

# Modèle linéaire

## M1 Maths IA

Yannig Goude <sup>1, 2</sup>

<sup>1</sup> EDF R&D, EDF Lab Saclay

<sup>2</sup> Laboratoire de Mathématiques d'Orsay

Janvier 2025

# Sommaire

**1** Régression linéaire

**2** Régression linéaire pénalisée

**3** Régression quantile

# Régression linéaire

# Cadre général

## Formulation

La régression recouvre plusieurs méthodes d'analyse statistique permettant d'approcher une variable aléatoire  $Y$  partir un ensemble d'autres variables aléatoires  $X_1, X_2, \dots, X_p$ , regroupées dans un vecteur aléatoire  $\mathbf{X}$ , qui lui sont corrélées

Soient

- $\ell : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}^+$  une **fonction de coût**
- $\mathcal{F}$  un **espace de fonctions** dans lequel la solution est recherchée

L'objectif est de résoudre le problème de minimisation suivant :

$$\min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\mathbf{X}))]$$

# Cadre général

## En pratique

- $\mathbb{E}[\ell(Y, f(\mathbf{X}))]$  n'est pas connue
- $\mathcal{F}$  et  $\ell$  peuvent être choisis

$\mathbb{E}[\ell(Y, f(\mathbf{X}))]$  est estimée grâce à un **échantillon de réalisations** de  $(Y, f(\mathbf{X}))$

Abus de notation :

- $\mathbf{Y} = (y_1, y_2, \dots, y_n)$  est un vecteur de  $n$  réalisations de la variable aléatoire  $Y$
- $\mathbf{X}$  est une matrice à  $n$  lignes et  $p$  colonnes de  $n$  réalisations  $\mathbf{X}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  du vecteur aléatoire  $\mathbf{X}$

# Cadre général

## Approximation

$\mathbb{E} [\ell(Y, f(\mathbf{X}))]$  est approximée par sa moyenne empirique :

$$\mathbb{E} [\ell(Y, f(\mathbf{X}))] \approx \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

Nous chercherons donc le meilleur modèle  $\hat{f} : \mathbb{R}^p \rightarrow \mathbb{R}$  tel que

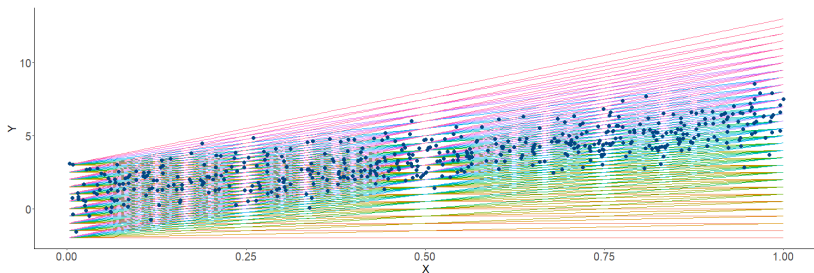
$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

# Exemple

## Régression linéaire uni-variée

$$\mathcal{F} = \{f_{\alpha, \beta} : x \mapsto \alpha + x\beta\}$$

$$\ell : \begin{array}{ccc} \mathbb{R} \times \mathbb{R} & \longrightarrow & \mathbb{R} \\ (y, x) & \longmapsto & (y - x)^2 \end{array}$$

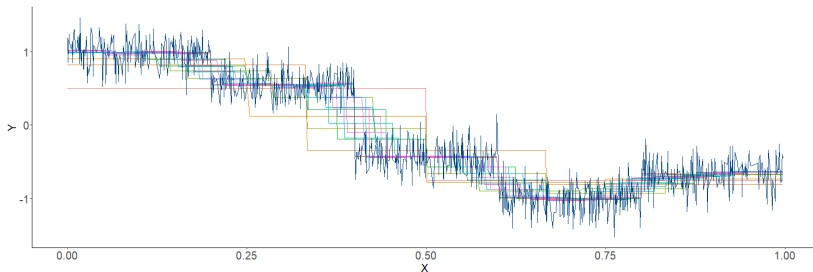


# Exemple

## Détection de ruptures

$$\mathcal{F} = \left\{ f_{x_0, a_0, \dots, x_K, a_K} : x \mapsto \sum_{k=1}^K a_k \mathbb{1}_{x_{k-1} \leq x < x_k}(x) \right\}$$

$$\ell : \begin{array}{ccc} \mathbb{R} \times \mathbb{R} & \longrightarrow & \mathbb{R} \\ (y, x) & \longmapsto & (y - x)^2 \end{array}$$





