

# Multicollinearity

## (1) Transformation:

### ① Linear Model:

Thm. For  $Y = X\beta + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2 I_n)$

If  $C(X_1) \perp C(X_2)$ . Then LSE of  $\beta_1, \beta_2 = \hat{\beta}_1, \hat{\beta}_2$

satisfies  $X_1\hat{\beta}_1 = M_1 Y$ ,  $X_2\hat{\beta}_2 = M_2 Y$ ,  $M_1 = P_{C(X_1)}$ ,  $M_2 = P_{C(X_2)}$

Moreover, if  $C(X)$  is full rank, then  $\beta_1, \beta_2$  are estimable.  $\hat{\beta}_i = (X_i^T X_i)^{-1} X_i^T Y$ .

Pf:  $M = M_1 + M_2^* = M_1^* + M_2$

$$C(X_1) \perp C(X_2) \Rightarrow C(M_2^*) = C((I - M_1)X_2) = C(X_2)$$

$$\therefore M = M_1 + M_2, \quad X_i \hat{\beta}_i = M_i X \hat{\beta} = M_i Y.$$

$C(X)$  is full rank  $\Rightarrow \beta$  is estimable.

So as  $(I \ 0)\beta$ ,  $(0 \ I_2)\beta$ .

Rmk: Apply on general case:  $C(X)$  full rank

$$Y = X\beta + \varepsilon = X_1\beta_1 + M_1 X_2\beta_2 + (I - M_1)X_2\beta_2 + \varepsilon$$

$$= X_1\beta_1 + (X_1^T X_1)^{-1} X_1^T X_2\beta_2 + (I - M_1)X_2\beta_2 + \varepsilon.$$

$$=: X_1\delta_1 + (I - M_1)X_2\beta_2 + \varepsilon.$$

satisfies the condition of Thm. So that:

$$\begin{cases} \hat{\delta}_1 = (X_1^T X_1)^{-1} X_1^T Y \\ \hat{\beta}_2 = (X_2^T (I - M_1)X_2)^{-1} X_2^T (I - M_1)Y \end{cases} \Rightarrow \text{solve } \hat{\beta}_i$$

$$\text{We obtain: } \hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} (X_1^T X_1)^{-1} X_1^T (Y - X_2 \hat{\beta}_2) \\ (X_2^T (I - P_1) X_2)^{-1} X_2^T (I - P_1) Y \end{pmatrix}$$

Prop. For  $Y = (J_n \ X) \begin{pmatrix} \beta_0 \\ \beta \end{pmatrix} + \varepsilon = J_n \beta_0 + X\beta + \varepsilon$ ,  $P_n = M J_n$

$\varepsilon \sim N(0, \sigma^2 I_n)$ ,  $r(J_n \ X) = p+1$  (full rank)

$$\Rightarrow \hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} \frac{1}{n} J_n^T Y - \frac{1}{n} J_n^T X (X^T (I - P_n) X)^{-1} X^T (I - P_n) Y \\ (X^T (I - P_n) X)^{-1} X^T (I - P_n) Y \end{pmatrix}$$

pf: Set  $X_1 = J_n$ ,  $X_2 = X$  in the remark above.

Cor. Set  $\bar{Y} = (Y - J_n \bar{c}_0) / \lambda_0$ ,  $\bar{X} = (X - J_n C^T) \Lambda^{-1}$ .

$C = (c_1, \dots, c_p)^T$ ,  $\Lambda = \text{diag} \{ \lambda_1, \dots, \lambda_p \}$ ,  $r(J_n \ \bar{X}) = p+1$ .

For  $\bar{Y} = J_n \bar{\beta}_0 + \bar{X} \bar{\beta} + \bar{\varepsilon}$ ,  $\bar{\varepsilon} \sim N(0, \bar{\sigma}^2 \bar{I})$ .

$$\Rightarrow \begin{pmatrix} \hat{\bar{\beta}}_0 \\ \hat{\bar{\beta}} \end{pmatrix} = \begin{pmatrix} \frac{\hat{\beta}_0 - c_0 + C^T \hat{\beta}}{\lambda_0} \\ \Lambda \hat{\beta} / \lambda_0 \end{pmatrix} \text{ where}$$

$\hat{\beta}_0, \hat{\beta}$  are LSE's of initial:  $Y = X\beta + J_n \beta_0 + \varepsilon$ .

