

# Régression linéaire

Pierre Gaillard

12 février 2019

Dans ce cours, nous allons étudier le problème simple mais encore abondamment utilisé de nos jours : la régression<sup>1</sup> linéaire. Ce cours est basé sur les livres RIVOIRARD et STOLTZ, 2012 et CORNILON et MATZNER-LØBER, 2011.

## 1 Introduction

Commençons par un exemple de problème pratique. Pour mieux optimiser sa production, un producteur s'intéresse à modéliser la consommation électrique en France en fonction de la température (cf. Figure 1).

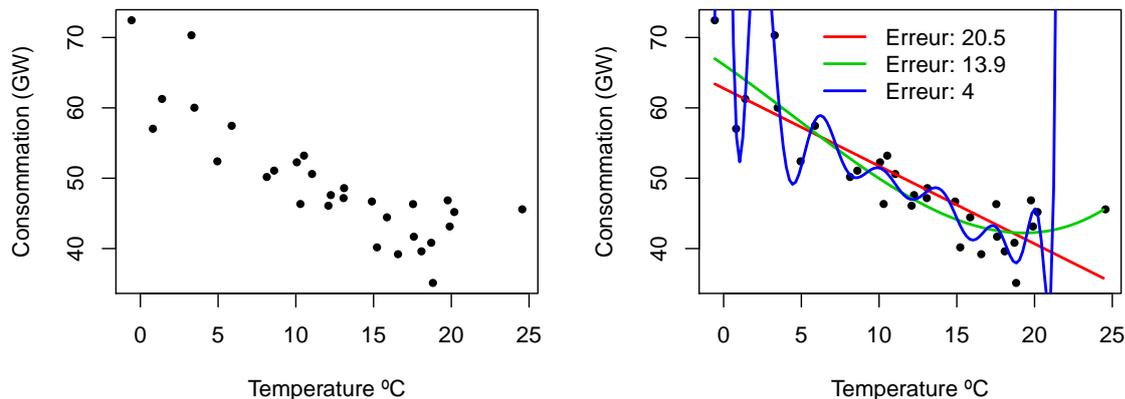


FIGURE 1 – Consommation électrique française (GW) en fonction de la température (°C). À droite, sont tracés fonctions minimisant l'erreur pour les espace de polynômes de degrés 1 (rouge), 3 (vert) et 30 (bleu)

L'objectif est de trouver une fonction  $f$  telle qui explique bien la consommation électrique  $(y_i)_{1 \leq i \leq n}$  en fonction de la température  $(x_i)_{1 \leq i \leq n}$ , soit  $y_i \approx f(x_i)$ . Pour cela, on peut choisir un espace de fonction  $\mathcal{F}$  et résoudre le problème de minimisation du risque empirique :

$$\hat{f}_n \in \arg \min_{f \in \mathcal{F}} \hat{\mathcal{R}}_n(f) := \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2. \quad (1)$$

---

1. Le mot « régression » aurait été introduit par Galton au XIXe siècle. En modélisant la taille des individus en fonction de celle de leurs pères, Galton a observé un retour (« regression » en anglais) vers la taille moyenne. Les pères plus grands que la moyenne ont tendance a avoir des enfants plus petits qu'eux et vice-versa pour les pères plus petits.

Il faut faire attention au choix de l'espace de fonction pour éviter le sur-apprentissage (cf. Figures 1). Bien que l'erreur quadratique moyenne diminue quand l'espace  $\mathcal{F}$  s'agrandit (degrés de polynômes plus grands), l'estimateur  $\hat{f}_n$  perd son pouvoir prédictif. La question qu'il faut se poser est :  $\hat{f}_n$  sera-t-il performant sur de nouvelles données ?

L'espace des fonctions linéaires de la forme  $f : x \mapsto ax + b$  est le plus simple. C'est celui que nous allons étudier dans ce cours.

