

Régression linéaire

Extension à la régression logistique pour la classification

Machine Learning

Cédric RICHARD

Université Nice Sophia Antipolis

RÉGRESSION LINÉAIRE

Exemple de modèle linéaire simple

Problème : On s'intéresse à la concentration d'ozone O_3 dans l'air.

On cherche à savoir s'il est possible d'expliquer la concentration maximale d'ozone de la journée par la température T à midi.

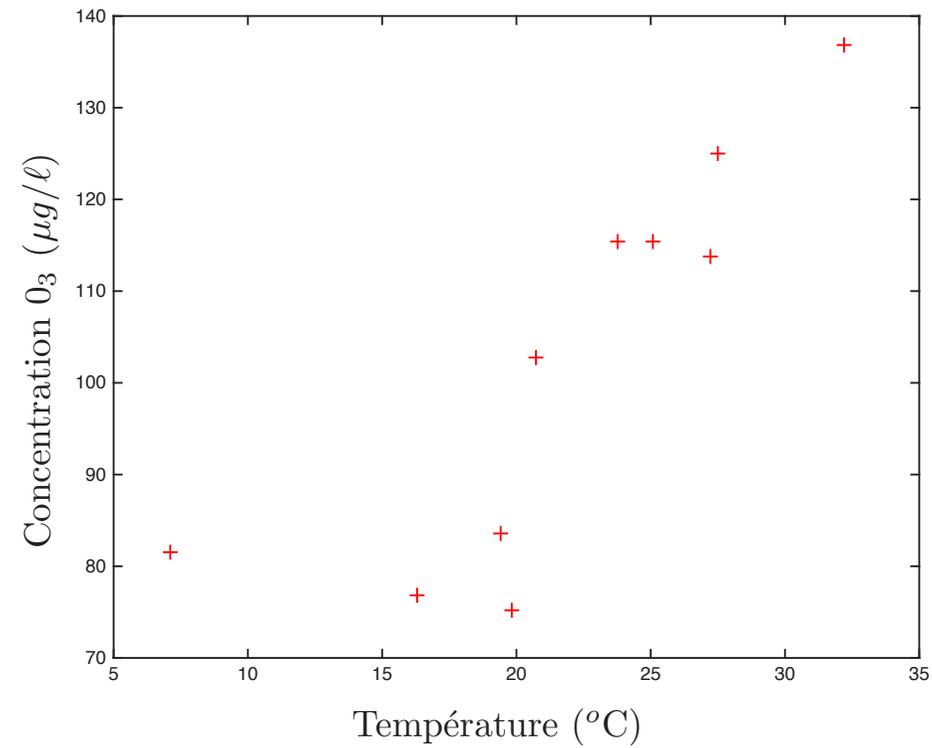
Objectif : Le but de la régression est :

- Ajuster les paramètres d'un modèle pour expliquer O_3 à partir de T
- Prédire O_3 à partir de nouvelles valeurs de T

T	23.8	16.3	27.2	7.1	25.1	27.5	19.4	19.8	32.2	20.7
O_3	115.4	76.8	113.8	81.6	115.4	125	83.6	75.2	136.8	102.8

RÉGRESSION LINÉAIRE

Exemple de modèle linéaire simple



RÉGRESSION LINÉAIRE

Exemple de modèle linéaire simple

Problème : Le botaniste Joseph D. Hooker a mesuré en 1849, en Himalaya, la pression atmosphérique p_i et la température d'ébullition de l'eau t_i .

Modèle : Selon les lois de la physique, $y_i = \ln p_i$ devrait être en première approximation proportionnel à t_i . On pose donc :

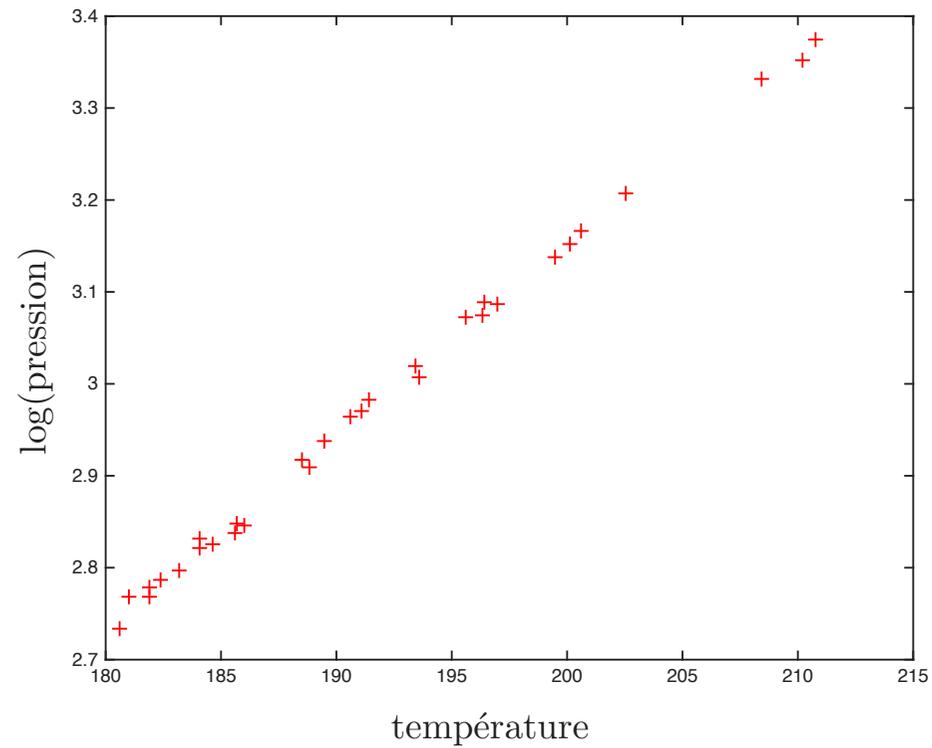
$$y_i = \beta_1 + \beta_2 t_i + \varepsilon_i$$

où ε_i est une erreur d'approximation de moyenne nulle et de variance σ^2 .