# Modèle linéaire uni-varié

## Formulation

Soient  $n$  observations  $(y_i, x_i)$ , avec  $y_i \sim Y_i$  où  $Y_i$  est une variable aléatoire

Le modèle linéaire suppose que :

$$y_i = \beta^* x_i + \varepsilon_i$$

avec  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  **indépendants** d'espérance nulle et de variance constante

Le processus  $(\varepsilon_i)_i$  est un bruit blanc

# Modèle linéaire uni-varié

## Moindres carrés ordinaires

L'objectif est de trouver  $\beta$  de façon à minimiser les erreurs d'estimation au carré. En dérivant l'erreur quadratique sur l'échantillon constitué des  $n$  observations  $(y_i, x_i)_{i=1, \dots, n}$

$$\frac{\partial \text{Err}(\beta)}{\partial \beta} = \frac{\partial (\sum_{i=1}^n (y_i - \beta x_i)^2)}{\partial \beta} = - \sum_{i=1}^n 2x_i (y_i - \beta x_i) = 0$$

L'estimateur des moindres carrés ordinaires est

$$\hat{\beta}^{\text{MCO}} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

# Modèle linéaire uni-varié

Distribution de l'estimateur MCO

Sous l'hypothèse de normalité de la variable aléatoire  $Y_i$  :

$$Y_i \sim \mathcal{N}\left(\beta x_i, \sigma^2\right)$$

L'estimateur des moindres carrés ordinaires est sans biais et vérifie

$$\hat{\beta}^{MCO} \sim \mathcal{N}\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n x_i^2}\right)$$

# Modèle linéaire multi-varié

## Formulation

Soient  $n$  observations  $(y_i, x_{i,1}, x_{i,2}, \dots, x_{i,p})$ , le modèle linéaire suppose que

$$y_i = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p} + \varepsilon_i$$

Avec

- $x_{i,j} \in \mathbb{R}$  considéré comme déterministe
- $\beta_1, \beta_2, \dots, \beta_p$  inconnus
- $\mathbb{E}[\varepsilon_i] = 0$  et  $\text{Var}(\varepsilon_i) = \sigma^2$

Remarque : une constante peut être incluse ou non dans modèle ( $x_{i,1} = 1$ )

# Modèle linéaire multi-varié

## Écriture matricielle

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

■  $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$  vecteur d'observation et  $\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \in \mathbb{R}^n$  vecteur de bruit

■  $X = \begin{bmatrix} x_{1,1} & \dots & x_{1,p} \\ \vdots & & \vdots \\ x_{n,1} & \dots & x_{n,p} \end{bmatrix} \in \mathbb{R}^{n \times p}$  matrice de *design* de rang  $p$

■  $\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} \in \mathbb{R}^p$  à estimer

# Modèle linéaire multi-varié

## Moindres carrés ordinaires

Minimisation de la perte quadratique :

$$\hat{\beta} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^p \beta_j x_{i,j} \right)^2 = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \| \mathbf{Y} - \mathbf{X}\beta \|^2$$

L'estimateur des moindres carrés est obtenu en annulant la différentielle en  $\beta$  la perte quadratique, qui est convexe

$$\hat{\beta}^{MCO} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

# Modèle linéaire multi-varié

## Biais et variance de l'estimateur MCO

L'estimateur des moindres carrés ordinaire est non biaisé

$$\begin{aligned}\mathbb{E} \left[ \hat{\beta}^{MCO} \right] &= \mathbb{E} \left[ (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \right] \\ &= \mathbb{E} \left[ (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \beta + \epsilon \right] \\ &= \beta\end{aligned}$$

et de variance

$$\begin{aligned}\text{Var} \left( \hat{\beta}^{MCO} \right) &= \text{Var} \left( (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \right) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var} (\mathbf{Y}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2\end{aligned}$$