## Formalisation

En apprentissage statistique, on suppose que les observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$  (ici  $\mathcal{X} = \mathbb{R}^p$  et  $\mathcal{Y} = \mathbb{R}$ ) sont les réalisations de variables aléatoires  $(X_i, Y_i)$ . Le jeu de données aléatoire sera souvent noté  $D_n := \{(X_i, Y_i), i = 1, \dots, n\}$ . Les  $X_i$  à valeur dans  $\mathcal{X}$  sont les variables d'entrée, les  $Y_i$  à valeur dans  $\mathcal{Y}$  sont les variables de sortie. On fera l'hypothèse dans ce cours que les données sont i.i.d.

$$(X_i, Y_i) \stackrel{i.i.d.}{\sim} \nu.$$

La loi  $\nu$  est inconnue du statisticien, il s'agit de l'apprendre à partir des données  $D_n$ . Une règle d'apprentissage  $\mathcal{A}$  est une fonction qui associe à des données d'entraînement  $D_n$  une fonction de prédiction  $\hat{f}_n$  (le chapeau sur  $f$  indique qu'il s'agit d'un estimateur) :

$$\begin{aligned} \mathcal{A} : \cup_{n \in \mathbb{N}} (\mathbb{R}^p \times \mathbb{R})^n &\rightarrow \mathbb{R}^{\mathbb{R}^p} \\ D_n &\mapsto \hat{f}_n \end{aligned}$$

La fonction estimée  $\hat{f}_n$  est construite en vue de prédire  $Y'$  à partir d'un nouveau  $X'$ , où  $(X', Y') \sim \nu$  est une paire de données test, c'est à dire non observée dans les données d'entraînement. La fonction  $\hat{f}_n$  est un estimateur car elle dépend des données  $D_n$  et d'aucun paramètre non observé (comme  $\nu$ ). Comme  $D_n$  est aléatoire, il s'agit d'une fonction aléatoire.

L'objectif est de trouver construire un estimateur  $\hat{f}_n$  qui prévoit bien des nouvelles données en minimisant le risque :

$$\mathcal{R}(\hat{f}_n) := \mathbb{E} \left[ (Y' - \hat{f}_n(X'))^2 \mid D_n \right] \quad \text{où} \quad (X', Y') \sim \nu. \quad (\text{Risque})$$

Cependant le statisticien ne peut pas calculer le risque car il ne connaît pas  $\nu$ . Une méthode courante en apprentissage statistique et intuitive consiste donc à remplacer le risque par le risque empirique

$$\hat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 = \frac{1}{n} \|Y - f(X)\|^2 \quad (\text{Risque empirique})$$

où  $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$  et  $X = (X_1, \dots, X_n)^\top \in \mathbb{R}^{n \times p}$ .

Il faut cependant faire attention au sur-apprentissage (cas où  $\hat{\mathcal{R}}_n(f)$  est beaucoup plus faible que  $\mathcal{R}(f)$ , cf. Figure 2). Dans ce cours, on va étudier la performance du minimiseur des moindres carrés dans le cas du modèle linéaire.

## 2 Le modèle linéaire

Dans cette partie, on pose un cadre stochastique qui va nous permettre d'analyser la performance de minimiseur du risque empirique dans le cas linéaire. Dans un modèle linéaire on suppose une équation de la forme :

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \dots + \beta_p X_{i,p} + \varepsilon_i, \quad (2)$$

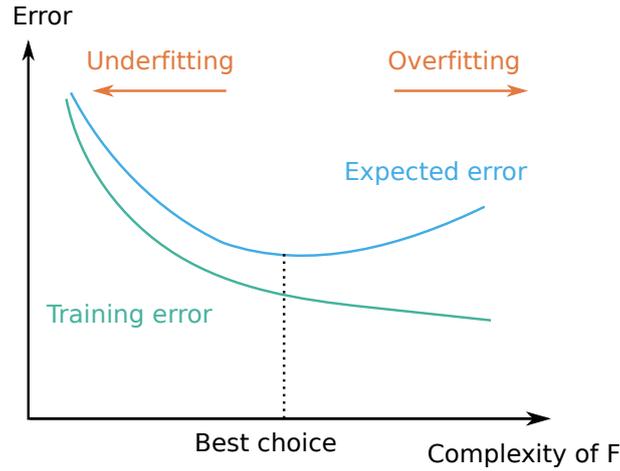


FIGURE 2 – Sur-apprentissage et sous-apprentissage en fonction de la complexité de  $\mathcal{F}$ . En bleu le vrai risque  $\mathcal{R}(f)$ , en vert le risque empirique  $\widehat{\mathcal{R}}_n(f)$ .

