

INTRODUCTION TO MACHINE LEARNING

GBM

Sara El Bouch

February 27, 2025

Rappels

- Données : $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d}) \in \mathbb{R}^d$
- Valeurs à prédire : $y_i \in \mathbb{R}$ (régression)
- Ensemble d'entraînement : $\mathcal{D}_m = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$

Le modèle linéaire est donné par :

$$y_i = \beta_0 + \sum_{j=1}^d \beta_j x_{i,j} + \quad i \in \{1, \dots, m\} \quad (1)$$

où :

- $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_d)^\top \in \mathbb{R}^{d+1}$ est le **vecteur des paramètres inconnus** à estimer.

Nous cherchons β en minimisant l'erreur quadratique :

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^m \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{i,j} \right)^2 \quad (2)$$

$$= \arg \min_{\beta} \frac{1}{2} RSS(\beta) \quad (3)$$

où $RSS(\beta)$ est la somme des erreurs quadratiques résiduelles.

En notation matricielle, on a :

$$RSS(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \quad (4)$$

avec :

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ 1 & x_{2,1} & \dots & x_{2,d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{m,1} & \dots & x_{m,d} \end{pmatrix} \in \mathbb{R}^{m \times (d+1)} \quad (5)$$

Nous voulons minimiser :

$$\min_{\boldsymbol{\beta}} RSS(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \quad (6)$$

Théorème (condition d'optimalité) : Si $RSS(\boldsymbol{\beta})$ est convexe, alors le minimum est obtenu lorsque :

$$\nabla RSS(\hat{\boldsymbol{\beta}}) = \mathbf{0}. \quad (7)$$

Le gradient est donné par :

$$\nabla RSS(\boldsymbol{\beta}) = -\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} \quad (8)$$

L'optimisation donne la solution analytique :

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (9)$$

sous la condition que $\mathbf{X}^\top \mathbf{X}$ soit inversible ($m > d + 1$).

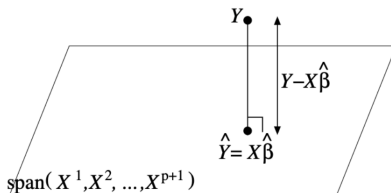
Remarque. Si d est très grand, inverser $\mathbf{X}^\top \mathbf{X}$ devient prohibitif. (La descente du gradient, stay tuned!)

Projection orthogonale sur l'espace des prédictions :

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \boldsymbol{\epsilon}$$

La minimisation revient à projeter \mathbf{y} sur le sous-espace engendré par les colonnes de \mathbf{X} :

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0$$



Mesures de performance :

- Erreur quadratique moyenne (MSE) :

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

- Coefficient de corrélation r :

$$r = \frac{\sum_{i=1}^m (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^m (y_i - \bar{y})^2 \sum_{i=1}^m (\hat{y}_i - \bar{\hat{y}})^2}}$$

Attention : Toujours évaluer sur un **jeu de test** indépendant.

Pour éviter le sur-apprentissage, on ajoute une pénalisation λ :

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \quad (10)$$

Effet de la régularisation :

- Réduit la variance du modèle.
- Implique un compromis biais-variance.
- Empêche la singularité de $\mathbf{X}^\top \mathbf{X}$.

Principe :

- Ajoute une pénalisation L_1 à la régression linéaire :

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (y_i - X_i \beta)^2 + \lambda \sum_{j=1}^d |\beta_j|$$

- Contrairement à Ridge, qui réduit les coefficients sans les annuler, **Lasso force certains coefficients à zéro.**

Pourquoi utiliser Lasso ?

- Sélectionne automatiquement les variables importantes.
- Produit des modèles plus interprétables.
- Utile quand beaucoup de variables sont inutiles.

Effet des régularisations :

Méthode Sélection de variables	Effet sur les coefficients
Ridge (L_2)	Réduit l'amplitude mais ne les annule pas
Lasso (L_1)	Peut annuler certains coefficients

Qu'est-ce que la régression logistique ?

- La régression linéaire est utilisée pour une variable dépendante **continue**.
- On va s'intéresser à **la classification** au travers de la régression logistique = Le modèle \rightarrow coût \rightarrow entraînement.

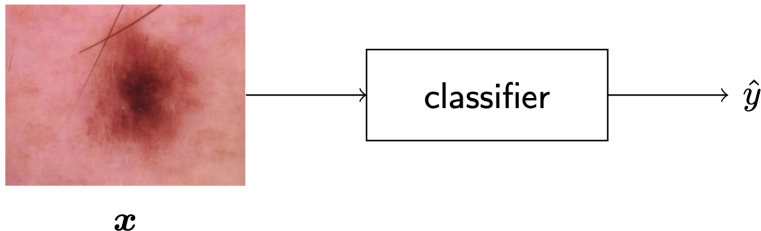


Figure: y and $\hat{y} \in \{benign, malignant\}$

Classification simple

- On suppose que x est un nombre et $y \in \{0, 1\}$
- On suppose que les différentes valeurs de x sont linéairement séparables.