Rmk:  $\bar{Y}^T (I - P_n) \bar{Y} = Y^T (I - P_n) Y / \lambda_0^2$

$$\bar{Y}^T (I - \tilde{M}) \bar{Y} = Y^T (I - M) Y / \lambda_0^2 \Rightarrow$$

For test  $H_0: \beta = 0$  or  $\tilde{H}_0: \bar{\beta} = 0$

on initial or transformed model. Then:

$$F = \tilde{F} \text{ (test statistics.)}$$

## ② Centralization:

For  $\bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$ ,  $\hat{\beta}_0 = \bar{Y} - \sum_k \hat{\beta}_k \bar{X}_k$



cross  $(\bar{X}_1, \dots, \bar{X}_p, \bar{Y})$ . Set  $\bar{X} = (X - J_n C^T) \Lambda^{-1}$  and

$$\bar{Y} = (Y - C_0 J_n) / \Lambda_0, \quad C^T = \frac{1}{n} J_n^T X, \quad \Lambda = I, \quad C_0 = \frac{1}{n} J_n^T Y, \quad \Lambda_0 = 1$$

$$\Rightarrow \hat{\bar{Y}} = \sum_{k=1}^p \hat{\bar{\beta}}_k \bar{X}_k, \quad \begin{pmatrix} \hat{\bar{\beta}}_0 \\ \hat{\bar{\beta}} \end{pmatrix} = \begin{pmatrix} \hat{\beta}_0 + C^T \hat{\beta} - C_0 / \Lambda_0 \\ \Lambda \hat{\beta} / \Lambda_0 \end{pmatrix} = \begin{pmatrix} 0 \\ \hat{\beta} \end{pmatrix}$$

### (3) Standardization:

For  $Y = X\beta + \varepsilon$ , set  $X^* = (X - J_n \frac{J_n^T X}{n}) L_x^{-1}$  and

$$Y^* = (Y - \frac{1}{n} J_n^T Y) / L_Y, \quad L_X = \text{diag}\{L_{11}, \dots, L_{pp}\}, \quad L_Y = Y^T (I - P_n) Y.$$

$$L_{ii} = X_j^T (I - P_n) X_j, \quad X = (X_1, \dots, X_p).$$

$$\text{i.e.} \quad \begin{cases} x_{ij}^* = x_{ij} - \bar{x}_j / \sqrt{L_{jj}} \\ y_i^* = y_i - \bar{y} / \sqrt{L_Y} \end{cases} \Rightarrow \begin{pmatrix} \hat{\beta}_0^* \\ \hat{\beta}^* \end{pmatrix} = \begin{pmatrix} 0 \\ \hat{\beta} \end{pmatrix}$$

$$\hat{\beta}_i^* = \frac{\sqrt{L_{ii}}}{\sqrt{L_Y}} \hat{\beta}_i, \quad \text{by the Formula.}$$

### (2) Background of Collinearity:

Def.  $\{X_k\}_1^p$  is multicollinear if  $\exists \vec{c} \neq \vec{0}_n$  s.t.

$$\sum_{j=1}^p c_j x_{ij} + c_0 = 0, \quad 1 \leq i \leq n.$$

Rmk: i) It's not common that multicollinearity

exists. But it's common:  $c_0 + \sum_{j=1}^p c_j x_{ij}$

$\approx 0$ . Which is called complex MCL.

ii) Note that  $Y = \beta_0 + \sum_{k=1}^p \beta_k X_k + \varepsilon$ . MCL of

$X$  means:  $r(X) < p$ .  $(X^T X)^{-1}$  doesn't exist.

If complex MCL happened  $\Rightarrow r(X) = p$ , but

$|X^T X| \approx 0 \Rightarrow D(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$  has large

elements  $\Rightarrow \hat{\beta}$  isn't accuracy estimation.

Ex. 1. ( $p=2$ )  $\hat{Y} = \hat{\beta}_1 \tilde{X}_1 + \hat{\beta}_2 \tilde{X}_2$ .  $\tilde{X}_i = X_i - \bar{X}_i$ .

$$\Rightarrow (X^T X)^T = \frac{1}{L_{11}L_{22}(1-r_{12}^2)} \begin{pmatrix} L_{22} & -L_{12} \\ -L_{12} & L_{11} \end{pmatrix}$$

$r_{12} = L_{12} / \sqrt{L_{11}L_{22}}$  if  $X_1, X_2$  are related

a lot, then  $r_{12} \uparrow 1$ .  $D(\hat{\beta}_1) + D(\hat{\beta}_2) \uparrow \infty$ .

