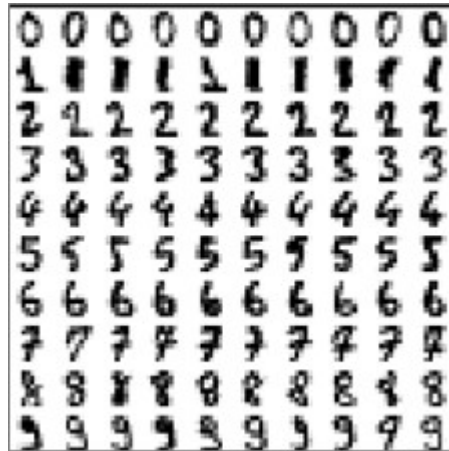


## Exercice : Reconnaissance des chiffres

Dans cet exercice, nous allons à nouveau utiliser Scikit-learn qui est une bibliothèque libre Python destinée à l'apprentissage automatique, et notamment le jeu de données `digits` : nous disposons de 1 797 images de  $8 \times 8$  pixels en niveaux de gris représentant les chiffres de 0 à 9.



Ces données nous sont fournies par l'intermédiaire de 2 tableaux :

- X est une matrice de :

- o 1797 lignes pour les images,

- o 64 colonnes pour les caractéristiques des images.

Pour tout  $k \in [0, 1796]$ ,  $X[k]$  est un vecteur de  $\mathbb{R}^{64}$  qui représente l'image digitalisée d'un chiffre.

- Y est un vecteur de taille 1 797.

Pour tout  $k \in [0, 1796]$ ,  $Y[k]$  correspond à l'étiquette de  $X[k]$ , ie l'entier compris entre 0 et 9 représenté par  $X[k]$ .

Le fichier `IA_Chiffres_Eleves.py` est à votre disposition : il comporte des fonctions permettant d'afficher les chiffres, ainsi que le squelette des fonctions que vous aurez à développer

Entrée[5]: ▶

```
1 from math import *
2 from random import *
3 import matplotlib.pyplot as plt
4 import numpy as np
5 from copy import deepcopy
6
7 from sklearn.datasets import load_digits
8
9 # Chargement du jeu de données des chiffres
10 chiffres=load_digits()
11
12 X=chiffres.data
13 Y=chiffres.target # recuperation des classes
14
15 # Affichage des chiffres
16 def afficherChiffres( X, Y, limit_max=10 ):
17     """ afficherChiffres( X: list, Y: list, limit_max=10 : int ):
18         entrees : X, liste du jeu de données
19             : Y, liste des etiquettes associees
20             : limit_max, entier, nombre de caractères affichés de chaque t
21     """
22     classes, nombres = np.unique( Y, return_counts = True )
23     nombre_max = min( np.max(nombres), limit_max )
24     img = np.zeros( ( 100, nombre_max*10 ) )
25     for i in range( 10 ) :
26         index_classe = np.where(Y == i)[0][:limit_max]
27         for j, echantillon in enumerate( index_classe ):
28             img[i*10+1:i*10+9,j*10+1:j*10+9] = X[echantillon].reshape((8, 8))
29     plt.imshow( img, cmap='binary' )
30     plt.xticks([])
31     plt.yticks( 5 + 10*np.arange(10), np.arange(10) )
32     plt.show()
33 # tests affichage
34 # afficherChiffres( X, Y)
35 # afficherChiffres( X, Y, 14)
36
37 def obtenir_echantillons( X , Y, nbEchantillons = 10 ):
38     '''Donne une sélection aléatoire d indices d'échantillons pour chaque classe
39     index = []
40     for classe in np.unique(Y):
41         index_classe = np.where(Y == classe)[0]
42         replace = len(index_classe) > nbEchantillons
43         index += list( np.random.choice(index_classe, size=nbEchantillons, repla
44     return index
45
46 index = obtenir_echantillons( X, Y )
47 afficherChiffres( X[index], Y[index] )
48 plt.show()
```

## Apprentissage supervisé

On envisage l'algorithme des k plus proches voisins.

1. Écrire la fonction `distance( im1, im2 )` qui renvoie la distance euclidienne entre les images `im1` et `im2`.

Entrée[6]: ▶

```
1 def distance( im1 , im2 ):
2     d=0
3     for i in range( len( im1 ) ):
4         d = d + ( im2[i] - im1[i] )**2
5     d = np.sqrt( d )
6     return d
7
8
```

2. Écrire la fonction `PlusProchesVoisins(X, Y, im, k)` qui renvoie la ou les étiquettes les plus présentes parmi les `k` voisins de `im` (image de  $8 \times 8$  pixels représentée sous la forme d'une liste).

