

Formation Machine Learning

Atelier Pratique AP-ML2

Atelier Pratique AP-ML2 : Exercice ML2.5

Objectif : Calculer la chance de survivre sur le Titanic

- 1) Charger les colonnes *pclass*, *survived*, *sex*, *age* du fichier *Titanic.xls* dans un Dataframe df , puis afficher le début du Dataframe
- 2) ENCODAGE : Remplacer les mots "*male*" et "*female*" par des valeurs numériques *0* et *1*
syntaxe : `df['sex'].replace(['male', 'female'], [0,1], inplace=True)`
- 3) Afficher les statistiques de base : `df.describe()`
- 4) Supprimer les lignes où il y a des données manquantes
- 5) Séparer le dataset en 2 parties : *Y (target)* et *X (features)*, où la target est la colonne *survived*

Atelier Pratique AP-ML2 : Exercice ML2.5

6) Créer et entraîner le modèle

7) Définir une fonction *predire (model, pclass, sex, age)*

Entrée : les information sur un passager (pclass, sex, age)

Sortie : l'information que le passager *survivrait (1) ou pas (0)*

8) Afficher le *% de chance de survivre d'un passager* (la probabilité que predict retourne 1)

en utilisant la fonction de Sklearn : *model.predict_proba(passager)*



SOLUTION

Atelier Pratique AP-ML2 : Exercice ML2.6

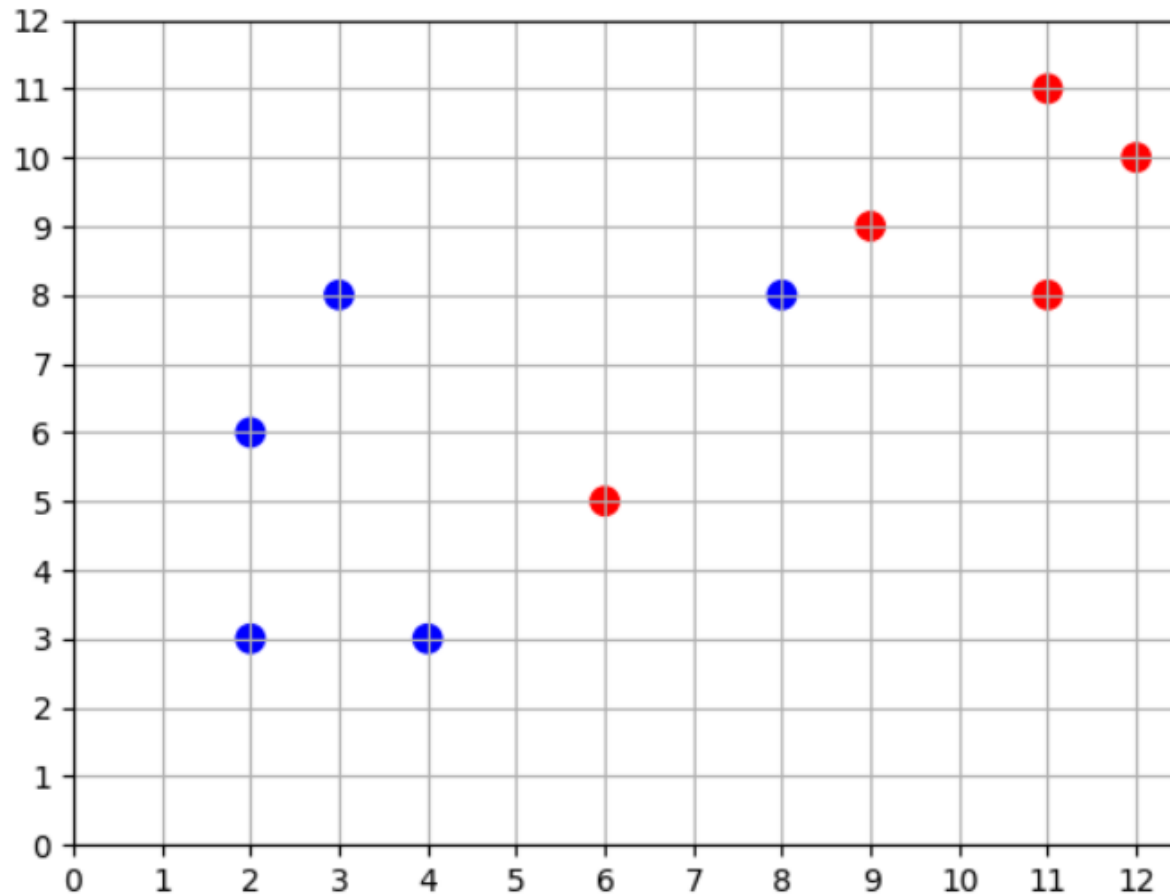
Objectif : Programmer l'algorithme KNN

1) Charger le dataset à partir du fichier *dataset-KNN.csv* dans un dataframe *df* puis afficher *df*