# Modèle linéaire multi-varié

MCO et maximum de vraisemblance

La vraisemblance de  $\beta$  au vu des  $n$  observations ( $\sim$  probabilité d'observer ces observations si celles si sont bien distribuées selon le modèle défini par  $\beta$ ) dans le cas où le bruit est gaussien s'écrit

$$L(\mathbf{X}, \beta, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\|\mathbf{Y} - \mathbf{X}\beta\|^2}{2\sigma^2}\right)$$

Maximiser la vraisemblance revient à minimiser  $\|\mathbf{Y} - \mathbf{X}\beta\|^2$  dans le cas gaussien, l'estimateur du maximum de vraisemblance donc égal à l'estimateur des moindres carrés ordinaires

Lorsque les données ne respectent plus l'hypothèse d'indépendance ou de variance constante :  $\mathbf{Y} \sim (\mathbf{X}\beta, \mathbf{V}\sigma^2)$  où  $\mathbf{V}$  est une matrice définie positive,

$$L(\mathbf{X}, \beta, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2|\mathbf{V}|}} \exp\left(-\frac{(\mathbf{Y} - \mathbf{X}\beta)^\top \mathbf{V} (\mathbf{Y} - \mathbf{X}\beta)}{2\sigma^2}\right)$$

et les deux estimateurs ne sont plus égaux.



# Modèle linéaire généralisé

## Formulation

Avec les mêmes notation, un modèle linéaire généralisé suppose

$$g(\mathbb{E}[Y]) = \mathbf{X}\beta$$

où  $g$  est une fonction de lien monotone et régulière. Les observations sont toujours supposées indépendantes et suivent une distribution exponentielle.

Une variable aléatoire  $Y$  est dans la famille exponentielle si sa densité de probabilité dépend de trois fonctions  $a$ ,  $b$  et  $c$  d'un paramètre d'échelle  $\phi$  et d'un paramètre canonique  $\theta$  de sorte que :

$$f_{\theta}(y) = \exp \left( \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right)$$

# Modèle linéaire généralisé

Famille exponentielle

	Gaussienne ( $\mu, \sigma^2$ )	Poisson ( $\lambda$ )	Binomiale ( $n, p$ )	Gamma ( $\alpha, \beta$ )
$\theta$	$\mu$	$\log \lambda$	$\log \frac{p}{1-p}$	$-\frac{\alpha}{\beta}$
$\phi$	$\sigma^2$	1	1	$\frac{1}{\alpha}$
$a(\phi)$	$\phi$	$\phi$	$\phi$	$\phi$
$b(\theta)$	$\frac{\theta^2}{2}$	$e^\theta$	$n \log(1 + e^\theta)$	$-\log -\theta$
$c(y, \phi)$	$\frac{1}{2} \left[ \frac{y^2}{\phi} + \log 2\pi\phi \right]$	$-\log y !$	$\log \binom{n}{y}$	$\frac{1}{\phi} \log \left( \frac{y}{\phi} \right) - \log \left( y \Gamma \left( \frac{1}{\phi} \right) \right)$
$f(y)$	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$	$\frac{\lambda^y e^{-\lambda}}{y !}$	$\binom{n}{y} p^y (1-p)^{n-y}$	$\frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y}$

Exemples :

- Gaussienne : modélisation de la consommation électrique
- Poisson : modélisation d'arrivée/départ à des stations de charge de véhicules électriques

# Modèle linéaire généralisé

Maximum de vraisemblance et Algorithme IRLS

Si la variable aléatoire  $Y$  est dans la famille exponentielle alors

$$\mathbb{E}[Y] = b'(\theta) \text{ et } \text{Var}(Y) = b''(\theta)a(\phi)$$

Dans le modèle linéaire généralisé,  $g(\mathbb{E}[Y]) = \mathbf{X}\beta$  et la vraisemblance de  $\beta$  au vu des  $n$  observations s'écrit

$$L(\mathbf{X}, \beta) = \prod_{i=1}^n f_{a_i, b_i, c_i, \theta_i, \phi_i}(y_i)$$

Comme il est alors difficile de maximiser de manière exacte la vraisemblance, la méthode de Newton (méthode numérique avec étape de calcul du gradient et de la Hessienne de la log-vraisemblance) est utilisée pour estimer itérativement  $\beta$

À chaque itération  $k$ ,  $\beta^{[k]}$  est solution d'un problème des moindres carrés pondérés

