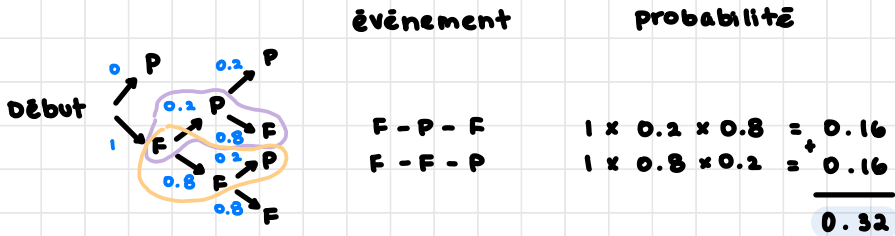


Q1 - Rappels de probabilités : probabilité conditionnelle et règle de Bayes

(a)
$$P(X|Y) = \frac{P(X \cap Y)}{P(Y)}$$

(b) $P(\text{face}) = 4/5 = 0.8$
 $P(\text{pile}) = 1/5 = 0.2$



(c) Expressions équivalentes de $P(X, Y)$

(i)
$$P(Y|X) = \frac{P(X, Y)}{P(X)}$$

$\Leftrightarrow P(X, Y) = P(Y|X) \cdot P(X)$

(ii)
$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

$\Leftrightarrow P(X, Y) = P(X|Y) \cdot P(Y)$

(d) Théorème de Bayes

Selon la définition de la probabilité conditionnelle, on a :

$$P(X|Y) = \frac{P(X \cap Y)}{P(Y)} \quad \text{et} \quad P(Y|X) = \frac{P(X \cap Y)}{P(X)}$$

En réarrangeant la deuxième équation, on a :

$$P(X \cap Y) = P(Y|X) \cdot P(X)$$

En substituant ce dernier dans la première équation, on a :

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$

Ainsi, on obtient la formule du théorème de Bayes :

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)} \quad \blacksquare$$

$$(e) P(UdeM) = 0.56$$

$$\begin{aligned} (i) P(McGill) &= 1 - P(UdeM) \\ &= 1 - 0.56 \\ &= 0.44 \end{aligned}$$

$$\begin{aligned} (ii) P(UdeM | bilingue) &= 0.7 \\ P(McGill | bilingue) &= 0.4 \end{aligned}$$

$$\begin{aligned} ① P(bilingue) &= P(bilingue | UdeM) \cdot P(UdeM) + \\ &\quad P(bilingue | McGill) \cdot P(McGill) \\ &= 0.7 \cdot 0.56 + 0.44 \cdot 0.4 \\ &= 0.568 \end{aligned}$$

$$\begin{aligned} ② P(McGill | bilingue) &= \frac{P(bilingue | McGill) \cdot P(McGill)}{P(bilingue)} \\ &= \frac{0.44 \cdot 0.4}{0.568} \\ &= 0.309859 \\ &\approx 0.31 \end{aligned}$$

Q2 - Bag of words & modèle de sujet unique

$$P(\text{sport}) = 1/3$$

$$P(\text{politique}) = 2/3$$

$$(a) P(\text{mot} = \text{goal} \mid \text{sujet} = \text{politique}) = 12/1000$$

$$\begin{aligned}(b) E &= P(\text{mot} = \text{congress} \mid \text{sujet} = \text{sport}) \cdot \text{nb total de mots} \\ &= 4/1000 \cdot 1000 \\ &= 4\end{aligned}$$

$$\begin{aligned}(c) P(\text{mot} = \text{goal}) &= P(\text{mot} = \text{goal} \mid \text{sujet} = \text{sport}) \cdot P(\text{sport}) + \\ &\quad P(\text{mot} = \text{goal} \mid \text{sujet} = \text{politique}) \cdot P(\text{politique}) \\ &= 4/100 \cdot 1/3 + 12/1000 \cdot 2/3 \\ &= 8/375 \\ &\approx 0.02\end{aligned}$$

$$\begin{aligned}(d) \textcircled{1} P(\text{mot} = \text{kick}) &= P(\text{mot} = \text{kick} \mid \text{sujet} = \text{sport}) \cdot P(\text{sport}) + \\ &\quad P(\text{mot} = \text{kick} \mid \text{sujet} = \text{politique}) \cdot P(\text{politique}) \\ &= 15/100 \cdot 1/3 + 2/1000 \cdot 2/3 \\ &= 0.051\bar{3}\end{aligned}$$

