

Théor

OLS

* Par hypothèse $\tilde{\beta}$ est non biaisé donc

$$\beta = E(\tilde{\beta})$$

$$= E(Hy) + E(Dy)$$

$$= E(\beta^*) + DE(y), \text{ or } \beta^* \text{ est aussi non biaisé.}$$

De plus dans le cadre du modèle gaussien $y \sim N(X\beta, \sigma^2 I_n)$, $\sigma^2 > 0$.

D'où $DX\beta = 0$ on a ici alors $DX = 0$.

$$* \text{Var}(\tilde{\beta}) = \text{Var}(Cy) = \sigma^2 CCT$$

$$= \sigma^2 (H+D)(H^T+D^T)$$

$$= \sigma^2 HHT + \sigma^2 DDT + \sigma^2 HD^T + \sigma^2 DH^T$$

$$\text{Or: } DH^T = DX(X^T X)^{-1} = 0$$

$$* HTD = (DH^T)^T = 0^T = 0$$

$$* \text{Var}(\beta^*) = \sigma^2 X^T X^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} = \sigma^2 HHT$$

$$\text{D'où: } \underline{\text{Var}(\tilde{\beta}) = \text{Var}(\beta^*) + \sigma^2 DDT}$$

$$\text{Ainsi } \text{Var}(\tilde{\beta}) - \text{Var}(\beta^*) = \sigma^2 DDT$$

$$\text{Or } \forall z \in \mathbb{R}^d \setminus \{0\}, z^T DDT z = \|Dz\|^2 > 0 \text{ car } D \neq 0$$

$$\text{D'où } \underline{\text{Var}(\tilde{\beta}) > \text{Var}(\beta^*)}$$

On a ici supposé qu'on se place dans le cas du modèle gaussien tel que (X) soit aussi inversible.

Ridge regression

$$\text{On pose pour } \beta \in \mathbb{R}^d, f(\beta) = \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

$$\text{Alors } \beta_{\text{ridge}}^{\text{opt}} = \underset{\beta}{\text{argmin}} f(\beta).$$

On a donc $f'(\beta) = -2x_c^T(y_c - x_c\beta) + 2\lambda\beta$
 $f''(\beta) = 2x_c^T x_c + 2\lambda \geq 0$.

Donc β_{ridge}^* est solution de $f'(\beta) = 0$.

D'où $\beta_{ridge}^* = (x_c^T x_c + \lambda I_m)^{-1} x_c^T y_c$

• Alors $E(\beta_{ridge}^*) = (x_c^T x_c + \lambda I_m)^{-1} x_c^T x_c \beta \neq \beta$.

The ridge estimator is biased.

• Soit (μ_i) les valeurs singulières de x_c . $D = \text{diag}(\mu_i)$.

On a $x_c^T x_c = V D U^T U D V^T = V D^2 V^T$ car U est orthogonale.

D'où $\beta_{ridge}^* = (V D^2 V^T + \lambda V V^T)^{-1} x_c^T y_c$ car V est aussi orthogonale.

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

$$= V \text{diag}(\mu_i^2 + \lambda)^{-1} V^T x_c^T y_c$$

$$= V \text{diag}\left(\frac{1}{\mu_i^2 + \lambda}\right) V^T x_c^T y_c$$

L'avantage d'utiliser la décomposition en valeurs singulières est qu'au lieu de devoir inverser une matrice compliquée, il suffit d'inverser des coefficients.

• $\text{Var}(\beta_{ridge}^*) = \sigma^2 V (D^2 + \lambda I_n)^{-1} V^T x_c^T x_c V (D^2 + \lambda I_n)^{-1} V^T$

$$= \sigma^2 V (D^2 + \lambda I_n)^{-1} D^2 (D^2 + \lambda I_n)^{-1} V^T$$

$$= \sigma^2 V \cdot \text{diag}\left(\frac{1}{\mu_i^2 + \lambda}\right) \text{diag}(\mu_i^2) \text{diag}\left(\frac{1}{\mu_i^2 + \lambda}\right) V^T$$

$$= \sigma^2 V \cdot \text{diag}\left(\frac{\mu_i^2}{(\mu_i^2 + \lambda)^2}\right) V^T$$

Et $\text{Var}(\beta_{OLS}^*) = \sigma^2 \text{diag}\left(\frac{1}{\mu_i^2}\right) V^T$

Donc $\text{Var}(\beta_{OLS}^*) - \text{Var}(\beta_{ridge}^*) = \sigma^2 V \text{diag} \left(\frac{1}{\mu_i^2} - \frac{\mu_i^2}{(\mu_i^2 + \lambda)^2} \right) V^T$

Or $\forall i, \frac{1}{\mu_i^2} - \frac{\mu_i^2}{(\mu_i^2 + \lambda)^2} = \frac{1}{\mu_i^2} \left(1 - \underbrace{\frac{1}{(1 + \lambda)^2}}_{< 1} \right) \geq 0$

D'où le résultat $\text{Var}(\beta_{ridge}^*) < \text{Var}(\beta_{OLS}^*)$

• On a $E(\beta_{ridge}^*) = V \text{diag} \left(\frac{\mu_i^2}{\mu_i^2 + \lambda} \right) V^T \beta$ et $\text{Var}(\beta_{ridge}^*) = \sigma^2 V \text{diag} \left(\frac{\mu_i^2}{(\mu_i^2 + \lambda)^2} \right) V^T$

Donc si λ augmente alors le biais augmente mais la variance diminue et inversement si λ diminue le biais diminue mais la variance augmente.

• Si $x_c^T x_c = Id$ alors on a $\beta_{OLS}^* = x_c^T y_c$ et

$\beta_{ridge}^* = ((1 + \lambda) Id)^{-1} x_c^T y_c = \frac{\beta_{OLS}^*}{1 + \lambda}$

Elastic net

On pose $g(\beta) = \|y_c - x_c \beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$

Selon le signe de β_i , le sous gradient de $\beta_i \mapsto |\beta_i|$ vaut ± 1 d'où on peut noter :

$g'(\beta) = -2x_c^T(y_c - x_c \beta) + 2\lambda_2 \beta \pm \lambda_1$

Et $g''(\beta) = 2x_c^T x_c + 2\lambda_2 > 0$

Ainsi $\beta_{Elastic}^*$ est solution de $g'(\beta) = 0$

D'où $\beta_{Elastic}^* = (x_c^T x_c + \lambda_2 Id)^{-1} (x_c^T y_c \pm \frac{\lambda_1}{2})$

Donc avec $x_c^T x_c = Id$ on a bien

$\beta_{Elastic}^* = \frac{\beta_{OLS}^* \pm \frac{\lambda_1}{2}}{1 + \lambda_2}$