

# Régression linéaire

## "Ridge" et "Lasso"

C. Frenck

$Y, \{X_1, \dots, X_p\}$  centrés

$$Y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \Sigma$$

$$\Sigma \sim N(0, \sigma^2)$$

$n$  - observations

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} X_{i1} & \dots & X_{ip} \\ \vdots & & \vdots \\ X_{i1} & \dots & X_{ip} \\ \vdots & & \vdots \\ X_{n1} & \dots & X_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \Sigma_1 \\ \vdots \\ \Sigma_i \\ \vdots \\ \Sigma_n \end{pmatrix}$$

En forme matricielle

$$Y = X \cdot \beta + \Sigma$$

↑  
matrice de design.

Estimateur de MCO (ou de maximum de vraisemblance)

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$
$$V(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

si  $X^T X$  est singulière

$$\det(X^T X) = 0$$

alors  $(X^T X)^{-1}$  n'existe pas.

Si  $\lambda_1, \dots, \lambda_p$  sont v. p. de  $(X^T X)$  (3)  
alors

$$\det(X^T X) = 0 \Leftrightarrow \exists q \in [1:p] + 1 \dots p$$

$$\underline{\lambda_q = \lambda_{q+1} = \dots = \lambda_p = 0}$$

Or, on a vu que (lecture 3, p. 47) :

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{n-p-1} \sum_{k=1}^p u_{kj}^2 \cdot \frac{1}{\lambda_k}$$

$$\text{Donc si } \lambda_{q+1} = 0 \Rightarrow \text{Var}(\hat{\beta}_j) \rightarrow \infty$$

\*  $j = 1 \dots p$ .

Solutions :

PCR

PLS

Ridge

Lasso

# Idee de la regression Ridge :

Considerer au lieu de

$$\hat{\beta}_{MCO} = (X^T X)^{-1} X^T Y$$

l'estimateur

$$\hat{\beta}_{ridge} = (X^T X + \lambda I_p)^{-1} X^T Y$$

avec  $\lambda > 0$  - parametre de regularisation  
et  $I_p$  - la matrice identite'  $p \times p$

$$I_p = \begin{pmatrix} 1 & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & 1 \end{pmatrix}$$

Remarque :

$(X^T X + \lambda I_p)$  est toujours inversible  
car si  $\lambda_i$  est valeur propre de  $(X^T X)$   
alors  $\mu_i = \lambda_i + \lambda$  est v.p de  $X^T X + \lambda I_p$

En effet, soit  $(\lambda_i, u_i)$  un couple  
vp - valeur propre et vecteur propre de  $X^T X$ .

$$X^T X u_i = \lambda_i u_i$$

Alors

$$\begin{aligned} (X^T X + \lambda I_p) \cdot u_i &= X^T X u_i + \lambda I_p u_i \\ &= \lambda_i u_i + \lambda \cdot u_i \\ &= \underbrace{(\lambda_i + \lambda)}_{\mu_i} \cdot u_i = \mu_i \cdot u_i \end{aligned}$$

Donc  $X^T X + \lambda I_p$  a les mêmes vecteurs  
propres que  $X^T X$  associés aux  
valeurs propres

$$\mu_i = \lambda_i + \lambda.$$

Comme  $\lambda_i \geq 0$  et  $\lambda > 0$  alors

$$\mu_i > 0.$$

et donc  $X^T X + \lambda I_p$  est inversible

On sait que  $\hat{\beta}_{MCO}$  est sans biais

$$E(\hat{\beta}_{MCO}) = \beta$$

Proposition: L'estimateur ridge est biaisé

En effet:

$$\begin{aligned}\hat{\beta}_{\text{ridge}} &= (X^T X + \lambda I_p)^{-1} X^T Y \\ &= (X^T X + \lambda I_p)^{-1} \underbrace{(X^T X)}_{I_p} \cdot \underbrace{(X^T X)^{-1} X^T Y}_{\hat{\beta}_{MCO}} \\ &= (X^T X + \lambda I_p)^{-1} X^T X \cdot \hat{\beta}_{MCO}\end{aligned}$$