Objectif : Estimer $\beta = [\beta_1 \ \beta_2]^\top$ permettant d'expliquer y_i à partir de t_i .

RÉGRESSION LINÉAIRE

Exemple de modèle linéaire simple



RÉGRESSION LINÉAIRE

Exemple de modèle linéaire multiple

Fonction de production : La fonction de Cobb-Douglas (1928) est largement utilisée en économie comme modèle de fonction de production. Elle exprime le niveau de production p_i d'un bien en fonction du capital utilisé k_i et de la quantité de travail t_i :

$$p_i = \alpha_1 \cdot k_i^{\alpha_2} \cdot t_i^{\alpha_3}$$

où α_1 , α_2 et α_3 sont déterminés par la technologie.

Modèle : Le modèle de Cobb-Douglas peut être linéarisé ainsi :

$$\ln p_i = \ln \alpha_1 + \alpha_2 \ln k_i + \alpha_3 \ln t_i + \varepsilon_i$$

où ε_i est une erreur d'approximation de moyenne nulle et de variance σ^2 .

Objectif : Estimer $\beta = [\ln(\alpha_1) \alpha_2 \alpha_3]^\top$.

RÉGRESSION LINÉAIRE

Exemple de modèle linéaire multiple

Données : Cobb et Douglas (1928) disposent des données suivantes :

Année	p_i	k_i	t_i	Année	p_i	k_i	t_i	Année	p_i	k_i	t_i
1899	100	100	100	1907	151	176	138	1915	189	266	154
1900	101	107	105	1908	126	185	121	1916	225	298	182
1901	112	114	110	1909	155	198	140	1917	227	335	196
1902	122	122	118	1910	159	208	144	1918	223	366	200
1903	124	131	123	1911	153	216	145	1919	218	387	193
1904	122	138	116	1912	177	226	152	1920	231	407	193
1905	143	149	125	1913	184	236	154	1921	179	417	147
1906	152	163	133	1914	169	244	149	1922	240	431	161

Résultat : Dans le cas de rendements d'échelle constants ($\alpha_2 + \alpha_3 = 1$),
Cobb et Douglas trouvent :

$$\alpha_2 = \frac{1}{4} \quad \alpha_3 = \frac{3}{4}$$

RÉGRESSION LINÉAIRE

Notations

▷ Les données consistent en des variables observées y_i (réponses) et des variables explicatives (ou régresseurs) $\mathbf{x}_i = [x_{i1}, \dots, x_{ip}]^\top$, avec $i = 1, \dots, n$. Chaque paire (y_i, \mathbf{x}_i) représente une expérience (un individu).

▷ On présume l'existence d'une relation de la forme :

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i$$

où ε_i est une erreur d'approximation.

▷ On suppose que $\boldsymbol{\varepsilon} = [\varepsilon_1, \dots, \varepsilon_n]^\top$ est un n -échantillon d'espérance nulle. Souvent, il est supposé distribué selon une loi Gaussienne $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$.

RÉGRESSION LINÉAIRE

Notations

▷ On arrange les n individus de p variables explicatives ainsi :

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & & \\ & & \ddots & & \\ \vdots & & & x_{ij} & \vdots \\ & & & & \ddots \\ 1 & x_{n1} & \cdots & & x_{np} \end{pmatrix}$$

On convient de mettre le régresseur constant (associé au paramètre β_0), s'il y a lieu, dans la première colonne de \mathbf{X} .

▷ Pour simplifier les notations, et éviter de devoir distinguer les cas avec et sans composante continue β_0 , on supposera que $\mathbf{X} \in \mathbb{R}^{n \times p}$ et $\boldsymbol{\beta} \in \mathbb{R}^p$.

RÉGRESSION LINÉAIRE

Modèle linéaire

Définition 1. *Un modèle linéaire se définit par une équation de la forme :*

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

où :

- ▷ $\mathbf{y} \in \mathbb{R}^n$ est le vecteur des observations
- ▷ $\mathbf{X} \in \mathbb{R}^{n \times p}$ est la matrice des variables explicatives (régresseurs)
- ▷ $\boldsymbol{\beta} \in \mathbb{R}^p$ le vecteur des coefficients à estimer.

Hypothèses. On suppose que la matrice \mathbf{X} est de rang plein p , et que le vecteur de bruit $\boldsymbol{\varepsilon}$ est un n -échantillon centré de variance σ^2 .

