

# INTRODUCTION AU MACHINE LEARNING

## MÉTHODES LINÉAIRE

**Théo Lopès-Quintas**

BPCE Payment Services,  
Université Paris Dauphine

2022 - 2025

Régression Linéaire

1    Comment prédire un nombre? . . . . . 1

2    Comment mesurer la performance d’une régression? . . . . . 6

3    Régressions pénalisées . . . . . 9

Régression Logistique

1    Comment prédire une classe? . . . . . 11

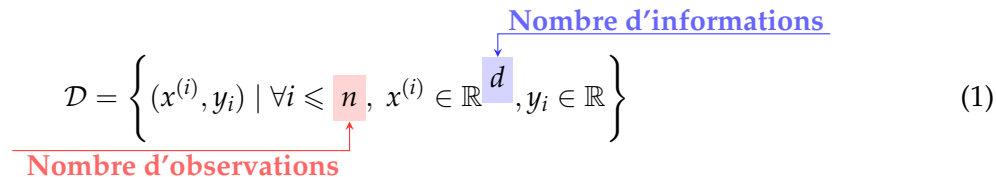
2    Comment mesurer la performance d’une classification? . . . . . 16

# COMMENT PRÉDIRE UN NOMBRE ?

## FORMALISATION D'UN PROBLÈME DE RÉGRESSION

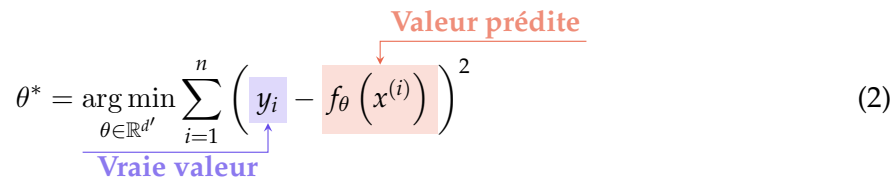
On considère le problème de régression associé au dataset décrit par (1) : on cherche à prédire un nombre.

$$\mathcal{D} = \left\{ (x^{(i)}, y_i) \mid \forall i \leq n, x^{(i)} \in \mathbb{R}^d, y_i \in \mathbb{R} \right\} \quad (1)$$



On suppose qu'il existe un lien entre la cible et les informations contenue dans  $\mathcal{D}$ . Ce lien est une fonction notée  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  de paramètre  $\theta \in \mathbb{R}^{d'}$ . On cherche donc à résoudre le problème d'optimisation :

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^{d'}} \sum_{i=1}^n \left( y_i - f_{\theta}(x^{(i)}) \right)^2 \quad (2)$$



# COMMENT PRÉDIRE UN NOMBRE ?

## FORMALISATION D'UN PROBLÈME DE RÉGRESSION LINÉAIRE

Dans le cadre d'une régression **linéaire** on suppose que le lien entre la cible et les informations disponibles est **linéaire** :

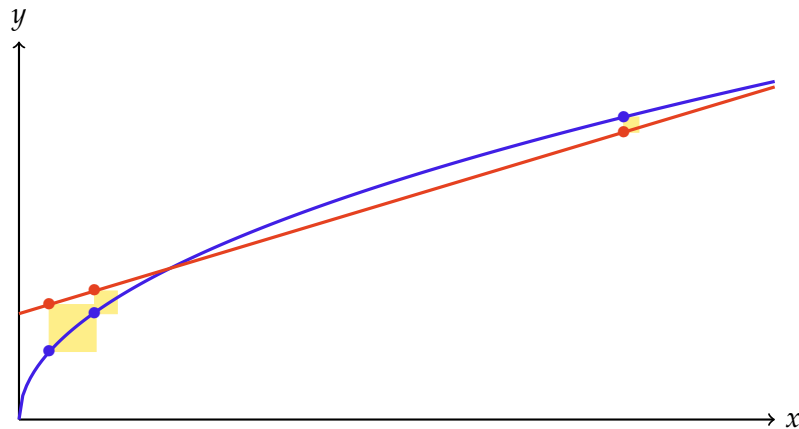
$$\hat{y} = \theta_0 + \sum_{j=1}^d \theta_j x_j$$

On peut réécrire notre problème (2) comme :

