

# Theoretical Questions

①

OLS:

$$(X^T X)^{-1} X^T X = Id \text{ donc } ((X^T X)^{-1} X^T)^{-1} = X$$

$$\begin{aligned} \bullet E[\tilde{\beta}] &= E[cy] = E[(H + D)y] = E\left[\left((X^T X)^{-1} X^T + D\right)y\right] \\ &= E\left[(Id + DX)(X^T X)^{-1} X^T y\right] \\ &= (Id + DX) E[\beta^*] \text{ comme } \beta^* \text{ non biaisé} \\ &= (Id + DX) \beta \end{aligned}$$

Si on veut que  $\tilde{\beta}$  soit non biaisé, Il faut  $DX = 0$ .

$$\begin{aligned} \bullet \text{var}(\tilde{\beta}) &= \text{var}(cy) = C \text{var}(y) C^T = \sigma^2 C C^T \\ &= \sigma^2 (X^T X)^{-1} + \sigma^2 (X^T X)^{-1} (DX)^T + \sigma^2 DX (X^T X)^{-1} + \sigma^2 DD^T \end{aligned}$$

$$\begin{aligned} \text{OR } DX = 0 \text{ donc } \text{var}(\tilde{\beta}) &= \sigma^2 (X^T X)^{-1} + \sigma^2 DD^T \\ &= \text{var}(\beta^*) + \underbrace{\sigma^2 DD^T}_{\sigma^2 \|D\|_2^2} > 0 \text{ car } D \text{ non-nulle} \end{aligned}$$

$$\text{donc } \text{var}(\tilde{\beta}) > \text{var}(\beta^*)$$

On doit utiliser l'assomption que  $\tilde{\beta}$  et  $\beta^*$  sont non biaisés  
c'est-à-dire que  $X$  est de rang complet et  $y = BX$ .  
Donc l'OLS est l'estimateur de plus petite variance.

Ridge Regression:

• On a

$$\beta^*_{\text{ridge}} = (X_c^T X_c + \lambda Id)^{-1} X_c^T y_c$$

$$\text{car } f: z \mapsto (y_c - x_c z)^T (y_c - x_c z) + \lambda \|z\|_2^2$$

$$\text{est différentiable et } f'(\beta^*_{\text{ridge}}) = 0$$

$$\text{avec } f'(z) = X_c^T (y_c - X_c z) + \lambda z$$

$$(X_c^T X_c + \lambda Id) \text{ inversible pour } \lambda > 0$$

Ainsi  $E[\beta_{ridge}^*] = (X_c^T X_c + \lambda I_d)^{-1} X_c^T X_c E[\beta]$  (2)

$$= (X_c^T X_c + \lambda I_d)^{-1} (X_c^T X_c + \lambda I_d - \lambda I_d) \beta$$

$$= \beta + \lambda (X_c^T X_c + \lambda I_d)^{-1} \beta \neq \beta$$

L'estimateur est biaisé

•  $X_c = U D V^T$  avec  $\begin{cases} U, V \text{ orthogonales} \\ D = \begin{pmatrix} d_1 & & & 0 \\ & \ddots & & \\ & & d_r & \\ 0 & & & 0 \end{pmatrix} \end{cases}$

$$\beta_{ridge}^* = (V D^2 V^T + \lambda I_d)^{-1} V D U^T y_c$$

$$= V (D^2 + \lambda I_d)^{-1} D U^T y_c$$

$$= V \begin{pmatrix} \frac{d_1}{d_1^2 + \lambda} & & & 0 \\ & \ddots & & \\ & & \frac{d_r}{d_r^2 + \lambda} & \\ & & & 0 \end{pmatrix} U^T y_c$$

Cette décomposition est utile pour calculer le pseudo-inverse

•  $\beta_{ridge}^* = (X_c^T X_c + \lambda I_d)^{-1} X_c^T \text{var}(X_c \beta - \varepsilon)$  avec  $y_c = X_c \beta - \varepsilon$ .

