

# TP: Supervised Learning - Partie théorique

## 1. OLS

On a:  $\beta^* = (X^T X)^{-1} X^T y = H y$  où on suppose  $y = X\beta + \epsilon$  avec  $\epsilon \sim N(0, \sigma^2 I)$

on considère  $\tilde{\beta} = (H + D) y$  où  $D = (I - H)$  vient

$E(\tilde{\beta}) = (H + D) X\beta = (I + DX)\beta$  où  $I$  est l'identité  
On suppose que  $E(\tilde{\beta}) = \beta$  on en déduit donc que  $DX = 0$

on calcule ensuite:  $Var(\tilde{\beta}) = Var(Cy) = C Var(y) C^T = \sigma^2 C C^T$  d'après notre hypothèse  
d'où  $Var(\tilde{\beta}) = \sigma^2 (X^T X)^{-1} X^T + D (X^T X)^{-1} X^T + D$

En développant et en utilisant  $DX = 0$  on a  
 $Var(\tilde{\beta}) = \underbrace{\sigma^2 (X^T X)^{-1}}_{Var(\beta^*)} + \sigma^2 (D D^T)$

on a ainsi prouvé que  $Var(\tilde{\beta}) \geq Var(\beta^*)$

## 2. Ridge

2.1) Cette fois-ci on prend  $\beta^* = (X_c^T X_c + \lambda I)^{-1} X_c^T y_c$   
 $\beta^*$  ainsi défini est, par rapport, l'estimateur de Ridge

En remarquant que  $E(\beta^*) = X_c^T \beta$  on a par linéarité:

$$E(\beta^*) = (X_c^T X_c + \lambda I)^{-1} X_c^T X_c \beta$$

Par  $\lambda \neq 0$ ,  $E(\beta^*) \neq \beta$   
et dans ces cas l'estimateur est biaisé.

2.2. on écrit  $X_c = U D V^T$  de sorte que

$$\begin{aligned} \beta^* &= ((U D V^T)^T (U D V^T) + \lambda I)^{-1} (U D V^T)^T y_c \\ &= ((V D^T U^T) (U D V^T) + \lambda I)^{-1} (U D V^T)^T y_c \\ &= V (D^T D + \lambda I)^{-1} V^T U^T y_c = V (D^T D + \lambda I)^{-1} D^T U^T y_c \end{aligned}$$

Avec cette forme, on n'a plus besoin de recourir à l'inverse d'une matrice puisque on peut utiliser le fait que  $(D^T D + \lambda I)^{-1} D^T$  est une matrice de la forme  $\begin{cases} \frac{1}{\lambda^2 + 1} & \text{sur la diagonale, où les } \lambda_i \text{ sont les valeurs propres} \\ 0 & \text{ailleurs} \end{cases}$

2-3. On a que

$$\begin{aligned} \text{Var}(\beta^k)_{\text{ridge}} &= \left( (x_c^T x_c + \lambda I)^{-1} x_c^T \right) \text{Var}(y_c) \left( (x_c^T x_c + \lambda I)^{-1} x_c^T \right)^T \\ &= \sigma^2 \left( x_c^T x_c + \lambda I \right)^{-1} x_c^T x_c \left( x_c^T x_c + \lambda I \right)^{-1} \end{aligned}$$

Comme  $\text{Var}(\beta^k_{OLS}) = \sigma^2 (x_c^T x_c)^{-1}$ , on a :  $\text{Var}(\beta^k_{OLS}) \rightarrow \text{Var}(\beta^k_{\text{ridge}})$

2-4. Plus  $\lambda$  augmente, plus le biais augmente et la variance diminue.