

## Exercice 5

### Classification et prédiction

#### Distances

|                         |  |
|-------------------------|--|
| Distance de Manhattan : | $d(A, B) = \sum_{i=1}^n  A_i - B_i $   |
| Distance euclidienne :  | $d(A, B) = \sqrt{\sum_{i=1}^n (A_i - B_i)^2}$  |
| Distance de Minkowski : | $d(A, B) = (\sum_{i=1}^n  A_i - B_i ^p)^{1/p}$   |
| Distance de cosinus :   | $d(A, B) = 1 - \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$ |
| Distance de Hamming :   | $d(A, B) = \# \{ i : A_i \neq B_i \}$  |

#### Algorithme

|  |
|--|
| <b>Algorithme : k-NN</b>   |
| <p><b>Entrées :</b> <math>D</math> : Dataset ; <math>k</math> : nombre de voisins à considérer ; <math>Inst</math> : l'instance à classifier ;</p> <p><b>Sorties :</b> <math>y</math> : la classe de l'instance à classifier</p> <p><b>Var</b></p> <p style="padding-left: 20px;"><math>dist</math> : tableau de <math>[1..N]</math> de paire (instance , distance) ; //avec <math>N</math> la taille de <math>D =  D </math></p> <p style="padding-left: 20px;"><math>knn</math> : tableau de <math>[1..k]</math> d'instances ;</p> <p><b>Début</b></p> <p style="padding-left: 20px;"><b>Pour chaque</b> instance <math>X</math> de <math>D</math> <b>faire</b></p> <p style="padding-left: 40px;"><math>dist[X] \leftarrow</math> Calculer la distance entre <math>X</math> et <math>Inst</math> ;</p> <p style="padding-left: 20px;"><b>Fait;</b></p> <p style="padding-left: 20px;"><math>dist \leftarrow</math> Trier <math>dist</math> dans l'ordre croissant des distances;</p> <p style="padding-left: 20px;"><math>knn \leftarrow</math> les <math>k</math> premières instances de <math>dist</math> ;</p> <p style="padding-left: 20px;"><math>y \leftarrow</math> La classe dominante dans <math>knn</math> ;</p> <p style="padding-left: 20px;"><b>Retourner</b> <math>y</math>;</p> <p><b>Fin.</b></p> |

#### Questions :

- 1- Écrire une fonction Python permettant de calculer la distance entre deux instances du dataset (implémenter les distances applicables parmi ceux listés).
- 2- Écrire une fonction Python permettant de trier les instances du dataset selon la valeur d'une distance calculée "d".
- 3- Écrire une fonction Python permettant de retourner la classe dominante parmi un ensemble de  $K$  classes.
- 4- Implémenter en Python l'algorithme KNN ayant comme paramètres :  $k$ , la métrique de distance à employer.
- 5- Déduire la classe de l'instance  $\langle 5.2, 3.5, 1.41, 0.25 \rangle$  avec  $K = 3$  puis avec  $K = 5$ .