$$= (X_c^T X_c + \lambda I_d)^{-1} X_c^T \text{var}(y) ((X_c^T X_c + \lambda I_d)^{-1} X_c^T)^T$$

$$= \sigma^2 (X_c^T X_c + \lambda I_d)^{-1} X_c^T X_c (X_c^T X_c + \lambda I_d)^{-1}$$

on utilise la décomposition SVD  
vue à la question précédente  $X_c = U D V^T$

$$\text{var}(\beta_{ridge}^*) = \sigma^2 V \left( D^2 + \lambda I_d \right)^{-1} D^2 (D^2 + \lambda I_d)^{-1} V^T$$

$$= \sigma^2 V \begin{pmatrix} \frac{d_1^2}{(d_1^2 + \lambda)^2} & & & 0 \\ & \ddots & & \\ & & \frac{d_r^2}{(d_r^2 + \lambda)^2} & \\ & & & 0 \end{pmatrix} V^T$$

$$\leq \sigma^2 V \begin{pmatrix} \frac{1}{d_1^2} & & & 0 \\ & \ddots & & \\ & & \frac{1}{d_r^2} & \\ & & & 0 \end{pmatrix} V^T = \text{var}(\beta_{OLS}^*)$$

$$= \sigma^2 X_c^T X_c$$

• Bias ( $\beta^*_{ridge}$ ) =  $\lambda (X_c^T X_c + \lambda Id)^{-1} \beta$

$$= V \begin{pmatrix} \frac{\lambda}{d_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\lambda}{d_r^2 + \lambda} & & 1 \dots 1 \end{pmatrix} V^T \beta$$

$\xrightarrow{\lambda \rightarrow +\infty}$

$$V \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \dots 1 \end{pmatrix} V^T \beta = \beta$$

$Var(\beta^*_{ridge}) = \sigma^2 V \begin{pmatrix} \frac{d_1^2}{(d_1^2 + \lambda)^2} & & \\ & \ddots & \\ & & \frac{d_r^2}{(d_r^2 + \lambda)^2} & & 0 \dots 0 \end{pmatrix} V^T \xrightarrow{\lambda \rightarrow +\infty} 0$

•  $\beta^*_{ridge} = (X_c^T X_c + \lambda Id)^{-1} X_c^T y_c$  avec  $X_c^T X_c = Id$

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

avec  $X_c^T y_c = \beta^*_{OLS}$

$$= \frac{\beta^*_{OLS}}{1 + \lambda}$$

car  $((1 + \lambda) Id)^{-1} = \frac{1}{1 + \lambda} Id$ .

### Elastic Net

Avec  $X_c^T X_c = Id$ ,  $\beta^*_{elNet} = \argmin_{\beta} (y_c - X_c \beta)^T (y_c - X_c \beta) + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$

$$= \argmin_{\beta} - y_c^T X_c \beta - \beta^T X_c^T y_c + \beta^T \beta + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$$

$$= \argmin_{\beta} - 2\beta^T \beta^*_{OLS} + \beta^T \beta + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$$

avec  $\beta^*_{OLS} = X_c^T y_c$ ,

on adonc  $\beta^*_{elNet} = \argmin_{\beta} - 2\beta^T \beta^*_{OLS} + (\lambda_2 + 1) \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$

$$= \argmin_{\beta} - 2\beta^T \beta^*_{OLS} + (\lambda_2 + 1) \|\beta\|_2^2 + \int \lambda_1 \beta \text{ if } \beta \geq 0$$

$$- \lambda_1 \beta \text{ else}$$

on pose  $f: \beta \mapsto - 2\beta^T \beta^*_{OLS} + (\lambda_2 + 1) \|\beta\|_2^2 + \int \lambda_1 \beta \text{ if } \beta \geq 0$

post dérivable sur  $\mathbb{R}^+$  et  $\mathbb{R}^-$ :

$$f'(\beta_{elNet}) = 0 \quad (\Rightarrow) 0 = - 2\beta^*_{OLS} + 2(\lambda_2 + 1) \beta_{elNet} + \begin{cases} \lambda_1 & \text{if } \beta_{elNet} \geq 0 \\ -\lambda_1 & \text{if } \beta_{elNet} < 0 \end{cases}$$

Donc  $\beta_{elNet} = \frac{\beta^*_{OLS} \pm \frac{\lambda_1}{2}}{\lambda_2 + 1}$

## LDA:

- Supposons que chaque classe possède sa propre matrice de covariance  $\Sigma_k$ .

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

$$= \arg \min_{c_k} x^T C_k x - x^T b_k + a_k$$

$$\text{avec } \begin{cases} a_k = \mu_k^T \Sigma_k^{-1} \mu_k + \log |\Sigma_k| - 2 \log(\pi_{c_k}) \\ b_k = -2 \Sigma_k^{-1} \mu_k \\ C_k = \Sigma_k^{-1} \end{cases}$$

- On obtient bien une solution quadratique de  $x$ , à cause du terme en  $x^T C_k x$ .