→ **Algorithme IRLS** : iteratively re-weighted least square (cf. Wood)

# Régression linéaire pénalisée

# Compromis Biais-Variance

## Modèle linéaire pour la prévision

L'hypothèse de modèle linéaire permet d'estimer  $\hat{f}(\mathbf{X}) = \mathbf{X}\hat{\beta}$  à partir d'un échantillon de  $n$  observations  $(y_i, \mathbf{X}_i)_{i=1, \dots, n}$  grâce la méthode des moindres carrés ordinaires (estimateur sans biais et de variance minimale parmi les estimateurs sans biais - Théorème de Gauss-Markov)

Objectif en prévision :

- Prévision de  $y_{n+1}$  à partir des variables explicatives  $\mathbf{X}_{n+1}$

$$\hat{y}_{n+1} = \mathbf{X}_{n+1}\hat{\beta}$$

- Observation de  $y_{n+1}$  et calcul de l'erreur quadratique (MSE Mean Squared Error)

$$(\hat{y}_{n+1} - y_{n+1})^2$$

# Compromis Biais-Variance

## Erreur de prévision

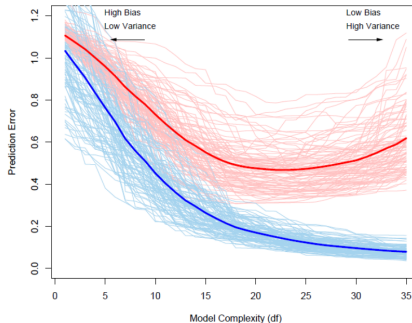
Sous les hypothèses du modèle linéaire, calcul de l'espérance de l'erreur quadratique :

$$\begin{aligned}
 \mathbb{E} \left[ (y_{n+1} - \hat{y}_{n+1})^2 \right] &= \mathbb{E} \left[ (\mathbf{X}_{n+1} \boldsymbol{\beta} + \varepsilon_{n+1} - \mathbf{X}_{n+1} \hat{\boldsymbol{\beta}})^2 \right] \\
 &= \sigma^2 + \mathbb{E} \left[ (\mathbf{X}_{n+1} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}))^2 \right] \\
 &= \sigma^2 + \mathbf{X}_{n+1} \text{Var}(\hat{\boldsymbol{\beta}}) \mathbf{X}_{n+1} + \left( \mathbb{E} [\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}] \right)^2
 \end{aligned}$$

Erreur en prévision = Erreur irréductible + Variance + Biais<sup>2</sup>

# Compromis Biais-Variance

## Illustration



**FIGURE 7.1.** Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error  $\text{err}_T$ , while the light red curves show the conditional test error  $\text{Err}_T$  for 100 training sets of size 50 each, as the model complexity is increased. The solid curves show the expected test error  $\text{Err}$  and the expected training error  $E[\text{err}]$ .

The elements of statistical learning : Data Mining, Inference, and Prediction, Trevor Hastie, Robert Tibshirani and Jerome Friedman, Springer series in statistics - 2001

# Régression Ridge

## Motivation

Cas extrême :

- Modèle uni-varié :  $\mathbf{Y} = \mathbf{X}_1\beta_1 + \varepsilon$
- Ajout d'une variable explicative :  $\mathbf{X}_2 = \mathbf{X}_1 + \text{bruit}$

$$\forall a \in \mathbb{R}, \quad \beta_a = \begin{bmatrix} (a+1)\beta_1 \\ -a\beta_1 \end{bmatrix}$$

est un estimateur qui donne une prévision non biaisée

$$\mathbb{E} [\hat{\mathbf{Y}}] = \mathbb{E} [(a+1)\mathbf{X}_1\beta_1 - a\mathbf{X}_2\beta_1] = \mathbf{X}_1\beta_1 = \mathbb{E} [\mathbf{Y}]$$

mais de variance

$$\begin{aligned} \text{Var}(\hat{\mathbf{Y}}) &= \mathbb{E} \left[ \left( (a+1)\mathbf{X}_1\beta_1 - a\mathbf{X}_2\beta_1 - \mathbf{X}_1\beta_1 \right)^2 \right] \\ &= a^2 \beta_1^2 \text{Var}(\text{bruit}) \end{aligned}$$



