

TD : Algorithme des K plus proches voisins

Spé NSI - Lycée du parc

Année 2020-2021

Introduction

On se propose ici d'écrire un algorithme dérivé de l'algorithme de tri par insertion dont le but est d'obtenir les k plus petits éléments d'une liste. On verra ensuite comment adapter et utiliser cet algorithme pour résoudre un problème de classification.

I Les k plus petits éléments d'une liste

Question 1

Compléter la fonction suivante après avoir bien lu sa spécification dans la docstring.

Si vous n'y arrivez pas, vous pouvez ouvrir votre cours au chapitre qui traite des tris et consulter le paragraphe sur le tri par insertion.

```
1 def insertion_croissante(L, x):  
2     """  
3     L est une liste de flottants qui est triée par ordre croissant  
4     et x est un flottant.  
5     Ajoute l'élément x à la liste L en l'insérant à sa place.  
6     """
```


Question 2

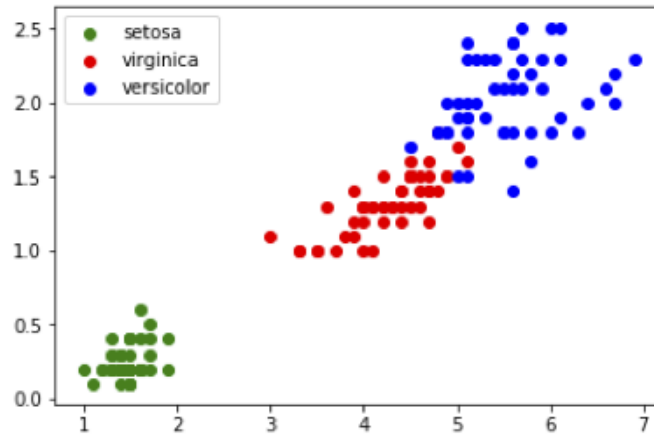
Pour déterminer les k plus petits éléments de la liste L, on procède en deux temps :

On range un à un les k premiers éléments de L dans une nouvelle liste `plus_petits` (ce qui revient à faire un tri par insertion de cette partie de L).

On parcourt le reste de la liste et pour chaque élément x, on le compare avec le plus grand de la liste `plus_petits`. Si x est plus petit, on retire le dernier élément de `plus_petits` et on y insère x à l'aide de `insertion_croissante`.

- la largeur des pétales (en cm)
- l'espèce d'iris : Iris setosa, Iris virginica ou Iris versicolor

On ne peut pas représenter simultanément quatre données numériques dans le plan mais, si on se restreint par exemple à la longueur des pétales en abscisse et la largeur en ordonnée, on obtiendrait l'ensemble des points suivants :



Les données sont incluses dans le fichier présent sur le réseau : C'est la définition de la liste `BDIris`

(base de données) comportant 150 éléments qui sont des tuples de 5 éléments sous la forme :
`(long_sep, larg_sep, long_pet, larg_pet, esp)`

où les quatre premiers éléments sont des flottants qui correspondent aux mesures et le cinquième est un entier entre 0 et 2 qui code l'espèce (0 pour Iris setosa, 1 pour Iris virginica et 2 pour Iris versicolor)

On va utiliser ces données pour essayer de déterminer à quelle espèce appartient une fleur connaissant les quatre mesures (longueur et largeur des sépales et pétales). On considère une fleur dont on a les mesures mais pas l'espèce (code -1 pour l'espèce) :

`nouv_fleur = (long_sep, larg_sep, long_pet, larg_pet, -1)`

Question 3

Pour mesurer à quel point cette fleur ressemble à chacune des 150 de la base de données, on va utiliser une distance. Parmi les k plus proches de cette fleurs (c'est à dire les k plus petites distances à `nouv_fleur`), on va ensuite déterminer quelle est l'espèce la plus fréquente. La fonction suivante calcule la distance entre deux fleurs.

Compléter la fonction suivante qui calcule la distance entre deux fleurs. n utilisera la distance euclidienne :

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2 + (t_2 - t_1)^2}$$

```

1 from math import sqrt
2 def distance(fleur_1, fleur_2):
3     """Renvoie la distance entre les deux fleurs."""
4
5
6
7
8

```

Question 4

Adapter la fonction `insertion_croissante` pour en déduire la fonction `insertion_fleur` en utilisant la fonction `distance` pour comparer les distances de deux fleurs avec `nouv_fleur`.

```
1 def insertion_fleur(nouv_fleur, L, fleur):
2     """
3     nouv_fleur est la fleur de référence. L est une liste de fleurs
4     qui est triée par ordre croissant des distances à nouv_fleur.
5     Ajoute fleur à la liste L en l'insérant à sa place pour la
6     distance à nouv_fleur.
7     """
8
9
```

Question 5

Adapter la fonction `k_plus_petits` pour en déduire la fonction `k_NN` (de l'anglais k-nearest neighbors) à l'aide de la fonction `insertion_fleur`. On utilisera la liste `BDIris` comme une variable globale.

```
1 def k_NN(nouv_fleur, k):
2     """
3     Renvoie la liste de k fleurs les plus proches de nouv_fleur.
4     """
5
6     assert type(k) == int and 1 <= k <= 150,
7     "k n'est pas entier entre 1 et 150"
8
9     plus_proches = []
10
11     #
12     # Partie à compléter
13     #
14
15     return plus_proches
16
```

Question 6

Lire et compléter les codes des fonctions pour en déduire une fonction qui renvoie le code de l'espèce que la méthode des k plus proches voisins attribue à `nouv_fleur`.

```
1 def indice_max(L):
2     """
3     L est une liste non vide d'entiers. Renvoie l'indice de la
4     première occurrence du maximum de la liste.
5     """
6
7     #
8     # Partie à compléter.
9     #
10
11
12 def prediction_espece(nouv_fleur, k):
13     """
14     Renvoie le code de l'espèce que la méthode des k plus proches
15     voisins attribue à nouv_fleur d'après les données de BDIris.
16     Dans le cas d'égalité, n'importe laquelle des réponses possibles
17     est bonne.
18     """
19
20     plus_proches = k_NN(nouv_fleur, k)
21     compteurs = [0,0,0] # compteurs[i] va compter les
22     # fleurs d'espèce codée par i.
23
24     #
25     # Partie à compléter.
26     #
27
28
29     return indice_max(compteurs)
```