Entrée[7]: ▶

```

1 def PlusProchesVoisins( X , Y , im , k ):
2     # calcul des distances entre chaque image du jeu de données et l'image im
3     lesDistances=[]
4     nbClasses = 10
5     for i in range( len(X) ):
6         d = distance( X[i] , im )
7         lesDistances.append( [d,Y[i]] )
8
9     # tri par ordre croissant de distance
10    lesDistances.sort()
11
12    # on cherche la classe majoritaire parmi les voisins les plus présents :
13    # on compte le nombre de voisins plus proches dans chaque catégorie
14    voisins=[0]*nbClasses
15    for i in range( k ): #on regarde les k plus proches voisins de im
16        voisins[ lesDistances[i][1] ]+=1
17    # on cherche la ou les catégories plus représentées
18    rep = [] # liste des chiffres les plus représentés
19    maxi = 0 # compte le maximum
20    for i in range( nbClasses ):
21        if voisins[ i ] > maxi:
22            maxi = voisins[ i ]
23            rep = [ i ]
24        elif voisins[ i ]== maxi :
25            maxi = voisins[ i ]
26            rep.append( i )
27    return rep

```

3. Proposer des instructions pour partager les données en 2 groupes : un groupe pour l'apprentissage et l'autre pour le test.

Entrée[8]: ▶

```

1 # Séparation des données d'apprentissage et de tests
2 AX,AY=[],[]
3 TX,TY=[],[]
4 dA={} # dictionnaire des clés déjà choisies comme donnée d'apprentissage
5 dT={}
6 i=0 # compteur
7 while i<= len( X )-1:
8     j=randint( 0, len( X )-1 )
9     if j not in dA:
10        AX.append( X[j] )
11        AY.append( Y[j] )
12        dA[j]=1
13        i=i+1
14    elif j not in dT:
15        TX.append( X[j] )
16        TY.append( Y[j] )
17        dT[ j ]=1
18        i=i+1
19 # après separation puis echantillonnage
20 index = obtenir_echantillons(AX, AY)
21 #afficherChiffres(AX, AY)
22 afficherChiffres( X[index], Y[index] )
23
24
25
26

```

4. Écrire des instructions pour calculer le taux de bonnes prédictions sur les données de tests.

```

Entrée[9]: ▶ 1 # Précision
2 nb = 0 # compteur des bonnes prédictions
3 for i in range(len(TX)):
4     prediction=PlusProchesVoisins( AX, AY, TX[i], 3)
5     if len( prediction ) == 1: # pas d'indécision
6         if prediction[0] == TY[i]:
7             nb = nb + 1
8     # print( "Prédiction : ", prediction, " Réponse : ", TY[i] )

```

## Apprentissage non supervisé

On s'intéresse à nouveau à l'ensemble des images de chiffres de 0 à 9, mais cette fois-ci sans prendre en compte leur étiquette.

Nous allons appliquer l'algorithme des k moyennes avec  $k = 10$ , de manière à obtenir un regroupement de ces images en 10 classes, que l'on espère correspondre aux différents chiffres représentés. Pour rappel, ces images sont regroupées dans une liste X. Chaque image  $X[i]$  est représentée par un vecteur de R64.

5. Expliquer l'algorithme des k-moyennes.

Entrée[ ]: ▶

6. Écrire une fonction `barycentre( s )` qui calcule le barycentre d'un ensemble d'images représenté par la liste s. Ce dernier est un vecteur de R64 (il représente lui aussi une image  $8 \times 8$ )

```

Entrée[10]: ▶ 1 def barycentre(s):
2     bary=[0]*64
3     for k in range(len(s)):
4         for i in range(64):
5             bary[i]=bary[i]+s[k][i]#/len(s)
6             bary[i] = bary[i] / len( s )
7     return bary
8
9 def dist( X1, X2 ):
10     d=0
11     for i in range(len(X1)):
12         d=d+(X1[i]-X2[i])**2
13     d= np.sqrt(d)
14     return d
15
16

```

7. Écrire la fonction `kmoyennes(X,k)` qui prend en argument une liste X d'images et un entier k et renvoie une liste de k classes où sont regroupées les images. Pour cela on choisit aléatoirement k centres initiaux `lesCentres`. Pour chaque image  $X_i$ , on calcule la valeur de `lesCentres[j]` la plus proche de  $X_i$  pour ranger ce dernier dans la classe. On calcule les barycentres de ces classes. On s'arrête quand ces derniers ne sont pas modifiés.

Entrée[11]: ▶

```
1 def kmoyennes( X, k ):
2     # choix au hasard de k images
3     changement = True
4     lesCentres=np.array([ X[randint(0,len(X)-1)] for i in range(k)] )
5
6     while changement: # on boucle tant que les centres bougent
7         lesClasses=[ [] for i in range(k) ] # k classes
8         for i in range(len(X)):
9             # recherche du centre le plus proche de Xi
10            dmin = np.inf
11            jmin = -1
12            for j in range(k):
13                d = dist(X[i],lesCentres[j])
14                if d < dmin:
15                    dmin = d
16                    jmin = j
17
18            # on ajoute Xi à la classe lesClasses[jmin]
19            lesClasses[ jmin ].append( X[i] )
20
21            # calcul des nouveaux barycentres pour chaque classe
22            nouveauxCentres=np.array([barycentre(lesClasses[i]) for i in range(k)])
23            if nouveauxCentres.all() == lesCentres.all():
24                changement = False
25            else :
26                lesCentres=deepcopy( nouveauxCentres )
27        return lesClasses , lesCentres
28
29 lesClasses , lesCentres = kmoyennes(X,10)
30 index=[i for i in range(10)]
31
32 afficherChiffres( lesCentres, Y[index] )
33 plt.show()
```

Entrée[ ]: ▶