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^{d+1}} \sum_{i=1}^n \left[ y_i - \left( \theta_0 + \sum_{j=1}^d \theta_j x_j^{(i)} \right) \right]^2$$

# COMMENT PRÉDIRE UN NOMBRE ?

## FORMALISATION D'UN PROBLÈME DE RÉGRESSION LINÉAIRE : VISUALISATION



**Figure** – Visualisation de l'erreur entre la régression linéaire et la vraie courbe

# COMMENT PRÉDIRE UN NOMBRE ?

## TROUVER LES PARAMÈTRES OPTIMAUX : CAS PARTICULIER

### Exercice 1 (Régression linéaire avec une seule information)

On suppose que l'on dispose d'un dataset  $\mathcal{D} = \{(x^{(i)}, y_i) \mid \forall i \leq n : x^{(i)} \in \mathbb{R}, y_i \in \mathbb{R}\}$ . On a donc une seule information pour prédire la valeur  $y$ .

1. Écrire le problème (2) dans le cadre de l'exercice.

2. Donner le meilleur vecteur de paramètre  $\theta$ .

On note  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$ . On rappelle avec cette convention que pour  $u, v \in \mathbb{R}^n$  :

$$\begin{aligned}\text{Cov}(u, v) &= \overline{uv} - \bar{u} \times \bar{v} \\ \mathbb{V}[u] &= \overline{u^2} - \bar{u}^2\end{aligned}$$

3. Montrer que  $\theta_0^*$  et  $\theta_1^*$  les deux paramètres optimaux peuvent s'écrire :

$$\begin{aligned}\theta_0^* &= \bar{y} + \theta_1^* \times \bar{x} \\ \theta_1^* &= \frac{\text{Cov}(x, y)}{\mathbb{V}[x]}\end{aligned}$$

## COMMENT PRÉDIRE UN NOMBRE ?

TROUVER LES PARAMÈTRES OPTIMAUX : EN GÉNÉRAL

Puisqu'on suppose un lien linéaire, on peut exploiter l'écriture matricielle :

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \cdots & x_d^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & x_3^{(n)} & \cdots & x_d^{(n)} \end{pmatrix} \times \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$\Longleftrightarrow$

$$Y = X\theta + \varepsilon, \text{ avec } \varepsilon \text{ un vecteur de bruit.}$$

On peut réécrire notre problème (2) comme  $\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \|Y - X\theta\|^2$

### Proposition 1

Si la matrice  $X$  est de rang plein, alors  $\theta^* = ({}^tXX)^{-1} {}^tXY$

# COMMENT MESURER LA PERFORMANCE D'UNE RÉGRESSION ?

## ERREUR QUADRATIQUE MOYENNE


Une première manière de mesurer la performance est de considérer l'erreur quadratique moyenne (MSE) :

$$\text{MSE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Cette métrique peut être décomposée en deux quantités :

- ▶ Bias  $[\hat{f}(x)] = \mathbb{E} [\hat{f}(x)] - f(x)$  : l'écart moyen entre la valeur prédite et la vraie valeur
- ▶  $\mathbb{V} [\hat{f}(x)] = \mathbb{E} \left[ \left( \mathbb{E} [\hat{f}(x)] - \hat{f}(x) \right)^2 \right]$  : la dispersion moyenne des valeurs prédites autour de la moyenne

$$\text{MSE}(y, \hat{f}(x)) = \left( \text{Bias} [\hat{f}(x)] \right)^2 + \mathbb{V} [\hat{f}(x)] + \sigma^2 \quad (3)$$

 Erreur incompressible



## COMMENT MESURER LA PERFORMANCE D'UNE RÉGRESSION ?

### RMSE

Une deuxième manière de mesurer la performance est de calculer la racine carrée de la MSE :

$$\text{RMSE}(y, \hat{y}) = \sqrt{\sum_{i=1}^n \frac{1}{n} (y_i - \hat{y}_i)^2}$$

