

INTRODUCTION AU MACHINE LEARNING

RÉGRESSION LOGISTIQUE

Théo Lopès-Quintas

BPCE Payment Services,
Université Paris Dauphine

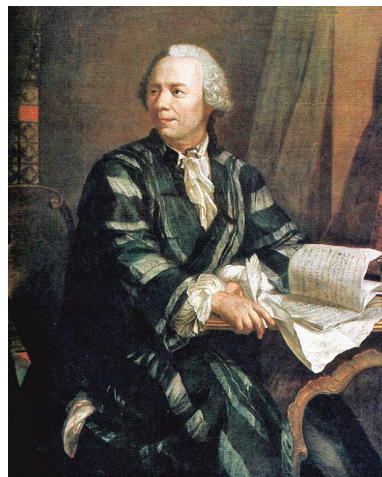
2023

MOTIVATION

COMPLÉTUDE

Quoique la doctrine des logarithmes soit si solidement établie, que les vérités qu'elle renferme semblent aussi rigoureusement démontrées que celles de la Géométrie; les Mathématiciens sont pourtant encore fort partagés sur la nature des logarithmes négatifs et imaginaires.

— Leonhard Euler (1749)



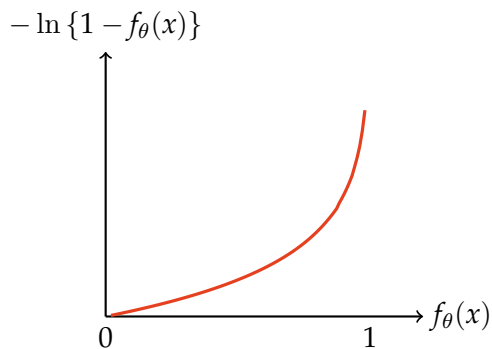
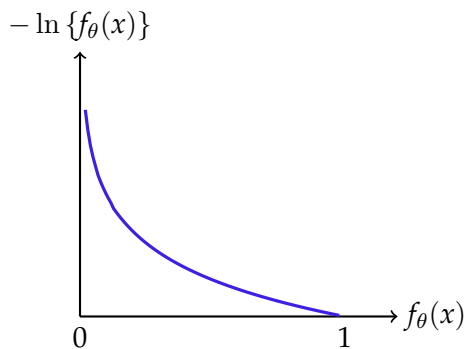
PROBLÈME

MODÉLISATION

Observation positive

$$\mathcal{L}(\theta; x, y) = - \left[y \ln \{f_{\theta}(x)\} + (1 - y) \ln \{1 - f_{\theta}(x)\} \right]$$

Observation négative

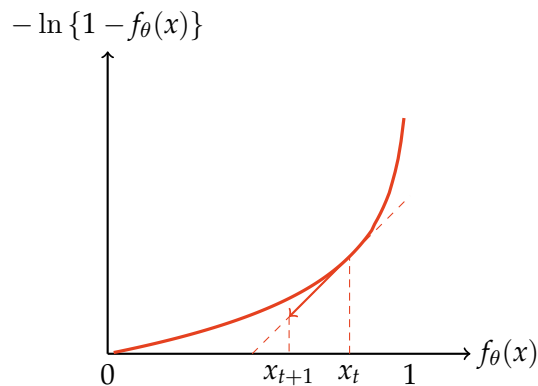
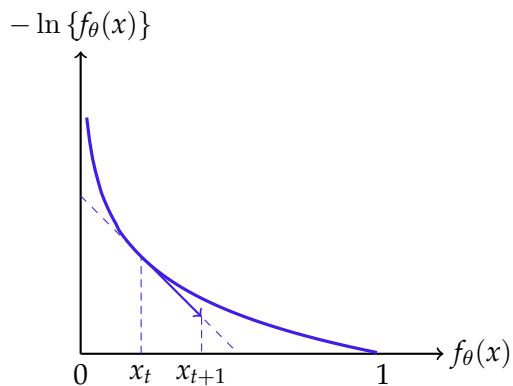


PROBLÈME

DESCENTE DE GRADIENT

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

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



PROBLÈME

SOLUTION

Exercice 1

On rappelle que $f_\theta(x) = \frac{1}{1 + e^{-\langle \theta, x \rangle}}$. 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 le problème avec la fonction de coût est :

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

MESURER LA PERFORMANCE D'UNE RÉGRESSION

ACCURACY

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

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

Exercice 2

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.

MESURER LA PERFORMANCE D'UNE RÉGRESSION

PRÉCISION, RECALL ET F1-SCORE

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

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

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

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

MESURER LA PERFORMANCE D'UNE RÉGRESSION

MÉTRIQUES

Exercice 3

Vous avez trop de mails, et vous demandez à votre data scientist de concevoir un algorithme qui va prioriser les mails en essayant de prédire les mails qui sont les plus importants. Vous lui donnez un dataset d'entraînement et un dataset de test. 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 :

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

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

Réal	Prédit	
	Classe 0	Classe 1
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.