RÉGRESSION LINÉAIRE

Estimateur des moindres carrés ordinaires

Objectif. Les points (y_i, \mathbf{x}_i) étant observés, il s'agit d'estimer la fonction affine définie par $\hat{f} : \mathbf{x} \mapsto \mathbf{x}^\top \hat{\boldsymbol{\beta}}$ minimisant l'erreur quadratique moyenne :

$$\hat{f} \triangleq \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

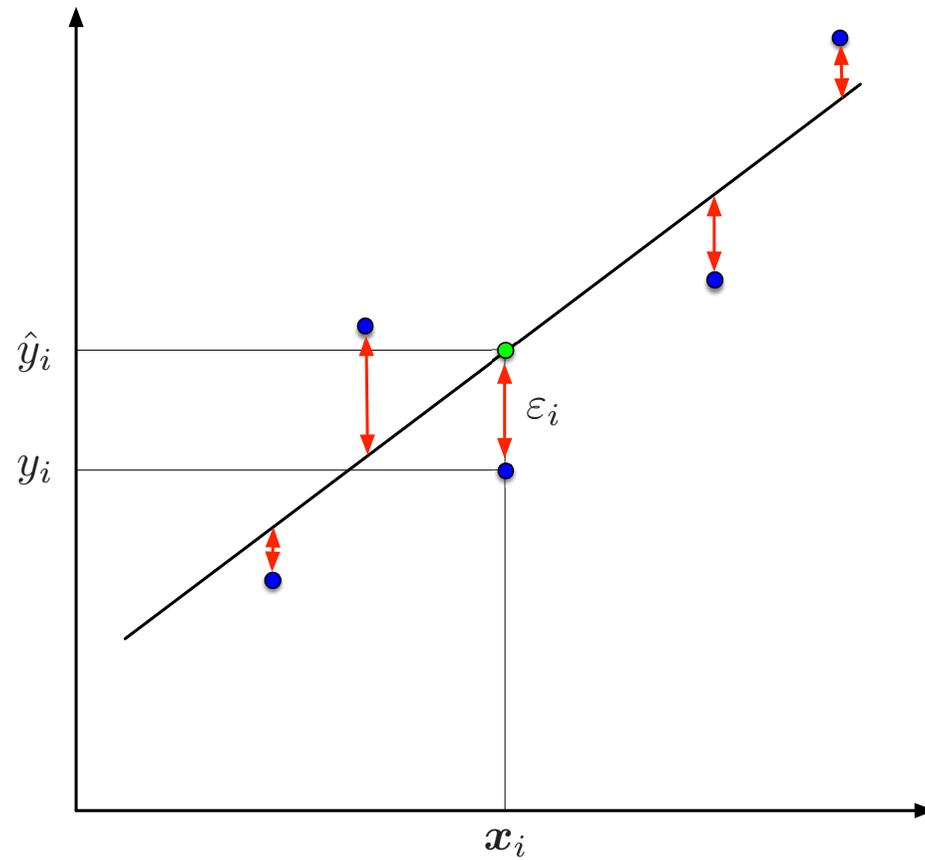
Définition 2 (Estimateur des moindres carrés ordinaire (OLS)). *On appelle estimateur des moindres carrés ordinaire, le minimiseur $\hat{\boldsymbol{\beta}}$ du risque empirique :*

$$\begin{aligned} \hat{f} &\triangleq \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \\ &\triangleq \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \end{aligned}$$

où $\|\cdot\|^2$ désigne la norme euclidienne telle que $\|\boldsymbol{\varepsilon}\|^2 \triangleq \sum_{i=1}^n \varepsilon_i^2$.

RÉGRESSION LINÉAIRE

Estimateur des moindres carrés ordinaires



RÉGRESSION LINÉAIRE

Estimateur des moindres carrés ordinaires

Proposition 1 (Estimateur des moindres carrés ordinaire (OLS)). *Sous les hypothèses et les notations de la Définition 1, l'OLS existe et est unique. Il s'écrit :*

$$\hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

RÉGRESSION LINÉAIRE

Estimateur des moindres carrés ordinaires

Démonstration :

On considère le problème à résoudre :

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} J(\boldsymbol{\beta}) \triangleq \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

Une condition nécessaire d'optimalité est : $\nabla J(\boldsymbol{\beta}) = 0$. Or :

$$\nabla J(\boldsymbol{\beta}) = -\frac{2}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

On note que $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{p \times p}$ est inversible puisque \mathbf{X} est de rang p par hypothèse. Afin que $\hat{\boldsymbol{\beta}}_{\text{OLS}}$ soit optimum au sens de $J(\boldsymbol{\beta})$, il doit donc nécessairement vérifier :

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Cette condition est suffisante car $J(\boldsymbol{\beta})$ est une fonction strictement convexe. En effet, sa matrice Hessienne $\nabla^2 J(\boldsymbol{\beta}) \triangleq \frac{2}{n} (\mathbf{X}^\top \mathbf{X})$ est définie positive (strictement).

