1 \Rightarrow (\hat{\beta}_j - \beta_j) / (\hat{\sigma} \|\mathbf{w}_j\|_2) \rightarrow N(0, 1)$$

- Sample size requirement:  $n \gg (s^* \log p)^2$  under “proper conditions”
- Bonferroni adjustment:  $\eta_j \|\hat{\beta}^{(init)} - \boldsymbol{\beta}\|_1 \ll \sqrt{\log p}$
- Low-dimensional projection estimator (LDPE, Z-Zhang, 14; arXiv 11)

## Bias correction (details)

- Error decomposition:

$$\hat{\beta}_j - \beta_j = \mathbf{w}_j^\top \boldsymbol{\varepsilon} - (\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j)^\top (\hat{\boldsymbol{\beta}}^{(init)} - \boldsymbol{\beta})$$

- A relaxed LS property is sufficient
  - Noise factor:  $\tau_j = \|\mathbf{w}_j\|_2$ , hopefully  $\tau_j \asymp n^{-1/2}$
  - Bias factor:  $\eta_j = \|\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j\|_\infty / \|\mathbf{w}_j\|_2$ , hopefully  $\eta_j \leq \sqrt{2 \log p}$
  - Asymptotic theory based on  $\ell_\infty$ - $\ell_1$  split:

$$\eta_j \|\hat{\boldsymbol{\beta}}^{(init)} - \boldsymbol{\beta}\|_1 \ll 1 \Rightarrow (\hat{\beta}_j - \beta_j) / (\hat{\sigma} \|\mathbf{w}_j\|_2) \rightarrow N(0, 1)$$

- Sample size requirement:  $\lambda_{univ} = \sigma \sqrt{(2/n) \log p}$ 
  - Sparsity assumption:  $\sum_j \min(|\beta_j| / \lambda_{univ}, 1) \leq s^*$
  - $\ell_1$  error bound:  $\|\hat{\boldsymbol{\beta}}^{(init)} - \boldsymbol{\beta}\|_1 \leq C_{\text{est},1} s^* \lambda_{univ}$
  - $\eta_j \|\hat{\boldsymbol{\beta}}^{(init)} - \boldsymbol{\beta}\|_1 \leq C_{\text{est},1} s^* \sigma (\log p) \sqrt{2/n}$
  - Sample size requirement:  $n \gg (s^* \log p)^2$
- Bonferroni adjustment for simultaneous interval estimation of all  $\beta_j$ 
  - $\eta_j \|\hat{\boldsymbol{\beta}}^{(init)} - \boldsymbol{\beta}\|_1 \ll \sqrt{\log p}$
  - Sample size requirement:  $n \gg (s^*)^2 \log p$

## Finding a “score vector” for LD projection

- Under the sample size and proper regularity conditions,
  - $(\hat{\beta}_j - \beta_j)/\tau_j \rightarrow N(0, \sigma^2)$  with  $\tau_j = \|\mathbf{w}_j\|_2$
  - when  $\eta_j = \|\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j\|_\infty / \|\mathbf{w}_j\|_2 \leq C\sqrt{2\log p}$
- When  $\text{rank}(\mathbf{X}) = p$ ,
  - $\mathbf{w}_j = \mathbf{X}_j^\perp / \|\mathbf{X}_j^\perp\|_2^2$  yields the LSE
  - $\eta_j = 0$  and  $\tau_j = 1/\|\mathbf{X}_j^\perp\|_2$
- When  $p > n$ ,
  - $\mathbf{X}_j^\perp = 0$  when  $\mathbf{X}$  is in “general position”
  - $\mathbf{z}_j$  = a relaxed version of  $\mathbf{X}_j^\perp$
  - $\mathbf{w}_j = \mathbf{z}_j / \mathbf{z}_j^\top \mathbf{X}_j$  to ensure  $\mathbf{w}_j^\top \mathbf{X}_j = 1$ ; Alternatively,  $\mathbf{w}_j = \mathbf{z}_j / \|\mathbf{z}_j\|_2^2$
- For random designs with  $\mathbb{E} \mathbf{X}^\top \mathbf{X} / n = \boldsymbol{\Sigma}$ 
  - Regression model

$$\mathbf{X}_j = \mathbf{X}_{-j} \boldsymbol{\gamma}_{-j} + \mathbf{z}_j^\circ, \quad \sigma_j^2 = \mathbb{E} \|\mathbf{z}_j^\circ\|_2^2 / n = 1 / (\boldsymbol{\Sigma}^{-1})_{jj}$$

- Population version of  $\mathbf{X}_j^\perp$ :  $\mathbf{z}_j^\circ = \mathbf{X} \mathbf{u}^\circ \propto \mathbf{X} \boldsymbol{\Sigma}^{-1} \mathbf{e}_j$ ,  $\mathbf{u}_j^\circ = 1$ ,  $\mathbf{u}_{-j}^\circ = -\boldsymbol{\gamma}_{-j}$
- When  $\mathbf{z}_j^\circ$  is unknown:  $\mathbf{z}_j \approx \mathbf{z}_j^\circ$ , or  $\mathbf{z}_j = \mathbf{X} \mathbf{u}$  with  $\mathbf{u} \approx \mathbf{u}^\circ$
- Expected:  $\tau_j = (1 + o(1)) n^{-1/2} / \sigma_j$
- Algorithms:

Choice 1:  $\mathbf{z}_j$  = residual vector of PLSE( $\mathbf{X}_{-j}, \mathbf{X}_j$ )

Choice 2:  $\mathbf{z}_j = \arg \min_{\mathbf{z}} \left\{ \|\mathbf{z}\|_2^2 : \mathbf{z}^\top \mathbf{X}_j = n, \|\mathbf{z}^\top \mathbf{X}_{-j} / n\|_\infty \leq \lambda' \right\}$

## Finding a score vector with Lasso/PLSE

- Under the sample size and proper regularity conditions,
  - $(\hat{\beta}_j - \beta_j)/\tau_j \rightarrow N(0, \sigma^2)$  with  $\tau_j = \|\mathbf{w}_j\|_2$
  - when  $\eta_j = \|\mathbf{X}^T \mathbf{w}_j - \mathbf{e}_j\|_\infty / \|\mathbf{w}_j\|_2 \leq C\sqrt{2 \log p}$
  - $\mathbf{w}_j = \mathbf{z}_j / \mathbf{z}_j^T \mathbf{X}_j$ ,  $\mathbf{X}_j^T \mathbf{w}_j = 1$
- Lasso:  $\mathbf{z}_j = \mathbf{X}_j - \mathbf{X}_{-j} \hat{\gamma}_{-j}$

$$\hat{\gamma}_{-j} = \arg \min_{\gamma} \left\{ \|\mathbf{X}_j - \mathbf{X}_{-j} \gamma\|_2^2 / (2n) + \lambda' \|\gamma\|_1 \right\}$$