où  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top \in \mathbb{R}^n$  est un vecteur d'erreurs (ou de bruit). Il vient du fait qu'en pratique les observations  $y_i$  ne collent jamais complètement à la prévision linéaire. On suppose qu'il s'agit d'un  $n$ -échantillon (i.e., les  $\varepsilon_i$  sont i.i.d.) d'espérance nulle  $\mathbb{E}[\varepsilon_i] = 0$ . Souvent, il est supposé de loi Gaussienne  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . Pour simplifier les notations, nous supposons que le premier vecteur des variables explicatives est le vecteur constant  $X_{\cdot 1} = (1, \dots, 1)^\top$  de sorte que le coefficient  $\beta_0$  de l'équation (2) n'est plus nécessaire. Nous définissons la matrice des variable explicatives

$$X := \begin{bmatrix} 1 & X_{1,2} & \dots & X_{1,p} \\ \vdots & \vdots & \dots & \vdots \\ 1 & X_{n,2} & \dots & X_{n,p} \end{bmatrix}$$

composée des lignes  $(1, X_i) \in \mathbb{R}^p$ . Nous pouvons réécrire la définition du modèle linéaire (2) de façon matricielle.

**Definition 2.1** (Modèle linéaire). *Un modèle linéaire se définit par une équation de la forme :*

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

où  $Y \in \mathbb{R}^n$  est le vecteur des observations,  $X \in \mathbb{R}^{n \times p}$  est la matrice de variables explicatives et  $\beta \in \mathbb{R}^p$  le vecteur de coefficients à estimer. On suppose que la matrice  $X$  est injective (i.e.,  $\text{rg}(X) = p$ ) et que le vecteur de bruit  $\varepsilon \in \mathbb{R}^n$  est un  $n$ -échantillon centré ( $\mathbb{E}[\varepsilon] = 0$  et  $\text{Var}[\varepsilon] = \sigma^2 I_n$ ).

**Design déterministe vs design aléatoire** Dans toute la suite, bien que  $X$  puisse être aléatoire, il s'agit en pratique de variables explicatives observées (ou calculées) par le prévisionniste avant d'effectuer une prévision. On supposera donc qu'il s'agit de variables déterministes. Le cas aléatoire peut se retrouver en considérant des espérances conditionnelles.

### 3 L'estimateur des moindres carrés ordinaire (OLS)

Le jeu de données  $(X_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}$  étant observé, l'objectif est de trouver la fonction affine  $\widehat{f} : x \mapsto x^\top \widehat{\beta}$  minimisant l'erreur quadratique moyenne (1). On suppose que l'intercepte est inclus

dans les variables explicatives de sorte que  $X_{i,1} = 1$ . Dans le cas linéaire, le risque empirique s'écrit

$$\mathcal{R}_n(\beta) := \frac{1}{n} \sum_{i=1}^n \left( Y_i - \beta_1 + \sum_{j=2}^p \beta_j X_{i,j} \right)^2 = \frac{1}{n} \|Y - X\beta\|^2,$$

où  $\|\cdot\|$  est la norme euclidienne.

**Définition 3.1** (Estimateur des moindres carrés ordinaires (OLS)). *On appelle estimateur des moindres carrés ordinaire, le minimiseur  $\hat{\beta} \in \mathbb{R}^p$  du risque empirique :*

$$\hat{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \widehat{\mathcal{R}}_n(\beta) = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|. \quad (3)$$

**Proposition 3.1.** *Sous les hypothèses et les notations de la Définition 2.1, l'OLS existe et est unique. Il s'écrit  $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ .*

*Démonstration.* Par définition, l'OLS est le minimiseur du risque empirique  $\widehat{\mathcal{R}}_n : \mathbb{R}^p \rightarrow \mathbb{R}_+$  définie par : pour tout  $\beta \in \mathbb{R}^p$

$$\widehat{\mathcal{R}}_n(\beta) = \frac{1}{n} \|Y - X\beta\|_2^2 = \frac{1}{n} (\beta^\top (X^\top X) \beta - 2\beta^\top X^\top Y + \|Y\|^2).$$

Pour être un minimum, il est nécessaire que le gradient s'annule. Celui-ci s'écrit :

$$\nabla \widehat{\mathcal{R}}_n(\hat{\beta}) = \frac{1}{n} (\hat{\beta}^\top (X^\top X) + (X^\top X) \hat{\beta} - 2X^\top Y) = \frac{2}{n} ((X^\top X) \hat{\beta} - X^\top Y).$$

où la dernière égalité est parce que la matrice  $X^\top X \in \mathbb{R}^{p \times p}$  est symétrique. Comme  $X$  est de rang  $p$  par hypothèse,  $X^\top X$  est même symétrique définie positive et inversible (laissé en exercice). On en déduit donc qu'un optimum de  $\widehat{\mathcal{R}}_n$  vérifie nécessairement :