### (3) Diagnosis:

#### ① Variance Inflation Factor:

$X^* = X - \frac{1}{n} J_n J_n^T X$ .  $(X^{*T} X^*) = (r_{ij}) = R > 0$  is the correlated matrix of  $X$ . Denote:  $C = (c_{ij}) = R^{-1}$

Def:  $VIF_i = c_{ii}$  is the variance inflation factor of variable  $X_i$ .

prop. i)  $\text{Var}(\hat{\beta}_i) = c_{ii} \sigma^2 / L_{ii}$ .  $L_{ii} = X_i^T (I - P_n) X_i$ .

ii)  $c_{ii} = 1 / (1 - R_{i|(i-1)}^2)$ .  $R_{i|(i-1)}^2 = \frac{X_{ij}^T (P_n - P_{n-1}) X_{ij}}{X_{ij}^T (I - P_n) X_{ij}}$

Pf: i) Directly  $\hat{\beta} = (X^T (I - P_n) X)^{-1} X^T (I - P_n) Y$

$$\Rightarrow \text{Var}(\hat{\beta}) = \sigma^2 (X^T (I - P_n) X)^{-1} = \sigma^2 (\Lambda (X^{*T} X^*) \Lambda)^{-1}$$

ii)  $c_{ii}^{-1} = X_{i(i)}^{*T} (I - P_{n-1}) X_{i(i)}^*$

Criteria: If  $\exists i$ .  $VIF_i \geq 10$ . Then it means:

$X_i$  is heavily collinear with other  $X_k$ 's.

Or if  $\overline{VIF} = \frac{1}{p} \sum VIF_i / p \geq 1 \Rightarrow \text{problem exists.}$



② By eigenvalues:

Note that if  $|X^T X| \approx 0$ . Then it  $(X^T X)$  has at least one eigenvalue  $\lambda_0 \approx 0$  (corresp. c.)

$$\Rightarrow X^T X c = \lambda_0 c \quad \therefore c^T X^T X c \approx 0 \Rightarrow X c \approx 0$$

which implies: multicollinearity exists!

Rmk: Denote eigenvalues of  $X^T X: \lambda_1 \geq \dots \geq \lambda_{p+1}$ .

$k_i = \frac{\lambda_1}{\lambda_i}$  is called condition index of  $\lambda_i$ .

Set  $k = \max_i k_i$ . if  $k > 100 \Rightarrow$  Problem exists.

(4) Correction:

① Ridge Estimate:

The idea: Since  $|X^T X| \approx 0$ . Add  $kI_p$  on  $X^T X$

$\Rightarrow X^T X + kI$  may be far away from singular.

Def:  $\hat{\beta}(k) = (X^T X + kI)^{-1} X^T Y$  is ridge estimator.

Rmk:  $\{\hat{\beta}(k)\}_{k \in \mathbb{R}}$  is a family with para.  $k$ .

Denote:  $\phi = (l_1 \dots l_p)$  is matrix of orthonormal eigenfun's of  $X^T X$ . i.e.  $\phi^T X^T X \phi = \text{diag}[\lambda_i]_1^p \equiv \Lambda$ . Set  $Z = X \phi$ .  $\alpha = \phi^T \beta$ .

Rmk:  $Y = Z\alpha + \varepsilon$ .  $\hat{\alpha}(k) = (Z^T Z + kI)^{-1} Z^T Y$ .

$$\hat{\beta}(k) = \phi \hat{\alpha}(k) \Rightarrow \|\hat{\beta}(k)\| = \|\hat{\alpha}(k)\|$$

$$\|\cdot\| \text{ is } : \|\vec{x}\| = (\sum_i x_i^2)^{\frac{1}{2}}.$$

i) properties:

i)  $\hat{\beta}(k)$  is biased estimator of  $\beta$  if  $k \neq 0$

ii) If  $k$  is indep't with  $Y$ . Then  $\hat{\beta}(k)$  is linear transform of  $\hat{\beta}$  and  $Y$

Pf:  $\hat{\beta}(k) = (X^T X + kI)^{-1} X^T X \hat{\beta}$ .