- KKT:  $\mathbf{X}_{-j}^T (\mathbf{X}_j - \mathbf{X}_{-j} \hat{\gamma}_{-j}) / n = \mathbf{X}_{-j}^T \mathbf{z}_j / n \in \lambda' \partial \|\hat{\gamma}_{-j}\|_1$
- $\mathbf{X}_j^T \mathbf{z}_j / n = \|\mathbf{z}_j\|_2^2 / n + (\mathbf{X}_{-j} \hat{\gamma}_{-j})^T \mathbf{z}_j / n = \|\mathbf{z}_j\|_2^2 / n + \lambda' \|\hat{\gamma}_{-j}\|_1$
- $\tau_j = \|\mathbf{w}_j\|_2 = \|\mathbf{z}_j\|_2 / |\mathbf{z}_j^T \mathbf{X}_j| \leq 1 / \|\mathbf{z}_j\|_2$
- $\eta_j = \|\mathbf{X}_{-j}^T \mathbf{w}_j\|_\infty / \|\mathbf{w}_j\|_2 = \|\mathbf{X}_{-j}^T \mathbf{z}_j\|_\infty / \|\mathbf{z}_j\|_2 = \lambda' n / \|\mathbf{z}_j\|_2 \leq \sqrt{2 \log p}$
- $\lambda' = \|\mathbf{z}_j / n^{1/2}\|_2 \sqrt{(2/n) \log p} = \hat{\sigma}_j \lambda_0$
- This suggests  $\eta_j = \sqrt{2 \log p}$  for the scaled Lasso
- $\eta_j = \|\mathbf{X}_{-j}^T \mathbf{z}_j\|_\infty / \|\mathbf{z}_j\|_2 \downarrow \lambda'$
- $\tau_j \leq 1 / \|\mathbf{z}_j\|_2 \uparrow \lambda'$

## Finding a score vector by quadratic programming

- Under proper regularity conditions,
  - $(\hat{\beta}_j - \beta_j)/\tau_j \rightarrow N(0, \sigma^2)$  with  $\tau_j = \|\mathbf{w}_j\|_2$
  - when  $\eta_j = \|\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j\|_\infty / \|\mathbf{w}_j\|_2 \leq \sqrt{2 \log p}$
  - $\mathbf{w}_j = \mathbf{z}_j / \mathbf{z}_j^\top \mathbf{X}_j$
  - $\eta_j = \|\mathbf{X}_{-j}^\top \mathbf{z}_j\|_\infty / \|\mathbf{z}_j\|_2$  as  $\mathbf{X}_j^\top \mathbf{w}_j = 1$
  - $\tau_j = \|\mathbf{z}_j\|_2 / |\mathbf{z}_j^\top \mathbf{X}_j| = \|\mathbf{z}_j\|_2 / n$  when  $\mathbf{X}_j^\top \mathbf{z}_j = n$
- Quadratic programming: minimize  $\tau_j$  given  $\eta_j \leq \sqrt{2 \log p}$

$$\begin{aligned} \mathbf{z}_j &= \arg \min_{\mathbf{z}} \left\{ \tau_j : \mathbf{z}^\top \mathbf{X}_j = n, \eta_j \leq \sqrt{2 \log p} \right\} \\ &= \arg \min_{\mathbf{z}} \left\{ \|\mathbf{z}\|_2^2 : \mathbf{z}^\top \mathbf{X}_j = n, \|\mathbf{X}_{-j}^\top \mathbf{z}\|_\infty \leq \|\mathbf{z}\|_2 \sqrt{2 \log p} \right\} \end{aligned}$$

- Z-Zhang (14; arXiv 11), Javanmard-Montanari (14)
- Near estimability:
  - If  $\eta_j \leq \sqrt{2 \log p}$  is attainable, or  $C\sqrt{2 \log p}$ ,  $\beta_j$  is “nearly estimable”
  - Otherwise,  $\beta_j$  is not nearly estimable
  - For Gaussian designs,  $\eta_j \leq \sqrt{2 \log p}$  is attainable with large probability

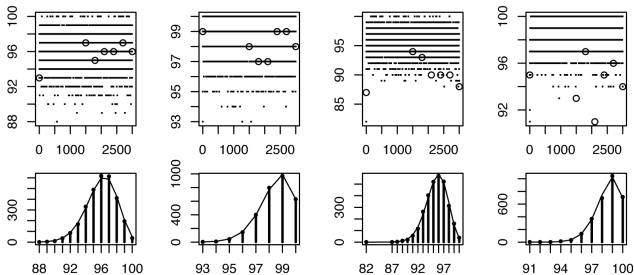
## Some simulation results:

- $n = 200$ ,  $p = 3000$ ,  $\sigma = 1$ ,  $\lambda_{univ} = \sqrt{(2/n) \log p} = 0.283$ ,  
 $\beta_j = 3\lambda_{univ} = 0.849$  for  $j = 1500, 1800, \dots, 3000$ , and  $\beta_j = 3\lambda_{univ}/j^\alpha$   
 otherwise;  $\beta_j \neq 0$  for all  $j$
- $(s, s * (\log p)/\sqrt{n}) = (8.93, 5.05)$  and  $(29.24, 16.55)$  respectively for  
 $\alpha = 1$  and  $2$ , while the theory requires  $s(\log p)/\sqrt{n} \rightarrow 0$ , where  
 $s = \sum_j \min(|\beta_j|/\lambda_{univ}, 1)$ .
- Generate  $(\tilde{\mathbf{X}}, \mathbf{X}, \varepsilon)$  in each replication, where  $\tilde{\mathbf{X}}$  has iid  $N(0, \mathbf{\Sigma})$  rows with  
 $\mathbf{\Sigma} = (\rho^{|j-k|})_{p \times p}$  and  $\mathbf{X}$  is the column normalized version of  $\tilde{\mathbf{X}}$
- Four settings, labeled (A), (B), (C), and (D), respectively, with  
 $(\alpha, \rho) = (2, 1/5)$ ,  $(1, 1/5)$ ,  $(2, 4/5)$ , and  $(1, 4/5)$
- Case (D) is most difficult

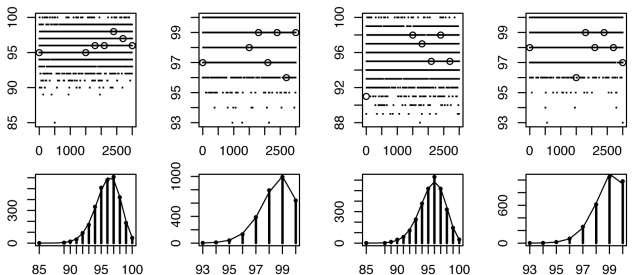
		(A)	(B)	(C)	(D)
all $\beta_j$	LDPE	0.9597	0.9845	0.9556	0.9855
	R-LDPE	0.9595	0.9848	0.9557	0.9885
maximal $\beta_j$	LDPE	0.9571	0.9814	0.9029	0.9443
	R-LDPE	0.9614	0.9786	0.9414	0.9786

**Table:** Mean coverage probability of LDPE and R-LDPE.

## LDPE



## Restricted LDPE





## Bias correction (summary)

- Linear bias correction from an initial estimator, e.g.  $\hat{\beta}^{(init)} = \hat{\beta}^{(lasso)}$ :

$$\hat{\beta}_j = \hat{\beta}_j^{(init)} + \mathbf{w}_j^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}^{(init)})$$

- Error decomposition:

$$\hat{\beta}_j - \beta_j = \mathbf{w}_j^\top \boldsymbol{\varepsilon} - (\mathbf{X}^\top \mathbf{w}_j - \mathbf{e}_j)^\top (\hat{\beta}^{(init)} - \boldsymbol{\beta})$$

- Asymptotic normality  $(\hat{\beta}_j - \beta_j)/(\hat{\sigma} \|\mathbf{w}_j\|_2) \rightarrow N(0, 1)$  when

- $\mathbf{w}_j = \mathbf{z}_j / \mathbf{z}_j^\top \mathbf{X}_j$
- $\eta_j = \|\mathbf{X}_{-j}^\top \mathbf{z}_j\|_\infty / \|\mathbf{z}_j\|_2 \leq \sqrt{2 \log p}$
- $\sum_j \min(|\beta_j| / \lambda_{univ}, 1) \leq s^*$
- $s^* \log p \ll n^{1/2}$

- Algorithms:

Choice 1:  $\mathbf{z}_j$  = residual vector of PLSE( $\mathbf{X}_{-j}$ ,  $\mathbf{X}_j$ )

Choice 2:  $\mathbf{z}_j = \arg \min_{\mathbf{z}} \{ \|\mathbf{z}\|_2^2 : \mathbf{z}^\top \mathbf{X}_j = n, \|\mathbf{X}_{-j}^\top \mathbf{z}\|_\infty \leq \|\mathbf{z}\|_2 \sqrt{2 \log p} \}$

- Approximate confidence interval for  $\beta_j$  is available if either Choice 2 is feasible or Choice 1 yields a feasible solution for Choice 2
- Statistical inference for  $\beta_j$  is “difficult” conditionally on  $\mathbf{X}$  if the algorithms are infeasible, i.e. failing to be nearly estimable

## Connection to semiparametric inference (Z, 11)

- Semiparametric inference:

Parametric component + NP component

- Low-dimensional statistical inference with high-dimensional data, or “semi-LD inference”:

LD component + HD component

- General methodology:

HD estimation  $\Rightarrow$  semi-LD inference  
is parallel to

NP estimation  $\Rightarrow$  semiparametric inference

- We borrow ideas from Engle et al (81), Bickel (82), Wahba (1984), Chen (85,88), Rice (86), Heckman (86), Bickel et al (90), ...

