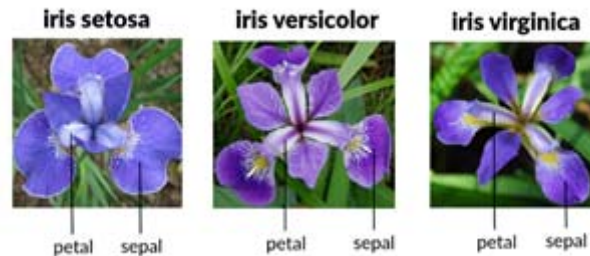


Exercice guidé sur la classification d'iris - Apprentissage automatique non supervisé

Algorithme des k moyennes (k-means)

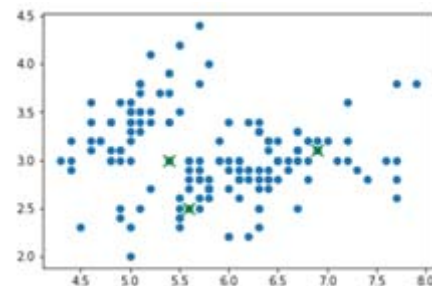
Pour mettre en place une classification à l'aide de l'algorithme des k-moyennes. Nous allons repartir du jeu de données des iris de Fischer. Mais cette fois-ci nous ne tiendrons pas compte de la classe définie dans la liste target...puisque nous cherchons à la définir.



Dans cet exercice, nous ne travaillerons que sur les 2 1ères caractéristiques (sur les 4).

sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
5.1	3.5	1.4	0.2
4.9	3.0	1.4	0.2
4.7	3.2	1.3	0.2
4.6	3.1	1.5	0.2
5.0	3.6	1.4	0.2

Le script suivant permet de tracer un nuage de points et de déterminer 3 centres initiaux aléatoires :



Analyser le code fourni puis définir et tester les fonctions suivantes :

Entrée[20]:



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import ListedColormap
4 from sklearn.datasets import load_iris
5 import random
6
7 #Chargement du jeu de données
8 iris = load_iris()
9 X= iris.data
10 X = X[ : , :2 ] # extraction des 2 1eres colonnes de la matrice des iris => uniquement les 2 1eres caracteristiques
11
12 # détermination aléatoire des centres parmi l'ensemble des données fournies
13 k= 3
14 lesCentres = []
15 for i in range(k):
16     ind = random.randint(0,len(X)-1)
17     lesCentres.append(X[ind].tolist())
18
19
20 # affichage de la situation initiale
21 plt.scatter([x[0] for x in X], [x[1] for x in X]) # affichage des points pour les iris
22 plt.scatter([x[0] for x in lesCentres], [x[1] for x in lesCentres], marker='x', s=100, c='g') #affichage des centres x vert
23 plt.title("k-means iris de Fischer")
24 plt.show()
25 #plt.closeAll()
26
27 """ Code à utiliser lorsque la fonction kmeans sera codée pour la tester
28 plt.close('all')
29 classes = kmeans(X, k, lesCentres) # test de kmeans
30 cmap = ListedColormap(['g', 'r', 'b', 'purple'])
31 centres = calculer_centres( classes, k )
32
33 plt.scatter([x[0] for x in X], [x[1] for x in X], c=[plus_proche(x, centres) for x in X], cmap=cmap)
34 plt.scatter([x[0] for x in centres], [x[1] for x in centres], marker='x', s=100, c=range(k), cmap=cmap)
35 plt.title("Inertie = " + str(inertie( classes, centres,k )))
36 plt.show()
37 """
```

Sortie[20]: ' Code à utiliser lorsque la fonction kmeans sera codée pour la tester\nplt.close('\all')\nclassees = kmeans(X, k, lesCentres)\ns) # test de kmeans\nncmap = ListedColormap(['g', 'r', 'b', 'purple'])\nncentres = calculer_centres(classees, k)\n\nplt.scatter([x[0] for x in X], [x[1] for x in X], c=[plus_proche(x, centres) for x in X], cmap=cmap)\nplt.scatter([x[0] for x in centres], [x[1] for x in centres], marker='x', s=100, c=range(k), cmap=cmap)\nplt.title("Inertie = " + str(inertie(classees, centres,k)))\nplt.show()\n'

1- Définir la fonction `dist(x, y)` renvoyant la distance euclidienne de 2 vecteurs x, y passés en paramètres.