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

Il reste cependant à vérifier qu'il s'agit bien d'un minimum et donc que la Hessienne est définie positive. Ce qui est le cas car :  $\nabla^2 \widehat{\mathcal{R}}_n(\hat{\beta}) = \frac{2}{n} (X^\top X)$ .  $\square$

**Une interprétation géométrique** Le modèle linéaire cherche à modéliser le vecteur  $Y \in \mathbb{R}^n$  des observations par une combinaison linéaire de la forme  $X\beta \in \mathbb{R}^n$ . L'image de  $X$  est l'espace des solutions, noté  $\text{Im}(X) = \{Z \in \mathbb{R}^n : \exists \beta \in \mathbb{R}^p \text{ tq } Z = X\beta\} \subseteq \mathbb{R}^n$ . C'est le sous-espace vectoriel de  $\mathbb{R}^n$  engendré par les  $p < n$  colonnes de la matrice de variables explicatives. Comme  $\text{rg}(X) = p$ , il est de dimension  $p$ . Selon la définition du modèle linéaire,  $Y$  est la somme d'un élément de  $\text{Im}(X)$  et d'un bruit. En minimisant  $\|Y - X\beta\|$  (cf. Définition 3.1), on cherche donc l'élément de  $\text{Im}(X)$  le plus proche de  $Y$ . Il s'agit du projeté orthogonal de  $Y$  sur  $\text{Im}(X)$ , noté  $\hat{Y}$ . Par définition de l'OLS et par la Proposition 3.1, on a :

$$\hat{Y} \stackrel{\text{Déf. 3.1}}{=} X\hat{\beta} \stackrel{\text{Prop. 3.1}}{=} X(X^\top X)^{-1} X^\top Y.$$

En particulier,  $P_X := X(X^\top X)^{-1} X^\top$  est la matrice de projection orthogonal sur  $\text{Im}(X)$ .

On peut montrer que l'OLS a de bonnes propriétés statistiques.

**Proposition 3.2.** *Sous les hypothèses de la Définition 2.1, l'estimateur  $\hat{\beta}$  est sans biais,  $\mathbb{E}[\hat{\beta}] = \beta$  et de variance  $\text{Var}(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$ .*

*Démonstration.* On peut calculer son espérance en utilisant  $\mathbb{E}[\varepsilon] = 0$  :

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}[(X^\top X)^{-1} X^\top Y] = \mathbb{E}[(X^\top X)^{-1} X^\top X \beta + \cancel{(X^\top X)^{-1} X^\top \varepsilon}] = \beta.$$

et sa variance en utilisant  $\text{Var}(Y) = \text{Var}(\varepsilon) = \sigma^2 I_n$ .

$$\text{Var}(\hat{\beta}) = \text{Var}((X^\top X)^{-1} X^\top Y) = (X^\top X)^{-1} X^\top \text{Var}(Y) X (X^\top X)^{-1} = \sigma^2 (X^\top X)^{-1}.$$

□

On peut même montrer que l'OLS vérifie la propriété de Gauss-Markov. Il est optimal parmi les estimateurs sans biais de  $\beta$ , dans le sens où il a une matrice de variance-covariance minimale (pour une certaine relation d'ordre partiel). Cependant, la Proposition 3.2 est peu pratique pour estimer la fiabilité de  $\hat{\beta}$  en pratique car elle fait intervenir la variance du bruit, qui est inconnue. Il peut également être intéressant de l'estimer aussi. Un estimateur naturel et de regarder l'erreur moyenne

$$\hat{\sigma}_{\text{nat}}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - X_{i,\cdot} \hat{\beta})^2 = \frac{\|Y - X \hat{\beta}\|^2}{n}.$$

Cependant cet estimateur est biaisé. On utilise donc l'estimateur non biaisé (exercice) :

$$\hat{\sigma}^2 = \frac{\|Y - X \hat{\beta}\|^2}{n - p}.$$

On peut de même estimer le risque associé à l'OLS.

