

Modern Statistical Methods  
Review questions

$$1. \quad Y \in \mathbb{R}^n \quad X \in \mathbb{R}^{n \times p} \quad \text{rank}(X) = p$$

$$Y = X\beta^* + \varepsilon - \bar{\varepsilon}\mathbf{1} \quad E(\varepsilon) = 0 \quad \text{Var}(\varepsilon) = \sigma^2 I \quad \bar{\varepsilon} = \mathbf{1}^\top \varepsilon / n$$

$$\hat{\beta}_{OLS} := \arg_{\beta} \min \|Y - X\beta\|_2^2$$

$$\frac{\partial}{\partial \beta} (Y - X\beta)^\top (Y - X\beta) = -2X^\top Y + 2X^\top X\beta = 0$$

$$\Rightarrow \hat{\beta}_{OLS} = (X^\top X)^{-1} X^\top Y \quad (X^\top X)^{-1} \text{ exists or } X \text{ full column rank}$$

$$\hat{\beta}_\lambda^R := \arg_{\beta} \min \left( \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right)$$

$$\frac{\partial}{\partial \beta} \left( (Y - X\beta)^\top (Y - X\beta) + \lambda \beta^\top \beta \right) = -2X^\top Y + 2X^\top X\beta + 2\lambda\beta = 0$$

$$\Rightarrow \hat{\beta}_\lambda^R = (X^\top X + \lambda I)^{-1} X^\top Y \quad \text{unless always exists}$$

$$\mathbb{E} \left\{ (\hat{\beta}_{OLS} - \beta^*) (\hat{\beta}_{OLS} - \beta^*)^\top \right\} - \mathbb{E} \left\{ (\hat{\beta}_\lambda^R - \beta^*) (\hat{\beta}_\lambda^R - \beta^*)^\top \right\} \quad \text{if } \text{pos indefinit?}$$

$$\hat{\beta}_{OLS} - \beta^* = (X^\top X)^{-1} X^\top (X\beta^* + \varepsilon - \bar{\varepsilon}\mathbf{1}) - \beta^* = (X^\top X)^{-1} X^\top (\varepsilon - \bar{\varepsilon}\mathbf{1})$$

$$\begin{aligned} \mathbb{E} \left\{ (\hat{\beta}_{OLS} - \beta^*) (\hat{\beta}_{OLS} - \beta^*)^\top \right\} &= \mathbb{E} \left\{ (X^\top X)^{-1} X^\top (\varepsilon - \bar{\varepsilon}\mathbf{1}) (\varepsilon - \bar{\varepsilon}\mathbf{1})^\top X (X^\top X)^{-1} \right\} \\ &= \sigma^2 (X^\top X)^{-1} \quad \text{as } X^\top \mathbf{1} = 0 \quad (X \text{ cols centred}) \end{aligned}$$