x1	x2	label
2	3	0
2	6	0
4	3	0
3	8	0
8	8	0
6	5	1
9	9	1
11	8	1
11	11	1
12	10	1

Atelier Pratique AP-ML2 : Exercice ML2.6

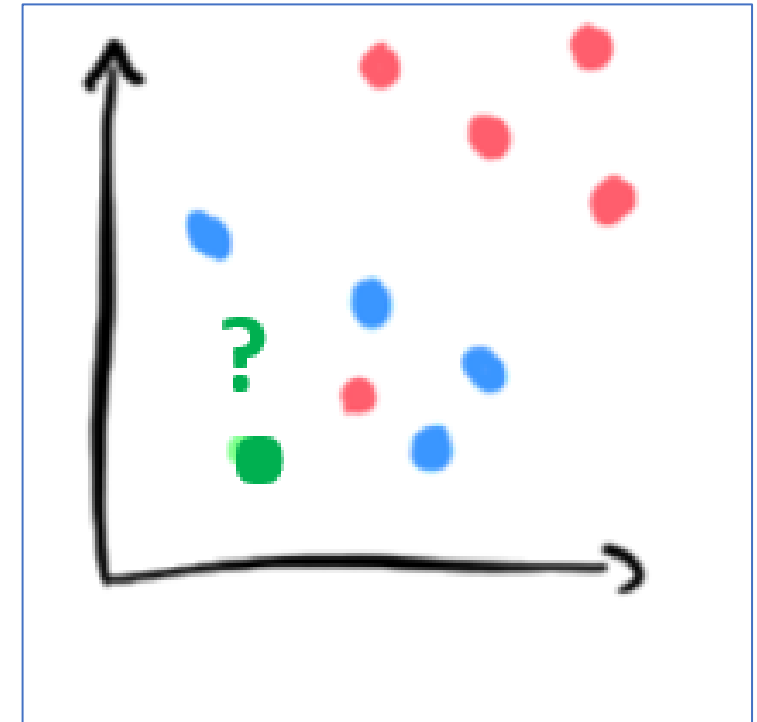
2) Afficher les données sur un plan en utilisant la fonction *scatter*. La figure suivante devrait s'afficher :



Atelier Pratique AP-ML2 : Exercice ML2.6

3) Ajouter une colonne '*distance*':

x1	x2	label	distance
2	3	0	
2	6	0	
4	3	0	
3	8	0	
8	8	0	
6	5	1	
9	9	1	
11	8	1	
11	11	1	
12	10	1	



Atelier Pratique AP-ML2 : Exercice ML2.6

4) isoler les features dans **X**, et les colonnes **label + distance** dans **Y**

X

x1	x2
2	3
2	6
4	3
3	8
8	8
6	5
9	9
11	8
11	11
12	10

Y

label	distance
0	
0	
0	
0	
0	
1	
1	
1	
1	
1	

Atelier Pratique AP-ML2 : Exercice ML2.6

5) Définir une fonction *distance (pt1, pt2)* qui calcule la *distance Euclidienne* entre deux points pt1 et pt2.

pt1 et pt2 sont des **tuples** :

- pt1 (a1, a2, ..., an)
- pt2 (b1, b2, ..., bn)

$$d(A,B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

Atelier Pratique AP-ML2 : Exercice ML2.6

6) Définir une fonction *myKNN* (*newpt*, *K*) qui retourne le label qu'elle va attribuer au nouveau point *newpt* :

Démarche :

a. Pour chaque point *pt* dans la dataset :

- calculer la distance entre *newpt* (nouveau point passé en paramètre) et ce point *pt*
- stocker la *distance* calculée dans la colonne *distance* de *Y*

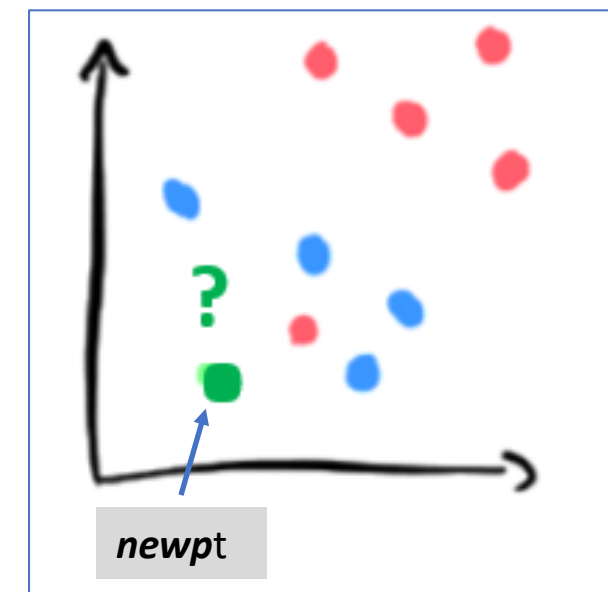
X

x1	x2
2	3
2	6
4	3
3	8
8	8
6	5
9	9
11	8
11	11
12	10

pt

Y

label	distance
1	8.60
0	2.41
1	9.67
0	4.32
0	5.39
1	1.01
1	7.53
1	8.60
1	10.42
0	3.25