Fisher information for  $\theta = \tau(\beta)$  in small submodels (Stein, 56)

- Suppose  $\mathbf{X}$  has iid rows with  $\Sigma = \mathbb{E}(\mathbf{X}^T \mathbf{X} / n)$
- Let

$$\mathbf{a} = (\partial / \partial \beta) \tau(\beta)$$

- The minimum Fisher information in the submodel  $\beta = \mathbf{u}\theta$ :

$$F_\theta = \min_{\mathbf{u}^T \mathbf{a} = 1} \mathbf{u}^T \Sigma \mathbf{u} / \sigma^2 = (\mathbf{a}^T \Sigma^{-1} \mathbf{a})^{-1} / \sigma^2$$

- The least favorable submodel:

$$\mathbf{u}^o = \Sigma^{-1} \mathbf{a} (\mathbf{a}^T \Sigma^{-1} \mathbf{a})^{-1}$$

- One-step correction: With  $\mathbf{u} \approx \mathbf{u}^o$  and  $\mathbf{u}^T \mathbf{a} = 1$ ,

$$\hat{\theta} = \tau(\beta^{(init)}) + \arg \max_{\phi} \log\text{-lik}(\beta^{(init)} + \mathbf{u}\phi)$$

- Le Cam (69), Bickel (82), Schick (86), Klaassen (87), Bickel et al (90); Z(11, Oberwolfach), Z-Zhang (14), van de Geer et al (2014), van de Geer (2014); Sun-Z (12), Ren et al (13)

The LDPE attains the marginal minimum Fisher information:

- Random  $\mathbf{X}$  with iid sub-Gaussian rows;  $\mathbf{\Sigma} = \mathbb{E}\mathbf{X}^\top \mathbf{X}/n$ ,  $\mathbf{D} = \text{diag}(\mathbf{\Sigma})$ ,

$$s = \sum_{j=1}^p \min(1, \mathbf{D}_j \beta_j^2 / (\sigma^2 (2/n) \log p))$$

- Assume that  $\mathbf{X} \hat{\mathbf{D}}^{-1/2}$  is subGaussian and

$$\lambda_{\min}(\mathbf{D}^{-1/2} \mathbf{\Sigma} \mathbf{D}^{-1/2}) \geq c_*, \quad s \log p \ll n^{1/2}$$

- Then, the minimum Fisher information is  $F_j = 1/\{\sigma^2(\mathbf{\Sigma}^{-1})_{jj}\}$ , and

$$(n \hat{F}_j)^{1/2}(\hat{\beta}_j - \beta_j) \rightarrow N(0, 1), \quad \hat{F}_j / F_j \rightarrow 1$$

- Interpretation: the conditional inference “works” with high marginal probability
- Z (11), van de Geer-Bühlmann-Ritov (14)

The sample size condition is rate optimal

- Suppose  $\mathbf{X}$  has iid  $N(0, \mathbf{\Sigma})$  rows under  $\mathbb{P}$ . Let

$$\mathcal{P} = \left\{ \mathbb{P} : c_* \leq \text{eigen}(\mathbf{D}^{-1/2} \mathbf{\Sigma} \mathbf{D}^{-1/2}) \leq c^*, s \log p \leq c'_* n \right\}$$

- For certain  $c_*'' > 0$  depending on  $\{c_*, c^*, c'_*\}$  only,

$$\inf_{\hat{\theta}} \sup_{\mathcal{P}} \mathbb{P} \left\{ n^{1/2} |\hat{\theta} - \theta| \geq c_*'' s (\log p) / n^{1/2} \right\} > 0$$

- Le Cam (73), Ren et al (13)
- When  $n^{1/2} \ll s \log p \ll n$ ,  $n^{-1/2}$  rate is impossible without some additional assumption
- Javanmard-Montanari (14):  $\mathbf{\Sigma} = \mathbf{I}$  (known) with  $n \gg s \log p$

## Some additional references

- Leeb-Pötscher (06)
- Laber-Murphy (11)
- Bühlmann (13), Belloni et al (14)
- Meinshausen-Bühlmann (10)
- Lockhart et al (14), Lee et al (14)

## De-biasing regularized estimators

- Low-dimensional projection from  $\hat{\beta}^{(init)}$ , e.g.  $\hat{\beta}^{(lasso)}$ :

$$\hat{\beta}_j = \hat{\beta}_j^{(init)} + (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \mathbf{z}_j^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}^{(init)}),$$

i.e.  $\mathbf{w}_j = \mathbf{z}_j / \mathbf{z}_j^\top \mathbf{X}_j$  with  $\mathbf{w}_j^\top \mathbf{X}_j = 1$

- Asymptotic theory:

$$\hat{\beta}_j - \beta_j = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \mathbf{z}_j^\top \boldsymbol{\varepsilon} - (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \sum_{k \neq j} \mathbf{z}_j^\top \mathbf{X}_k (\hat{\beta}_k^{(init)} - \beta_k)$$

Bias after de-biasing:  $\|\mathbf{z}_j\|_2^{-1} |\sum_{k \neq j} \mathbf{z}_j^\top \mathbf{X}_k (\hat{\beta}_k^{(init)} - \beta_k)| = o_P(1)$

- Standard error for homoscedastic  $\boldsymbol{\varepsilon}$ :  $\text{s.e.}_j = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|\mathbf{z}_j\| \sigma$
- Studentization:  $T_j = (\hat{\beta}_j - \beta_j) / \widehat{\text{s.e.}}_j$
- Conservative simultaneous interval/inference (Bonferroni):

$$\max_{j \in G} |T_j| \leq \Phi^{-1}(1 - \alpha/(2|G|))$$

- Bootstrap: more accurate simultaneous inference, heteroscedastic  $\boldsymbol{\varepsilon}$

## Some references

- Joint with Ruben Dezeure and Peter Bühlmann
- Residual bootstrap: Efron (79)
- Paired bootstrap and heteroscedasticity: Efron (79), Liu-Singh (92)
- Wild bootstrap and heteroscedasticity: Wu (86), Liu-Singh (92), Mammen (93), Chernozhukov et al. (13), Zhang-Cheng (16)
- Chatterjee-Lahiri (11)



**Bootstrapping the LDPE:**  $\hat{\beta}_j = \hat{\beta}_j^{(init)} + (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \mathbf{z}_j^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}^{(init)})$