Rmk: Commonly,  $k$  should depend on data  $Y$ .

iii) For  $\|\hat{\beta}\| \neq 0$ ,  $k > 0 \Rightarrow \|\hat{\beta}(k)\| < \|\hat{\beta}\|$ .

Pf:  $\|\hat{\beta}(k)\| = \|\hat{\gamma}(k)\| = \|(I + kI)^{-1} \Lambda \hat{\gamma}\|$   
 $< \|\hat{\gamma}\| = \|\hat{\beta}\|$ .

Rmk:  $R_k(\beta) = \beta(k)$  is a contraction

$k \rightarrow \infty \rightarrow \hat{\beta}(k) \rightarrow \vec{0}$ .

iv)  $\exists k > 0$ ,  $MSE(\hat{\beta}(k)) < MSE(\hat{\beta})$  where

$MSE_{\theta}(\hat{\theta}) = E \|\hat{\theta} - \theta\|^2 = E((\hat{\theta} - \theta)^T (\hat{\theta} - \theta))$

Pf:  $MSE(\hat{\beta}(k)) = MSE_{\gamma}(\hat{\gamma}(k))$

$= \text{tr}(Var(\hat{\gamma}(k))) + \|E(\hat{\gamma}(k)) - \gamma\|^2$

$$\begin{cases} Var(\hat{\gamma}(k)) = \sigma^2 (I + kI)^{-1} \Lambda (kI + \Lambda)^{-1} \\ E(\hat{\gamma}(k)) = (I + kI)^{-1} \Lambda \gamma \end{cases}$$

$$\Rightarrow MSE(\hat{\beta}(k)) = \sigma^2 \sum_{i=1}^p \frac{\lambda_i}{(\lambda_i + k)^2} + k^2 \sum_{i=1}^p \frac{\tau_i^2}{(\lambda_i + k)^2}$$

ii) Choice for  $k$ :

i) By trace of  $\hat{\beta}(k)$  in Plot:



choose  $k$  to make sign of  $\hat{\beta}(k)$  reasonable.  
and SSE won't increase so much. (since  $\hat{\beta}(k)$  deviates  $\hat{\beta}$  a lot for large  $|k|$ .)

ii) By VIF:

$$\text{Var}(\hat{\beta}(k)) = \sigma^2 (X^T X + kI)^{-1} X^T X (X^T X + kI)^{-1} \\ \triangleq \tilde{\sigma}^2(c_{ij}(k))$$

When  $k \uparrow$ ,  $c_{ii}(k) = \text{VIF}_i(k) \downarrow$ .

Choose  $k$  st.  $\text{VIF}_i(k) \leq 10, \forall i$ .

iii) By SSE:

Set a const:  $c$  st.  $\text{SSE}(k) < c \text{SSE}$ , ( $c > 1$ ).

iv) By Hoerl-Kennard Formula:

Denote  $f(k) = \text{MSE}(\hat{\beta}(k))$ ,  $f'(k) = 2 \sum \frac{\lambda_i}{(\lambda_i + k)^2} (k q_i^2 - \sigma^2)$

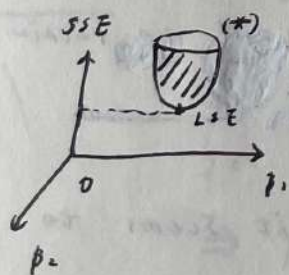
Choose  $\hat{k} = \lfloor \sigma^2 / \max_i q_i^2 \rfloor$ . Since  $f(k) \downarrow$  when  $k \uparrow$   
in  $[0, \hat{k}]$ .  $\Rightarrow$  minimize  $\text{MSE}_p(\hat{\beta}(k))$

iii) Geometric Interpretation:

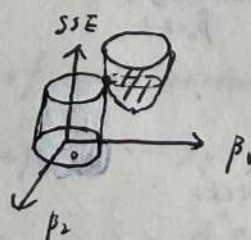
Generally, fix  $c > 1$  st.  $(Y - X \hat{\beta}(k))^T (Y - X \hat{\beta}(k))$