~~EXPLANATION OF THE DERIVATION~~

$$\begin{aligned}
\hat{\beta}_\lambda^R - \beta^0 &= (X^T X + \lambda I)^{-1} X^T (X \beta^0 + \varepsilon - \bar{\varepsilon} \mathbb{I}) - \beta^0 \\
&= (X^T X + \lambda I)^{-1} (X^T X \beta^0 + \lambda I \beta^0 - \lambda I \beta^0 + X^T \varepsilon - \bar{\varepsilon} X^T \mathbb{I}) - \beta^0 \\
&= (X^T X + \lambda I)^{-1} (-\lambda \beta^0 + X^T \varepsilon) \\
\mathbb{E}\{(\hat{\beta}_\lambda^R - \beta^0)(\hat{\beta}_\lambda^R - \beta^0)^T\} &= \mathbb{E}\left\{ (X^T X + \lambda I)^{-1} (-\lambda \beta^0 + X^T \varepsilon)(\varepsilon^T X - \lambda \beta^0 T) \right. \\
&\quad \left. (X^T X + \lambda I)^{-1} \right\} \\
&= (X^T X + \lambda I)^{-1} \lambda^2 \beta^0 \beta^0 T (X^T X + \lambda I)^{-1} + \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \\
&= (X^T X + \lambda I)^{-1} (\lambda^2 \beta^0 \beta^0 T - \sigma^2 \lambda I) (X^T X + \lambda I)^{-1} + \sigma^2 (X^T X + \lambda I)^{-1} \\
\mathbb{E}\{(\hat{\beta}_{OLS} - \beta^0)(\hat{\beta}_{OLS} - \beta^0)^T\} - \mathbb{E}\{(\hat{\beta}_\lambda^R - \beta^0)(\hat{\beta}_\lambda^R - \beta^0)^T\} \\
&= \sigma^2 (X^T X)^{-1} - \sigma^2 (X^T X + \lambda I)^{-1} + \lambda (X^T X + \lambda I)^{-1} (\sigma^2 \mathbb{I} - \lambda \beta^0 \beta^0 T / (X^T X + \lambda I))^{-1} \\
&= (X^T X + \lambda I)^{-1} ((X^T X + \lambda I) \sigma^2 (X^T X)^{-1} - \sigma^2) + \dots \\
&= (X^T X + \lambda I)^{-1} (\sigma^2 \lambda (X^T X)^{-1} (X^T X + \lambda I)) (X^T X + \lambda I)^{-1} + \dots \\
&= \lambda (X^T X + \lambda I)^{-1} \left( \sigma^2 [I + \lambda (X^T X)^{-1} + I] - \lambda \beta^0 \beta^0 T \right) (X^T X + \lambda I)^{-1} \\
&= \lambda (X^T X + \lambda I)^{-1} \left[ \sigma^2 \{2I + \lambda (X^T X)^{-1}\} - \lambda \beta^0 \beta^0 T \right] (X^T X + \lambda I)^{-1}
\end{aligned}$$

Since  $\lambda > 0$ ,  $(X^T X + \lambda I)^T = (X^T X + \lambda I) \neq 0$ , then this is positive definite iff  
 $\sigma^2 \{2I + \lambda (X^T X)^{-1}\} - \lambda \beta^0 \beta^0 T \succ 0$  which is always true for sufficiently small  $\lambda > 0$ .

$$\begin{cases} X^T X \succ 0 & \text{by pos def of Euclidean norm} \\ X^T X + \lambda I \succ 0 & \text{for } \lambda > 0 \end{cases}$$

$$A \succ 0 \Leftrightarrow A^{-1} \succ 0 \quad \text{from eigenvalues}$$

Therefore,

$$\mathbb{E}\{(x^T \hat{\beta}_\lambda^R - x^T \beta^0)^2\} < \mathbb{E}\{(x^T \hat{\beta}_{OLS} - x^T \beta^0)^2\} \quad \text{for all } x^* \in \mathbb{R}^P$$

(and indeed all  $\|x^*\|_2 = 1$ )

$$\lambda \left( \frac{\|x^T x + \lambda I\|^{-1}}{\sigma^2} (2I + \lambda(x^T x)^{-1}) - \lambda \rho^\circ \rho^{\circ T} \right) (x^T x + \lambda I)^{-1} = M$$

Want to show that  $\exists x^*, \rho^0$  with

$$-x^{*T} M x^* > 5$$

$$\text{Take } x^* = \frac{(x^T x + \lambda I) \rho^0}{\|(x^T x + \lambda I) \rho^0\|_2} \quad \text{let } c_1 \text{ be the max eigenval}$$

$$\text{then } -x^{*T} M x^* = \lambda \frac{\lambda \| \rho^0 \|_2^2 - \rho^{\circ T} (\sigma^2 2I + \sigma^2 \lambda (x^T x)^{-1}) \rho^0}{\|(x^T x + \lambda I) \rho^0\|_2^2}$$