- Notation:  $(\mathbf{v})_{\text{cent}} = \mathbf{v} - (\mathbf{1}\mathbf{1}^\top/n)\mathbf{v}$ , the Hadamard product  $\mathbf{u} \circ \mathbf{v}$
- Asymptotic theory
  - $\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}^{(init)}$ ,  $\hat{\boldsymbol{\varepsilon}}_{\text{cent}} = (\hat{\boldsymbol{\varepsilon}})_{\text{cent}}$
  - $\widehat{\mathbf{s.e.}}_j = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|\mathbf{z}_j\|_2 \|\hat{\boldsymbol{\varepsilon}}_{\text{cent}}\|_2 / \sqrt{n}$
  - $T_j = (\hat{\beta}_j - \beta_j) / \widehat{\mathbf{s.e.}}_j$
  - $\widehat{\mathbf{s.e.}}_{j,\text{robust}} = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|(\mathbf{z}_j \circ \hat{\boldsymbol{\varepsilon}})_{\text{cent}}\|_2$  for heteroscedastic  $\varepsilon$
  - $T_{j,\text{robust}} = (\hat{\beta}_j - \beta_j) / \widehat{\mathbf{s.e.}}_{j,\text{robust}}$

**Bootstrapping the LDPE:**  $\hat{\beta}_j = \hat{\beta}_j^{(init)} + (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \mathbf{z}_j^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}^{(init)})$

- Notation:  $(\mathbf{v})_{\text{cent}} = \mathbf{v} - (\mathbf{1}\mathbf{1}^\top/n)\mathbf{v}$ , the Hadamard product  $\mathbf{u} \circ \mathbf{v}$

- Asymptotic theory

- $\hat{\varepsilon} = \mathbf{Y} - \mathbf{X} \hat{\beta}^{(init)}$ ,  $\hat{\varepsilon}_{\text{cent}} = (\hat{\varepsilon})_{\text{cent}}$
- $\widehat{\text{s.e.}}_j = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|\mathbf{z}_j\|_2 \|\hat{\varepsilon}_{\text{cent}}\|_2 / \sqrt{n}$
- $T_j = (\hat{\beta}_j - \beta_j) / \widehat{\text{s.e.}}_j$
- $\widehat{\text{s.e.}}_{j,\text{robust}} = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|(\mathbf{z}_j \circ \hat{\varepsilon})_{\text{cent}}\|_2$  for heteroscedastic  $\varepsilon$
- $T_{j,\text{robust}} = (\hat{\beta}_j - \beta_j) / \widehat{\text{s.e.}}_{j,\text{robust}}$

- Residual bootstrap

- $\varepsilon^*$  iid uniform distribution from elements of  $\hat{\varepsilon}_{\text{cent}}$
- $\mathbf{Y}^* = \mathbf{X} \hat{\beta}^{(init)} + \varepsilon^*$
- $\hat{\beta}^{*(init)}$  = the Lasso estimator based on  $(\mathbf{X}, \mathbf{Y}^*)$
- $\hat{\beta}_j^* = \hat{\beta}_j^{*(init)} + (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \mathbf{z}_j^\top (\mathbf{Y}^* - \mathbf{X} \hat{\beta}^{*(init)})$
- $\widehat{\text{s.e.}}_j^* = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|\mathbf{z}_j\|_2 \|(\mathbf{Y}^* - \mathbf{X} \hat{\beta}^{*(init)})_{\text{cent}}\|_2 / \sqrt{n}$ ,
- $T_j^* = (\hat{\beta}_j^* - \hat{\beta}_j^{(lasso)}) / \widehat{\text{s.e.}}_j^*$
- $\widehat{\text{s.e.}}_{j,\text{robust}}^* = (\mathbf{z}_j^\top \mathbf{X}_j)^{-1} \|(\mathbf{z}_j \circ (\mathbf{Y}^* - \mathbf{X} \hat{\beta}^{*(init)}))_{\text{cent}}\|_2$
- $T_{j,\text{robust}}^* = (\hat{\beta}_j^* - \hat{\beta}_j^{(lasso)}) / \widehat{\text{s.e.}}_{j,\text{robust}}^*$

## Wild bootstrap and xyz-paired bootstrap

- Wild bootstrap
  - Draw iid  $W_i$  with  $\mathbb{E}W_i = 0$  and  $\mathbb{E}W_i^2 = 1$
  - Set  $\varepsilon^* = \mathbf{W} \circ \hat{\varepsilon}_{\text{cent}}, \mathbf{Y}^* = \mathbf{X}\hat{\beta}^{(init)} + \varepsilon^*$
  - Proceed as in the residual bootstrap, i.e. plug-in with  $(\mathbf{X}, \mathbf{Y}^*, \mathbf{z}_j)$

## Wild bootstrap and xyz-paired bootstrap

- Wild bootstrap
  - Draw iid  $W_i$  with  $\mathbb{E}W_i = 0$  and  $\mathbb{E}W_i^2 = 1$
  - Set  $\varepsilon^* = \mathbf{W} \circ \hat{\varepsilon}_{\text{cent}}, \mathbf{Y}^* = \mathbf{X}\hat{\beta}^{(init)} + \varepsilon^*$
  - Proceed as in the residual bootstrap, i.e. plug-in with  $(\mathbf{X}, \mathbf{Y}^*, \mathbf{z}_j)$
- xyz-paired bootstrap
  - Paired bootstrap: sample rows of  $(\mathbf{X}, \mathbf{Y})$  with replacement

## Wild bootstrap and xyz-paired bootstrap

- Wild bootstrap
  - Draw iid  $W_i$  with  $\mathbb{E}W_i = 0$  and  $\mathbb{E}W_i^2 = 1$
  - Set  $\varepsilon^* = \mathbf{W} \circ \hat{\varepsilon}_{\text{cent}}, \mathbf{Y}^* = \mathbf{X}\hat{\beta}^{(init)} + \varepsilon^*$
  - Proceed as in the residual bootstrap, i.e. plug-in with  $(\mathbf{X}, \mathbf{Y}^*, \mathbf{z}_j)$
- xyz-paired bootstrap
  - Paired bootstrap: sample rows of  $(\mathbf{X}, \mathbf{Y})$  with replacement
  - Construction of a proper linear model
    - $\hat{\varepsilon}_{\text{cent}} = (\mathbf{Y} - \mathbf{X}\hat{\beta}^{(init)})_{\text{cent}},$
    - $\hat{\mathbf{X}}_k = \mathbf{X}_k - (\hat{\varepsilon}_{\text{cent}}\hat{\varepsilon}_{\text{cent}}^\top / \|\hat{\varepsilon}_{\text{cent}}\|_2^2)\mathbf{X}_k; \hat{\mathbf{X}}_k^\top \hat{\varepsilon}_{\text{cent}} = 0$
    - $\hat{\mathbf{Y}} = \hat{\mathbf{X}}\hat{\beta}^{(init)} + \hat{\varepsilon}_{\text{cent}}$
    - $\hat{\mathbf{z}}_j = \mathbf{z}_j - (\hat{\varepsilon}_{\text{cent}}\hat{\varepsilon}_{\text{cent}}^\top / \|\hat{\varepsilon}_{\text{cent}}\|_2^2)\mathbf{z}_j; \hat{\mathbf{z}}_j^\top \hat{\varepsilon}_{\text{cent}} = 0$
  - $(\mathbf{X}^*, \mathbf{Y}^*, \mathbf{z}_j^*, j \in G)$ : iid sample of rows of  $(\hat{\mathbf{X}}, \hat{\mathbf{Y}}, \hat{\mathbf{z}}_j, j \in G)$
  - $T_j^*$  = the plug-in version based on  $(\mathbf{X}^*, \mathbf{Y}^*, \mathbf{z}_j^*)$
  - $T_{j,\text{robust}}^*$  = the plug-in version based on  $(\mathbf{X}^*, \mathbf{Y}^*, \mathbf{z}_j^*)$