Donc

$$E(\hat{\beta}_{\text{ridge}}) = (X^T X + \lambda I_p)^{-1} (X^T X) \cdot \beta$$

$$\text{Biais: } \beta - E(\hat{\beta}_{\text{ridge}}) = \underline{\underline{\left( I - (X^T X + \lambda I_p)^{-1} X^T X \right) \beta}}$$

7

Par contre, la variance des estimateurs ridge est plus petite que ceux MCO

$$\text{Var}(\hat{\beta}_{\text{ridge}, j}) = \frac{\sigma^2}{n-p-1} \sum_{k=1}^p u_{kj}^2 \cdot \frac{1}{\lambda_k + \lambda}$$

et, de manière générale,

$$V(\hat{\beta}_{\text{ridge}}) = \sigma^2 \cdot W X^T X W$$

avec  $W = (X^T X + \lambda I_p)^{-1}$

et

$$\text{Biais}(\hat{\beta}_{\text{ridge}}) = -\lambda W \beta$$

Remarque importante

① On peut montrer (calcul direct) que l'estimateur ridge est le minimiseur du problème :

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \underbrace{\sum_{i=1}^n (Y_i - (\beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2}_{\text{MCO}} + \lambda \underbrace{\sum_{j=1}^p \beta_j^2}_{\text{pénalité}} \right\}$$

$\lambda$  est donc un paramètre que "joue" (arbitre) entre l'adéquation aux données (MCO) et la "grandeur" des paramètres (stabilité)

$$\lambda \rightarrow 0 \Rightarrow \hat{\beta}_{\text{ridge}} \rightarrow \hat{\beta}_{\text{MCO}}$$

$$\lambda \rightarrow \infty \Rightarrow \hat{\beta}_{\text{ridge}} \rightarrow 0$$

② En format matriciel

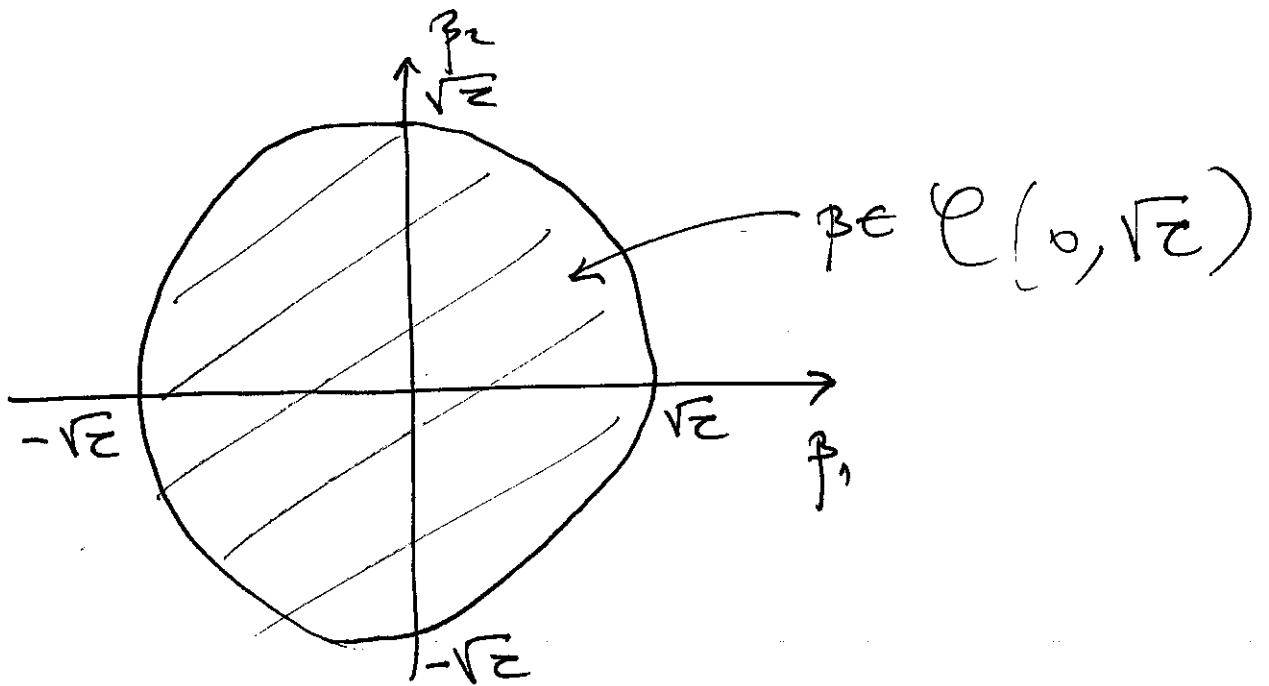
$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|^2 + \lambda \|\beta\|_2^2$$



