

## Retour sur les moindres carrés et la régression Ridge

Cours: SALMON Joseph

Scribes: ABOUQATEB Mouad et KANDOUCI Walid

## 1 Ridge (Régression)

Considérons le modèle linéaire :

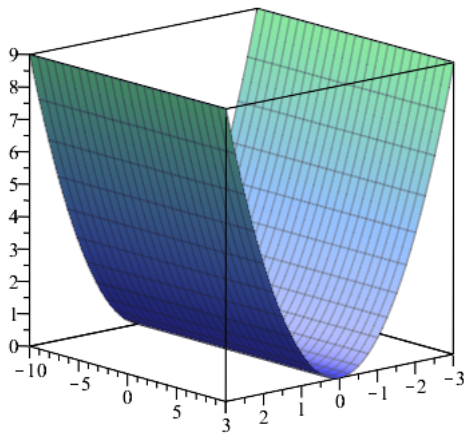
$$Y = X\beta^* + \varepsilon$$

Où  $Y$  est la matrice observations définies sur  $\mathbb{R}^n$ ,  $X$  la matrice des covariables définie sur  $\mathbb{R}^{n \times p}$ ,  $B^*$  dans  $\mathbb{R}^p$  est l'opérateur linéaire et  $\varepsilon$  le bruit.

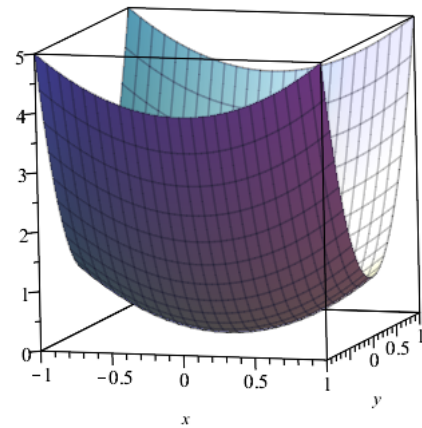
Considérons la matrice  $(X^T X)$ , c'est la matrice de covariance des covariables, nommée **matrice de Gram**.

Le but est d'exprimer  $Y$  en fonction des variables explicatives  $X$  sous forme de l'opérateur linéaire  $\beta$  dans les paramètres inconnus du modèle.

L'approche naturelle pour traiter ce problème est d'estimer  $\beta$  en utilisant la méthode des moindres carrés. (LS)



(a)  $(x_1, x_2) \rightarrow x_1^2$  de forme matricielle donnée par  $f_1(x) = x^T A x = x^T \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} x$  avec  $A$  non-inversible



(b)  $(x_1, x_2) \rightarrow x_1^2 + x_2^2$  de forme matricielle  $f_2(x) = x^T B x = x^T \begin{pmatrix} 1 + \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix} x$  avec  $B$  inversible.

FIGURE 1 – Moindres carrés : cas non unique, avant et après régularisation

On cherche  $\beta$  minimisant l'erreur :

$$\hat{\beta}^{LS} \in \text{Argmin} \|y - x\beta\|^2 \quad (1)$$

Pour cela il suffit de trouver  $\beta$  vérifiant l'équation du gradient au minimum, qui s'écrit :

$$X^T(X\hat{\beta}^{LS} - y) \quad (2)$$

Notre  $\hat{\beta}$  vérifie donc :  $(X^T X)\hat{\beta}^{LS} = X^T y$

Remarquons que notre **matrice de Gram** ( $X^T X$ ) est symétrique donc diagonalisable dans une base orthonormée, et semi-définie, (i.e.,  $v^T X^T X v = \|Xv\|^2 \geq 0$ ), pas forcément définie, ce qui pourra nous empêcher dans un premier temps d'utiliser l'estimateur de la méthode de moindres carrés, qui correspond à  $(X^T X)^{-1} X^T y$  néanmoins toutes ses valeurs propres sont donc positives, ou nulles.

Il existe donc ( $U$ ) matrice orthogonale (son inverse est sa transposée), et  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$  dans  $\mathbb{R}^+$ , tel que :

$$(X^T X) = U^T \text{diag}(\lambda_1, \dots, \lambda_p) U \quad (3)$$

C'est la décomposition spectrale de  $X^T X$ , pour résoudre notre problème d'inversibilité, on en ajoute une pénalité sous forme  $\lambda Id$  avec  $\lambda > 0$ , notre nouvelle matrice  $\text{diag}(\lambda_1, \dots, \lambda_n)$  sera définie positive et par suite inversible, d'inverse  $U^T \cdot \text{diag}(\frac{1}{\lambda_1 + \lambda}, \dots, \frac{1}{\lambda_p + \lambda}) \cdot U$ .