**Proposition 3.3.** Soit  $(X', Y') \in \mathbb{R}^p \times \mathbb{R}$  satisfaisant le modèle linéaire  $Y' = \beta X' + \varepsilon'$ , avec  $\varepsilon'$  bruit centré indépendant de  $\varepsilon \in \mathbb{R}^n$  et de variance  $\text{Var}(\varepsilon') = \sigma^2$ . Alors

$$\mathbb{E}[\mathcal{R}(\hat{\beta})] = \sigma^2 (1 + X'^\top (X^\top X)^{-1} X').$$

*Démonstration.*

$$\begin{aligned} \mathbb{E}[\mathcal{R}(\hat{\beta})] &:= \mathbb{E}[(Y' - X'^\top \hat{\beta})^2] \\ &= \mathbb{E}[(X'^\top \hat{\beta} - \mathbb{E}[X'^\top \hat{\beta}])^2] + \mathbb{E}[(\mathbb{E}[X'^\top \hat{\beta}] - \mathbb{E}[Y'])^2] + \mathbb{E}[(Y' - \mathbb{E}[Y'])^2] \\ &= X'^\top \text{Var}(\hat{\beta}) X' + \sigma^2 \\ &= \sigma^2 (1 + X'^\top (X^\top X)^{-1} X'). \end{aligned}$$

□

*Remarque :*  $X$  et  $X'$  sont supposés déterministes et connus ici. On parle de “deterministic design”. On peut facilement étendre les résultats au cas de design aléatoire, dans quel cas  $X$  et  $X'$  sont aussi aléatoires. Il faut alors expliciter le fait que les espérances  $\mathbb{E}[Y]$  et  $\mathbb{E}[X^\top \beta]$  dans les calculs ci-dessus sont conditionnelles à  $X$  et  $X'$ .

**Le cas Gaussien** Un cas particulier très considéré est celui du bruit Gaussien :  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ . Ce choix vient non seulement du fait qu'il permet de calculer de nombreuses propriétés statistiques supplémentaires sur  $\hat{\beta}$  et de faire des tests (intervalles de confiance, significativité des variables, ...). Il est en pratique motivé par le théorème central limite et le fait que le bruit est souvent une addition de nombreux phénomènes non expliqués par la combinaison linéaire des variables explicatives.

**Proposition 3.4.** *Dans le cas Gaussien, les estimateurs du maximum de vraisemblance de  $\beta$  et  $\sigma$  vérifient respectivement :*

$$\widehat{\beta}_{MV} = (X^\top X)^{-1}XY \quad \text{et} \quad \widehat{\sigma}_{MV}^2 = \frac{\|Y - X\widehat{\beta}\|^2}{n}.$$

On retrouve donc l'estimateur des moindres carrés obtenu par la minimisation du risque empirique. L'estimateur de la variance est biaisé.

**Problème non-linéaire : régression polynomiale, spline, à noyau** L'hypothèse que les observations  $Y_i$  s'explique comme une combinaison des variables explicatives  $X_{i,j}$  peut paraître forte. Cependant, on peut appliquer la théorie précédente sur des transformations des variables  $x_{i,j}$ . Par exemple, en ajoutant les puissance des variables  $X_{i,j}^k$  ou leurs produits  $X_{i,j}X_{i,j'}$ . Cela permet de se comparer aux espaces polynomiaux. Faire une régression linéaire sur des transformation polynomiales des variables revient à faire une régression polynomiale.

Bien sûr d'autres bases (transformations) existent comme des régression sur des bases de splines (polynômes par morceaux avec contraintes sur les bords). C'est le modèle qu'utilise par exemple EDF en opérationnel pour prévoir la consommation électrique en fonction de variables comme l'heure, le jour de la semaine, la température ou la couverture nuageuse. Une autre forme de régression que nous aborderons dans la suite sera la régression à noyau. De façon générale, on peut considérer des transformations  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$  avec le modèle

$$Y_i = \phi(X_i)^\top \beta + \varepsilon_i,$$

pour un paramètre  $\beta \in \mathbb{R}^d$  inconnu.

Dans la Figure 1, on a de cette façon minimisé le risque empirique sur les espaces polynomiaux de degré 1 (modèle linéaire), 3 et 30. On voit qu'il faut faire attention à ne pas considérer des espaces trop grands, au risque que le modèle soit mal posé ( $X$  non injective). Inversement, pour que la Proposition 3.3 soit vérifiée, il faut se trouver dans le vrai modèle  $Y = \phi(X)\beta + \varepsilon$ . Il faut donc s'assurer que  $\phi(X)$  contienne suffisamment de descripteurs pour que la dépendance entre  $Y$  et  $\phi(X)$  soit effectivement linéaire. Autrement on paye un terme de biais supplémentaire.