## Bootstrap methods, a summary

- Residual bootstrap
  - $\varepsilon^*$  iid from elements of  $\widehat{\varepsilon}_{\text{cent}} = (\mathbf{Y} - \mathbf{X}\widehat{\beta}^{(\text{init})})_{\text{cent}}$
  - $\mathbf{Y}^* = \mathbf{X}\widehat{\beta}^{(\text{init})} + \varepsilon^*$
  - The plug-in estimates of  $T_j^*$  and  $T_{j,\text{robust}}^*$  based on  $(\mathbf{X}, \mathbf{Y}^*, \mathbf{z}_j)$
- Wild bootstrap
  - Draw iid  $W_i$  with  $\mathbb{E}W_i = 0$  and  $\mathbb{E}W_i^2 = 1$
  - $\mathbf{Y}^* = \mathbf{X}\widehat{\beta}^{(\text{init})} + \mathbf{W} \circ \widehat{\varepsilon}_{\text{cent}}$
  - The plug-in estimates of  $T_j^*$  and  $T_{j,\text{robust}}^*$  based on  $(\mathbf{X}, \mathbf{Y}^*, \mathbf{z}_j)$
- The xyz-paired bootstrap
  - $\widehat{\mathbf{X}} \perp \widehat{\varepsilon}_{\text{cent}}, \widehat{\mathbf{Y}} = \widehat{\mathbf{X}}\widehat{\beta}^{(\text{init})} + \widehat{\varepsilon}_{\text{cent}}, \widehat{\mathbf{z}}_j \perp \widehat{\varepsilon}_{\text{cent}}$
  - $(\mathbf{X}^*, \mathbf{Y}^*, \mathbf{z}_j^*, j \in G)$ : iid sample of rows of  $(\widehat{\mathbf{X}}, \widehat{\mathbf{Y}}, \widehat{\mathbf{z}}_j, j \in G)$
  - The plug-in estimates of  $T_j^*$  and  $T_{j,\text{robust}}^*$  based on  $(\mathbf{X}^*, \mathbf{Y}^*, \mathbf{z}_j^*)$
- No re-computation of  $\mathbf{z}_j$  in bootstrap replications

## Theoretical assumptions for simultaneous inference of $\beta_j, j \in G$ :

- (A1):  $\|\hat{\beta}^{(init)} - \beta\|_1 = o_P(1)/\sqrt{(\log p) \log(1 + |G|)}$
- (A2):  $\|\mathbf{z}_G^\top \mathbf{X}_{-j}/n\|_{\max} \lesssim \sqrt{(\log p)/n}$ ,  $\|\mathbf{z}_j\|_2^2/n \geq L_z$ ,  $\|\mathbf{z}_j\|_{2+\delta}^{2+\delta} \ll \|\mathbf{z}_j\|_2^{2+\delta}$
- (A3):  $\varepsilon_i$  independent,  $\mathbb{E}\varepsilon_i = 0$ ,  $\mathbb{E}\varepsilon_i^2 = \sigma_i^2 \geq L$ ,  $\mathbb{E}|\varepsilon_i|^{2+\delta} \leq C$
- (A4)  $\|\hat{\beta}^{*(init)} - \hat{\beta}^{(init)}\|_1 = o_{P^*}(1)/\sqrt{(\log p) \log(1 + |G|)}$  in probability
- (A5)  $\|\mathbf{X}\|_{\max} \leq C$
- (A6)  $\max_{j \in G} \|\mathbf{z}_j\|_{\infty} \leq K$ ,  $\delta = 2$ ,  $\log(|G|) = o(n^{1/7})$

## Theoretical assumptions for simultaneous inference of $\beta_j, j \in G$ :

- (A1):  $\|\hat{\beta}^{(init)} - \beta\|_1 = o_P(1)/\sqrt{(\log p) \log(1 + |G|)}$
- (A2):  $\|\mathbf{z}_G^\top \mathbf{X}_{-j}/n\|_{\max} \lesssim \sqrt{(\log p)/n}$ ,  $\|\mathbf{z}_j\|_2^2/n \geq L_z$ ,  $\|\mathbf{z}_j\|_{2+\delta}^{2+\delta} \ll \|\mathbf{z}_j\|_2^{2+\delta}$
- (A3):  $\varepsilon_i$  independent,  $\mathbb{E}\varepsilon_i = 0$ ,  $\mathbb{E}\varepsilon_i^2 = \sigma_i^2 \geq L$ ,  $\mathbb{E}|\varepsilon_i|^{2+\delta} \leq C$
- (A4)  $\|\hat{\beta}^{*(init)} - \hat{\beta}^{(init)}\|_1 = o_{P^*}(1)/\sqrt{(\log p) \log(1 + |G|)}$  in probability
- (A5)  $\|\mathbf{X}\|_{\max} \leq C$
- (A6)  $\max_{j \in G} \|\mathbf{z}_j\|_{\infty} \leq K$ ,  $\delta = 2$ ,  $\log(|G|) = o(n^{1/7})$

For proper PLSE as  $\hat{\beta}^{(init)}$  and under regularity conditions on  $\mathbf{X}$  (RE or weaker)

- (A3) and (A5) imply  $\|\mathbf{X}^\top \varepsilon/n\|_{\infty} \lesssim \lambda_{univ}$
- (A1) and (A4) hold when  $n \gg (s \log p)^2 \log(1 + |G|)$
- (A2) and (A6) hold if  $\mathbf{X}$  has iid rows with  $\text{Var}(X_{ij}|X_{i,-j}) \geq v_0 > 0$



## Consistency of the residual bootstrap

- Homoscedastic case:  $\mathbb{E}\varepsilon_i^2 = \sigma^2$  for all  $i \leq n$ 
  - Suppose conditions (A1)-(A5) holds. If  $|G| = O(1)$ , then

$$\sup_{t_j, j \in G} \left| \mathbb{P}^* \{ T_j^* \leq t_j, j \in G \} - \mathbb{P} \{ T_j \leq t_j, j \in G \} \right| = o_P(1)$$

with  $T_j \rightarrow N(0, 1)$  for each  $j \in G$

- If in addition (A6) holds, then

$$\sup_t \left| \mathbb{P}^* \{ \max_{j \in G} h(T_j^*) \leq t \} - \mathbb{P} \{ \max_{j \in G} h(T_j) \leq t \} \right| = o_P(1)$$

for  $h(t) = t$ ,  $h(t) = -t$  or  $h(t) = |t|$ .

- Heteroscedastic case: Suppose (A1)-(A5). Then,

$$\sup_t \left| \mathbb{P}^* \{ T_{j,\text{robust}}^* \leq t \} - \mathbb{P} \{ T_{j,\text{robust}} \leq t \} \right| = o_P(1)$$

with  $T_{j,\text{robust}} \rightarrow N(0, 1)$  for each  $j \in G$

## Consistency of the wild bootstrap and xyz-paired bootstrap

- Suppose conditions (A1)-(A5) holds. If  $|G| = O(1)$ , then

$$\sup_{t_j, j \in G} \left| \mathbb{P}^* \{T_j^* \leq t_j, j \in G\} - \mathbb{P} \{T_j \leq t_j, j \in G\} \right| = o_P(1)$$

with  $T_j \rightarrow N(0, 1)$  for each  $j \in G$

- If in addition (A6) holds and  $\log p \ll n^{1/2}$ , then

$$\sup_t \left| \mathbb{P}^* \{ \max_{j \in G} h(T_j^*) \leq t \} - \mathbb{P} \{ \max_{j \in G} h(T_j) \leq t \} \right| = o_P(1)$$

for  $h(t) = t$ ,  $h(t) = -t$  or  $h(t) = |t|$ .