Nous définissons ainsi notre nouvel estimateur de l'opérateur linéaire :  $\hat{B}^{ridge} = (x^T x + \lambda Id_p)^{-1} x^T y$

**Remark.** Remarquons que notre estimateur Ridge  $\hat{\beta}^{ridge(\lambda)}$  tends vers l'estimateur des moindres carrés  $\hat{\beta}^{LS}$  quant  $\lambda \rightarrow 0$ . En effet pour un modèle de plein rang ( $X^T X$  inversible) l'estimateur des moindres carrés correspond à un  $\hat{\beta}^{ridge(0)}$

**Theorem.**  $\hat{\beta}^{ridge(\lambda)} \in \text{Argmin}_{\lambda_1, \dots, \lambda_T} \|y - x\beta\|^2 + \lambda \|\beta\|^2$

*Démonstration.* On définit  $g := \|y - xa\|^2 + \lambda \|\beta\|^2$ ,  $f := \|y - xa\|^2$  et  $\nabla g(\hat{\beta}^\lambda) = 0$

Notons  $\hat{\beta}^\lambda$  l'argument minimal de  $g$ , et montrons que  $\hat{\beta}^\lambda = \hat{\beta}^{ridge}$ .

Notre  $\hat{\beta}^\lambda$  vérifie l'équation du gradient au minimum de  $g$  (condition du premier ordre) est :

$$\begin{aligned} \nabla g(\hat{\beta}^\lambda) &= \nabla f(\hat{\beta}^\lambda) + 2\lambda\beta = 0 \\ (\text{avec}) \quad \nabla f(\hat{\beta}^\lambda) &= 2x^T x\beta - 2(x^T y)^T \beta \end{aligned}$$

En sommant nos deux équations, on obtient :

$$x^T x\beta\lambda - x^T y + \lambda\beta^\lambda = 0$$

Et puisque  $(X^T X + \lambda Id_p)$  est inversible comme on l'a déjà vue (1<sup>re</sup> partie), on déduit le résultat :

$$\hat{\beta}^{ridge(\lambda)} = (x^T x + \lambda Id_p)^{-1} x^T y$$

□

**Remark.** On peut constater à travers notre théorème que le rôle du coefficient  $\lambda$ , sera de privilégier la minimisation des erreurs au profit du biais (régularisation) et vice versa. En effet un  $\lambda$  de plus en plus grand impliquera des prédictions de moins en moins sensibles à la variation des variables, mathématiquement

parlant, un opérateur linéaire de plus en plus négligeable  $\hat{\beta} \xrightarrow{\lambda \rightarrow \infty} 0$ . Respectivement un  $\lambda$  négligeable nous donnera un estimateur proche des moindres carrés (**Remarque 1**).

La question qui se pose est quelle sera le choix optimale de  $\lambda$ .

$$Sp(x^T x + \lambda Id_p)^{-1} = \left( \frac{1}{\lambda_1 + \lambda} \cdots \frac{1}{\lambda_p + \lambda} \right) \quad (4)$$

## 2 Validation croisé ("Hold out")

Cette méthode consiste à décomposer notre matrice d'observations de taille  $n$  en deux sous-matrice,  $n_1$  observations d'"apprentissage", et  $n_2$  observations de "validation", (tel que  $n_1 + n_2 = n$ )

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}; x = \begin{pmatrix} x^A \\ \vdots \\ x^V \end{pmatrix}; y = \begin{pmatrix} y^A \\ \vdots \\ y^V \end{pmatrix}$$

et d'appliquer notre régression Ridge, pour différent choix de  $\lambda_i$ , c sous  $\lambda = (\lambda_1, \dots, \lambda_T)$

Pour des raisons pratiques, on considère une grille géométrique

Le choix du  $\lambda$  optimal est :

$$\hat{\lambda}^{cv} = \arg \min_{\lambda_1, \dots, \lambda_T} \|y_v - X_v \hat{\beta}^\lambda\|^2 \quad (5)$$

Et pour des raisons computationnelles, on procède par la méthode "Leave on out", qui comme son nom l'indique, consiste a valider le modèle sur la  $n^e$  observation après avoir appris sur  $n-1$  observations et l'on répète cette opération  $n$  fois.  $n_1 = n - 1$   $n_2 = 1$

(Paramètre par défaut courant,  $n_1 \approx 80\%n$ ;  $n_2 \approx 20\%n$ )

**Proposition.**  $(x^T + \lambda Id_p)$  est symétrique définie positif

*Démonstration.* Notre matrice est symétrique comme somme de deux matrices symétriques.