# Régression Ridge

## Motivation

Si les coefficients de  $\beta$  ne sont pas contraints, ils peuvent

- exploser
- être soumis à une très grande variance

S'il existe des variables corrélées ( $\beta_1$  et  $\beta_2$  dans l'exemple précédent) dans un modèle de régression linéaire, leurs coefficients peuvent être mal déterminés et présenter une variance élevée : Un coefficient extrêmement élevé sur une variable peut être annulé par un coefficient négatif également élevé sur son cousin corrélé

→ Imposer une contrainte de taille sur les coefficients

# Régression Ridge

## Pénalisation

Pour éviter ces écueils, les coefficients de  $\beta$  doivent être contraints, et le problème de régression Ridge devient

$$\operatorname{argmin}_{\beta \in \mathbb{R}^p} \|Y - X\beta\|^2 \quad \text{avec} \quad \|\beta\|^2 \leq \text{constante}$$

Ce problème équivaut à résoudre, avec  $\lambda > 0$

$$\operatorname{argmin}_{\beta \in \mathbb{R}^p} \|Y - X\beta\|^2 + \lambda \|\beta\|^2 = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \left( y_i - \sum_{j=1}^p \beta_j x_{i,j} + \lambda \sum_{j=1}^p \beta_j^2 \right)$$

# Régression Ridge

## Estimateur, Biais et Variance

L'expression explicite de l'estimateur Ridge s'obtient par différenciation

$$\frac{\partial \left( \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2 \right)}{\partial \beta} = 2\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta) + 2\lambda\beta$$

$$\rightarrow \hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{Y}$$

L'estimateur Ridge est biaisé :

$$\mathbb{E} [\hat{\beta}_\lambda] = \mathbb{E} \left[ (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \varepsilon) \right] = \beta - \lambda (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \beta$$

et de variance

$$\text{Var} (\hat{\beta}_\lambda) = \sigma^2 (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$$

# Régression Lasso

## Motivation et Formulation

La régression Lasso (*least absolute shrinkage and selection operator*) a été introduite dans une optique de sélection de variables et sous l'hypothèse que le vecteur  $\beta$  est parcimonieux (grand nombre de ses coefficients sont nuls)

$$\operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|^2 \quad \text{avec} \quad \|\beta\|_1 \leq \text{constante}$$

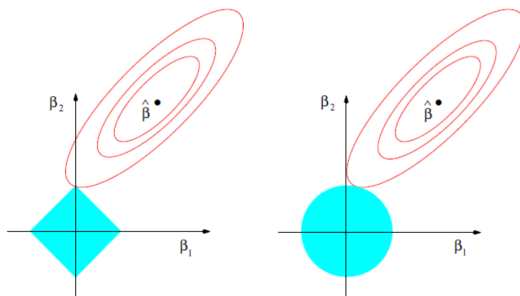
Ce problème équivaut à résoudre, avec  $\lambda > 0$

$$\operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \left( y_i - \sum_{j=1}^p \beta_j x_{i,j} + \lambda \sum_{j=1}^p |\beta_j| \right)$$

Remarque : la norme 1 est préférée à la norme 0 pour garder un problème sous-différentiel (cf. Introduction to High-Dimensional Statistics, Christophe Giraud, 2014)

# Régression Ridge et Lasso

## Illustration



**FIGURE 3.11.** Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.

The elements of statistical learning : Data Mining, Inference, and Prediction, Trevor Hastie, Robert Tibshirani and Jerome Friedman, Springer series in statistics - 2001

# Choix de la pénalité

## Comportements extrêmes

Les estimateurs Ridge et Lasso dépendent fortement du paramètre  $\lambda$

- Chaque  $\lambda$  donne une unique solution
- $\lambda$  est un paramètre de régularisation

- $\lambda = 0 : \hat{\beta}_{\lambda}^{Ridge} = \hat{\beta}_{\lambda}^{Lasso} = \hat{\beta}^{MCO}$

- $\lambda \rightarrow \infty : \hat{\beta}_{\lambda}^{Ridge} = \hat{\beta}_{\lambda}^{Lasso} = \mathbf{0}$

Le paramètre  $\lambda$  gère le compromis biais variance :