$\uparrow$  let  $c_2$  be the max eigenval

$$\geq \lambda \frac{\lambda \| \rho^0 \|_2^2 (1/\| \rho^0 \|_2^2 - c_1)}{\| \rho^0 \|_2^2 c_2} > 5 \quad \text{for some } \rho^0$$

$$2. \text{ FWER} = P(N_{01} \geq 1) \quad N_{01} - \text{rejected true hypothesis}$$

Bonferroni correction: reject  $H_i$  if  $p_i \leq \alpha/m$ .

then

$$\text{FWER} = P(N_{01} \geq 1) \leq E(N_{01}) \quad \text{Markov inequality}$$

$$\begin{aligned} &= E\left(\sum_{i \in I_0} \mathbb{1}_{\{p_i \leq \alpha/m\}}\right) \quad \text{where } I_0 \text{ are true null hypotheses} \\ &\quad |I_0| = m_0 \\ &= \sum_{i \in I_0} E(\mathbb{1}_{\{p_i \leq \alpha/m\}}) \\ &= \sum_{i \in I_0} P(p_i \leq \alpha/m) \leq \sum_{i \in I_0} \frac{\alpha}{m} = \frac{\alpha m_0}{m} \leq \frac{\alpha}{m} \end{aligned}$$

Interaction hypothesis  $H_I = \cap_{i \in I} H_i \quad (\text{all } H_i : i \in I \text{ are true})$

Closure of  $\{H_i\}_{i=1}^m$  w.r.t.  $\{H_I : I \subseteq \{1, \dots, m\}, I \neq \emptyset\}$

Closed testing procedure:

Reject  $H_I$  iff  $\forall J \subseteq I$

$H_J$  is rejected by the local test  $\phi_J$

Local test  $\phi_I$ :  $\alpha$ -level test for each  $I$  taking values in  $\{0, 1\}$ , under  $H_I$

$$P_{H_I}(\phi_I = 1) \leq \alpha$$

If  $I_0$  are true null hypotheses, in order for any false rejection  $H_{I_0}$  must have been rejected

$$\text{FWER} = P(N_{01} \geq 1) = P(\text{at least one false rejection}) \leq P(\phi_{I_0} = 1) \leq \alpha$$

$\neg H_I : I \in \mathcal{I}$ ,  $\forall I, J \in \mathcal{I}$  either  $I \cap J = \emptyset$ ,  $I \subseteq J$  or  $J \subseteq I$

$$P_I^{\text{adj}} := \max_{J: J \in \mathcal{I}, J \supseteq I} \frac{m}{|J|} P_J \quad : \text{Reject } H_I : P_I^{\text{adj}} \leq \alpha. \quad \text{Prob at least } 1-\alpha, \text{ no false rejections}$$

$$\begin{aligned} \text{FWER} &= P(N_{01} \geq 1) \leq \sum_{I \in \mathcal{I}} P(\phi_I = 1) = \sum_{I \in \mathcal{I}} P(P_I^{\text{adj}} \leq \alpha) = \sum_{I \in \mathcal{I}: \forall J \supseteq I} P(P_I^{\text{adj}} \leq \alpha) \quad \text{or otherwise all } J \subseteq I \text{ not rejected} \\ &= \sum_{I \in \mathcal{I}: \forall J \supseteq I} P\left(\frac{m}{|I|} P_I \leq \alpha\right) \leq \sum_{I \in \mathcal{I}: \forall J \supseteq I} \frac{\alpha |I|}{m} = \alpha \sum_{I \in \mathcal{I}} \frac{|I|}{m} \leq \alpha. \end{aligned}$$

3.  $n, p > 1$ ,  $k \in \{1, \dots, p\}$ ,  $j \in \{1, \dots, n\}$

$z \in \mathbb{R}^p$ ,  $z_{-k} \in \mathbb{R}^{p-1}$        $x \in \mathbb{R}^{n \times p}$ ,  $x_k$   $\text{ } k^{\text{th}} \text{ col}$ ,  $x_{-j,k}$   $\text{ } k^{\text{th}} \text{ col}, j^{\text{th}} \text{ row}$

$$Z \sim N_p(\mu, \Sigma), \Sigma > 0$$

$C_1 G_0$  is the minimal ~~Markov~~ graph satisfying the pairwise Markov property w.r.t.  $Z$