$$\begin{aligned}\textcircled{2} P(\text{sujet} = \text{sport} \mid \text{mot} = \text{kick}) &= \frac{P(\text{mot} = \text{kick} \mid \text{sujet} = \text{sport}) \cdot P(\text{sujet} = \text{sport})}{P(\text{mot} = \text{kick})} \\ &= \frac{15/100 \cdot 1/3}{0.051\bar{3}} \\ &= 0.974689 \\ &\approx 0.97\end{aligned}$$

$$\begin{aligned}(e) \textcircled{1} P(\text{mot} = \text{congress}) &= P(\text{mot} = \text{congress} \mid \text{sujet} = \text{sport}) \cdot P(\text{sujet} = \text{sport}) \\ &\quad + P(\text{mot} = \text{congress} \mid \text{sujet} = \text{politique}) \cdot P(\text{sujet} = \text{congress}) \\ &= 4/1000 \cdot 1/3 + 6/100 \cdot 2/3 \\ &= 31/750\end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \quad P(A) &= P(\text{sujet} = \text{sport} \mid \text{mot} = \text{congress}) \\
 &= \frac{P(\text{mot} = \text{congress} \mid \text{sujet} = \text{sport}) \cdot P(\text{sujet} = \text{sport})}{P(\text{mot} = \text{congress})} \\
 &= \frac{4/1000 \cdot 1/3}{31/750} \\
 &= 1/31
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{3} \quad P(A') &= P(\text{sujet} = \text{politique} \mid \text{mot} = \text{congress}) \\
 &= 1 - P(A) \\
 &= 1 - 1/31 \\
 &= 30/31
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{4} \quad P(\text{goal} \mid \text{congress}) &= P(\text{mot} = \text{goal} \mid \text{sujet} = \text{sport}) \cdot P(\text{sujet} = \text{sport} \mid \text{mot} = \text{congress}) + P(\text{mot} = \text{goal} \mid \text{sujet} = \text{politique}) \cdot P(\text{sujet} = \text{politique} \mid \text{mot} = \text{congress}) \\
 &= 4/100 \cdot 1/31 + 12/1000 \cdot 30/31 \\
 &= 2/155 \\
 &\approx 0.013
 \end{aligned}$$

(f) Estimation des probabilités des sujets

Il suffit de compter le nombre de documents associés à chaque sujet. Prenons N comme le nombre total de documents :

- N_{sport} : le nombre de documents de sport.
- $N_{\text{politique}}$: le nombre de documents de politique.

Alors la probabilité de chaque sujet :

$$\begin{aligned}
 \bullet \quad P(\text{sujet} = \text{sport}) &= \frac{N_{\text{sport}}}{N} \\
 \bullet \quad P(\text{sujet} = \text{politique}) &= \frac{N_{\text{politique}}}{N}
 \end{aligned}$$

Estimation des probabilités conditionnelles

Premièrement, il faut compter le nombre d'occurrences de chaque mot dans les documents de sport et de politique afin de construire une table de fréquence pour chaque mot et chaque sujet.

Ensuite, on peut estimer les probabilités à partir des fréquences.

- $C_{\text{goal, sport}}$: le nombre de fois que le mot goal apparaît dans les documents de sport.
- $C_{\text{total, sport}}$: le nombre total de mots dans les documents de sport.

Alors, la probabilité conditionnelle pour que le mot soit goal sachant que le sujet était sport :

$$P(\text{mot} = \text{goal} \mid \text{sujet} = \text{sport}) = \frac{C_{\text{goal, sport}}}{C_{\text{total, sport}}}$$