### Exercice 2 (Ordre de grandeur)

*Montrer que :*

$$\text{RMSE}(y, \bar{y}) = \sqrt{\overline{y^2} - \bar{y}^2}$$

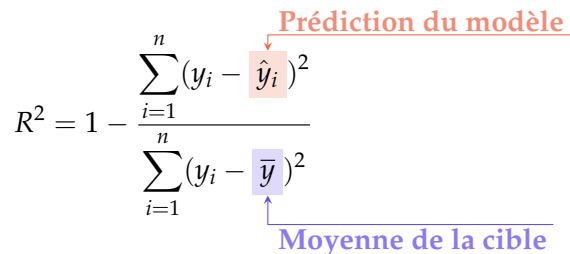
*En déduire une interprétation de la RMSE et un critère de performance d'une régression.*

# COMMENT MESURER LA PERFORMANCE D'UNE RÉGRESSION ?

## COEFFICIENT DE DÉTERMINATION $R^2$

Une troisième manière de mesurer la performance est d'étudier le coefficient de détermination  $R^2$  :

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



### Exercice 3 (Interprétation de $R^2$ )

On suppose que l'on dispose des vecteurs  $y$  et  $\hat{y}$ .

1. Comment interpréter la valeur 1 pour le  $R^2$  ? Et la valeur 0 ?
2. Le  $R^2$  peut-il être négatif ?

# RÉGRESSIONS PÉNALISÉES

## RÉGRESSION RIDGE

On dit qu'on **pénalise** un modèle quand on modifie le problème d'optimisation :

Apprentissage

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \left( y_i - f_{\theta} \left( x^{(i)} \right) \right)^2 + \mathcal{P}_{\lambda}(\theta) \quad \text{(Pénalisation)}$$

Pénalisation

On définit la régression Ridge comme une régression pénalisée :

$$\theta_{\text{Ridge}}^* = \arg \min_{\theta \in \mathbb{R}^d} \|Y - X\theta\|^2 + \lambda \|\theta\|^2 \quad \text{(Ridge)}$$

Dans ce cas, on obtient l'estimateur optimal à l'aide de la formule :

$$\theta_{\text{Ridge}}^* = ({}^tXX + n\lambda\mathbb{I}_d)^{-1} {}^tXY$$

# RÉGRESSIONS PÉNALISÉES

## RÉGRESSION LASSO

On dit qu'on **pénalise** un modèle quand on modifie le problème d'optimisation :

Apprentissage

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \left( y_i - f_{\theta} \left( x^{(i)} \right) \right)^2 + \mathcal{P}_{\lambda}(\theta) \quad (\text{Pénalisation})$$

**Pénalisation**

On définit la régression LASSO comme une régression pénalisée :

$$\theta_{\text{LASSO}}^* = \arg \min_{\theta \in \mathbb{R}^d} \|Y - X\theta\|^2 + \lambda \|\theta\|_1 \quad (\text{LASSO})$$

### Exercice 4 (Biais/Variance pour Ridge et LASSO)

*Pour la régression Ridge, puis la régression LASSO, comment évolue le biais quand  $\lambda$  augmente ? Même question pour la variance.*

# COMMENT PRÉDIRE UNE CLASSE ?

## FORMALISATION D'UN PROBLÈME DE CLASSIFICATION

On considère le problème de classification associé au dataset décrit par (4) : on cherche à prédire un nombre.

$$\mathcal{D} = \left\{ (x^{(i)}, y_i) \mid \forall i \leq n, x^{(i)} \in \mathbb{R}^d, y_i \in \mathcal{Y} \subset \mathbb{N} \right\} \quad (4)$$

Nombre d'observations (pointing to  $n$ )

Nombre d'informations (pointing to  $d$ )

On suppose qu'il existe un lien entre la cible et les informations contenue dans  $\mathcal{D}$ . Ce lien est une fonction notée  $f : \mathbb{R}^d \rightarrow \mathcal{Y}$  de paramètre  $\theta \in \mathbb{R}^{d'}$ . On cherche à résoudre le problème d'optimisation :

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^{d'}} - \left[ y \ln \{f_{\theta}(x)\} + (1 - y) \ln \{1 - f_{\theta}(x)\} \right]$$

Observation positive (pointing to  $y \ln \{f_{\theta}(x)\}$ )

Observation négative (pointing to  $(1 - y) \ln \{1 - f_{\theta}(x)\}$ )