## 4 Régularisation

Si  $X$  n'est pas injective (i.e.,  $\text{rg}(X) = p$ ), la matrice  $(X^\top X)$  n'est plus inversible et le problème admet plusieurs solutions qui minimisent le risque empirique. On dit que le problème est mal posé ou non identifiable.

La Proposition 3.3 nous rappelle que le risque augmente proportionnellement à la variance de  $\widehat{\beta}$  qui elle-même dépend du conditionnement de la matrice  $(X^\top X)^{-1}$ . Plus les colonnes de cette dernière risquent d'être dépendantes et moins  $\widehat{\beta}$  sera précis. Plusieurs solutions permettent de traiter le cas où  $\text{rg}(X) < p$  :

- *contrôle explicite de la complexité* en réduisant l'espace des solutions  $\text{Im}(X)$ . Cela peut se faire en retirant des colonnes de la matrice  $X$  jusqu'à ce qu'elle deviennent injective (par exemple, en réduisant le degré des polynômes). On peut aussi poser des contraintes d'identifiabilité de la forme  $\beta \in V$  un sous-espace vectoriel de  $\mathbb{R}^p$  tel que tout élément  $y \in \text{Im}(X)$  a un unique antécédent  $\beta \in V$  avec  $y = X\beta$ . On peut par exemple choisir  $V = \text{Ker}(X)^\perp$ .
- *contrôle implicite de la complexité* en régularisant le problème de minimisation du risque empirique comme on le voit rapidement ci-après avec les régression Ridge et Lasso.

**Régression Ridge** La régularisation la plus courante est la *régression Ridge*, où on utilise la norme 2 :

$$\widehat{\beta}_{(\text{ridge})} \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\}.$$

Le paramètre de régularisation  $\lambda > 0$  règle le compromis entre la variance de  $\widehat{\beta}$  et son biais. En refaisant les calculs de la Proposition 3.1, on voit que la solution est unique (même si  $X$  n'est pas injective) donnée par :

$$\widehat{\beta}_{(\text{ridge})} = (X^\top X + n\lambda I_n)^{-1} X^\top Y.$$

On voit qu'il n'y a plus le problème d'inversion de  $X^\top X$  puisque la régression Ridge revient à remplacer  $(X^\top X)^{-1}$  par  $(X^\top X + n\lambda I_n)^{-1}$  dans la solution de l'OLS. Le calibration du paramètre  $\lambda$  est cependant essentielle en pratique. Elle peut par exemple se faire de façon analytique, par *validation croisée (généralisée)* ou par *l'heuristique de pente* (cf. cours sur la sélection de modèle).

**Proposition 4.1** (Biais et variance de  $\widehat{\beta}_{(\text{ridge})}$ ). *L'estimateur Ridge vérifie :  $\widehat{\beta}_{(\text{ridge})} = (X^\top X + \lambda I_p)^{-1} X^\top X \beta$ ,  $\mathbb{E}[\widehat{\beta}_{(\text{ridge})}] = \beta - \lambda (X^\top X + \lambda I_p)^{-1} \beta$  et  $\text{Var}(\widehat{\beta}_{(\text{ridge})}) = \sigma^2 (X^\top X + \lambda I_p)^{-1} X^\top X (X^\top X + \lambda I_p)^{-1}$ .*

L'estimateur Ridge est donc biaisé au contraire de l'OLS ce qui constitue un handicap. Cependant, sa variance ne fait pas intervenir l'inverse de  $X^\top X$  mais de  $X^\top X + \lambda I_p$  qui est mieux conditionnée. Il aura donc une plus faible variance. Exercice : calculer le risque de l'estimateur Ridge.