- $\lambda = 0 : \mathbb{E} \left[ \hat{\beta}_{\lambda}^{Ridge} \right] = \mathbb{E} \left[ \hat{\beta}_{\lambda}^{Lasso} \right] = \mathbb{E} \left[ \hat{\beta}^{MCO} \right] = \beta$  mais la variance des estimateurs peut exploser (cf. exemple motivation Ridge)
- $\lambda \rightarrow \infty : \text{Var} \left( \hat{\beta}_{\lambda}^{Ridge} \right) = \text{Var} \left( \hat{\beta}_{\lambda}^{Lasso} \right) = \mathbf{0}$  mais les estimateurs sont forcement biaisés (de biais  $-\beta$ )

**Comment choisir  $\lambda$  ?**

# Choix de la pénalité

## Sélection de modèle

Choisir le meilleur  $\lambda$  de façon à avoir des performances optimales en prévision fait parti d'une gamme de problème plus large : **la sélection de modèles**

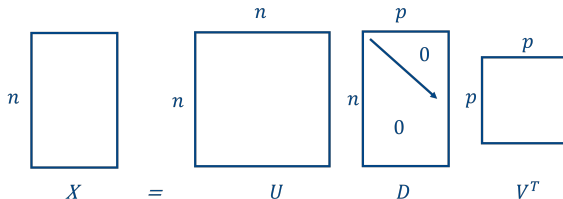
Afin d'avoir un modèle performant en prévision, il convient de construire un ensemble de données dites d'**estimation** (ou d'apprentissage) et d'un second ensemble de données dites de **validation** (ou de test) sur lesquelles le modèle se doit d'être performant

Des techniques de **validation croisée** (qui consiste à tirer plusieurs ensembles de validation d'une même base de données pour obtenir une estimation plus robuste) permettent de trouver les **paramètres de régularisation optimaux**

Avant d'explicitier le critère de validation pour le choix du paramètre de régularisation de la régression Ridge, nous allons définir les degrés de liberté effectifs d'un modèle (qui nécessite un rappel sur la décomposition en valeurs singulières)

# Critère de validation croisée

Décomposition en valeurs singulières (SVD) et nouvelle expression de  $\hat{\beta}_{\lambda}^{Ridge}$



Les matrices  $U$  et  $V$  sont orthogonales ( $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_n$ ) et  $D$  contient dans ses coefficients diagonaux les valeurs singulières de  $\mathbf{X}$  (elles correspondent aux racines des valeurs propres).

$$\hat{\beta}_{\lambda}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{V} \text{Diag} \left( \frac{d_j^2}{d_j^2 + \lambda} \right) \mathbf{U}^T \mathbf{Y}$$



# Critère de validation croisée

Matrice de lissage et *Effective Degrees of Freedom*

Une matrice de lissage  $\mathbf{A}$  est un opérateur linéaire tel que :

$$\hat{\mathbf{Y}} = \mathbf{A} \mathbf{Y}$$

- Moindres Carrés Ordinaires :  $\mathbf{A}^{MCO} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

Notons que la trace de  $\mathbf{A}^{MCO}$ ,  $\text{Tr}(\mathbf{A}^{MCO}) = p$ , le nombre de paramètres à estimer. Par analogie, les *Effective Degrees of Freedom* d'une matrice de lissage  $\mathbf{A}$  sont définis par

$$\text{df}(\mathbf{A}) = \text{Tr}(\mathbf{A})$$

- Régression Ridge :  $\mathbf{A}_\lambda^{\text{Ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top$

$$\text{df}(\mathbf{A}_\lambda^{\text{Ridge}}) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

# Critère de validation croisée

## Principe

$\forall i = 1, \dots, n$

- Retirer l'observation  $(y_i, \mathbf{X}_i)$  de l'échantillon
- Estimer  $\hat{\beta}_{\lambda}^{-i}$  à l'aide de toutes les autres données
- Mesurer l'erreur de prévision  $(y_i - \hat{\beta}_{\lambda}^{-i} X_i)^2$

Le critère de validation croisée s'écrit

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_{\lambda}^{-i} X_i)^2$$

Il est possible de montrer qu'en espérance, le critère de validation croisée est égal à l'erreur de prévision ; il convient donc de le minimiser

# Critère de validation croisée

## Généralisation

Un tel critère prend du temps à être calculé ( $n$  estimateurs pour chaque  $\lambda$ ) !

