

CHAPITRE 3

Apprentissage

Hugo SALOU MPI*

Dernière mise à jour le 31 octobre 2022

0 Motivation

L'intelligence artificielle est vu comme un « objet magique » mais ce n'est pas le cas : c'est ce que nous allons étudier dans ce chapitre. Il existe plusieurs méthodes permettant l'apprentissage : descente de gradient, ?, ...

La base de donnée la plus utilisée est MNIST : elle contient 60 000 images de 28×28 pixels représentant un chiffre, et le chiffre correspondant. L'idée de l'apprentissage est de « deviner » le chiffre dessiné en connaissant l'image.

1 Vocabulaire

Définition: On appelle *signature de données* un n -uplet de paires nom, ensemble ; on le typographie

$$(\text{nom}_1 : S_1, \text{nom}_2 : S_2, \dots, \text{nom}_n : S_n).$$

Définition: Étant donné une signature de données $S = (\text{nom}_1 : S_1, \dots, \text{nom}_n : S_n)$, on appelle *donnée* un vecteur

$$\bar{v} = (v_1, v_2, \dots, v_n) \in S_1 \times S_2 \times \dots \times S_n.$$

Définition: Étant donné une signature de données S , on appelle *jeu de données* un ensemble fini de vecteurs de signature S .

Définition: Étant donnée une signature de données S et un ensemble de classes \mathcal{C} , on appelle *jeu de données classifié* la donnée

- d'un jeu de données S ,
- d'une fonction $f : S \rightarrow \mathcal{C}$ de classification.

2 Apprentissage supervisé

L'objectif de cette section est de construire des fonctions de classification, à partir d'un jeu de données classifié.

Définition: Étant donné une signature de données S , et un ensemble de classes \mathcal{C} , on appelle *fonction de classification* une fonction des données de signature S dans \mathcal{C} .

REMARQUE:

On discutera de la « qualité » d'une fonction de classification en fonction de ses résultats sur les données d'un jeu de données et sur des exemples de tests.

2.1 k plus proches voisins

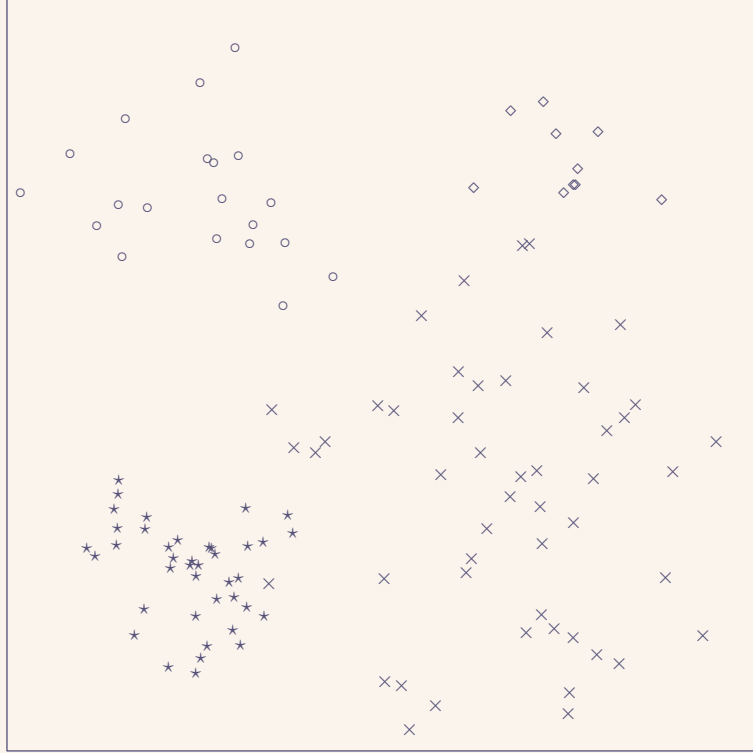


FIGURE 1 – Représentation de l'algorithme des k plus proches voisins

Algorithme 1 k -NN (k nearest neighbors)

Entrée Un jeu de données classifié (S, c) , un vecteur d'entrée \bar{v}

- 1: On trie S par distance à croissante de v en $d_1, d_2, \dots, d_k, d_{k+1}, \dots$
 - 2: Soit D un dictionnaire de \mathbb{C} vers \mathbb{N} initialisé à 0¹
 - 3: **pour** $j \in \llbracket 1, k \rrbracket$ **faire**
 - 4: $D[c(d_j)] \leftarrow D[c(d_j)] + 1$
 - 5: **retourner** $\operatorname{argmax}_{d \in \mathbb{C}} D[d]$
-

REMARQUE:

On doit avoir $k \leq n$, et l'espace doit être muni d'une distance. Les résultats de l'algorithme dépendent fortement du jeu de données, du paramètre k et de la distance choisie.