Atelier Pratique AP-ML2 : Exercice ML2.6

6) Définir une fonction *myKNN* (*newpt*, *K*) qui retourne le label qu'elle va attribuer au nouveau point *pt* :

Démarche (cont):

b. Copier Y dans *Y_temp*, puis trier *Y_temp* sur la colonne *distance*, puis afficher le *Y_temp*

Y_temp

label	distance
1	1.01
0	2.41
0	3.25
0	4.32
0	5.39
0	6.46
1	7.53
1	8.60
1	9.67
1	10.74

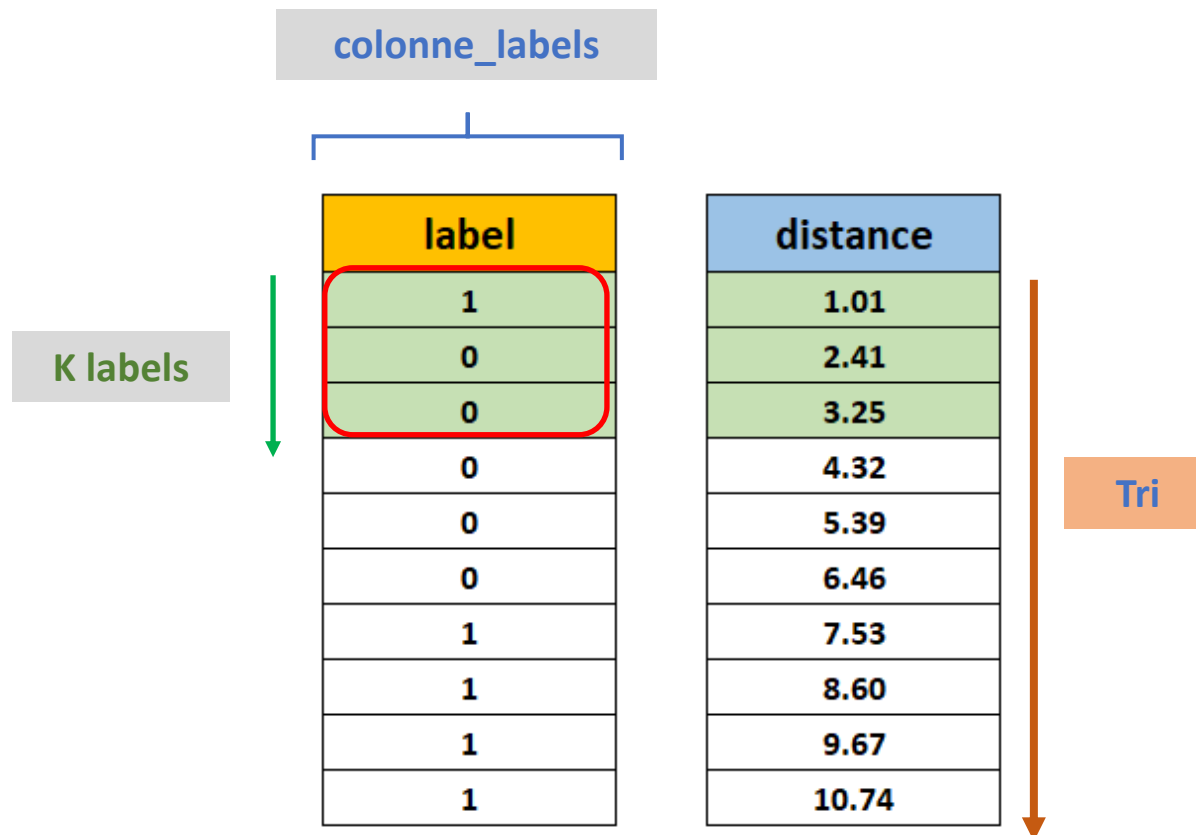
Tri

Atelier Pratique AP-ML2 : Exercice ML2.6

6) Définir une fonction *myKNN* (*pt*, *K*) qui retourne le label qu'elle va attribuer au nouveau point *pt* :

Démarche (cont):

c. Isoler la colonne 'label' avec uniquement les **K premiers labels** (en créant une série *colonne_labels*)



Atelier Pratique AP-ML2 : Exercice ML2.6

Démarche : (cont)

d. Sur ces K premiers labels, calculer le **nombre d'occurrence** de chaque label (*value_counts()*)

Le résultat du *value_counts()* est une série qu'on appellera *lab_count*

e. Le label recherché correspond à l'index du plus grand des nombres d'occurrences.