**Le Lasso** Une deuxième régularisation courante est la régularisation par la norme 1. On parle alors de l'estimateur Lasso.

$$\widehat{\beta}_{(\text{lasso})} \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}.$$

L'estimateur Lasso n'a pas de formule analytique dans le cadre général. Cependant, dans le cas où  $X^\top X = I_p$ , le Lasso a une solution explicite. L'estimateur Lasso correspond alors à un seuillage doux de  $\widehat{\beta}$  la solution des moindres carrés (OLS) : pour tout  $j = 1, \dots, p$

$$\widehat{\beta}_{\text{lasso}} = \text{sgn}(\widehat{\beta}) \max(0, |\widehat{\beta}| - \lambda).$$

L'estimateur Lasso est en particulier très utilisé quand le nombre de variables explicatives est très grand voir supérieur au nombre d'observations :  $p \gg n$ . Le caractère "anguleux" de la norme 1 va avoir tendance à tronquer (annuler) de nombreux coefficients de  $\widehat{\beta}_{\text{lasso}}$  pour proposer des solutions parcimonieuses. Ceci se voit mieux avec la formulation équivalente du Lasso sous la forme d'un problème de régularisation sous contrainte : pour chaque  $\lambda > 0$ , il existe  $U = U(X) > 0$  tel que :

$$\widehat{\beta}_{(\text{lasso})} \in \arg \min_{\beta \in \mathbb{R}^p: \|\beta\|_1 \leq U} \|Y - X\beta\|^2.$$

Ridge se reformule de la même façon en remplaçant la norme 1 par la norme 2 au carré.

## 5 Calcul de $\widehat{\beta}$

La formule close  $\widehat{\beta} = (X^\top X)^{-1} X^\top Y$  de l'OLS est utile pour l'analyser. Cependant, la calculer de façon naïve peut s'avérer prohibitif. En particulier quand  $p$  est grand, on préfère éviter l'inversion de la matrice  $X^\top X$  qui coûte  $\mathcal{O}(p^3)$  par la méthode de Gauss-Jordan et qui peut être très instable quand la matrice est mal conditionnée.

**Décomposition QR** Cependant, pour améliorer la stabilité, on peut utiliser la décomposition QR. Rappelons que  $\hat{\beta}$  est solution de l'équation :

$$X^T X \hat{\beta} = X^T Y .$$

On écrit  $X \in \mathbb{R}^{n \times p}$  de la forme  $X = QR$ , où  $Q \in \mathbb{R}^{n \times p}$  est une matrice orthogonale (i.e.,  $QQ^T = I_n$ ) et  $R \in \mathbb{R}^{p \times p}$  triangulaire supérieure. Les matrices triangulaires supérieures sont très utiles pour résoudre les systèmes linéaires. En substituant dans l'équation précédente, on obtient :

$$\begin{aligned} R^T (Q^T Q) R \hat{\beta} = R^T Q^T Y &\Leftrightarrow R^T R \hat{\beta} = R^T Q^T Y \\ &\Leftarrow R \hat{\beta} = Q^T Y . \end{aligned}$$

Il ne reste alors plus qu'à résoudre un système linéaire avec une matrice triangulaire supérieure, ce qui est facile.

**Descente de gradient.** Pour éviter totalement le coût de l'inversion de matrice, une autre solution est la descente de gradient. Qui consiste à résoudre le problème de minimisation pas à pas en s'approchant du minimum grâce à des pas de gradient. On initialise par exemple  $\hat{\beta}_0 = 0$ , puis on met à jour :

$$\begin{aligned} \hat{\beta}_{i+1} &= \hat{\beta}_i - \eta \nabla \widehat{\mathcal{R}}_n(\hat{\beta}_i) \\ &= \hat{\beta}_i - \frac{2\eta}{n} ((X^T X) \hat{\beta}_i - Y^T X) , \end{aligned}$$

où  $\eta > 0$  est un paramètre d'apprentissage. On voit que si l'algorithme converge, alors il converge vers un point annulant le gradient, donc vers la solution de l'OLS. Pour avoir convergence, il faut que le paramètre  $\eta$  soit bien calibré, mais on verra cela plus en détails dans le cours sur la descente de gradient.

Si le jeu de données est beaucoup trop grand  $n \gg 1$ . Il peut également être prohibitif de charger toutes les données pour faire le calcul de  $\nabla \widehat{\mathcal{R}}_n(\hat{\theta}_i)$ . La solution courante est alors de faire une descente de gradient stochastique, où on ne fait des pas de gradients que sur des estimés de  $\nabla \widehat{\mathcal{R}}_n(\hat{\theta}_i)$ , calculés sur une sous-partie aléatoire des données.

## Références

- CORNILLON, P.-A. et E. MATZNER-LØBER (2011). "La régression linéaire simple". In : *Régression avec R*, p. 1-28.
- RIVOIRARD, V. et G. STOLTZ (2012). *Statistique mathématique en action*.