

Définition 5.1 - *apprentissage*

L'*apprentissage* en informatique permet l'approche de différents problèmes :

- la *classification*, par exemple déterminer un objet sur une image, ou un son sur un flux audio ;
- la *régression*, par exemple prévoir la valeur du cours de la bourse.

Définition 5.2 - *apprentissage supervisé*

L'*apprentissage supervisé* consiste à entraîner un modèle ou un algorithme à l'aide d'un *ensemble d'apprentissage* :

$$S = \{(x_i, y_i), i \in \llbracket 1, n \rrbracket\} \subset X \times Y \quad \text{avec} \quad \begin{cases} X & \text{l'ensemble des } \textit{objets} \text{ manipulés par le modèle ou l'algorithme} \\ Y & \text{l'ensemble des } \textit{classes} \text{ ou } \textit{valeurs} \text{ associées aux objets de } X \end{cases}$$

$(x, y) \in S$ signifie que *l'objet x est dans la classe y ou a pour valeur y .*

On cherche pour un objet inconnu x à déterminer une classe ou une valeur y convenable en s'appuyant sur l'ensemble d'apprentissage S .

Définition 5.3 - *fonction de prédiction*

À tout modèle ou algorithme d'apprentissage supervisé on peut associer une *fonction* $f : X \rightarrow Y$ dite de *prédiction* qui à un objet associe la *classe estimée raisonnable par le modèle ou l'algorithme*.

Définition 5.4 - *fonction de perte*

À tout modèle ou algorithme d'apprentissage supervisé on peut associer une fonction $L : Y^2 \rightarrow \mathbb{R}_+$ qui à une prédiction associe une *valeur mesurant son inexactitude*. On a :

$$\forall y \in Y, L(y, y) = 0$$

Définition 5.5 - fonction de risque empirique

Pour tout modèle ou algorithme d'apprentissage supervisé on peut associer à toute fonction f de prédiction une espérance appelée *risque* R par :

$$\begin{aligned} R(f) &= \mathbb{E}_{X,Y} \left(L(Y, f(X)) \right) \\ &= \sum_{(x,y) \in X \times Y} L(y, f(x)) \mathbb{P}_{X,Y}(x, y) \end{aligned}$$

En pratique, on n'a jamais accès à $\mathbb{P}_{X,Y}$. On définit alors le *risque empirique* comme étant la moyenne des pertes sur l'ensemble d'apprentissage. Lui est calculable :

$$R_{\text{emp}}(f) = \frac{1}{|S|} \sum_{(x,y) \in S} L(y, f(x))$$

Dès lors, un algorithme d'apprentissage supervisé mettra en œuvre des algorithmes d'optimisation afin de trouver une fonction f qui minimise le risque empirique.

Implémentation - algorithme de classification des k plus proches voisins - classique

- **Entrée :**
 - un ensemble d'apprentissage $S \subset X \times Y$ indexé sur $\llbracket 0, n \rrbracket$
 - une distance $d : X^2 \rightarrow \mathbb{R}_+$
 - un objet x de classe inconnue
 - k le nombre de voisins à considérer
- **Sortie :** la classe y du voisin majoritairement présent parmi les k plus proches de x

```
1 | file = file de priorité max vide
2 | pour tout  $i \in \llbracket 0, n \rrbracket$  :
3 |     si  $|file| < k$  :
4 |         ajouter  $x_i$  à file avec la priorité  $d(x, x_i)$ 
5 |     sinon :
6 |         si  $d(x, x_i) < file[0]$  :
7 |             supprimer le maximum de file
8 |             ajouter  $x_i$  à file avec la priorité  $d(x, x_i)$ 
9 | renvoyer la classe majoritaire parmi les éléments de file
```

Implémentation - construction d'un arbre d -dimensionnel

Cet algorithme constitue un prétraitement de l'ensemble d'apprentissage, pour "faciliter" la recherche des k plus proches voisins en aval. Attention, si d est trop grand, on fait face au fléau de la dimension qui rend la méthode trop peu efficace.

- **Entrée initiale :**
 - l'ensemble $X_S \subset X$ pour lequel on connaît la classe de chaque objet
 - d la dimension de X_S
- **Sortie :** un arbre d -dimensionnel décrivant les positions relatives de chaque objet de X_S

```
1 | Construit( $i$ , X_courant) :  
2 |      $n = |X\_courant|$   
3 |     si  $n == 0$  :  
4 |         renvoyer Vide  
5 |     sinon :  
6 |         X_courant = tri de X_courant dans l'ordre croissant de  
7 |             la  $i$ -ème coordonnée  
8 |         val = X_courant[ $\lfloor \frac{n}{2} \rfloor$ ]  
9 |         X_g, X_d = séparer strictement X_courant en deux listes,  
10 |             en l'indice  $\lfloor \frac{n}{2} \rfloor$   
11 |         g = Construit( $d$ ,  $i+1 \bmod d$ , X_g)  
12 |         d = Construit( $d$ ,  $i+1 \bmod d$ , X_d)  
13 |         renvoyer Noeud(val, g, d)
```

Implémentation - *algorithme de classification des k plus proches voisins - avec arbre d -dimensionnel (1/2)*

- **Entrée initiale :**

- un arbre d -dimensionnel décrivant les positions relatives de chaque objet de X_S
- un objet x de classe inconnue
- k le nombre de voisins à considérer

- **Sortie :** une pile contenant le chemin de la racine vers une feuille de \mathcal{A} modélisant le plus proche voisin de x .

```
1 | explore( $\mathcal{A}$ ,  $i$ ,  $x$ ) :
2 |     si  $\mathcal{A}$  est Vide :
3 |         renvoyer []
4 |     sinon :
5 |         on identifie  $\mathcal{A} = \text{Noeud}(\text{val}, g, d)$ 
6 |          $x_i = i$ -ème coordonnée de  $x$ 
7 |          $v_i = i$ -ème coordonnée de val
8 |         si  $x_i \leq v_i$  :
9 |             renvoyer  $\mathcal{A} :: \text{explore}(g, i+1 \bmod d, x)$ 
10 |        sinon :
11 |            renvoyer  $\mathcal{A} :: \text{explore}(d, i+1 \bmod d, x)$ 
```

Implémentation - *algorithme de classification des k plus proches voisins - avec arbre d -dimensionnel (2/2)*

• **Entrée initiale :**

- une distance $d : X^2 \rightarrow \mathbb{R}_+$
- l la pile du chemin parcouru vers le plus proche de x , contenant des sommets de \mathcal{A} un arbre d -dimensionnel
- un objet x de classe inconnue
- k le nombre de voisins à considérer

• **Sortie :** la classe y du voisin majoritairement présent parmi les k plus proches de x

```
1 | file = file de priorité max vide
2 | kPPV(1, i, x, k) :
3 |     Noeud(val, g, d), q = depiler l
4 |      $x_i$  =  $i$ -ème coordonnée de  $x$ 
5 |      $v_i$  =  $i$ -ème coordonnée de val
6 |     si file[0].prio < d( $x_i E_i, v_i E_i$ ) // distance unidimensionnelle :
7 |         // dépend de d mais souvent égal à  $|x_i - v_i|$ 
8 |         si besoin, ajouter/remplacer  $x$  à file avec la priorité d( $x, val$ )
9 |     sinon :
10 |         chemin = explore(Noeud(val, g, d),  $i \bmod d$ ,  $x$ )
11 |         kPPV(chemin, |file| +  $i \bmod d$ ,  $x$ ,  $k$ )
12 |     renvoyer la classe majoritaire parmi les éléments de file
```