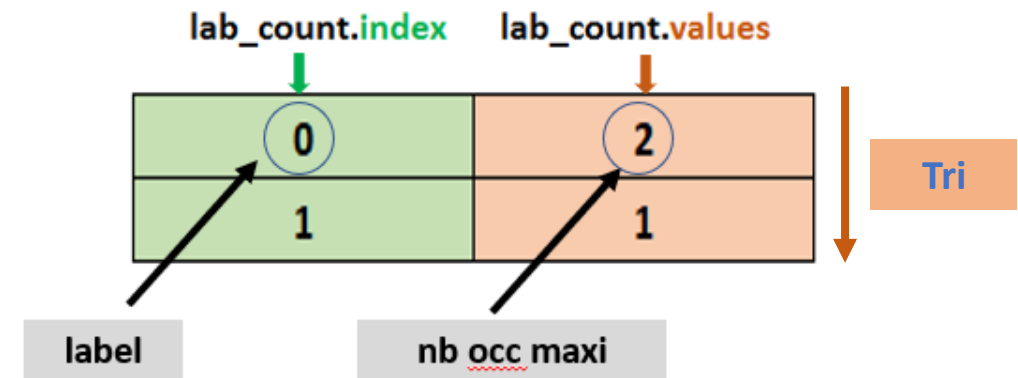
Pour y accéder :

. Méthode-1:

- utiliser `argmax()`

. Méthode-2:

- trier *lab_count* dans l'ordre décroissant
- récupérer le label , c'est *l'index* de la ligne 0



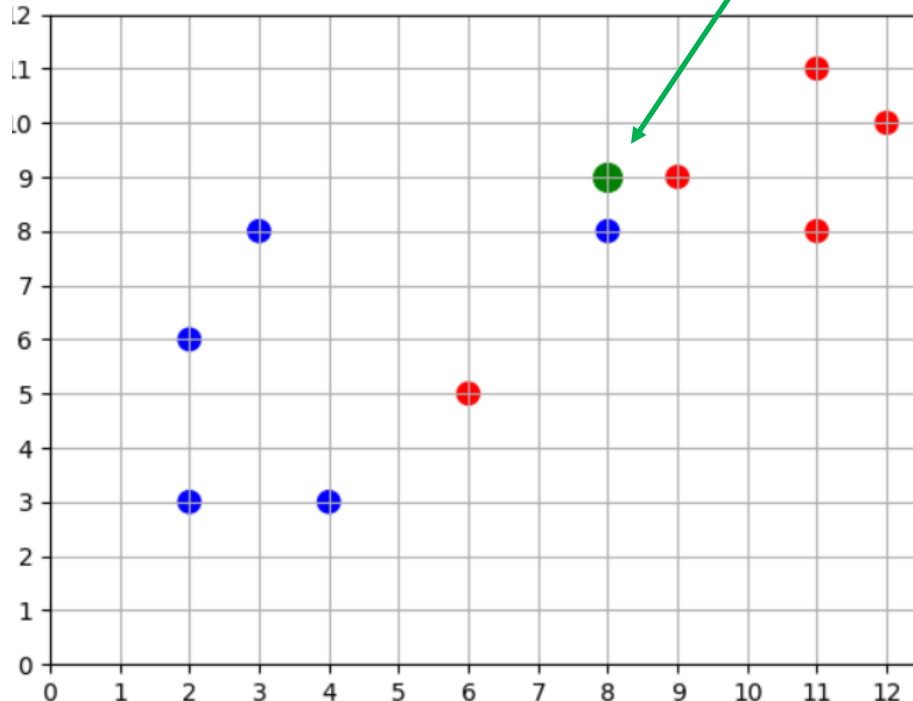
Atelier Pratique AP-ML2 : Exercice ML2.6

Démarche : (cont)

7) Faire des prédictions en entrant des points (pt) avec des coordonnées x1 et x2

- donner les features (x1,x2) d'un nouveau point pt
- afficher le label attribué au nouveau point pt

8) Afficher le Dataset + le nouveau point en couleur = 'green'



SOLUTION

Atelier Pratique AP-ML2 : Exercice ML2.7

Objectif : Tester l'algorithme *MyKNN* avec le dataset des images (chiffres manuscrits)

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	←	Classe 0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	←	Classe 1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	←	Classe 2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	←	Classe 3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	←	Classe 4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	←	Classe 5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	←	Classe 6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	←	Classe 7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	←	Classe 8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	←	Classe 9



SOLUTION