Mais, pour la régression Ridge, il est possible de montrer que

$$\text{CV}(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_{\lambda}^{-i} x_i)^2 = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{\beta}_{\lambda} x_i)^2}{(1 - \mathbf{A}_{\lambda_{i,i}})^2}$$

Il n'y a donc plus qu'un seul estimateur à calculer !

L'approximation  $\mathbf{A}_{\lambda_{i,i}} \approx \frac{\text{Tr}}{n}$

( $\mathbf{A}_{\lambda}$ ) permet de proposer le critère de validation croisée généralisé :

$$\text{GCV}(\lambda) = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{\beta}_{\lambda} x_i)^2}{\left(1 - \frac{\text{df}(\mathbf{A}_{\lambda})}{n}\right)^2}$$

C'est ce dernier qui est généralement utilisé par les algorithmes pour choisir le meilleur paramètre de régularisation

# Régression quantile

# Régression quantile

## Motivation

Alors que la méthode des moindres carrés fournit une estimation de l'espérance (conditionnellement aux variables explicative) de la variable réponse, la régression quantile cherche à approcher la médiane ou d'autres quantiles de la variable réponse

Elle s'avère utile pour prévoir notamment des seuils et lorsque plusieurs régressions sont réalisées, il est possible d'avoir une bonne idée de la distribution générale de la variable réponse

La régression quantile est moins sensible aux points aberrants. Elle peut être définie comme une régression avec une norme 1

# Régression quantile

## Formulation

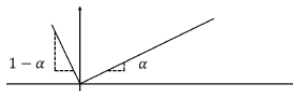
Le quantile  $q_\alpha$  de la variable aléatoire  $Y$  de densité  $f$  et de fonction de répartition  $F$  vérifie

$$F(q_\alpha) = \int_{-\infty}^{q_\alpha} f(y)dy = \mathbb{P}(Y \leq q_\alpha) = \alpha$$

En introduisant la **pinball loss**

$$\ell_\alpha(y - q) = \alpha|y - q|^+ + (1 - \alpha)|y - q|^-,$$

$$\text{avec } |x|^+ = \max(x, 0) \text{ et } |x|^- = \max(-x, 0)$$



Le quantile  $q_\alpha$  minimise la quantité  $\mathbb{E}[\ell_\alpha(Y - q)]$

# Régression quantile

## Preuve

Résolvons le problème de minimisation convexe

$$\operatorname{argmin}_q \mathbb{E} [\ell_\alpha(Y - q)]$$

par différentiation :

$$\begin{aligned} 0 &= \mathbb{E} \left[ \frac{\partial \ell_\alpha(Y - q)}{\partial q} \right] \\ &= \int_{-\infty}^{+\infty} \frac{\partial \ell_\alpha(y - q)}{\partial q} f(y) dy \\ &= -(1 - \alpha) \int_{-\infty}^q f(y) dy + \alpha \int_q^{+\infty} f(y) dy \\ &= (\alpha - 1)F(q) + \alpha(1 - F(q)) = \alpha - F(q) \end{aligned}$$

La solution vérifie donc

$$F(\hat{q}) = \alpha$$

# Régression quantile

## Estimateur et Algorithme

Lorsque l'on dispose de l'échantillon  $(y_1, \dots, y_n)$ , l'estimateur de la régression quantile est obtenu en minimisant le critère

$$\hat{\beta}^\alpha \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \ell_\alpha(y_i, \mathbf{x}_i \beta)$$

Il est possible d'utiliser une méthode de descente de gradient puisque la fonction à minimiser est presque partout dérivable. L'algorithme Iteratively reweighted least squares permet aussi de calculer l'estimateur (plusieurs implémentations existent dans différents packages R ou python)



# Références

- Wood, Simon N. Generalized additive models: an introduction with R. chapman and hall/CRC, 2006.
- Hastie, Trevor, Robert Tibshirani, Jerome H. Friedman, and Jerome H. Friedman. The elements of statistical learning: data mining, inference, and prediction. Vol. 2. New York: springer, 2009.
- Giraud, Christophe. Introduction to high-dimensional statistics. Chapman and Hall/CRC, 2021.