

Introduction à l'apprentissage supervisé

①

I. Des données et des prédictions

Observations: $(x_i, y_i) \in X \times Y, i=1, \dots, n$
(Données d'entraînement)
(Training Data)

- Inputs:
- Images
 - Sons
 - Vidéos
 - Textes
 - ...

- Outputs:
- Classification binaire
 $\in \{0, 1\}$ ou $\in \{-1, 1\}$
 - Classification multiclasse
 $\in \{0, 1, \dots, k-1\}$ ou $\in \left\{ \begin{array}{l} (1, 0, 0, \dots, 0) \\ (0, 1, 0, \dots, 0) \\ \vdots \\ (0, \dots, 0, 1) \end{array} \right\}$
 - Valeurs $\in \mathbb{R}^d$

Objectif:

À partir d'un x observé ou non pendant l'entraînement, être capable de "prédire" le y qui lui correspond

le mieux

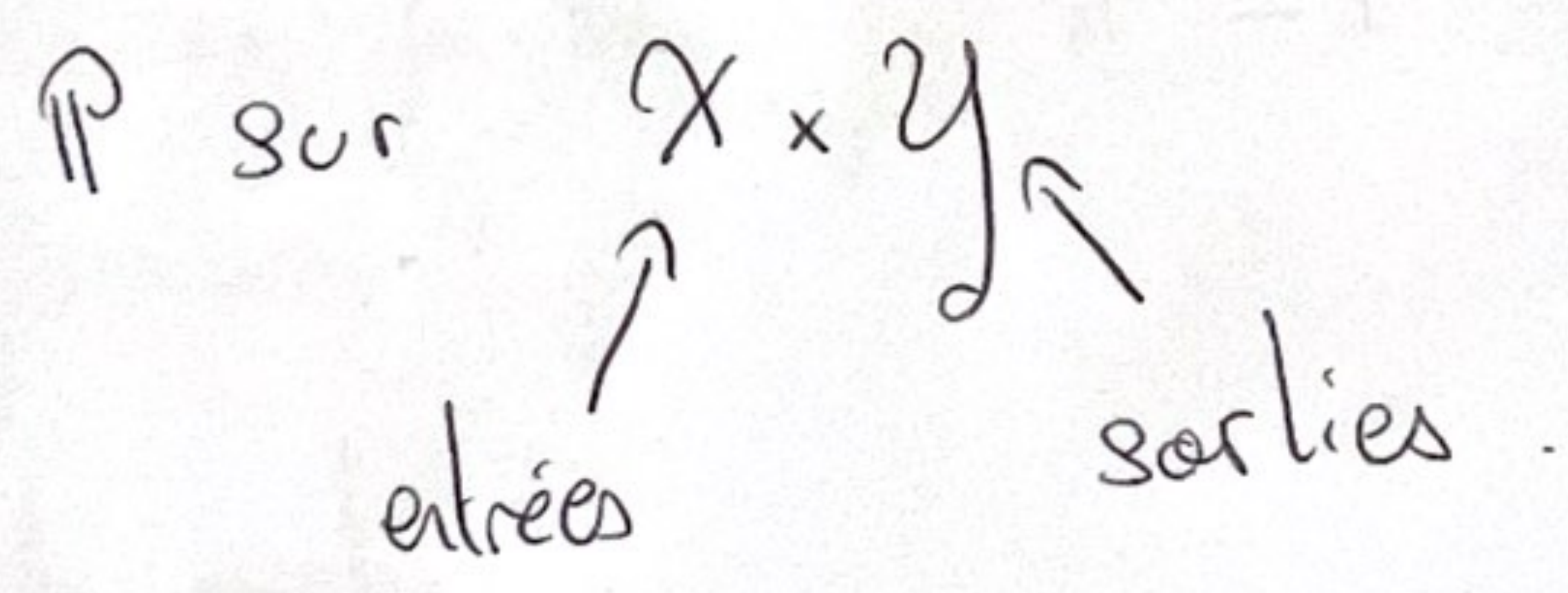
$$y \approx f(x)$$

- y peut être une fonction aléatoire de x
- f peut être complexe
- On n'observe qu'un nombre fini de données d'échantillon
- La dimension peut être élevée

II. Formalisation Mathématique



La modélisation se fait via la Théorie des probabilités



Hypothèse d'indépendance

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \stackrel{i.i.d}{\sim} \mathbb{P}$

Objet:

Estimer $P(y|x)$ à partir des conséquences conditionnellement aux causes.