$< c Y^T (I - M) Y$ . (Guarantee SSE(k) won't be large)

$$\Rightarrow \hat{\beta}_{\text{ridge}} =: \underset{\sum \beta_i^2 \leq t}{\text{argmin}} \sum_i (y_i - \beta_0 - \sum \beta_j x_{ij})^2$$



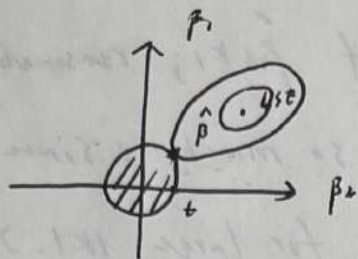
Restrict  
 $\sum \beta_i^2 \leq t$



(\*) = SSE is  
quadratic  
func. of  $\hat{\beta}$ .

consider the cylinder  $\sum \beta_i^2 \leq t$  intersects  $f(\vec{\beta})$

i.e.



$\hat{\beta}$  is the ridge estimate under restrict of  $\sum \beta_i^2 = t$ .

If we see  $\hat{\beta}(k)$  is contraction of  $\hat{\beta}$ .

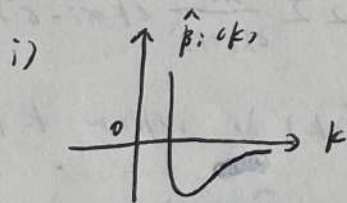
$$\Rightarrow \hat{\beta}^{\text{ridge}} = \arg \min_{\|\beta\| = c\|\hat{\beta}\|} \|Y - X\beta\|^2 \quad ( \|\cdot\| = \|\cdot\|_{L^2} )$$

$$= \arg \min_{\|\beta\| = c\|\hat{\beta}\|} \|X\hat{\beta} - X\beta\|^2 \quad ( \text{By Lagrange} )$$

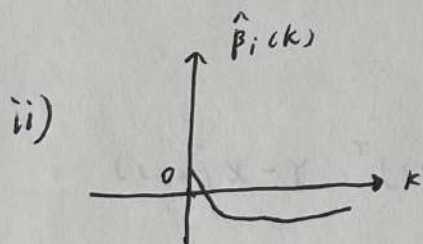
Rank: Another form:  $\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} ( \|Y - X\beta\|^2 + \lambda \|\beta\|^2 )$

iv) Solve variables:

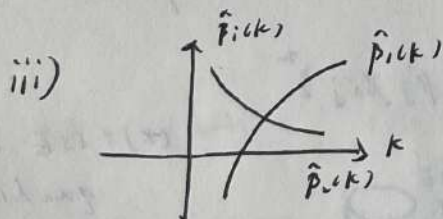
Analysis by trace:



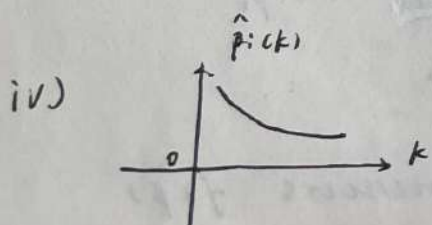
When  $k \uparrow$ ,  $\hat{\beta}_1(k) \downarrow$  rapidly and changes its sign. It's unstable. Besides,  $\hat{\beta}_1(k) \rightarrow 0$  means it loses predictable ability.



When  $k \uparrow$ , Its influential ability  $\uparrow$ .



It means: There's a strong relation between  $\hat{\beta}_1(k)$  and  $\hat{\beta}_2(k)$ . We only need to retain one of them.



It's stable. So it seems to be reasonable.



## V) Generalization:

Consider  $\hat{\beta}(k) = \phi \hat{\alpha}(k)$ .  $\hat{\alpha}(k) = (Z^T Z + k)^{-1} Z^T Y$ .

for general matrix  $k = \text{diag}\{k_1, \dots, k_p\}$ .  $k_i > 0$ .

It retains the properties of  $\hat{\beta}(\text{ridge})$ .

Remark:  $\hat{k}_i = \hat{\sigma}^2 / \hat{\sigma}_i^2 - \hat{\sigma}^2 / \lambda_i$  is common choice.