Q3 - Estimateur du maximum de vraisemblance

$$\begin{aligned} (a) \quad f_{\theta}(x_1, x_2, \dots, x_n) &= f_{\theta}(x_1) \cdot f_{\theta}(x_2) \cdot \dots \cdot f_{\theta}(x_n) \\ &= \prod_{i=1}^n f_{\theta}(x_i) \end{aligned}$$

$$\begin{aligned} (b) \quad f_{\theta}(x_1, x_2, \dots, x_n) &= f_{\theta}(x_1) \cdot f_{\theta}(x_2) \cdot \dots \cdot f_{\theta}(x_n) \\ &= \prod_{i=1}^n f_{\theta}(x_i) \\ &= \prod_{i=1}^n \frac{1}{\theta} \quad \text{pour } x_i \in]0, \theta] \\ &= \begin{cases} \left(\frac{1}{\theta}\right)^n & \text{si } 0 < x_i \leq \theta \\ 0 & \text{sinon} \end{cases} \\ &= \begin{cases} \frac{1}{\theta^n} & \text{pour } \theta \geq \max(x_1, x_2, \dots, x_n) \\ 0 & \text{sinon} \end{cases} \end{aligned}$$

La valeur de $\frac{1}{\theta^n}$ est inversement proportionnelle à θ . Ainsi,

puisque l'on veut maximiser notre fonction, on cherche à avoir

le plus petit θ tout en respectant que $\theta \geq \max(x_1, x_2, \dots, x_n)$.

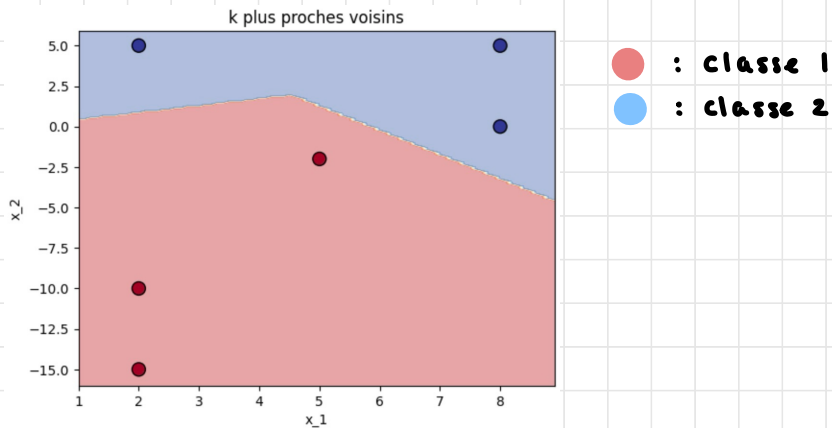
L'estimateur du maximum de vraisemblance θ est donc égal

à $\max(x_1, \dots, x_n)$.



Q4 - k plus proches voisins

(a)



(b) ① Moyenne classe 1 = (3, -9)

$$x_1 = \frac{5 + 2 + 2}{3} = 3$$

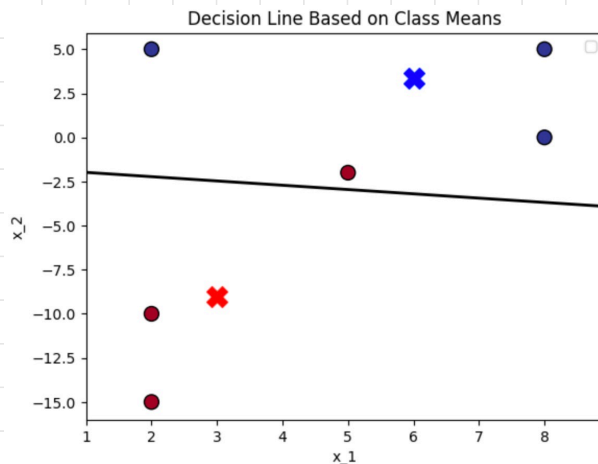
$$x_2 = \frac{-2 + (-10) + (-15)}{3} = -9$$

② Moyenne classe 2 = (6, 3.3)