$Z$  satisfies the pairwise Markov property w.r.t.  $G$  if for any pair  $j, k \in V$  with  $j \neq k$  and  $(j, k), (k, j) \notin E$

$$Z_j \perp\!\!\!\perp Z_k \mid Z_{-jk}$$

$$z \in \mathbb{R}^p$$

$$Z_k \mid Z_{-k} = z_k$$

Consider general  $Z_A \mid Z_B = z_B$ . Write

$$Z_A = M Z_B + (Z_A - M Z_B) \text{ s.t. }$$

$$\text{Cov}(Z_B, Z_A - M Z_B) = \Sigma_{B,A} - \Sigma_{B,B} M^T = 0$$

$$\Rightarrow M^T = \Sigma_{B,B}^{-1} \Sigma_{B,A} \quad M = \Sigma_{A,B} \Sigma_{B,B}^{-1}$$

then  $Z_B \perp\!\!\!\perp Z_A - M Z_B$  and

$$E(Z_A - M Z_B) = \mu_A - \Sigma_{A,B} \Sigma_{B,B}^{-1} \mu_B$$

$$\text{Var}(Z_A - M Z_B) = \Sigma_{A,A} + \Sigma_{A,B} \Sigma_{B,B}^{-1} \Sigma_{B,B} \Sigma_{B,B}^{-1} \Sigma_{B,A}$$

$$- 2 \Sigma_{A,B} \Sigma_{B,B}^{-1} \Sigma_{B,A}$$

$$= \Sigma_{A,A} - \Sigma_{A,B} \Sigma_{B,B}^{-1} \Sigma_{B,A}$$

$$\therefore Z_A \mid Z_B = z_B \sim N_{|A|}(\mu_A + \Sigma_{A,B} \Sigma_{B,B}^{-1} (z_B - \mu_B), \Sigma_{A,A} - \Sigma_{A,B} \Sigma_{B,B}^{-1} \Sigma_{B,A})$$

$$\Rightarrow Z_k \mid Z_{-k} = z_k \sim N(\mu_k + \sum_{k,-k} \Sigma_{-k,-k}^{-1} (z_{-k} - \mu_{-k}), \Sigma_{k,k} - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \Sigma_{-k,k})$$

Nodewise regression given data  $X$ ,

conditioning on  $Z_{-k} = z_{-k}$  gives

$$Z_k = m_k + Z_{-k}^T \Sigma_{-k,-k}^{-1} \Sigma_{-k,k} + \varepsilon_k \quad (*)$$

where  $m_k = \mu_k - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \mu_{-k}$

$$\varepsilon_k | Z_{-k} = z_{-k} \sim N(0, \Sigma_{k,k} - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \Sigma_{-k,k})$$

In nodewise regression we estimate  $\Sigma_{-k,-k}^{-1} \Sigma_{-k,k}$  for each  $k$  by regressing  $X_k$  on  $X_{-k}$  (crossed out) using least squares.

$$\frac{1}{2n} \sum_{k=1}^p \|X_k - \mu_k - X_{-k} \Theta_{-k,k}\|_2^2 + \lambda \sum_{j < k} \sqrt{\Theta_{jk}^2 + \Theta_{kj}^2}$$

estimates the intercept and coefficient vector. The penalty is similar to Lasso but the  $L_2$ -norm for symmetric terms prefers  $\Theta$  to be symmetric.

$$\Theta_{-k,k} = \Sigma_{-k,-k}^{-1} \Sigma_{-k,k}$$

Hence

$$(Z_{-k}^T Z_{-k}) (\Theta_{-k,k})_j = 0 \Leftrightarrow \begin{cases} \Delta_{j,k} = 0 & : j < k \\ \Delta_{j+1,k} = 0 & : j \geq k \end{cases}$$

and

$$\Delta = \Sigma^{-1}, \quad Z_k \perp\!\!\!\perp Z_j \mid Z_{-k} \Leftrightarrow \Delta_{j,k} = 0$$