Remarque importante (continuation)

② On peut montrer que  $\hat{\beta}_{\text{ridge}}$  est solution du pb. de minimisation sous contrainte (multiplicateurs de Lagrange)

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - (\beta_1 X_{i1} + \dots + \beta_p X_{ip}))^2$$
$$\|\beta\|_2^2 < \tau, \quad \tau > 0.$$



$\lambda$  est alors fonction de  $\tau$

$$\underline{\lambda = g(\tau)}$$

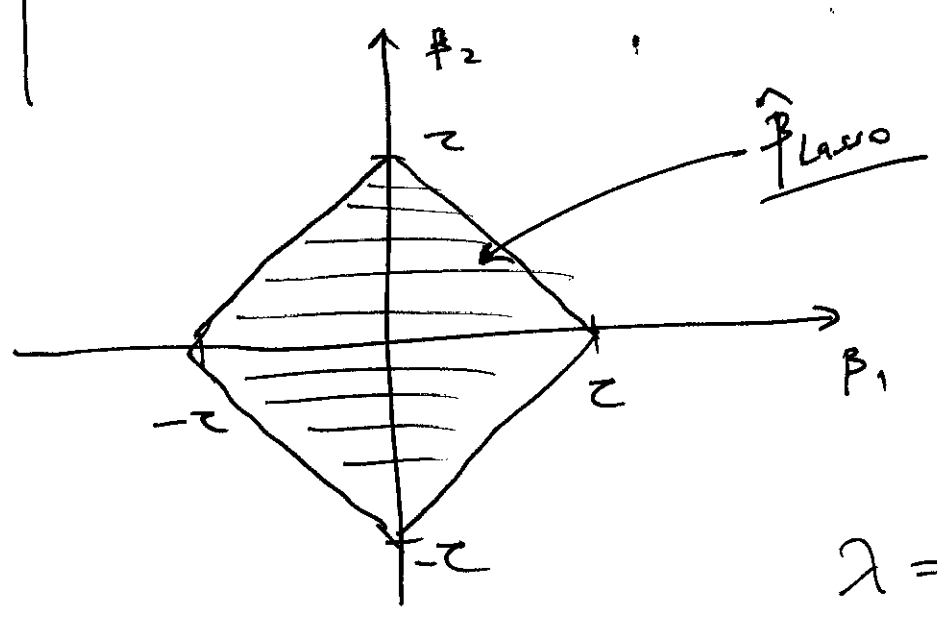
# Régression LASSO

$$\hat{\beta}_{\text{Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

avec  $\|\beta\|_1 = |\beta_1| + |\beta_2| + \dots + |\beta_p|$   
(la norme  $L_1$ )

Ce problème est équivalent à

$$\hat{\beta}_{\text{Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2$$
$$\|\beta\|_1 \leq \tau, \quad \tau > 0$$



$$\lambda = g(\tau)$$

! Sparsité :

La solution LASSO est telle que certains  $\hat{\beta}_{Lasso, j} = 0$

Donc Lasso régularise les coefficients et réalise également une sélection de variables

Elastic-net :

$$\hat{\beta}_{Elastic} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2 + \lambda \left[ (1-\alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1 \right]$$

$\alpha = 0 \Rightarrow$  ridge  
 $\alpha = 1 \Rightarrow$  Lasso

shrinkage + sélection des variables  
(réduction de variance) (Lasso)