The theorem is applicable in the heteroscedastic case

## Chernozukov et al (13)

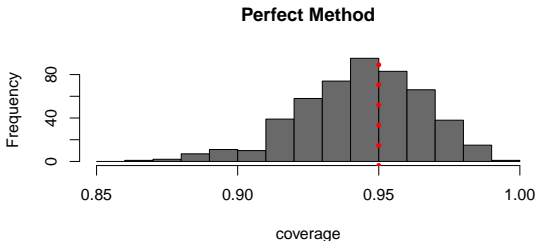
- $\xi \in \mathbb{R}^{n \times p}$  with independent rows,  $\mathbb{E}\xi = 0$ ,  $\mathbb{E}\xi^\top \xi / n = \Sigma^{(1)}$
- $\zeta \in \mathbb{R}^{n \times p}$ , Gaussian, with  $\mathbb{E}\zeta = 0$ ,  $\mathbb{E}\xi^\top \zeta / n = \Sigma^{(2)}$
- CLT for maxima:

$$\max_t \left| \mathbb{P}\left(\max_{j \leq p} \mathbf{1}^\top \xi \leq t\right) - \mathbb{P}\left(\max_{j \leq p} \mathbf{1}^\top \zeta \leq t\right) \right| = o(1)$$

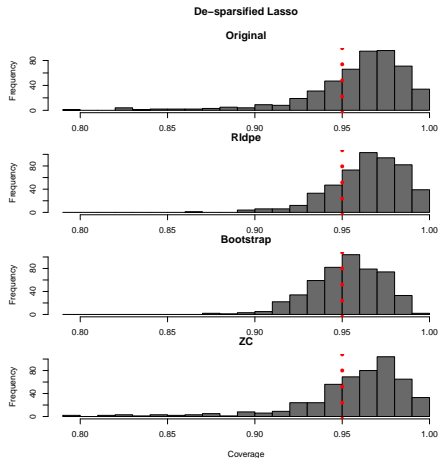
under the following conditions

- $c_1 \leq \Sigma_{jj}^{(k)} \leq c_2$ ,  $n^{-1} \sum_{i=1}^n \mathbb{E}\xi_{ij}^4 \leq M$
- $\log p \ll n^{1/7}$
- $\|\xi\|_{\max} + \|\zeta\|_{\max} = O_P(n^{1/2})/(\log(np))^{3/2}$
- $\|\Sigma^{(1)} - \Sigma^{(2)}\|_{\max} \ll 1/(\log p)^2$

## Some simulation results

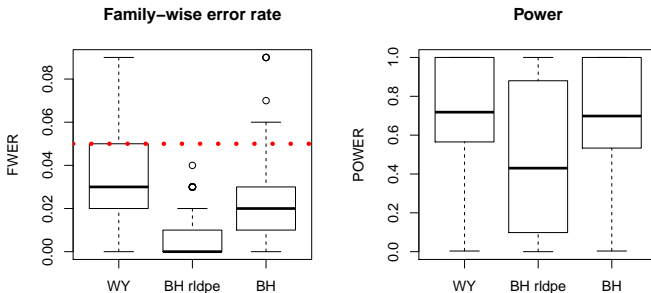


**Figure:** Histogram of the coverage probabilities of two sided 95% confidence intervals for 500 parameters. It illustrates how the results look like for a perfectly correct method for creating confidence intervals and one uses only 100 realizations to compute the probabilities.

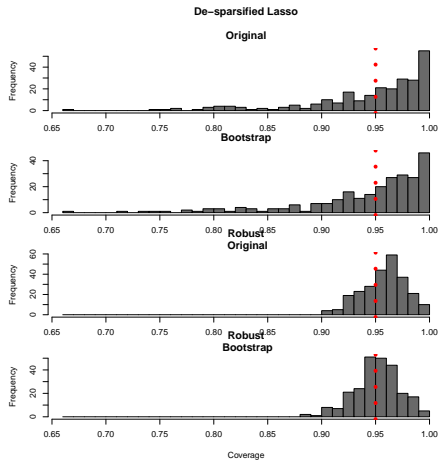


**Figure:** Histograms of the coverage probabilities of two-sided 95% confidence intervals for all 500 parameters in a linear model ( $n = 100, p = 500$ ), computed from 100 independent replications. Perfect performance would look like Figure 1. The fixed design matrix is of Toeplitz type, the single coefficient vector of type  $U(-2, 2)$  and **homoscedastic Gaussian errors**. The original estimator has more over-coverage and under-coverage than the bootstrapped estimator. The RLDPE estimator has little under-coverage, like the bootstrapped estimator, but it has too high coverage probabilities overall. The ZC approach to bootstrapping, which only bootstraps the linearized part of the estimator, doesn't show any improvements over the original de-sparsified Lasso.

## de-sparsified Lasso



**Figure:** Boxplot of the familywise error rate and the power for multiple testing for the de-sparsified Lasso. The target is controlling the FWER at level 0.05, highlighted by a red-dotted horizontal line. Two different approaches for multiple testing correction are compared, Westfall-Young (WY) and Bonferroni-Holm (BH). For Bonferroni-Holm, we make the distinction between the original method and the RLDPE approach. 300 linear models are investigated in total, where 50 Toeplitz design matrices are combined with 50 coefficient vectors for each of the 6 types  $U(0, 2)$ ,  $U(0, 4)$ ,  $U(-2, 2)$ , fixed 1, fixed 2, fixed 10. The variables belonging to the active set are chosen randomly. The errors in the linear model were chosen to be **homoscedastic Gaussian**. Each of the models has a data point for the error rate and the power in the boxplot. The error rate and power probabilities were calculated by averaging over 100 realizations.



**Figure:** The same plot as Figure 2 but for **heteroscedastic non-Gaussian errors** and without signal. The robust standard error estimation clearly outperforms the non-robust version. There seems to be hardly any difference between the bootstrap and the original estimator after choosing the standard error estimation.

## Summary

- Bootstrapping the de-biased Lasso turns out to improve the coverage of confidence intervals without increasing the confidence interval lengths
- The use of the conservative RLDPE (Zhang and Zhang, 2014) is not necessary: the bootstrap achieves reliable coverage, while for the original de-biased Lasso/LDPE, the RLDPE seems worthwhile to achieve reasonable coverage while paying a price in terms of efficiency
- Bootstrapping only the linearized part of the de-biased estimator as proposed by Zhang and Cheng (2016) is clearly sub-ideal in comparison to bootstrapping the entire estimator and using the plug-in principle as advocated here
- For multiple testing, the bootstrapped estimator had familywise error rates that were closer to the target level while Bonferroni-Holm adjustment is too conservative.
- The robust standard error turned out to be critical when dealing with heteroscedastic errors.



## Multivariate inference

- $\theta \in \mathbb{R}^d$
- The (minimum) Fisher information  $\mathbf{F} = \mathbf{F}_\theta$  for the estimation of  $\theta$
- $\hat{\theta}_n = n^{-1} \sum_{i=1}^n \psi(\text{data}_i) + n^{-1/2} \text{Rem}$
- $n^{1/2} \left\{ n^{-1} \sum_{i=1}^n \psi(\text{data}_i) - \theta \right\} \approx N(0, \mathbf{F}^{-1}) \in \mathbb{R}^d$
- Asymptotic  $\chi^2$ -confidence region:  $\text{Rem} \ll \sqrt{d} \asymp \|N(0, \mathbf{F}^{-1})\|_2$
- Efficient inference:  $\text{Rem} = o(1)$ 
  - Testing  $H_0 : \theta = 0$
  - $\chi_d \approx \sqrt{d} + N(0, 1/2)$
- Inference is much harder when  $p \gg n \gg d \rightarrow \infty$

Background: Inference about a complex effect

- Additive regression
- Varying coefficient models with longitudinal data
- Nonparametric graphical models

Fundamental problem: Inference about a large group of coefficients

$$\beta_G, \quad |G| \rightarrow \infty,$$

in the linear model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n), \quad \boldsymbol{\beta} \in \mathbb{R}^p.$$