```
Entrée[21]: ▶ 1 def dist( x, y ):
2     """ dist(x : list, y : list)-> float
3         entrees : x, y, 2 listes (vecteurs) de flottants dont on veut calculer la distance euclidienne
4         sortie : s, flottant, la distance entre les vecteur x et y
5     """
6     d = 0.
7     for i in range( len(x) ):
8         d+= (x[i] - y[i])**2
9     return d**.5
```


2- Écrire une fonction `plus_proche(x, lesCentres)` renvoyant le numéro i (indice) de la classe la plus proche de x parmi lesCentres, c'est-à-dire la classe telle que la distance de x à lesCentres[i] soit minimale.

Cette fonction fait appel à la fonction `dist` .

```
Entrée[22]: ▶ 1 #Q2
2 def plus_proche( x , lesCentres ):
3     """ plus_proche( x : list , lesCentres : list )-> int
4         entrees : x, liste, vecteur dont on veut connaitre le centre le plus proche
5                 : lesCentres, liste des centres
6         sortie : imin, entier, indice du centre le plus proche
7     """
8     imin = 0
9     for i in range(len(lesCentres)):
10         if dist( x, lesCentres[i] ) < dist( x, lesCentres[imin] ):
11             imin = i
12     return imin
```

3- Écrire une fonction `calculer_classes(X , lesCentres , k)` renvoyant une liste de classes telle que classes[i] soit la liste des données de X dont le centre le plus proche est lesCentres[i].

Cette fonction fait appel à la fonction `plus_proche` .

Entrée[23]: 

```


1 #Q3
2 def calculer_classes( X, lesCentres, k ):
3     """ calculer_classes( X:list, centres:list , k:int)->list
4         entrees : X, liste, liste des entrées
5             : lesCentres, liste des centres
6             : k, entier, nombre de classes
7         sortie : classes, liste de listes des points de chaque classe
8     """
9     classes = [ [] for i in range(k) ]
10    for x in X:
11        indPlusProche = plus_proche(x, lesCentres) # indice du centre le plus proche
12        classes[ indPlusProche ].append(x) # ajout de x à la classe possédant cet indice
13    return classes

```

4- Le centre (ou isobarycentre ou centroïde) d'un ensemble de vecteurs x_1, \dots, x_n est défini par le vecteur :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Écrire une fonction `centre(X)` renvoyant le centre de la liste de vecteurs X, passée en paramètre

Entrée[24]: 

```

1 #Q4
2 def centre( X ):
3     """ centre( X )
4         determination du barycentre d'un ensemble X de vecteurs
5         entree : X, liste, liste des donnees
6         sortie : c, liste(vecteur) du barycentre de X
7     """
8     dim=0
9     if len(X) != 0:
10        dim = len(X[0])
11        c = [0.]*dim
12        if len(X) == 0 :
13            return c
14        for x in X:
15            for i in range(len(x)): # pour chaque caracteristique
16                c[i] = c[i] + x[i]
17
18        for i in range(len(c)):
19            c[i] = round( c[i] /len(X) , 3) #arrondi à 3 decimales
20        return c
21

```

5- Écrire une fonction `calculer_centres(classes,k)` renvoyant la liste des centres de chaque classe. Cette fonction fait appel à la fonction

centre .