# COMMENT PRÉDIRE UNE CLASSE ?

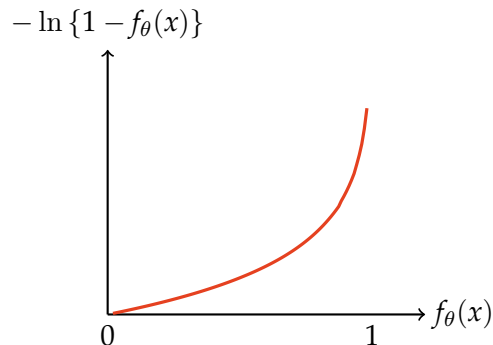
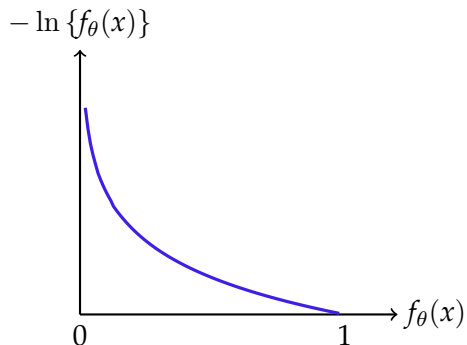
## FORMALISATION D'UN PROBLÈME DE CLASSIFICATION : VISUALISATION

Observation positive

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^{d'}} - \left[ y \ln \{f_{\theta}(x)\} + (1 - y) \ln \{1 - f_{\theta}(x)\} \right]$$

Observation négative

Pour le problème d'optimisation considéré, on peut visualiser les variations de chaque partie comme :



# COMMENT PRÉDIRE UNE CLASSE ?

## DESCENTE DE GRADIENT

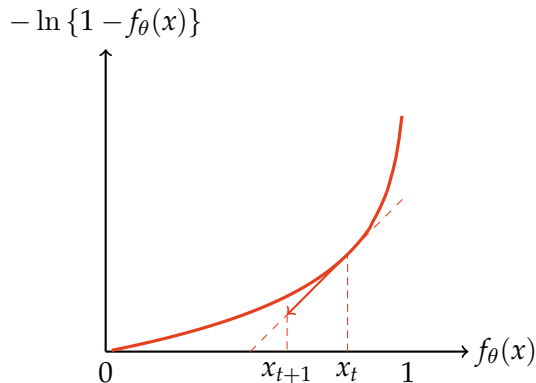
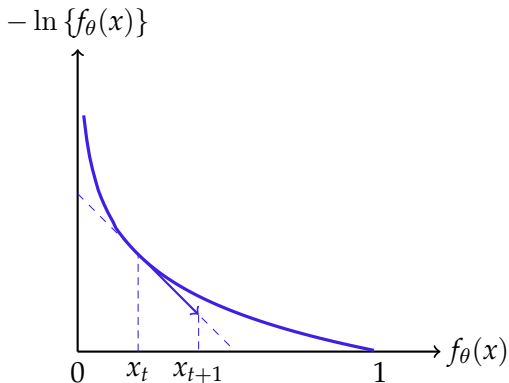
Soit  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  une fonction différentiable. Pour le problème d'optimisation :

$$x^* = \arg \min_{x \in \mathbb{R}^d} f(x)$$

La descente de gradient correspond au schéma :

$$x_{t+1} = x_t - \eta \nabla f(x_t) \quad \text{avec } \eta > 0$$

Visuellement pour notre problème cela se traduit par :



## COMMENT PRÉDIRE UNE CLASSE ?

### MODÉLISATION D'UN PROBLÈME DE CLASSIFICATION PAR LA RÉGRESSION LOGISTIQUE

La **régression logistique** suppose un lien *linéaire* entre les informations et la côte que l'observation soit de la classe d'intérêt. On modélise cela par la fonction  $f$  :

$$f_{\theta}(x) = \frac{1}{1 + e^{-(x_1\theta_1 + \dots + x_d\theta_d)}} = \frac{1}{1 + e^{-\langle x, \theta \rangle}} \quad (\text{Régression logistique})$$