## Direct extension from univariate bias correction

- We have

$$\hat{\beta}_j = \beta_j + \mathbf{w}_j^T \boldsymbol{\varepsilon} + \text{Rem}_j / \sqrt{n} \quad \text{Rem}_j \asymp (s \log p) / \sqrt{n}$$

- This yields asymptotic optimality in the sense of

$$\hat{\beta}_G = \beta_G + (\mathbf{w}_j, j \in G)^T \boldsymbol{\varepsilon} + \text{Rem}_G / \sqrt{n},$$

- For  $\chi^2$ -confidence regions, the sample size requirement is still

$$(s \log p) \ll \sqrt{n}$$

- However, for the approximate  $\chi^2$  test, direct application of the univariate bound with Cauchy-Schwarz leads to an undesirable requirement

$$\|\text{Rem}_G\|_2 \lesssim (s \log p) \sqrt{|G|/n} \rightarrow 0$$

## Debiasing under group sparsity

- For approximate  $\chi^2$ -confidence regions, it is sufficient to have

$$\frac{s + g \log M}{\sqrt{|G|n}} \rightarrow 0$$

where  $g = \#\{j : \beta_{G_j} \neq 0\}$ ,  $s = \sum_{\beta_{G_j} \neq 0} |G_j|$ ,  $M = \#$  of groups

- For approximate  $\chi^2$ -tests, it is sufficient to have

$$\frac{s + g \log M}{\sqrt{n}} \rightarrow 0$$

- Larger  $|G|$  the better under proper assumptions
- Twin benefits of group sparsity: factor  $\sqrt{|G|}$ ;  $(s + g \log M)$  vs.  $s \log p$

## Working Assumptions

- Group structure:

$$\cup_{j=1}^M G_j = \{1, \dots, p\}, \quad G_j \cap G_k = \emptyset$$

- Group sparsity:

$$\text{supp}(\beta^*) \subset G_{S^*}, \quad |S^*| \leq g, \quad |G_{S^*}| \leq s$$

- The existence of estimators satisfying

$$\left| \frac{\hat{\sigma}}{\sigma^*} - 1 \right| + \frac{1}{n^{1/2}} \sum_{j=1}^M \omega_j \left\| \mathbf{x}_{G_j} \hat{\beta}_{G_j}^{(init)} - \mathbf{x}_{G_j} \beta_{G_j}^* \right\|_2 \lesssim \left( \frac{s + g \log M}{n} \right)$$

with  $\omega_j \approx \sqrt{(|G_j| + 2 \log M)/n}$

- Group Lasso: Yuan-Lin (06), Bach (08), Koltchinskii and Yuan (08), Obozinski et al (08), Huang-Zhang (10), Lounici et al (11)

## Bias correction

- Model:

$$\mathbf{Y} = \mathbf{X}_G \boldsymbol{\beta}_G + \sum_{G_k \not\subseteq G} \mathbf{X}_{G_k \setminus G} \boldsymbol{\beta}_{G_k \setminus G} + \boldsymbol{\varepsilon}$$

- Bias correction with a projection  $\mathbf{P}_G$ :

$$\hat{\boldsymbol{\beta}}_G = (\mathbf{P}_G \mathbf{X}_G)^\dagger \mathbf{P}_G \left( \mathbf{Y} - \sum_{G_k \not\subseteq G} \mathbf{X}_{G_k \setminus G} \hat{\boldsymbol{\beta}}_{G_k \setminus G}^{(init)} \right)$$

- Error decomposition for confidence region:

$$\mathbf{P}_G (\mathbf{X}_G \hat{\boldsymbol{\beta}}_G - \mathbf{X}_G \boldsymbol{\beta}_G^*) = \mathbf{P}_G \boldsymbol{\varepsilon} - \sum_{G_k \not\subseteq G} \mathbf{P}_G \mathbf{X}_{G_k \setminus G} \left( \hat{\boldsymbol{\beta}}_{G_k \setminus G}^{(init)} - \boldsymbol{\beta}_{G_k \setminus G}^* \right)$$

- The LS property  $\mathbf{P}_G \mathbf{X}_{G_k \setminus G} = \mathbf{0}$  leads to

$$\|\mathbf{P}_G (\mathbf{X}_G \hat{\boldsymbol{\beta}}_G - \mathbf{X}_G \boldsymbol{\beta}_G^*)\|_2^2 / \sigma^2 \sim \chi_{\text{rank}(\mathbf{P}_G)}^2$$

## Relaxed LS property

- Error decomposition for confidence region:

$$\begin{aligned} \mathbf{P}_G(\mathbf{X}_G \hat{\boldsymbol{\beta}}_G - \mathbf{X}_G \boldsymbol{\beta}_G^*) &= \mathbf{P}_G \boldsymbol{\varepsilon} - \text{Rem}_G \\ \text{Rem}_G &= \sum_{G_k \not\subseteq G} \mathbf{P}_G \mathbf{X}_{G_k \setminus G} \left( \hat{\boldsymbol{\beta}}_{G_k \setminus G}^{(init)} - \boldsymbol{\beta}_{G_k \setminus G}^* \right) \end{aligned}$$

- Test statistic for  $H_0 : \mathbf{X}_G \boldsymbol{\beta}_G^* = 0$ :

$$T_G = \|\mathbf{P}_G \mathbf{X}_G \hat{\boldsymbol{\beta}}_G\|_2 / \hat{\sigma}$$

- Upper bound for the residual bias after bias correction:

$$\|\text{Rem}_G\| \leq \sum_{G_k \not\subseteq G} \|\mathbf{P}_G \mathbf{Q}_{G_k \setminus G}\|_s \left\| \mathbf{X}_{G_k \setminus G} \left( \hat{\boldsymbol{\beta}}_{G_k \setminus G}^{(init)} - \boldsymbol{\beta}_{G_k \setminus G}^* \right) \right\|_2$$

- Sufficient condition for approximate  $\chi^2$ -test:

$$\left| T_G - \frac{\|\mathbf{P}_G \boldsymbol{\varepsilon}\|_2}{\sigma} \right| \leq \frac{\|\text{Rem}_G\|_2}{\hat{\sigma}} + \left| \frac{\sigma}{\hat{\sigma}} - 1 \right| \frac{\|\mathbf{P}_G \boldsymbol{\varepsilon}\|_2}{\sigma} = o(1)$$

- Sufficient condition for approximate  $\chi^2$ -test:

$$\left| T_G - \frac{\|\mathbf{P}_G \boldsymbol{\varepsilon}\|_2}{\sigma} \right| \leq \frac{\|\text{Rem}_G\|_2}{\hat{\sigma}} + \left| \frac{\sigma}{\hat{\sigma}} - 1 \right| \frac{\|\mathbf{P}_G \boldsymbol{\varepsilon}\|_2}{\sigma} = o(1)$$

iff

$$\|\text{Rem}_G\|_2 = o(1), \quad \left| \frac{\hat{\sigma}}{\sigma^*} - 1 \right| = \frac{o(1)}{\sqrt{|G|}}$$

- Upper bound for the residual bias

$$\|\text{Rem}_G\| \leq \sum_{G_k \not\subseteq G} \|\mathbf{P}_G \mathbf{Q}_{G_k \setminus G}\|_s \left\| \mathbf{X}_{G_k \setminus G} \left( \hat{\boldsymbol{\beta}}_{G_k \setminus G}^{(init)} - \boldsymbol{\beta}_{G_k \setminus G}^* \right) \right\|_2$$