Allowing to estimate  $(1)$ .

$$\Sigma = \begin{pmatrix} \Sigma_{k,k} & \Sigma_{k,-k} \\ \Sigma_{-k,k} & \Sigma_{-k,-k} \end{pmatrix} \succ 0 \quad S = \Sigma_{k,k} - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \Sigma_{-k,k} \quad \text{Schur comp}$$

$$\Delta^* = \Sigma^{-1} = \begin{pmatrix} S^{-1} & -S^{-1} \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \\ -\Sigma_{-k,-k}^{-1} \Sigma_{-k,k} S^{-1} & \Sigma_{-k,-k}^{-1} + \Sigma_{-k,-k}^{-1} \Sigma_{-k,k} S^{-1} \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \end{pmatrix}$$

$$\therefore \Delta_{k,k}^{-1} = S = \Sigma_{k,k} - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \Sigma_{-k,k}, \quad -\Delta_{k,-k} \Delta_{-k,k}^{-1} - \Delta_{-k,-k} \Delta_{k,k}^{-1} = \Sigma_{-k,-k}^{-1} \Sigma_{-k,k} = \Theta_{-k,k}$$

and  $\Delta_{k,k} \in \mathbb{R} \setminus \{0\}$  as  $S \succ 0$ .

$$4. Y \in \mathbb{R}^n \quad X \in \mathbb{R}^{n \times p} \quad G_1, \dots, G_q \text{ partition of } \{1, \dots, p\}$$

Group Lasso,  $\lambda > 0$ ,  $m_1, \dots, m_q > 0$  weights

$$Q_\lambda(p) = \frac{1}{2n} \|Y - Xp\|_2^2 + \lambda \sum_{i=1}^q \|\beta_{G_i}\|_2 m_i$$

Multiphase balance different group well.

Penalty encourages either an entire group  $G_j$  to have  $\hat{\beta}_{G_j} = 0$  or  $\hat{\beta}_k \neq 0 \forall k \in G_j$ . Useful for categorical predictors or when expanding predictions using basis functions.

Blockwise coordinate descent: Starting with a guess  $p^{(0)}$ , notice for  $m \geq 1$ :

$$\beta_{G_1}^{(m)} = \underset{\beta_{G_1} \in \mathbb{R}^{1 \times p}}{\operatorname{argmin}} Q_\lambda(\beta_{G_1}, \beta_{G_2}^{(m-1)}, \dots, \beta_{G_q}^{(m-1)})$$

$$\beta_{G_j}^{(m)} = \underset{\beta_{G_j} \in \mathbb{R}^{1 \times p}}{\operatorname{argmin}} Q_\lambda(\beta_{G_1}^{(m)}, \dots, \beta_{G_j}^{(m)}, \dots, \beta_{G_q}^{(m-1)})$$

As other penalties do not change, for  $r$ th block the problem is

$$\min_b \left\{ \frac{1}{2n} \|Y - X_{Gr} \overset{\text{X has null}}{\beta_{Gr}} - X_{Gr} b\|_2^2 + \lambda m_r \|b\|_2 \right\}$$

$$= \min_b \left\{ \frac{1}{2} \|w - Ab\|_2^2 + \eta \|b\|_2 \right\}$$

$$\text{where } w = \frac{1}{\sqrt{n}}(Y - X_{Gr} \overset{\text{X has null}}{\beta_{Gr}}) \quad \text{and } \beta_{Gr}^* = \begin{cases} \beta_{Gj}^{(m)} & : j < r \\ \beta_{Gj}^{(m-1)} & : j > r \end{cases}$$

$$A = \frac{1}{\sqrt{n}} X_{Gr}, \quad \eta = \lambda m_r$$

$$\frac{\partial}{\partial b} \left\{ \frac{1}{2} (w - Ab)^T (w - Ab) + \eta \sqrt{b^T b} \right\} = -A^T A w + A^T A b + \frac{\eta b}{\sqrt{b^T b}} = 0$$