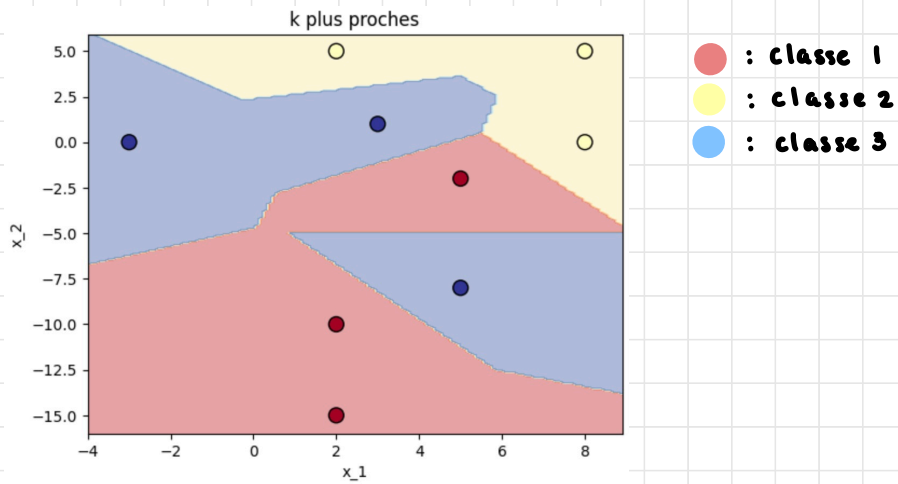
$$x_1 = \frac{8 + 2 + 8}{3} = 6$$

$$x_2 = \frac{0 + 5 + 5}{3} = 3.\bar{3}$$

③



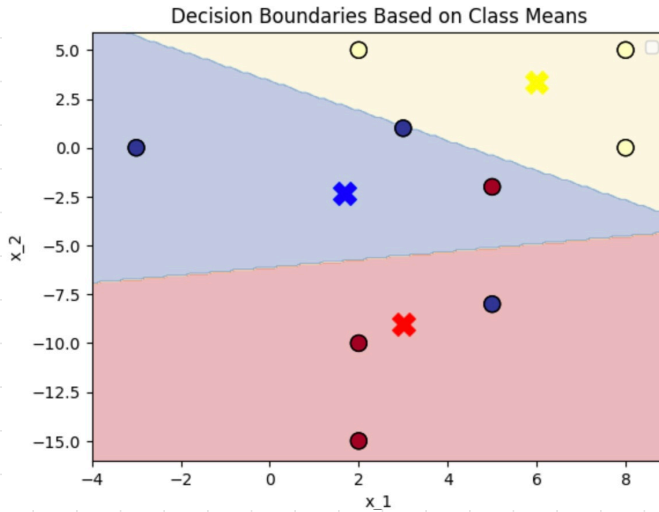
(c)



Moyenne classe 3 = $(1.\bar{6}, -2.\bar{3})$

$$x_1 = \frac{3 + (-3) + 3}{3} = 1.\bar{6}$$

$$x_2 = \frac{1 + 0 + (-8)}{3} = -2.\bar{3}$$



QS - Méthodes à base d'histogrammes

$$(a) I_{\{x \in S\}} = \begin{cases} 1, & \text{si } x \in S \\ 0, & \text{sinon} \end{cases}$$

$$\text{Ainsi, } E[I_{\{x \in S\}}] = 1 \cdot P(x \in S) + 0 \cdot (1 - P(x \in S)) \\ = P(x \in S)$$

■

$$(b) p_i = P(x \in V_i) \\ = \int_{V_i} f(x) dx$$

où p_i est la vraie probabilité d'être dans la région i , V_i est le volume de la région i et $f(x)$ est la fonction PDF inconnue

$$\hat{p}_i = \frac{k_i}{n}$$

où \hat{p}_i est la probabilité empirique d'être dans la région i , k_i est le nb de points dans la région i et n est le nb total de points

Ainsi, par la loi des grands nombres, on a :

$$\frac{k_i}{n} \rightarrow P(x \in V_i) = \int_{V_i} f(x) dx = p_i \text{ lorsque } n \rightarrow \infty$$

En effet, lorsque $n \rightarrow \infty$, la probabilité empirique \hat{p} s'approche de la vraie probabilité p_i .

■

$$(c) d = 464 \\ m = 2$$

$$\text{nb de régions} = m^d \\ = 2^{464}$$

$$(d) \text{ précision} = 90\% \\ \epsilon = 10\%$$

$$\textcircled{1} 90\% - 10\% = 80\%$$

$$\textcircled{2} (80\%) \div (5\%) = 16 \text{ (précision augmente en 16 intervalles de 5\%)}$$

$$\textcircled{3} 4 \times 16 = 64 \text{ nouveaux points / région}$$

(e) nb région = m^d

$$P(\text{point région}) = \frac{1}{m^d}$$

$$P(\text{point région vide}) = \left(1 - \frac{1}{m^d}\right)^n$$