RÉGRESSION LINÉAIRE

Interprétation géométrique

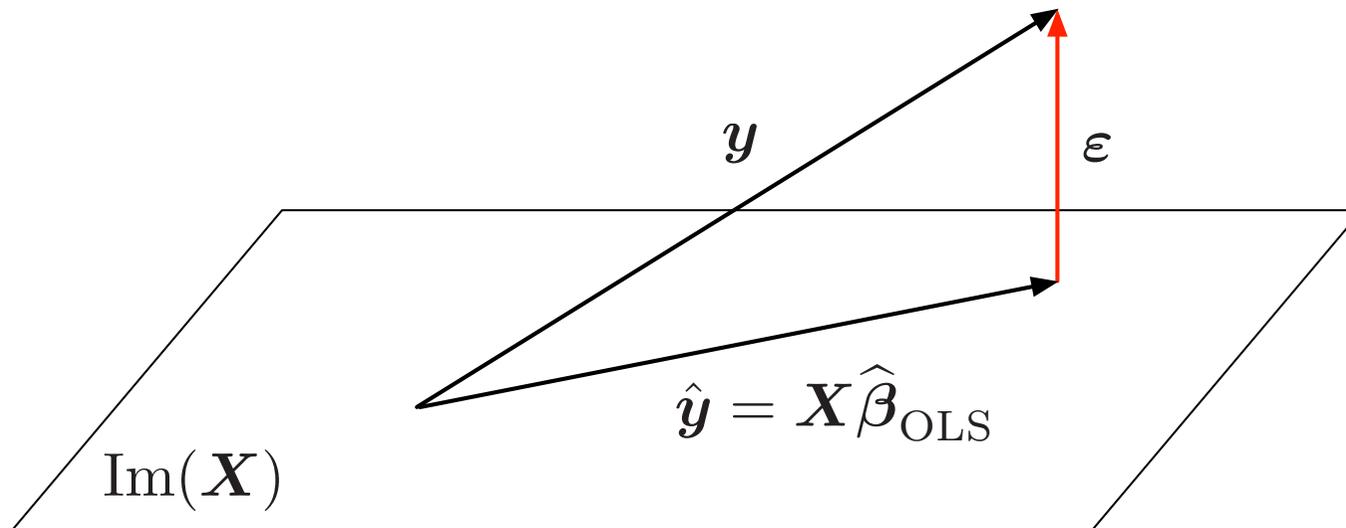
Le modèle linéaire vise à représenter \mathbf{y} par une combinaison linéaire $\mathbf{X}\boldsymbol{\beta}$ des colonnes de \mathbf{X} . Celles-ci constituent une famille libre de \mathbb{R}^n puisque $\text{rang}(\mathbf{X}) = p$.

En minimisant $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$, on recherche donc l'élément de $\text{Im}(\mathbf{X})$ le plus proche de \mathbf{y} au sens de la distance euclidienne. Il s'agit de la projection orthogonale de \mathbf{y} sur $\text{Im}(\mathbf{X})$, notée $\hat{\mathbf{y}} \triangleq \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}}$.

Notons que $\mathbf{P}_X \triangleq (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ est la matrice de projection orthogonale sur le sous-espace $\text{Im}(\mathbf{X})$.

RÉGRESSION LINÉAIRE

Interprétation géométrique



RÉGRESSION LINÉAIRE

Propriétés

Propriété 1. Soit $\varepsilon \triangleq \mathbf{y} - \hat{\mathbf{y}}$. On a :

$$\varepsilon \perp \hat{\mathbf{y}}$$

Démonstration :

$$\begin{aligned}\varepsilon^\top \hat{\mathbf{y}} &= (\mathbf{y} - \mathbf{P}_X \mathbf{y})^\top \mathbf{P}_X \mathbf{y} \\ &= \|\mathbf{P}_X \mathbf{y}\|^2 - \|\mathbf{P}_X \mathbf{y}\|^2 \\ &= 0\end{aligned}$$

RÉGRESSION LINÉAIRE

Mesures de performance du modèle linéaire

L'erreur quadratique moyenne $\frac{1}{n} \|\boldsymbol{\varepsilon}\|^2$ permet de caractériser la qualité d'un modèle linéaire, mais ce n'est pas une grandeur normalisée.

Il est préférable de recourir au coefficient de détermination, défini comme le rapport de la variance expliquée sur la variance totale :

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2} \quad \text{avec} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

- ▷ $R^2 = 1$: le modèle linéaire est parfait, les points (\boldsymbol{x}_i, y_i) sont alignées
- ▷ $R^2 \approx 0$: le modèle linéaire n'est pas approprié

RÉGRESSION LINÉAIRE

Mesures de performance du modèle linéaire

Propriété 2. *On a :*

$$R^2 = \rho_{y\hat{y}}$$

où ρ désigne le coefficient de corrélation

RÉGRESSION LINÉAIRE

Propriétés

Propriété 3. *L'estimateur $\hat{\beta}_{OLS}$ est sans biais, $\mathbb{E}\{\hat{\beta}_{OLS}\} = \beta$, et de covariance :*

$$\text{Cov}\{\hat{\beta}_{OLS}\} = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$$

Démonstration :

$$\begin{aligned}\mathbb{E}\{\hat{\beta}_{OLS}\} &= \mathbb{E}\left\{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}\right\} \\ &= \mathbb{E}\left\{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta) + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \varepsilon\right\} \\ &= \beta\end{aligned}$$

$$\begin{aligned}\text{Cov}\{\hat{\beta}_{OLS}\} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}\{\mathbf{y}\mathbf{y}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}\end{aligned}$$

RÉGRESSION LINÉAIRE

Propriétés

Pour caractériser l'écart entre $\hat{\beta}_{OLS}$ et β , on étudie l'écart quadratique suivant :

$$MSD(\hat{\beta}_{OLS}) \triangleq \mathbb{E}\{\|\hat{\beta}_{OLS} - \beta\|^2\}$$

Propriété 4.

$$MSD(\hat{\beta}_{OLS}) = \sigma^2 \text{trace}(\mathbf{X}^\top \mathbf{X})^{-1}$$

Démonstration :

$$\begin{aligned} MSD(\hat{\beta}_{OLS}) &\triangleq \mathbb{E}\{\|\hat{\beta}_{OLS} - \beta\|^2\} \\ &= \mathbb{E}\{(\hat{\beta}_{OLS} - \beta)^\top (\hat{\beta}_{OLS} - \beta)\} \\ &= \text{trace}\left\{\mathbb{E}\{(\hat{\beta}_{OLS} - \beta)(\hat{\beta}_{OLS} - \beta)^\top\}\right\} \\ &= \text{trace}\left\{\text{Cov}\{\hat{\beta}_{OLS}\}\right\} \end{aligned}$$

RÉGRESSION LINÉAIRE

Propriétés

La Propriété 4 est peu pratique pour caractériser la précision de $\hat{\boldsymbol{\beta}}$ car elle fait intervenir la variance du bruit σ^2 , inconnue dans un contexte applicatif.