Ce problème peut parfois être simplifié en estimant ③

$$f: X \rightarrow Y$$

telle que $P(y|x) \approx \delta_{y(f(x))}$, ou de manière simplifiée, $y \approx f(x)$.

III. Evaluation des performances

Étant donné une fonction $f: X \rightarrow Y$, comment évaluer ses performances ?

1) Fonction de perte

$$P: Y \times Y \rightarrow \mathbb{R}$$

avec pour interprétation $P(y, \hat{y}) =$ erreur induite par la prédiction de \hat{y} alors que la vérité est y .

Exemples

• Classification: $P(y, \hat{y}) = \mathbb{1}_{y \neq \hat{y}}$

• Régression: $P(y, \hat{y}) = (\hat{y} - y)^2$ ou $P(y, \hat{y}) = |\hat{y} - y|$

2) Risque

- Risque moyen

$$R(f) := \mathbb{E}_{(x,y)} (l(y, f(x)))$$

Exemple

En classification binaire

$$R(f) = 0 \times P(f(x) = y) + 1 \times P(f(x) \neq y) \\ = P(f(x) \neq y)$$

⚠ On aimerait choisir un f qui minimise $R(f)$. Le problème est que $R(\cdot)$ utilise la vraie distribution des données qui n'est pas connue ! Il faudra utiliser un estimateur.

3) Risque de Bayes - Prédicteurs de Bayes.

Proposition - Définition:

Le risque moyen $R(\cdot)$ est minimisé par le prédicteur de Bayes $f_* : X \rightarrow Y$ tq

$$\forall x' \in X, f_*(x') = \operatorname{argmin}_{y \in Y} \mathbb{E}(l(y, \cdot) | x = x') \\ =: r(y | x')$$

(5)

Le risque de Bayes est défini comme le risque du prédicteur de Bayes.

$$M^* := M(f_*)$$

Démonstration:

$$M(f) - M^* = M(f) - M(f_*)$$

$$= \int_{\mathcal{X}} \left(\underbrace{c(f(x') | x') - \min_{j \in \mathcal{Y}} c(j | x')}_{\geq 0} \right) dP(x')$$

≥ 0 .

□

Exemples:

- Classification

$$f_*(x') \in \underset{j}{\operatorname{argmin}} P(y \neq j | x = x')$$

$$= \underset{j}{\operatorname{argmax}} P(y = j | x = x').$$

- ~~Classification~~ Régression:

$$f_*(x') \in \underset{j}{\operatorname{argmin}} E\left((y - j)^2 | x = x'\right)$$

$$= \underset{j}{\operatorname{argmin}} E\left(\left(y - E(y | x = x')\right)^2 | x = x'\right) + \left(j - E(y | x = x')\right)^2$$

$$\text{donc } f_*(x') = E(y|x=x')$$

(6)

III. Paradigmes classiques d'apprentissage.

1) Minimisation du risque empirique.

$$\text{💡 } \forall f, R(f) \approx \frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i)) =: \hat{M}(f)$$

" donc "

$$\operatorname{argmin}_f R(f) \approx \operatorname{argmin}_f \hat{M}(f)$$

L'objet $\hat{M}(\cdot)$ s'appelle le risque empirique.

2) Méthodes génératives

Étape 1: Estimer $P(x|y)$ et $P(y)$

Étape 2: Dédurre $P(y|x)$ par la formule de Bayes

$$P(y|x) \propto P(x|y)P(y)$$

3) Méthodes par moyennes locales

💡 Approcher directement $f_*(x)$.