$$\Rightarrow b^* = (\theta^* I + A^T A)^{-1} A^T w \quad \text{where } \theta^* = \frac{\eta}{\|A^T A b\|_2} > 0.$$

$$A = UDV^T \quad , \quad \gamma = D U^T w$$

Rewriting from above

$$\eta^2 = w^T A (\theta^* I + A^T A)^{-1} A^T w \theta^{*3}$$

$$\eta^2 = w^T UDV^T (\theta^* I + VD^T U^T UDV^T)^{-1} VD^T U^T w \theta^{*3}$$

$$\eta^2 = \gamma^T V (\theta^* I + VD^2 V^T)^{-1} V \gamma \theta^{*3}$$

$$\eta^2 = \gamma^T V^T (I + \theta^{*-1} VD^2 V^T)^{-1} V \gamma$$

$$\eta^2 = \gamma^T (I + \theta^{*-1} D^2)^{-1} \gamma$$

$$\Rightarrow \sum_k \frac{\gamma_k^2}{(\theta_{kk}^* + \theta^*)^3} = \frac{\eta^2}{\theta^{*3}}$$

$$b^* = (\theta^* I + VD^2 V^T)^{-1} V \gamma$$

$$= (V (\theta^* I + D^2) V^T)^{-1} V \gamma$$

$$= V (\theta^* I + D^2)^{-1} \gamma$$

$$\theta^* = \frac{\eta^2}{\|b^*\|_2^2}$$

$$\|b^*\|_2^2 = \frac{\eta^2}{\theta^*} = \gamma^T (\theta^* I + D^2)^{-2} \gamma$$

$$= \sum_k \frac{\gamma_k^2}{(\theta^* + \theta_{kk}^*)^2}$$

$$FDP = \frac{N_{\alpha 1}}{\max(R, 1)} \quad E(FDP) = FDR \quad R = \max \{ i : p_{(i)} \leq \frac{i\alpha}{m} \}$$

Reject  $H_{(1)}, \dots, H_{(R)}$

All  $p$ -val independent

$$p_1, \dots, p_m \sim U[0, 1]$$

$$\{ p_i \leq \alpha r/m, R=r \} = \{ p_1 \leq \frac{\alpha r}{m}, p_{(r)} \leq \frac{\alpha r}{m}, p_{(s)} > \frac{\alpha s}{m} \forall s > r \}$$

Let  ~~$p_{(1)}, \dots, p_{(m-1)}$~~  be the order stats of  $p_2, \dots, p_m$

$$\text{Let } R^1 = \max \{ i : p_{(i)} \leq \frac{(i+1)\alpha}{m} \}$$

$$\dots = \{ p_1 \leq \frac{\alpha r}{m}, R^1 = r-1 \}$$

$p_{(1)}, \dots, p_{(m-2)}$  be the order stats ~~for~~  $p_3, \dots, p_m$

$$R'' = \max \{ i : p_{(i)} \leq \frac{(i+2)\alpha}{m} \}$$

$$\{ p_1 \leq \frac{\alpha r}{m}, p_2 \leq \frac{\alpha r}{m}, R=r \} = \{ p_1 \leq \frac{\alpha r}{m}, p_2 \leq \frac{\alpha r}{m}, R''=r-2 \}$$

$$\begin{aligned} E((FDP)^2) &= \sum_{r=1}^m E \left( \frac{N_{\alpha 1}^2}{r^2} \mathbb{1}_{\{R=r\}} \right) \\ &= \sum_{r=1}^m \frac{1}{r^2} E \left( \left( \sum_{j=1}^{m_0} \mathbb{1}_{\{p_j \leq \frac{\alpha r}{m}, R=r\}} \right)^2 \right) \\ &= \sum_{r=1}^m \frac{1}{r^2} \left( \sum_{j=1}^{m_0} P(p_j \leq \frac{\alpha r}{m}, R^1=r-1) + 2 \sum_{j=1}^{m_0} \sum_{k=j+1}^{m-1} P(p_j \leq \frac{\alpha r}{m}, \cancel{p_k \leq \frac{\alpha r}{m}}, \cancel{R''=r-2}) \right. \\ &\quad \left. \text{Independent events} \right) \end{aligned}$$