Un estimateur naturel est de considérer l'erreur moyenne :

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2$$

Propriété 5. *Dans le cas Gaussien, les estimateurs par maximum de vraisemblance de $\boldsymbol{\beta}$ et σ vérifient :*

$$\hat{\boldsymbol{\beta}}_{MV} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{y} \quad \text{et} \quad \hat{\sigma}_{MV}^2 = \frac{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2}{n}$$

RÉGRESSION LINÉAIRE

Identifiabilité

Si \mathbf{X} n'est pas de rang maximum p , la matrice $\mathbf{X}^\top \mathbf{X}$ n'est pas inversible et le problème admet plusieurs solutions qui minimisent le risque empirique.

\Rightarrow le problème est mal posé ou non-identifiable.

Plusieurs solutions permettent de traiter ce problème :

▷ **Par sélection de variables :**

Cette stratégie consiste à réduire la dimension de l'espace des solutions par sélection des colonnes de \mathbf{X} de sorte que celle-ci soit de rang maximum.

▷ **Par régularisation du problème :**

Cette stratégie consiste à régulariser le problème de minimisation du risque empirique de sorte que celui-ci admette une solution unique.

RÉGRESSION RIDGE

Problème et solution

La méthode de régression régularisée la plus couramment utilisée est certainement la régression Ridge. Elle vise à contraindre l'amplitude des composantes de $\boldsymbol{\beta}$:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} J(\boldsymbol{\beta}) \triangleq \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \text{ sous la contrainte } \|\boldsymbol{\beta}\|^2 \leq \tau, \quad \tau > 0$$

Ce problème peut être formulé de manière équivalente (Lagrangien) ainsi :

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} J(\boldsymbol{\beta}) \triangleq \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2, \quad \lambda > 0$$

où λ est lié à τ ...

RÉGRESSION RIDGE

Problème et solution

Le problème de régression Ridge s'écrit :

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} J(\boldsymbol{\beta}) \triangleq \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2, \quad \lambda > 0$$

En suivant le même cheminement que précédemment, on trouve :

Proposition 2. *La solution du problème de régression Ridge existe et est unique pour tout $\lambda > 0$. Elle s'écrit :*

$$\hat{\boldsymbol{\beta}}_{ridge} = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Remarque :

- ▷ Il n'y a plus de problème d'inversion relativement au rang de $\mathbf{X}^\top \mathbf{X}$ puisque la matrice $\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}$ est de rang d .
- ▷ Le choix du paramètre λ est essentiel en pratique (voir ci-après).

RÉGRESSION RIDGE

Propriétés

Propriété 6. L'estimateur $\hat{\beta}_{ridge}$ vérifie :

$$\mathbb{E}\{\hat{\beta}_{ridge}\} = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) \beta$$

$$Cov\{\hat{\beta}_{ridge}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1}$$

Remarque :

- ▷ L'estimateur Ridge est biaisé contrairement à l'estimateur OLS, ce qui constitue un inconvénient.
- ▷ La covariance de l'estimateur Ridge ne fait pas intervenir l'inverse de $\mathbf{X}^\top \mathbf{X}$, mais celui de $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$ qui est mieux conditionné.
- ▷ On note que $\lim_{\lambda \rightarrow \infty} Cov\{\hat{\beta}_{ridge}\} = \mathbf{0}$

RÉGRESSION RIDGE

Propriétés

Démonstration de la Propriété 5 :

$$\begin{aligned}\hat{\beta}_{\text{ridge}} &= (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) \hat{\beta}_{\text{OLS}}\end{aligned}$$

$$\begin{aligned}\text{Cov}\{\hat{\beta}_{\text{ridge}}\} &= \text{Cov}\left\{(\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) \hat{\beta}_{\text{OLS}}\right\} \\ &= (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) \text{Cov}\{\hat{\beta}_{\text{OLS}}\} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1}\end{aligned}$$

RÉGRESSION RIDGE

Propriétés

Propriété 7. *Les estimateurs $\hat{\beta}_{ridge}$ et $\hat{\beta}_{OLS}$ vérifient :*

$$\begin{aligned} & Cov\{\hat{\beta}_{OLS}\} - Cov\{\hat{\beta}_{ridge}\} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \left[2n\lambda \mathbf{I} + n^2 \lambda^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right] \left[(\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \right]^\top \end{aligned}$$

La matrice ci-dessus est définie positive, quel que soit $\lambda > 0$. Ceci signifie :

$$Cov\{\hat{\beta}_{OLS}\} \succeq Cov\{\hat{\beta}_{ridge}\}$$

l'inégalité étant stricte pour $\lambda > 0$.