Matrice de confusion

Définition: On appelle *matrice de confusion* d'un algorithme de prédiction \mathcal{A} sur un jeu de données classifié (T, c) , la matrice

$$\left(\operatorname{Card}\{t \in T \mid \mathcal{A}(t) = i \text{ et } c(t) = j\} \right)_{(i,j) \in \mathbb{C}^2}.$$

Dans le cas particulier dans le cas d'une classification (\mathbf{V}, \mathbf{F}) , on nomme

1. où toutes les valeurs sont initialisées à 0, pas un dictionnaire vide

$\mathcal{A} \backslash$	vrai	F	V
F		vrai négatif	faux négatif
V		faux positif	vrai positif

TABLE 1 – Matrice de confusion dans le cas d’une classification en V et F

Comment améliorer la performance de l'algorithme des k plus proches voisins ? En dimension 1, on peut utiliser une dichotomie. Mais, dans des dimensions plus grandes, l'ordre lexicographique, et l'ordre produit ne fonctionnent pas. Mais, on peut appliquer une "dichotomie" en changeant de dimension. Par exemple, en deux dimension, on a

À faire : Représenter le schéma

FIGURE 2 – Représentation de la “dichotomie” en dimension 2

Pour représenter cette structure de données, on utilise un arbre binaire comme montré ci-dessous. Cet arbre est appelé un arbre k -dimensionnels.

À faire : Faire l'arbre

FIGURE 3 – Arbre 2-dimensionnel représentant la “dichotomie” précédente

2.2 Arbres k -dimensionnels

REMARQUE (Notations):

Étant donné un jeu de données S , on note pour $v \in S$,

$$S^{\leq i v} = \{u \in S \mid u_i \leq v_i\} \quad \text{et} \quad S^{> i v} = \{u \in S \mid u_i > v_i\}.$$

Algorithme 2 “F” : Fabrication d’un arbre k -dimensionnel

Entrée \mathcal{V} un jeu de données et $i \in \llbracket 0, n - 1 \rrbracket$, où n est la dimension des données

1: **si** $\mathcal{V} = \emptyset$ **alors**

2: | retourner Vide

3: sinon

4: On cherche $v \in \mathcal{V}$ tel que v_i est la médiane de $\{u_i \mid u \in \mathcal{V}\}$

5: **retourner** $\mathbf{N}\mathbf{eud}\left((v, i), \mathbf{F}(\mathcal{V} \setminus \{v\})^{\leq i v}, i + 1 \bmod n), \mathbf{F}(\mathcal{V} \setminus \{v\})^{> i v}, i + 1 \bmod n)\right)$

Algorithme 3 “R” : Recherche du point le plus proche

Entrée Un arbre k -dimensionnel et un vecteur v

```
1: si  $T$  est vide alors
2: |   retourner  $\emptyset$ 
3: sinon
4: |    $\text{Nœud}((u, i), G, D) \leftarrow T$ 
5: |   si  $u_i \leq v_i$  alors
6: | |    $W \leftarrow \text{R}(D, v)$ 
7: | |   si  $W = \text{None}$  alors
8: | | |    $W' \leftarrow \text{R}(G, v)$ 
9: | | |   si  $W' = \emptyset$  alors
10: | | | |   retourner  $\text{Some}(u)$ 
11: | | |   sinon
12: | | | |    $\text{Some}(z) \leftarrow W'$ 
13: | | |   retourner le plus proche de  $v$  entre  $u$  et  $z$ 
14: |   sinon
15: | |    $\text{Some}(w) \leftarrow W$ 
16: | |   si  $v_i - u_i \leq d(w, v)$  alors  $\triangleright d(w, v)$  représente la distance entre  $w$  et  $v$ 
17: | | |    $W' \leftarrow \text{R}(G, v)$ 
18: | | |   si  $W' = \text{None}$  alors
19: | | | |   retourner  $\text{Some}(\text{plus proche de } v \text{ entre } u \text{ et } w)$ 
20: | | |   sinon
21: | | | |    $\text{Some}(z) \leftarrow w$ 
22: | | |   retourner  $\text{Some}(\text{plus proche de } v \text{ entre } u, v, \text{ et } w)$ 
23: |   sinon
24: | |   retourner  $W$ 
```