```
Entrée[25]: ▶ 1 #Q5
2 def calculer_centres( classes, k ):
3     """ calculer_centres(classes : lst , k : int) -> list
4         entrees : classes, liste, représentant les classes
5             : k, entier, nombre de classes à générer
6         sortie : centres, liste (de listes) des centres générés
7     """
8     centres = []
9     for classe in classes:
10         centrei = centre( classe )
11         centres.append( centrei )
12     return centres
13
```

6- Écrire une fonction `kmeans(X, k, lesCentres)` appliquant l'algorithme des k-moyennes à X en partant des centres et renvoyant la liste des classes obtenues.

Cette fonction fait appel à `calculer_classes` et `calculer_centres` .

```

Entrée[1]: ▶ 1 def affiche( X, classes,k):
2     cmap = ListedColormap(['g', 'r', 'b', 'purple'])
3     centres = calculer_centres( classes, k )
4     plt.close()
5     plt.scatter([x[0] for x in X], [x[1] for x in X], c=[plus_proche(x, centres) for x in X], cmap=cmap)
6     plt.scatter([x[0] for x in centres], [x[1] for x in centres], marker='x', s=100, c=range(k), cmap=cmap)
7     plt.title("Inertie = " + str(inertie( classes, centres,k )))
8     plt.show()
9
10 #Q6
11 def kmeans(X, k, lesCentres):
12     """ kmeans(X, k, lesCentres)->list
13         entrees : X, liste, liste des entrées
14                 : lesCentres, liste des centres
15                 : k, entier, nombre de classes
16         sortie : classes, liste des points de chaque classe apres traitement du kmeans
17     """
18     centres2 = []
19     while ( (sorted(lesCentres)) != (sorted(centres2)) ) : #tri sur les listes pour pouvoir utiliser ==
20     # while not ( np.array( sorted(lesCentres)==sorted(centres2) ).all() ):
21         centres2 = lesCentres
22         classes = calculer_classes( X, centres2, k )
23         # affiche(X,classes,k)
24         # plt.show()
25         lesCentres = calculer_centres( classes, k)
26
27     return classes

```

7- Créer une fonction `inertie(classes, centres, k)` qui calcule l'inertie

```

Entrée[35]: ▶ 1 def inertie( classes, centres, k ):
2     s = 0.
3     for i in range(k):
4         for x in classes[i]:
5             s += dist(x, centres[i])**2
6

```

Entrée[36]:



```
1
2 #_____ Test _____
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from matplotlib.colors import ListedColormap
6
7 from sklearn.datasets import load_iris
8 import numpy as np
9 import random
10 #Chargement du jeu de données
11 iris = load_iris()
12
13 k= 3
14 X= iris.data
15 X = X[ : , :2 ] # extraction des 2 1eres colonnes de la matrice des iris => uniquement les 2 1eres caracteristiques
16
17 # détermination aléatoire des centres
18 lesCentres = []
19 for i in range(k):
20     ind = random.randint(0,len(X)-1)
21     lesCentres.append(X[ind].tolist())
22
23
24 # affichage
25 plt.scatter([x[0] for x in X], [x[1] for x in X]) # affichage des points pour les iris
26 plt.scatter([x[0] for x in lesCentres], [x[1] for x in lesCentres], marker='x', s=100, c='g') #affichage des centres
27 plt.show()
28 #plt.closeAll()
29
30
31 classes = kmeans(X, k, lesCentres) # test de kmeans
32 cmap = ListedColormap(['g', 'r', 'b', 'purple'])
33 centres = calculer_centres( classes, k )
34 plt.close()
35 plt.scatter([x[0] for x in X], [x[1] for x in X], c=[plus_proche(x, centres) for x in X], cmap=cmap)
36 plt.scatter([x[0] for x in centres], [x[1] for x in centres], marker='x', s=100, c=range(k), cmap=cmap)
37 plt.title("Inertie = " + str(inertie( classes, centres,k )))
38 plt.show()
39
40
```