$$\begin{aligned} &= \sum_{r=1}^m \frac{1}{r^2} \left( \frac{\alpha r}{m} m_0 P(R^1=r-1) + \frac{m_0(m_0-1)\alpha^2 r^2}{m^2} P(R''=r-2) \right) \\ &= \frac{\alpha m_0}{m} E \left( \frac{1}{R^1+1} \right) + \frac{\alpha^2 m(m_0-1)}{m^2} \end{aligned}$$

$$b. \quad Y \in \mathbb{R}^n \quad X \in \mathbb{R}^{n \times p}$$

$$Y = X\beta^0 + z - \bar{z}\mathbf{1} \quad \varepsilon \sim N_n(0, \sigma^2 I)$$

Lasso estimator

$$\hat{\beta}_\lambda^L = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

$$S = \{k \in \{1, \dots, p\} : \beta_k^0 \neq 0\}, \quad N = \{1, \dots, p\} \setminus S, \quad s = |S|$$

$$\text{Arbitrary } A \subseteq \{1, \dots, p\}, \quad b \in \mathbb{R}^p, \quad b_A \in \mathbb{R}^{|A|}$$

Assume  $\exists \phi > 0$  s.t.

$\forall b \in \mathbb{R}^p$  with  $\|b_N\|_1 \leq 3\|b_S\|_1$  have

$$\|b_S\|_1^2 \leq \frac{s\|Xb\|_2^2}{n\phi^2}$$

Show:

If  $\lambda = A\sigma\sqrt{\frac{\log p}{n}}$  for some  $A > 0$ , then w.p. at least  $1 - 2p^{-(A^2/8 - 1)}$

$$\text{have } \frac{1}{n} \|X(\hat{\beta}_\lambda^L - \beta^0)\|_2^2 + \lambda \|\hat{\beta}_\lambda^L - \beta^0\|_1 \leq \frac{16A^2}{\phi^2} \frac{\sigma^2 s \log p}{n}.$$

Basic inequality

$$\frac{1}{2n} \|Y - X\hat{\beta}_\lambda^L\|_2^2 + \lambda \|\hat{\beta}_\lambda^L\|_1 \leq \frac{1}{2n} \|Y - X\beta^0\|_2^2 + \lambda \|\beta^0\|_1$$

$$-\frac{1}{n} Y^T X \hat{\beta}_\lambda^L + \frac{1}{2n} \|X \hat{\beta}_\lambda^L\|_2^2 + \cancel{\frac{1}{2n} \|Y^T X \beta^0 - \frac{1}{2n} \|X \beta^0\|_2^2} \leq \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}_\lambda^L\|_1$$

$$\frac{1}{2n} (X \hat{\beta}_\lambda^L - 2Y)^T X \hat{\beta}_\lambda^L - \frac{1}{2n} (X \beta^0 - 2Y)^T X \beta^0 \leq \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}_\lambda^L\|_1$$

$$\frac{1}{2n} (X \hat{\beta}_\lambda^L - 2X \beta^0 - 2\varepsilon + 2\bar{z}\mathbf{1})^T X \hat{\beta}_\lambda^L - \frac{1}{2n} (-X \beta^0 - 2\varepsilon)^T X \beta^0 \leq \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}_\lambda^L\|_1$$

$$\frac{1}{2n} \|X(\hat{\beta}_\lambda^L - \beta^0)\|_2^2 \leq \frac{1}{n} \varepsilon^T X(\hat{\beta}_\lambda^L - \beta^0) + \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}_\lambda^L\|_1$$

Work on the event  $\mathcal{S}_2 = \{\|X^T \varepsilon\|_\infty/n \leq \lambda\}$