Soit  $\beta$  non-nul dans  $\mathbb{R}^P$ , d'après **(1)**, et en conservant les mêmes notations :

$$X^T X + \lambda Id_p = U^T \text{diag}(\lambda_1 + \lambda, \dots, \lambda_n + \lambda) \beta^T (X^T X + \lambda Id_p) U + \lambda U^T U = \text{diag}(\lambda + \lambda_1, \dots, \lambda + \lambda_p)$$

Notre matrice est donc diagonalisable dans une base orthonormée avec des valeurs propres **strictement** positives, par conséquent, elle est définie positive.

Le risque(estimation) quadratique moyen est donnée par  $\mathbb{E}\|\hat{\beta} - \beta^*\|^2$  et le risque (prédiction) quadratique moyen avec  $\hat{\beta}$  vecteur de prédiction est donnée par  $\mathbb{E}\|X\hat{\beta} - X\beta^*\|^2$

(On rappelle que :  $\hat{\beta}^{LS} = (X^T X)^{-1} X^T y$ )

$$\begin{aligned}
R &= \mathbb{E} \|(X^T X)^{-1} X^T y - \beta^*\|^2 \quad (\text{Définie et de plein rang}) \\
&= \mathbb{E} \|(X^T X)^{-1} X^T (X\beta^* + \sigma\varepsilon) - \beta^*\|^2 \\
&= \mathbb{E} \|(X^T X)^{-1} X^T X\beta^* + \sigma(X^T X)^{-1} X^T \varepsilon - \beta^*\|^2 \\
&= \mathbb{E} \|\Sigma(X^T X)^{-1} X^T \varepsilon\|^2 \\
&= \sigma^2 \mathbb{E} \|(X^T X)^{-1} X^T \varepsilon\|^2 \\
&= \sigma^2 \mathbb{E} [\varepsilon^T X (X^T X)^{-1} (X^T X)^{-1} X^T \varepsilon] \\
&= \sigma^2 \text{tr}(\varepsilon^T X (X^T X)^{-1} (X^T X)^{-1} X^T \varepsilon) \\
R &= \sigma^2 \mathbb{E}(\text{tr}[\varepsilon^T X (X^T X)^{-2} X^T \varepsilon])
\end{aligned}$$

Or on a pour tout vecteur  $v$  de taille finie,  $v^T x = \|v\|^2 \in \mathbb{R}$  et pour tout réel  $a$ ,  $a = \text{tr}(a)$  en particulier  $\|v\|^2 = a$ , d'où :

$$R = \sigma^2 \text{tr}(\varepsilon^T X (X^T X)^{-1} (X^T X)^{-1} X^T \varepsilon) \quad (6)$$

Donc :

$$R = \sigma^2 \mathbb{E}(\text{tr}[\varepsilon^T X (X^T X)^{-2} X^T \varepsilon])$$

□

**Remark.**  $\text{tr}(AB) = \text{tr}(BA)$  même si  $A$  et  $B$  ne sont pas de taille identique. Par exemple :  $\text{tr}(xx^T) = \text{tr}(x^T x)$ , mais il faut que  $AB$  et  $BA$  existe.

$$\begin{aligned}
R &= \sigma^2 \text{tr}[\mathbb{E}((X^T X)^{-1} \varepsilon^T \varepsilon)] \\
&= \sigma^2 \text{tr}(X (X^T X)^{-2} X^T \mathbb{E}(\varepsilon \varepsilon^T)) \\
&= \sigma \text{tr}(X (X^T X)^{-2} X^T) \\
&= \sigma^2 \text{tr}((X^T X)^{-1}) \\
R &= \sigma^2 \sum_{j=1}^P \frac{1}{\lambda_j}
\end{aligned}$$

(Si  $Sp(X^T X) = (\lambda_1, \dots, \lambda_p)$ )

**Remark.** Le risque est proportionnel à  $\sigma^2$ , le risque se détériore si la structure de corrélation devient singulière.

### 3 Pour aller plus loin sur ce thème

Parmi les lectures intéressantes on notera : [Del15]<sup>1</sup>

## Références

[Del15] B. Delyon. Régression, 2015. 4

1. <https://perso.univ-rennes1.fr/bernard.delyon/regression.pdf>