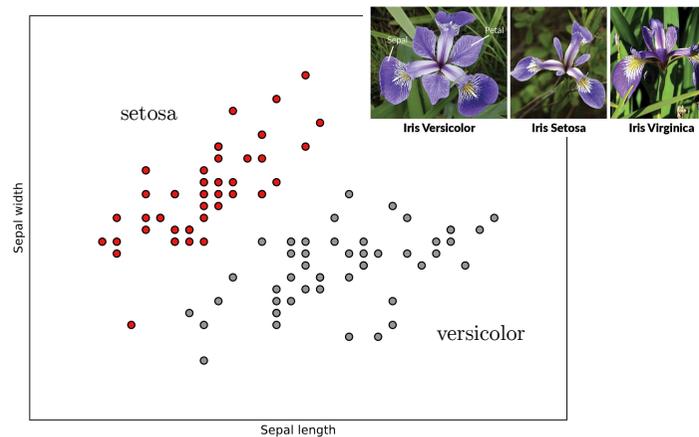
RÉGRESSION LOGISTIQUE

Contexte

Problème : On s'intéresse à la classification automatique d'iris.

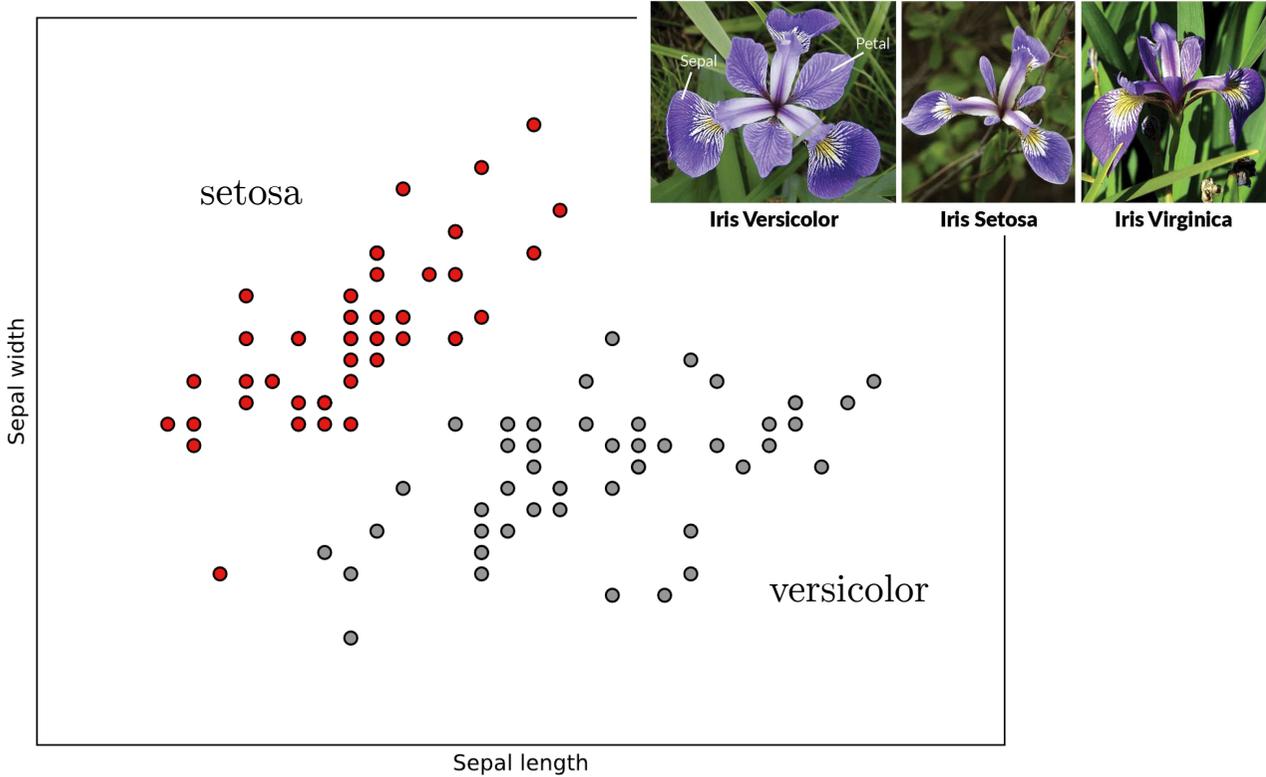
On dispose d'une base d'individus *Setosa* et *Versicolor*, caractérisés par les dimensions de leurs sépales et pétales.

Méthode : Le but de la régression logistique est d'effectuer une régression des individus x_i sur les étiquettes correspondantes $y_i \in \{-1, +1\}$.



RÉGRESSION LOGISTIQUE

Contexte



RÉGRESSION LOGISTIQUE

Spécification du modèle

Soit $y \in \{-1, +1\}$ la var. à prédire, et $\mathbf{x} \in \mathbb{R}^p$ le vecteur des var. explicatives.^a

Notations :

- ▷ $P(y = \pm 1)$: probabilités a priori des classes, notées $P(\pm 1)$
- ▷ $P(\mathbf{x} | y = \pm 1)$: distrib. conditionnelles des observations, notées $P(\mathbf{x} | \pm 1)$

Principe :

La régression logistique repose sur l'hypothèse fondamentale que :

$$\ln \frac{P(\mathbf{x} | + 1)}{P(\mathbf{x} | - 1)} = a_0 + a_1 x_1 + \cdots + a_{p-1} x_{p-1}$$

Contrairement à l'AFD, on ne s'intéresse pas ici aux distributions conditionnelles mais à leur rapport.

a. Les variables y et \mathbf{x} du problème sont aléatoires. Contrairement à l'usage, elles ne seront pas représentées par des lettres capitales (Y, \mathbf{X}) pour éviter toute confusion avec ce qui précède.