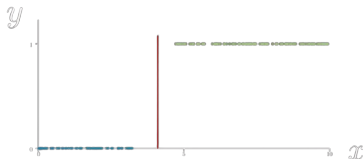
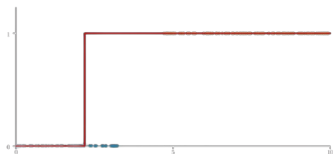


Figure: On choisit un seuil (ligne rouge) qu'on note s

$$\hat{y}_i = \begin{cases} 0 & \text{si } s \times x \leq 0 \\ 1 & \text{sinon} \end{cases} \quad (11)$$

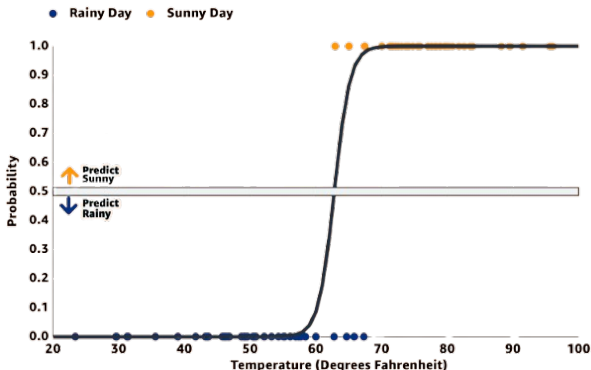
$$\hat{y}_i = \begin{cases} 0 & \text{si } a \times x + b \leq 0 \\ 1 & \text{sinon} \end{cases} \quad (12)$$



- On veut interpréter \hat{y}_i comme **une probabilité** On introduit la fonction **sigmoïde**:

$$\sigma(x) = \frac{1}{1 + \exp^{-x}} \quad (13)$$

$$\hat{y}_i = p(y = 1|x) = \sigma(ax + b) \in [0, 1]. \quad (14)$$



On définit le coût de l'entropie pour une donnée i comme:

$$J^i = -\log p_\beta(y = y_i|x_i) \quad (15)$$

- Pour une classification **binaire** on a 2 cas, $y_i = 0$ ou $y_i = 1$

$$J^i = \begin{cases} -\log p_\beta(y = 1|x_i), & y_i = 1 \\ -\log p_\beta(y = 0|x_i), & y_i = 0 \end{cases} \quad (16)$$

En notant $\hat{y}_i = p(y = 1|x)$ (c'est un choix arbitraire, on peut choisir $\hat{y}_i = p(y = 0|x)$, on peut réécrire J^i comme:

$$J^i = -y_i \log(\hat{y}_i) - (1 - y_i) \log(1 - \hat{y}_i) \quad (17)$$

Comment optimiser cette fonction coût ?

On peut utiliser **l'algorithme de descente du gradient** pour trouver les paramètres a et b optimaux.

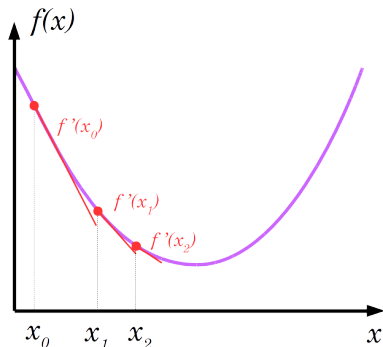


Figure: Descente du gradient

	Prédit Positif	Prédit Négatif
Réel Positif	VP (Vrais Positifs)	FN (Faux Négatifs)
Réel Négatif	FP (Faux Positifs)	VN (Vrais Négatifs)

Table: Matrice de confusion

$$\text{Précision} = \frac{VP}{VP + FP} \quad (18)$$

Interprétation : Proportion des prédictions positives qui sont réellement positives.

$$\text{Rappel} = \frac{VP}{VP + FN} \quad (19)$$

Interprétation : Proportion des cas positifs correctement identifiés.

Formule :

$$\text{Spécificité} = \frac{VN}{VN + FP} \quad (20)$$

Interprétation : Proportion des cas négatifs bien classés.

Formule :

$$F1 = 2 \times \frac{\text{Précision} \times \text{Rappel}}{\text{Précision} + \text{Rappel}} \quad (21)$$

Interprétation : Mesure équilibrée entre précision et rappel.

- Ce modèle peut être étendu au cas multi-classe, où $y \in \{1, \dots, K\}$, ce qui donne lieu à la régression logistique multinomiale (on utilise plus la fonction sigmoïde mais plutôt la fonction **softmax**).
- Bien que ce modèle soit utilisé pour la classification, on l'appelle un "régresseur" car la fonction $\sigma(ax + b)$ produit des valeurs réelles comprises entre 0 et 1. C'est l'interprétation de ces valeurs (par exemple, en appliquant un seuil à 0.5) qui permet de prendre des décisions de classification.
- On peut également utiliser la régularisation comme dans la régression linéaire pour éviter le sur-apprentissage.