- Classification : $f_*(x') = \underset{j}{\operatorname{argmax}} P(y=j | x=x')$

↳ Méthode des k-plus proches voisins.

Pour un x' , prédire la classe majoritaire parmi les k plus proches voisins de x' dans le jeu d'apprentissage.

- Régression : $f_*(x') = E(y | x=x')$

↳ Méthode des k-plus proches voisins.

Pour un x' , prédire la valeur moyenne des y des k plus proches voisins de x' dans le jeu d'apprentissage.

4) Méthodes Probabilistes - Maximum de Vraisemblance.

💡 Faire une hypothèse sur la classe de mesures de probabilités pertinentes pour le problème, et choisir celle qui maximise ~~de~~ les chances d'observer le jeu d'entraînement.

Si $P(x, y | \theta)$, $\theta \in \Theta$ est une paramétrisation de la P_i , Θ
 \uparrow
 paramètre de la P_i

$$\hat{\theta} \text{ e argmax}_{\theta \in \Theta} P((x_1, y_1), \dots, (x_n, y_n) | \theta)$$

$$= \underbrace{P(x_1, y_1 | \theta) \times \dots \times P(x_n, y_n | \theta)}_{\text{indépendance}}$$

$$=: \mathcal{L}((x_1, y_1), \dots, (x_n, y_n) | \theta)$$

\uparrow
 Fonction de vraisemblance (P.L.D.L.)

de manière équivalente,

$$\hat{\theta} \text{ e argmax}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ln(P(x_i, y_i | \theta))$$

$\underbrace{\hspace{10em}}_{\text{Log-vraisemblance (Log-P.L.D.L.)}}$

IV. Exemples:

1) Régression linéaire par minimisation du risque empirique

Modèle: $y = \theta^T x + \varepsilon \leftarrow \text{bruit}$.

(Perte quadratique + Modèle linéaire) $\Rightarrow \hat{\theta} \text{ e argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$

2) Méthodes génératives, cas gaussien (en 1D)

(9)

$$y \sim \mathcal{B}\left(\frac{1}{2}\right)$$

$$x|y \sim \mathcal{N}(\mu_y, 1)$$

$$\text{i.e. } p(x|y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-y)^2}{2}\right)$$

Dans ce cas simplifié, nous connaissons la loi de $x|y$ et la loi de y .
Nous n'avons donc pas à l'estimer.

Par la formule de Bayes,

$$P(y=0|x) \propto \frac{1}{2} \exp\left(-\frac{x^2}{2}\right)$$

$$P(y=1|x) \propto \frac{1}{2} \exp\left(-\frac{(x-1)^2}{2}\right)$$

$$\text{Ainsi, } \frac{P(y=0|x)}{P(y=1|x)} \geq 1 \quad \text{ssi} \quad \frac{(x-1)^2}{x^2} \geq 1 \quad \text{ssi} \quad \frac{1}{2} \geq x$$

3) Maximum de vraisemblance - Régression Logistique

$$\text{Fonction sigmoïde: } \sigma(x) := \frac{1}{1 + e^{-x}}$$

$$\text{Modèle Logistique: } p(y=1|x) = \frac{1}{1 + e^{-\theta^T x}} \quad y \in \{-1, 1\}$$
$$= \sigma(\theta^T x) = \sigma(y \theta^T x)$$

remarque

(10)

$$p(y = -1) = 1 - \sigma(\theta^T x)$$

$$= 1 - \frac{1}{1 + e^{-\theta^T x}} = \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{\theta^T x}}$$

$$= \sigma(y \theta^T x)$$

$$\text{donc } p(y|x) = \sigma(y \theta^T x)$$

Maximum de vraisemblance:

$$\hat{\theta} \text{ est argmax}_{\theta} \sum_{i=1}^n \ln \left(\frac{1}{1 + e^{-y_i \theta^T x_i}} \right)$$

$$\Rightarrow \text{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n \ln (1 + e^{-y_i \theta^T x_i})$$