- Working Assumption:

$$\left| \frac{\hat{\sigma}}{\sigma^*} - 1 \right| + \frac{1}{n^{1/2}} \sum_{j=1}^M \omega_j \left\| \mathbf{X}_{G_j} \hat{\boldsymbol{\beta}}_{G_j}^{(init)} - \mathbf{X}_{G_j} \boldsymbol{\beta}_{G_j}^* \right\|_2 \lesssim \left( \frac{s + g \log M}{n} \right)$$



Suppose

$$\begin{aligned}\|\mathbf{X}_{G_k \setminus G} \mathbf{b}_{G_k \setminus G}\|_2 &\leq M_k \|\mathbf{X}_{G_k} \mathbf{b}_{G_k}\|_2 \\ \|\mathbf{P}_G \mathbf{Q}_{G_k \setminus G}\|_s &\leq \omega'_k\end{aligned}$$

for all  $\mathbf{b}_{G_k}$  and all  $k$ . Suppose that the Working Assumptions hold with

$$\frac{|G|}{n} + \frac{s + g \log M}{n^{1/2}} \left( \frac{|G|^{1/2}}{n^{1/2}} + \max_{G_k \not\subseteq G} M_k \frac{\omega'_k}{\omega_k} \right) = o(1)$$

Then,

$$\mathbf{P}_G \mathbf{X}_G (\hat{\beta}_G - \beta_G^*) = \mathbf{P}_G \boldsymbol{\varepsilon} + o(1)$$

and

$$\left| T_G - \frac{\mathbf{P}_G \boldsymbol{\varepsilon}}{\sigma} \right| = o(1)$$

Recall that  $\omega_k \approx \sqrt{(|G_k| + 2 \log M)/n}$

Finding  $\mathbf{P}_G$ :

- An optimization scheme

$$\mathbf{P}_G = \arg \min_{\mathbf{P}} \left\{ \|\mathbf{P}\mathbf{Q}_G^\perp\|_S : \mathbf{P} = \mathbf{P}^T = \mathbf{P}^2, \|\mathbf{P}_G \mathbf{Q}_{G_k \setminus G}\|_S \leq \omega'_k \right\}$$

- An extension of an optimization scheme discussed in (Z-Zhang, 11)
- The conclusions in the previous slide hold for all feasible solutions
- Feasibility for sub-Gaussian designs with  $\mathbf{\Sigma} = \mathbb{E}\mathbf{X}^T\mathbf{X}/n$ :

$$\mathbf{X}_G = \mathbf{X}_{-G}\mathbf{\Gamma}_{-G,G} + \mathbf{Z}_G^\circ$$

$$\mathbf{P}_G^\circ = \text{projection to the column space of } \mathbf{Z}_G^\circ$$

$$\omega'_k = \xi \sqrt{|G|/n + |G_k \setminus G|/n + (\log M)/n}$$

$$\text{rank}(\mathbf{P}_G^\circ \mathbf{Q}_G) = |G|, \quad \|\mathbf{P}_G \mathbf{Q}_G^\perp\|_S \leq \sqrt{1 - a_n^2}$$

$$a_n \approx \lambda_{\min} \left( \mathbf{\Sigma}_{G,G}^{-1/2} (\mathbf{\Sigma}^{-1})_{G,G} \mathbf{\Sigma}_{G,G}^{-1/2} \right)$$

## Finding a feasible solution

- Recall that our theory holds for all feasible solutions of the following:

$$\mathbf{P}_G = \arg \min_{\mathbf{P}} \left\{ \|\mathbf{P} \mathbf{Q}_G^\perp\|_S : \mathbf{P} = \mathbf{P}^T = \mathbf{P}^2, \|\mathbf{P}_G \mathbf{Q}_{G_k \setminus G}\|_S \leq \omega'_k \right\}$$

with  $\omega'_k \asymp \sqrt{|G|/n + |G_k \setminus G|/n + (\log M)/n}$

- Finding a feasible solution

$$\hat{\Gamma}_{-G,G} = \arg \min \left\{ \frac{1}{2n} \|\mathbf{X}_G - \mathbf{X}_{-G} \Gamma_{-G,G}\|_F + R(\Gamma_{-G,G}) \right\}$$

- Duality:

$$R(\Gamma_{-G,G}) = \sum_{G_k \not\subseteq G} \frac{\xi \omega_k''}{n^{1/2}} \|\mathbf{X}_{G_k \setminus G} \Gamma_{G_k \setminus G, G}\|_N$$

- Group Lasso:

$$R(\Gamma_{-G,G}) = \sum_{G_k \not\subseteq G} \frac{\xi \omega_k''}{n^{1/2}} \|\mathbf{X}_{G_k \setminus G} \Gamma_{G_k \setminus G, G}\|_F$$

## Working Assumptions

- Group structure:

$$\cup_{j=1}^M G_j = \{1, \dots, p\}, \quad G_j \cap G_k = \emptyset$$

- Group sparsity:

$$\text{supp}(\beta^*) \subset G_{S^*}, \quad |S^*| \leq g, \quad |G_{S^*}| \leq s$$

- Initial estimator

$$\left| \frac{\hat{\sigma}}{\sigma^*} - 1 \right| + \frac{1}{n^{1/2}} \sum_{j=1}^M \omega_j \left\| \mathbf{X}_{G_j} \hat{\beta}_{G_j}^{(init)} - \mathbf{X}_{G_j} \beta_{G_j}^* \right\|_2 \lesssim \left( \frac{s + g \log M}{n} \right)$$

with  $\omega_j \approx \sqrt{(|G_j| + 2 \log M)/n}$

- Group Lasso:

$$\hat{\beta} = \arg \min_{\mathbf{b}} \left\{ \frac{\|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2}{2n} + \sigma \sum_{j=1}^M \omega_j \|\mathbf{b}_{G_j}\|_2 \right\}$$

- Scaled group Lasso:

$$\{\hat{\beta}, \hat{\sigma}\} = \arg \min_{\mathbf{b}, \sigma} \left\{ \frac{\|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2}{2n\sigma} + \frac{\sigma}{2} + \sum_{j=1}^M \omega_j \|\mathbf{b}_{G_j}\|_2 \right\}$$

- Extension from the scaled Lasso: Städler et al (10), Antoniadis (10), Sun-Z (10,12), Belloni et al (11)
- Sign-restricted cone invertibility factor;

$$SCIF = \inf_{\mathbf{u} \in \mathcal{C}_-} \frac{\max_j \omega_j^{-1} \|\mathbf{X}_{G_j} \mathbf{X} \mathbf{u}\|_2 \sum_{j \in S^*} \omega_j^2}{n \sum_{j \in S^*} \omega_j \|\mathbf{u}_{G_j}\|_2}$$

with

$$\mathcal{C}_- = \left\{ \mathbf{u} : \sum_{j \notin S^*} \omega_j \|\mathbf{u}_{G_j}\|_2 \leq \xi \sum_{j \in S^*} \omega_j \|\mathbf{u}_{G_j}\|_2 \neq 0, \right. \\ \left. \mathbf{u}_{G_j}^T \mathbf{X}_{G_j}^T \mathbf{X} \mathbf{u} \leq 0 \quad \forall j \notin S^* \right\}$$

Suppose  $\|\mathbf{X}_{G_j}/\sqrt{n}\|_s \leq 1$  and  $1/SCIF = O(1)$ . Suppose

$$\omega_j = A\sqrt{(|G_j| + 2\log M)/n}, \quad A > 1$$

Then,

$$\left| \frac{\hat{\sigma}}{\sigma^*} - 1 \right| + \frac{1}{n^{1/2}} \sum_{j=1}^M \omega_j \left\| \mathbf{X}_{G_j} \hat{\boldsymbol{\beta}}_{G_j}^{(init)} - \mathbf{X}_{G_j} \boldsymbol{\beta}_{G_j}^* \right\|_2 \lesssim \left( \frac{s + g \log M}{n} \right)$$

# Thanks!