$\theta = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$

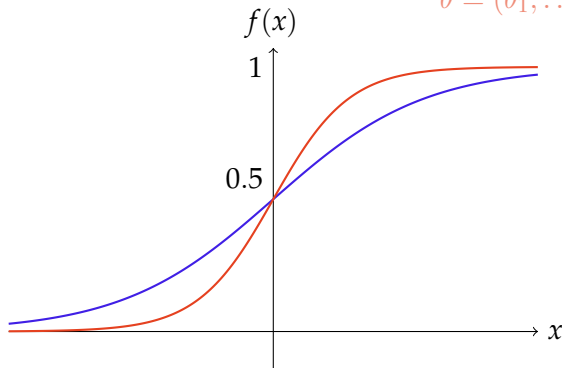


Figure –  $f(x) = \frac{1}{1 + e^{-x}}$  et  $f(x) = \frac{1}{1 + e^{-2x}}$



## COMMENT PRÉDIRE UNE CLASSE ?

### MODÉLISATION D'UN PROBLÈME DE CLASSIFICATION PAR LA RÉGRESSION LOGISTIQUE

#### Exercice 5 (Équation de la descente de gradient)

On rappelle que  $f_\theta(x) = \frac{1}{1 + e^{-\langle \theta, x \rangle}}$  et on note  $\mathcal{L}(\theta; x, y) = -[y \ln \{f_\theta(x)\} + (1 - y) \ln \{1 - f_\theta(x)\}]$ .

Montrer que :

1.  $f_\theta(x) = \frac{e^{\langle \theta, x \rangle}}{1 + e^{\langle \theta, x \rangle}}$

2.  $f_\theta(-x) = 1 - f_\theta(x)$

3.  $\frac{\partial \ln}{\partial \theta_j} (f_\theta(x)) = x_j (1 - f_\theta(x))$

4.  $\frac{\partial \ln}{\partial \theta_j} (1 - f_\theta(x)) = -x_j f_\theta(x)$

5.  $\frac{\partial \mathcal{L}}{\partial \theta_j} (\theta; x^{(i)}, y_i) = x_j^{(i)} (f_\theta(x^{(i)}) - y_i)$

6. Conclure que la descente de gradient pour notre problème est :

$$\theta_j^{t+1} = \theta_j^t - \eta \sum_{i=1}^n x_j^{(i)} (f_\theta(x^{(i)}) - y_i)$$

## COMMENT MESURER LA PERFORMANCE D'UNE CLASSIFICATION ?

### ACCURACY

À partir de la régression logistique, nous pouvons obtenir la **matrice de confusion** :

		Prédit	
		Classe 0 (baisse)	Classe 1 (hausse)
Réal	Classe 0	TN	FP
	Classe 1	FN	TP

On définit l'**accuracy** comme la proportion de bonne prédiction :

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

### Exercice 6

*On souhaite prédire une hausse exceptionnelle, et dans le dataset que l'on a à disposition, il y a 1% de classe 1 (hausse exceptionnelle). Construire un algorithme qui permet d'atteindre 99% d'accuracy.*

# COMMENT MESURER LA PERFORMANCE D'UNE CLASSIFICATION ?

## PRÉCISION, RECALL ET F1-SCORE

Réal	Prédit	
	Classe 0 (baisse)	Classe 1 (hausse)
	Classe 0	Classe 1
	TN	FP
	FN	TP

À partir de la matrice de confusion, on peut extraire d'autres métriques qui apportent des éclairages plus précis :

$$\text{Précision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{F1-score} = \frac{2}{\frac{1}{\text{Précision}} + \frac{1}{\text{Recall}}}$$

# COMMENT MESURER LA PERFORMANCE D'UNE CLASSIFICATION ?

## MÉTRIQUES

### Exercice 7

*Vous demandez à votre data scientist de concevoir un algorithme qui priorise les mails en essayant de prédire ceux qui sont les plus importants. Un dataset d'entraînement et un dataset de test lui sont fournis. Dans le dataset de test, il y a 1000 mails dont 200 sont importants. Il vous présente un premier modèle qui, pour un certain seuil (A), présente la matrice de confusion suivante :*

		Prédit	
		Classe 0	Classe 1
	Réel		
	Classe 0	700	100
	Classe 1	50	150

*Pour un autre seuil (B) il présente cette matrice de confusion :*

		Prédit	
		Classe 0	Classe 1
	Réel		
	Classe 0	760	40
	Classe 1	80	120

1. Calculer l'accuracy, la précision, le recall et le F1-score de chacun des seuils.
2. Conclure sur le seuil que vous souhaitez conserver.