Hölder inequality  $\|X^T X(\hat{\beta}_\lambda^L - \beta^*)\|_1 \leq \|X^T \varepsilon\|_\infty \|\hat{\beta}_\lambda^L - \beta^*\|_1$

$$\Rightarrow \frac{1}{n} \|X(\hat{\beta}_\lambda^L - \beta^*)\|_2^2 + 2\lambda \|\hat{\beta}_\lambda^L\|_1 \leq \lambda \|\hat{\beta}_\lambda^L - \beta^*\|_1 + 2\lambda \|\beta^*\|_1$$

Now  $\frac{1}{n} \|X(\hat{\beta}_\lambda^L - \beta^*)\|_2^2 \leq 3\lambda \|\hat{\beta}_\lambda^L - \beta^*\|_1$ , ~~(1)~~ and

$$\alpha + 2\|\hat{\beta}\|_1 \leq * \|\hat{\beta} - \beta^*\|_1 + 2\|\beta^*\|_1, \quad \text{where } \alpha = \frac{1}{n\lambda} \|\hat{\beta}_\lambda^L - \beta^*\|_1 \|\varepsilon\|_2^2$$

$$\alpha + 2(\|\hat{\beta}_N\|_1 + \|\hat{\beta}_S\|_1) \leq \|\hat{\beta}_S - \beta_S^*\|_1 + \|\hat{\beta}_N - \beta_N^*\|_1 + 2\|\beta_N^*\|_1 + 2\|\beta_S^*\|_1$$

$$\leq \|\hat{\beta}_S - \beta_S^*\|_1 + \|\hat{\beta}_N\|_1 + 2\|\beta_S^*\|_1$$

$$\alpha + 4\|\hat{\beta}_N\|_1 \leq \|\hat{\beta}_S - \beta_S^*\|_1 + 2\|\beta_S^*\|_1 - 2\|\hat{\beta}_S\|_1$$

$$\alpha + \|\hat{\beta}_N - \beta_N^*\|_1 \leq 3\|\hat{\beta}_S - \beta_S^*\|_1$$

$$\alpha + \|\hat{\beta} - \beta^*\|_1 \leq 4\|\hat{\beta}_S - \beta_S^*\|_1$$

$$\Rightarrow \frac{1}{n} \|X(\hat{\beta}_\lambda^L - \beta^*)\|_2^2 + \lambda \|\hat{\beta}_\lambda^L - \beta^*\|_1 \leq 4\lambda \|\hat{\beta}_S - \beta_S^*\|_1$$

Applying the assumption to  $\hat{\beta}_S - \beta_S^*$ ,  $\leq \frac{4\lambda}{\phi} \sqrt{\frac{s}{n}} \|\hat{\beta}_\lambda^L - \beta^*\|_2$

~~Noting that  $X^T X$  is full rank and~~ and  $\frac{1}{\sqrt{n}} \|X(\hat{\beta}_\lambda^L - \beta^*)\|_2 \leq \frac{4\lambda}{\phi} \sqrt{s}$

~~Thus  $\hat{\beta}_\lambda^L - \beta^* \leq \frac{4\lambda}{\phi} \sqrt{s}$~~

$$\therefore \frac{1}{n} \|X(\hat{\beta}_\lambda^L - \beta^*)\|_2^2 + \lambda \|\hat{\beta}_\lambda^L - \beta^*\|_1 \leq \frac{16\lambda^2 s}{\phi^2} \quad \text{the required result.}$$

Finally  $P(\mathcal{A}_2) \geq 1 - 2p^{-(A^2/8-1)}$  as  $P(2\|X^T \varepsilon\|_\infty/n > \lambda) \leq \sum_{j=1}^p P(2|X_j^T \varepsilon|/n > \lambda) \leq 2p^{-(A^2/8-1)}$   
 $X_j^T \varepsilon/n$  sub-Gaussian r.p.  $(\sigma^2 \|X_j\|_2^2/n^2)^{1/2} = \sigma/n$