## ① Stein estimate:

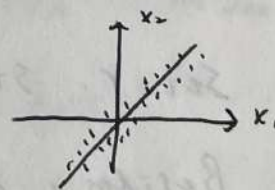
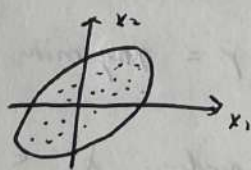
$\hat{\beta}_{SCC} = c \hat{\beta}$ .  $0 < c < 1$ .  $c \hat{\beta}$  is LSE. A contraction

$\exists 0 < c < 1$ . so,  $MSE_F(\hat{\beta}_{SCC}) \leq MSE_F(\hat{\beta})$

## ③ PCA:

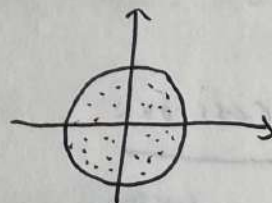
We want to get information from linear combination of  $X = (x_1, \dots, x_p)$ . so that reduces the dim of datas. (but retain most of information).

Remark: It only can be applied in the case  $\{x_i\}$  are correlated:



When  $\{x_i\}$  is uncorrelated

PCA doesn't work:



Denote:  $X$  has been standardized. (so  $X^T X = R$ )

$\phi = (\phi_1, \dots, \phi_p)$  matrix of orthonormal eigenvectors.

so,  $\phi^T X^T X \phi = \Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ .  $Z = X \phi$ .

$\Rightarrow Y = X \beta + \varepsilon = Z \alpha + \varepsilon$ .  $\alpha = \phi^T \beta$ .

Since  $|X^T X| \approx 0$ ,  $\exists r$  s.t.  $\lambda_{r+1} \dots \lambda_p \approx 0$

i.e.  $Z_{(p-r)}^T Z_{(p-r)} \approx 0$ ,  $Z_{(p-r)} = (Z_{r+1} \dots Z_p)$

$\Rightarrow$  Simplify the model:

$$Y = Z\alpha + \varepsilon \approx Z_{(r)}\alpha_{(r)} + \varepsilon, \quad \varepsilon_{(r)} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_r \end{pmatrix}$$

$$\Rightarrow \hat{\alpha}_{(r)} = (Z_{(r)}^T Z_{(r)})^{-1} Z_{(r)}^T Y = \begin{pmatrix} \sum z_{i1} y_i / \lambda_1 \\ \vdots \\ \sum z_{ir} y_i / \lambda_r \end{pmatrix}, \quad \phi \hat{\alpha}_{(r)} = \hat{\beta}_{(r)}$$

prop.  $\sigma^2 \text{tr}(\Lambda_{(p-r)}^{-1}) \geq \| \varepsilon_{(p-r)} \|^2 \Rightarrow \text{MSE}(\hat{\beta}) > \text{MSE}(\hat{\beta}_{(r)})$

$$\begin{aligned} \text{Pf: } E \| \hat{\beta}_{(r)} - \beta \|^2 &= E \left\| \begin{pmatrix} \hat{\alpha}_{(r)} \\ 0 \end{pmatrix} - \alpha \right\|^2 \\ &= \text{tr}(\text{Var}(\hat{\alpha}_{(r)})) + \| E \begin{pmatrix} \hat{\alpha}_{(r)} \\ 0 \end{pmatrix} - \alpha \|^2 \\ &= \text{tr}(\sigma^2 \Lambda_{(r)}^{-1}) + \| \varepsilon_{(p-r)} \|^2. \end{aligned}$$

$$E \| \hat{\beta} - \beta \|^2 = \text{tr}(\sigma^2 \Lambda^{-1}).$$

Rmk: choice of  $r$ : Fix  $c \in (0, 1)$ .

Set  $r$  s.t.  $r = \arg\min \frac{\sum_1^r \lambda_i}{\sum_1^p \lambda_i} > c$

Besides, discard  $\lambda_i < 0.01$ .

#### ④ Lasso Regression:

$$\hat{\beta} = \arg\min_{\beta} \sum_1^n (y_i - \beta_0 - \sum \beta_j x_{ij})^2 + \lambda \sum |\beta_j|, \text{ fix } \lambda.$$

Rmk: Modification:

$$\text{A-Lasso: } \hat{\beta} = \arg\min \|Y - X\beta\|_2^2 + \lambda \sum w_i |\beta_i|$$

$$\text{Elastic-net: } \hat{\beta} = \arg\min \|Y - X\beta\|_2^2 + \lambda_1 \sum |\beta_i| + \lambda_2 \sum \beta_i^2.$$