

Cours de Data Science

Marc Tommasi

19 octobre 2022

Outline

1 Naive Bayes

2 Arbres de décision

Outline

- 1 Naive Bayes
- 2 Arbres de décision

Probabilités conditionnelles et Bayes

- Notations :

- ▶ $P(A)$ la probabilité d'un événement A
- ▶ $P(A \cap B)$ ou $P(A, B)$ la probabilité d'avoir à la fois les événements A et B .
- ▶ $P(A | B)$ la probabilité d'avoir l'événement A sachant B .

Definition des probabilités conditionnelles

$$P(A|B) = \frac{P(A, B)}{P(B)} \text{ ou encore par symétrie } P(B|A) = \frac{P(A, B)}{P(A)}$$

Donc

$$P(A, B) = P(A|B)P(B) = P(B|A)P(A)$$

Théorème de Bayes

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}$$

Petit abus de notations

- Abus de notation : je note donc $P(X = \mathbf{x})$ simplement par $P(\mathbf{x})$. (Soit : $P(\mathbf{x})$ pour dire qu'une variable aléatoire X prend comme valeur \mathbf{x} donc la probabilité de l'événement $X = \mathbf{x}$)
 - Exemple : Soit X soit la variable sur deux attributs x_1 et x_2 qui peuvent prendre uniquement des valeurs binaires. Alors $P(X)$ désigne $P(X = (0, 0))$, $P(X = (0, 1))$, etc.. notés simplement $P(0, 0)$, $P(0, 1)$, etc...
 - On va noter aussi $P(\mathbf{x})$ ou $P(x_1, x_2)$ pour désigner l'un de ces 4 cas.

Revenant à l'apprentissage

Rappels

- Les données sont générées par une probabilité jointe fixée mais inconnue d'avoir une description des données \mathbf{x} et une classe y , écrite $P(\mathbf{x}, y)$.
- On veut chercher à résoudre le problème : trouver le meilleur y quand on observe un \mathbf{x}

La règle de Bayes

- C'est la meilleure règle qu'on puisse imaginer

$$\operatorname{argmax}_y P(y \mid \mathbf{x})$$

- **Erreur de Bayes** : Erreur de cette règle
- c'est la plus petite erreur qu'on puisse faire pour cet apprentissage si les exemples sont décrits par \mathbf{x} .

Difficile à calculer

- $P(y | \mathbf{x})$ ne peut être calculée car P est inconnue.
- Si on applique le principe ERM, par la règle de Bayes on a

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

- On cherche la valeur de y qui maximise cette quantité. Mais $P(\mathbf{x})$ ne dépend pas de y . Il suffit de résoudre

$$\operatorname{argmax}_y P(y)P(\mathbf{x} | y)$$

- Problème : le calcul ne peut être fait efficacement \mathbf{x} et des valeurs qu'il peut prendre. Exemple : Dans le cas binaire avec 5 attributs, on a 2^5 possibilités et donc 2 fois 2^5 quantités à estimer pour tous les cas de $P(\mathbf{x} | y)$.

La chain rule

$$\begin{aligned}P(x_1, x_2, \dots x_n) &= P(x_1 \mid x_2, \dots x_n)P(x_2, \dots x_n) \\&= P(x_1 \mid x_2, \dots x_n)P(x_2 \mid x_3, \dots x_n)P(x_3, \dots x_n) \dots\end{aligned}$$

- Obtenu en appliquant $P(A, B) = P(A|B)P(B)$ de façon répétée quand B est un événement qui peut être une conjonction d'événements
- Par récursion

$$P(\mathbf{x} \mid y) = \prod_{k=1}^n P(x_k \mid x_{k-1}, \dots, x_1, y)$$

Indépendance Conditionnelle et Naive Bayes

- Si A et B sont indépendants alors $P(A | B) = P(A)$.
- Si on fait l'hypothèse que tous les attributs sont indépendants, c'est-à-dire si x_i et x_j sont indépendants pour tout $i \neq j$, alors le produit s'écrit :

$$P(\mathbf{x} | y) = \prod_{k=1}^n P(x_k | y)$$

- Pour calculer chacun de ces $P(x_k | y)$ il suffit de compter !
- C'est une **approximation forte** !
- Mais le calcul est de **faible complexité**

$$\operatorname{argmax}_y P(y) \prod_{k=1}^n P(x_k | y)$$

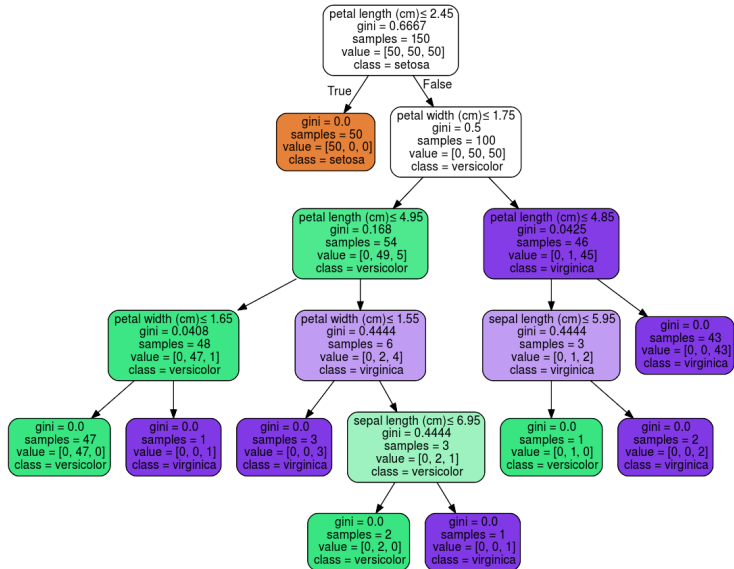
Outline

- 1 Naive Bayes
- 2 Arbres de décision

Principes

- Classifieur de \mathcal{X} dans \mathcal{Y}
- Modèle sous forme d'un ensemble de règles de décision successives, représentées dans un arbre
- Modèle simple à interpréter quand l'arbre est petit
- Exemple sur le jeu de données iris, 3 classes : setosa, virginica, versicolor ; 150 exemples ; attributs sepal length, sepal width, petal length, petal width.

Example



Fonctionnement

- On part de la racine, avec un exemple (e.g. 3, 2, 3, 2)
- On passe à travers les tests, chaque test coupe l'espace en 2 parties.
- On désigne donc un partitionnement récursif de l'espace de description
- Chaque feuille donne une étiquette à une partie.
- Bonne visualisation dans le livre de [Jake VandenPlas](#).