

## Theoretical questions

### OLS

- On a  $\tilde{\beta} = Cy = Hy + Dy$  et  $\beta^* = Hy$   
 Donc  $\text{Var}(\tilde{\beta}) = \text{Var}(Cy)$   

$$= C \text{Var}(y) C^T$$
  

$$= \sigma^2 C C^T$$
  

$$= \sigma^2 ((X^T X)^{-1} X^T + D) (X (X^T X)^{-1} + D^T)$$
  

$$= \sigma^2 (X^T X)^{-1} + \sigma^2 (X^T X)^{-1} (D X^T + \sigma^2 D X (X^T X)^{-1} + D D^T)$$
  

$$= \sigma^2 (X^T X)^{-1} + \sigma^2 D D^T \quad \text{car } D X = 0$$
  

$$\quad \text{car } \tilde{\beta} \text{ est non biaisé}$$
  

$$= \text{Var}(\beta^*) + \underbrace{\sigma^2 D D^T}_{> 0}$$
  

$$\quad \text{car } D D^T \text{ est semi-def positive}$$

$$E(\tilde{\beta}) = E[Cy]$$

$$= (I_d + D X) \beta \quad \text{donc} \quad \text{non biaisé} \Rightarrow D X = 0$$

Donc l'OLS est l'estimateur qui a la plus petite variance.

### Ridge Regression

- On a  $\beta_{\text{ridge}}^* = (X_c^T X_c + \lambda I_d)^{-1} X_c^T y_c$

$$\text{Donc } E[\beta_{\text{ridge}}^*] = (X_c^T X_c + \lambda I_d)^{-1} X_c^T \beta \neq \beta$$

et  $\beta_{\text{ridge}}^*$  est biaisé.



$$\bullet \quad \beta_{\text{ridge}}^* = (X_c^T X_c + \lambda I_p)^{-1} X_c^T y_c$$

on utilise la décomposition SVP de  $X_c$

$$X_c = U D V^T \quad \text{où } U, V \text{ sont orthogonales}$$

$$D \text{ est diagonale } = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & & \\ & & & & 0 & \dots \end{pmatrix}$$

$$\beta_{\text{ridge}}^* = (V D^T V^T + \lambda I_p)^{-1} V D U^T y_c$$

$$= \left( V \begin{pmatrix} \sigma_1^2 + \lambda & & \\ & \ddots & \\ & & \sigma_r^2 + \lambda & \\ & & & \lambda & \dots \end{pmatrix} V^T \right)^{-1} V D U^T y_c$$

$$= V \begin{pmatrix} \frac{1}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{1}{\sigma_r^2 + \lambda} & \\ & & & \frac{1}{\lambda} & \dots \end{pmatrix} \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & \\ & & & 0 & \dots \end{pmatrix} U^T y_c$$

$$= V \begin{pmatrix} \frac{\sigma_1}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_r}{\sigma_r^2 + \lambda} & \\ & & & 0 & \dots \end{pmatrix} U^T y_c$$

Utiliser cette décomposition est utile pour vérifier  
qu'il n'y ait pas de problème de colinéarité



(2)

On a

$$(X_c^T X_c + \lambda I_d) \beta_{\text{ridge}}^* = X_c^T y_c = X_c^T X_c \beta + X_c^T \varepsilon$$

Donc  $(X_c^T X_c + \lambda I_d) \text{Var}(\beta_{\text{ridge}}^*) (X_c^T X_c + \lambda I_d)^T = \sigma^2 X_c^T X_c$

et  $\text{Var}(\beta_{\text{ridge}}^*) = \sigma^2 (X_c^T X_c + \lambda I_d)^{-1} X_c^T X_c (X_c^T X_c + \lambda I_d)^{-1}$

On note  $X_c^T X_c = P D P^T$

Alors  $\text{Var}(\beta_{\text{ridge}}^*) = \sigma^2 P^T (D + \lambda I_d)^{-1} D (D + \lambda I_d)^{-1} P$

$$= \sigma^2 P^T \begin{pmatrix} \frac{\sigma_A^2}{(\sigma_A^2 + \lambda)^2} & & \\ & \ddots & \\ & & \frac{1}{(\sigma_1 + \lambda)^2} \\ & & & 0 \dots 0 \end{pmatrix} P$$

$$\leq \sigma^2 P^T \begin{pmatrix} \frac{1}{\sigma_1^2} & & \\ & \ddots & \\ & & \frac{1}{\sigma_1^2} \\ & & & 0 \dots 0 \end{pmatrix} P$$

$$= \text{Var}(\beta_{\text{OLS}}^*)$$

$\beta_{\text{ridge}}^*$  est bien sûr mieux à une plus petite variance !



- $E(\beta_{\text{ridge}}^*) = (X_c^T X_c + \lambda I_d)^{-1} X_c^T y$

Donc biais  $(\beta_{\text{ridge}}^*) = -P^T \begin{pmatrix} \frac{\lambda}{\sigma^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\lambda}{\sigma^2 + \lambda} \end{pmatrix} P \beta$

$\xrightarrow{\lambda \rightarrow +\infty} -P^T \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix} P \beta = -\beta$

et quand  $\beta \rightarrow +\infty$

$$\text{Var}(\beta_{\text{ridge}}^*) = \sigma^2 P^T \begin{pmatrix} \frac{\sigma^2}{(\sigma^2 + \lambda)^2} & & \\ & \ddots & \\ & & \frac{\sigma^2}{(\sigma^2 + \lambda)^2} \end{pmatrix} P$$

$\rightarrow 0$

C'est logique, quand on régularise beaucoup,

$E(\beta_{\text{ridge}}^*)$  tend vers 0, donc le biais vers  $\beta$ .

On fait des concessions sur l'écart à  $\beta$  mais on obtient une variance de plus en plus petite, qui tend vers 0.



3

$$\bullet \beta_{\text{ridge}}^* = (X_c^T X_c + \lambda I_d)^{-1} X_c^T y_c$$

assuming  $X_c^T X_c = I_d$ ,

$$\beta_{\text{ridge}}^* = ((1+\lambda) I_d)^{-1} X_c^T y_c$$

$$\beta_{\text{OLS}}^* = X_c^T y_c$$

$$= \frac{1}{1+\lambda} X_c^T y_c$$

$$\beta_{\text{ridge}}^* = \frac{\beta_{\text{OLS}}^*}{1+\lambda}$$

Elastic Net.



(4)

- Si chaque classe possède sa propre matrice de covariance  $\Sigma_k$ , on a

$$\begin{aligned}
 f^*(x_j) &= \arg \max_{c_k} P_{c_k}(x_j) \pi_{c_k} \\
 &= \arg \min_{c_k} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log(|\Sigma_k|) - 2 \log(\pi_{c_k}) \\
 &= \arg \min_{c_k} x^T \Sigma_k^{-1} x - 2x^T \Sigma_k^{-1} \mu_k + \mu_k^T \Sigma_k^{-1} \mu_k + \log(|\Sigma_k|) - 2 \log(\pi_{c_k}) \\
 &= \arg \min_{c_k} x^T a_k x + x^T b_k + c_k
 \end{aligned}$$

où

$$\begin{aligned}
 a_k &= \Sigma_k^{-1} \\
 b_k &= -2 \Sigma_k^{-1} \mu_k
 \end{aligned}$$

$$c_k = \mu_k^T \Sigma_k^{-1} \mu_k - 2 \log(\pi_{c_k}) + \log(|\Sigma_k|)$$

ce qui est bien une fonction quadratique de  $x$ .