RÉGRESSION LOGISTIQUE

Spécification du modèle

La spécification précédente peut être réécrite de manière équivalente :

$$\begin{aligned}\ln \frac{P(+1|\mathbf{x})}{1 - P(+1|\mathbf{x})} &= \beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1} \\ &= \boldsymbol{\beta}^\top \mathbf{x}\end{aligned}$$

On désigne par *logit* de $P(+1|\mathbf{x})$ l'expression ci-dessus.

Il s'agit d'une régression logistique car le modèle ci-dessus provient d'une loi logistique. En effet, par manipulation du modèle ci-dessus, on obtient :

$$P(+1|\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\beta}^\top \mathbf{x}}}$$

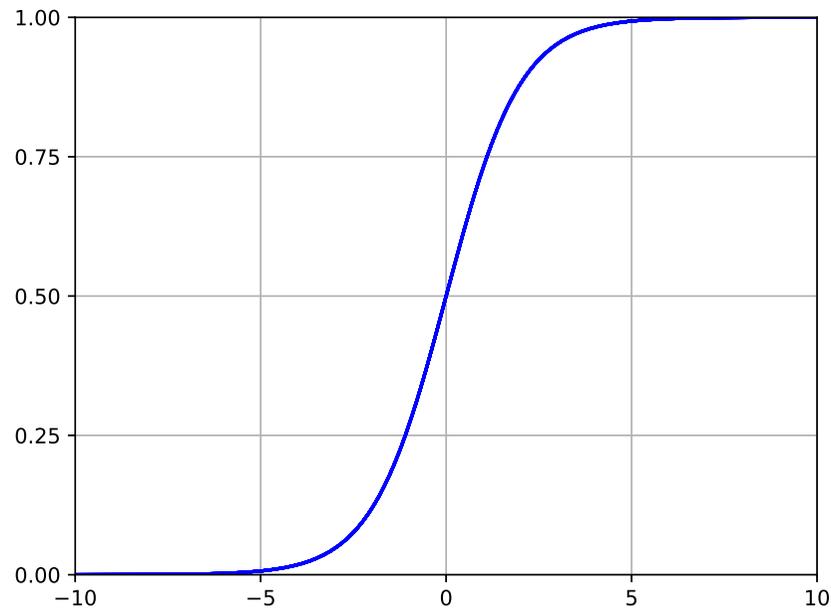
ainsi que :

$$P(-1|\mathbf{x}) = 1 - P(+1|\mathbf{x}) = \frac{1}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}}}$$

RÉGRESSION LOGISTIQUE

Fonction logistique

La loi logistique est définie par sa fonction caractéristique de même nom, appelée également sigmoïde.



RÉGRESSION LOGISTIQUE

Règle de décision

Etant donné β estimé comme ci-après, la règle de décision mise en œuvre est :

- ▷ choisir (+1) si $P(+1|\mathbf{x}) > P(-1|\mathbf{x})$
- ▷ choisir (-1) sinon

RÉGRESSION LOGISTIQUE

Estimation des paramètres

Loi de y :

La probabilité d'appartenance d'un individu \mathbf{x} à une classe $y = \pm 1$ est régit par une loi de Bernoulli :

$$P(y | \mathbf{x}) = P(+1 | \mathbf{x})^{\frac{1+y}{2}} \times P(-1 | \mathbf{x})^{\frac{1-y}{2}}, \quad y \in \{-1, +1\}$$

Estimation de β par maximum de vraisemblance :

Les variables $\{y_i\}_{i=1}^n$ sont supposées i.i.d. La vraisemblance de β est donnée par :

$$L(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n; \beta) = \prod_{i=1}^n P(y_i | \mathbf{x}_i)$$

qu'il faut maximiser par rapport à β .

RÉGRESSION LOGISTIQUE

Estimation des paramètres

Estimation de β par maximum de vraisemblance :

Afin de simplifier les calculs, on considère l'opposé de la log-vraisemblance, qu'il faut minimiser par rapport à β :

$$\begin{aligned} -\ln L(y_1, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n; \beta) &= -\sum_{i=1}^n \ln P(y_i \mid \mathbf{x}_i) \\ &= \sum_{i=1}^n \ln [1 + e^{-y_i \beta^\top \mathbf{x}_i}] = J(\beta) \end{aligned}$$

La fonction coût $J(\beta)$ est différentiable sur \mathbb{R}^d et strictement convexe. Elle admet un minimum global $\hat{\beta}$ satisfaisant la condition :

$$\nabla J(\hat{\beta}) = \mathbf{0}$$

RÉGRESSION LOGISTIQUE

Estimation des paramètres

$$\min_{\boldsymbol{\beta}} J(\boldsymbol{\beta}) \triangleq \sum_{i=1}^n \ln [1 + e^{-y_i \boldsymbol{\beta}^\top \mathbf{x}_i}]$$

Calcul du gradient :

La dérivée partielle de $J(\boldsymbol{\beta})$ par rapport à β_j est donnée par :

$$\frac{\partial J(\boldsymbol{\beta})}{\partial \beta_j} = - \sum_{i=1}^n \frac{y_i [\mathbf{x}_i]_j p_i(\boldsymbol{\beta})}{1 + p_i(\boldsymbol{\beta})} \quad \text{avec} \quad p_i(\boldsymbol{\beta}) = e^{-y_i \boldsymbol{\beta}^\top \mathbf{x}_i}$$

Ceci nous permet de réécrire le gradient de $J(\boldsymbol{\beta})$ sous la forme :

$$\nabla J(\boldsymbol{\beta}) = -\mathbf{X}^\top \mathbf{D}(\boldsymbol{\beta}) \mathbf{y}$$

où $\mathbf{D}(\boldsymbol{\beta})$ est la matrice diagonale de termes diagonaux $[\mathbf{D}(\boldsymbol{\beta})]_{ii} = \frac{p_i(\boldsymbol{\beta})}{1+p_i(\boldsymbol{\beta})}$

RÉGRESSION LOGISTIQUE

Estimation numérique des paramètres

La condition d'optimalité $\nabla J(\hat{\beta}) = \mathbf{0}$ n'admet pas de solution analytique. Il est nécessaire de recourir à une méthode numérique

Algorithme du gradient :

Choisir un point initial β_0 , un seuil ϵ , et un pas μ

Itérer les étapes suivantes à partir de $k = 0$

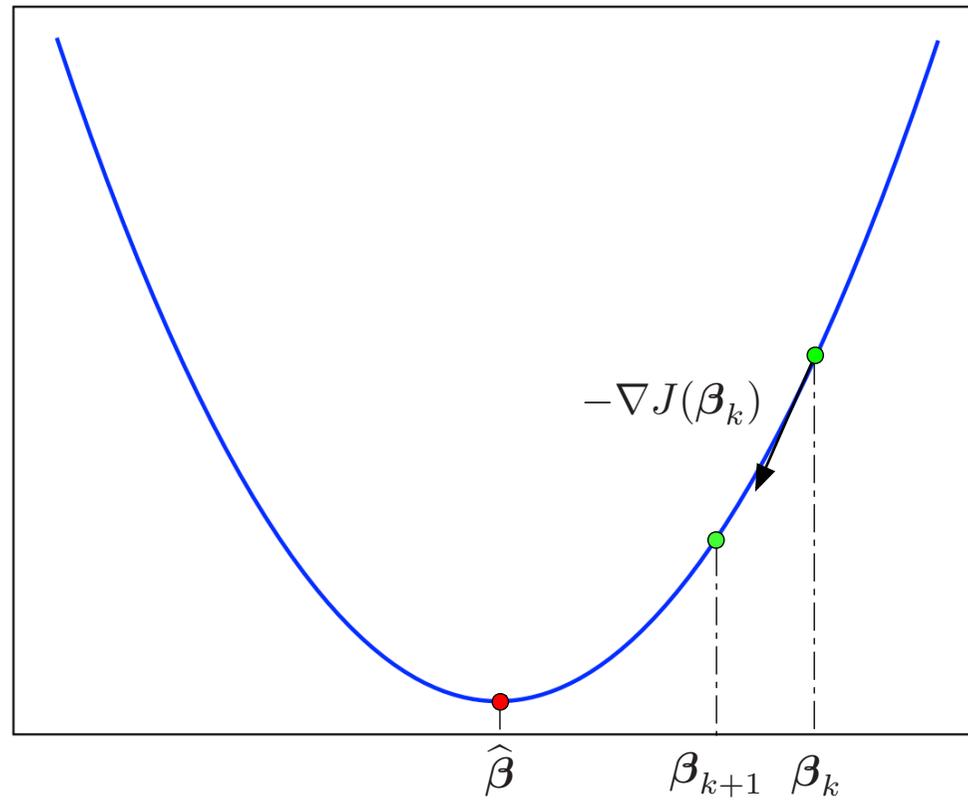
1. Calculer $\nabla J(\beta_k)$
2. Test d'arrêt : si $\|\nabla J(\beta_k)\| < \epsilon$, arrêt
3. Nouvel itéré : $\beta_{k+1} = \beta_k - \mu \nabla J(\beta_k)$

Régression logistique :

$$\beta_{k+1} = \beta_k + \mu \mathbf{X}^\top \mathbf{D}(\beta_k) \mathbf{y}$$

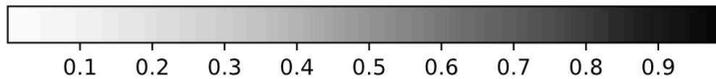
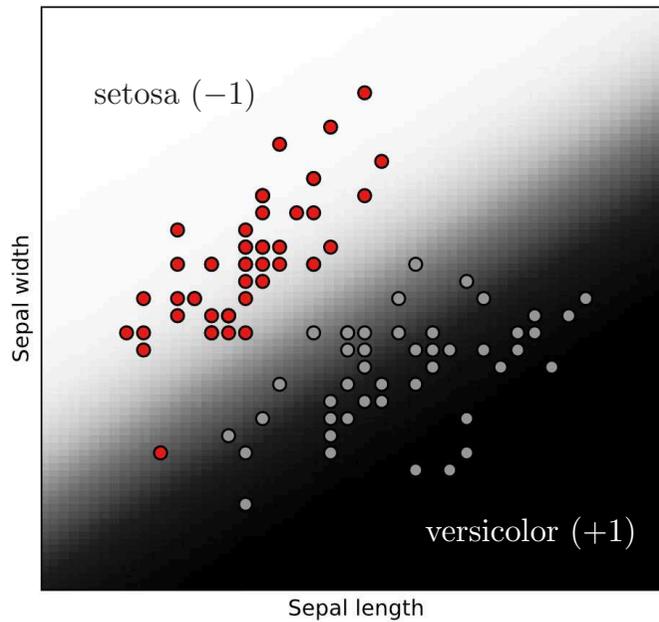
RÉGRESSION LOGISTIQUE

Estimation numérique des paramètres



RÉGRESSION LOGISTIQUE

Données iris



$P(+1|x)$