## TP4: Clustering de données Le but de ce TP est de créer une fonction Python capable de détecter des clusters de données homogènes dans un ensemble de données, puis d'analyser un jeu de données réelles. A) K-Moyennes:

1. Ecrivez en python l'algorithme des K-Moyennes sous la forme d'une fonction. Vous trouverez une description de l'algorithme dans le cours ou sur internet. La fonction prend en entrée deux paramètres : la matrice des données et le nombre de clusters que l'on la fonction Kmeans de sklearn.

souhaite. Testez sur les données Iris ou sur des données que vous générez. Comparez avec import numpy as np import matplotlib.pyplot as plt from sklearn import datasets, metrics from sklearn.metrics.pairwise import euclidean distances iris = datasets.load iris() X = iris.data Y = iris.target

In [40]: import random **def** dist(x,y): return [[ np.sqrt(np.sum((xi-yi)\*\*2)) for yi in y] for xi in x] def K Moyennes(X, K): try: k centroids = random.choices(X, k=K) cluster\_labels = []

while True: # distance = euclidean\_distances(X, k\_centroids) distance = dist(X, k\_centroids) cluster\_labels = [np.argmin(dist) for dist in distance] distance matrix = [ [ x for i, x in enumerate(X) if cluster labels[i] == k ] for k in range(K) ] for X\_K in distance\_matrix: if not X\_K : raise Exception("\nError") new\_k\_centroids = [ np.mean( a = X\_K, axis = 0 ) for X\_K in distance\_matrix ] if np.array\_equal(new\_k\_centroids, k\_centroids): break else : k\_centroids = new\_k\_centroids return cluster\_labels except: return K\_Moyennes(X, K) K\_Moyennes(X, 3)

0, 0, 0, 0, 0, Ο, Ο, 0, Ο, 0, 0, 0. 0, 0, 0, 0,

Out[40]: [0, 0, 0, 0, 0, 0, Ο, 0, 0, 0, Ο, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, Ο, 0, 0, 0, Ο, 0, 0, 0, Ο, 1, 2, 1, 2, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,

2, 2, 2, 1, 1, 2, 1, 1, 1, 2, 2, 1, 2, 1, 2, 1, 2, 1, 1, 1, 2, 1, 1, 2, 1, 2, 1, 1, 2, 1, 1, 2] In [41]: from sklearn.cluster import KMeans kmeans = KMeans(n clusters=3) kmeans.fit\_predict(X) 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 2, 2, 2, 2, 1, 2, 2, 2, 2, 2, 2, 1, 1, 2, 2, 2, 1, 2, 1, 2, 1, 2, 2, 1, 1, 2, 2, 2, 2, 2, 2, 1, 2, 2, 2, 1, 2, 2, 1, 2, 2, 1, 2, 2, 1]) 2. Expérimenter l'instabilité due à l'initialisation : les centres des clusters étant choisis au hasard lors de l'initialisation, le résultat obtenu peut varier d'une exécution à l'autre. Vérifiez que c'est le cas. In [42]: plt.scatter(X[:,2],X[:,3], c=K\_Moyennes(X, 3), cmap='rainbow')

2,

2.5

2.0

1.5

1.0

2.0

1.5

1.0

plt.scatter(X[:,2],X[:,3], c=K\_Moyennes(X, 3), cmap='rainbow')

3. Utiliser l'indice de Silhouette (qui est dans le package sklearn) pour stabiliser les résultats et sélectionner automatiquement le nombre de groupes. Pour ce faire, créez un script qui applique K-moyenne sur les données pour différents nombres de clusters allant de 2 à 10, 10 fois pour chaque nombre de clusters (soit 90 fois en tout) et qui renvoie la solution ayant

print('\n Le meilleur score de Silhouette pour 90 iteration est : ', np.max(scores), "Avec k = ", int(np.arg

3

2

1

0

 $^{-1}$ 

-2

-3

Le LDA et le PCA sont des techniques de transformation linéaire: le LDA est un supervisé tandis que le PCA non supervisé - le PCA ignore les étiquettes de classe. Nous pouvons imaginer l'ACP comme une technique qui trouve les directions de la variance maximale.

Les données choixprojetstab.csv contiennent des données anonymisées de voeux faits par les étudiant·e·s du master 2 pour des projets

jamais. Les projets non mentionnés peuvent être considérés comme « moyens » et on peut dire que très bien = 3, bien = 2, moyens = 1,

remplir deux variables : la liste des codes « C » représentant les étudiant (première colonne)

matrice M doit être de type array du package numpy. Faites attentions à ce que les valeurs dans M soient bien numériques (1, 2, 3) et non textuelle ('1', '2', '3'). Vous pouvez utiliser la

étudiant∙e ga.vTZVmBFaC. ga/mLSm4ai/6g ga04b5zeP48qY ga1ohlKbe4v9Y ga2f3zAu/j5w6 ga2tObQKD38MQ ga5EltwHBOBQU

2. Dans sklearn.cluster il existe différents algorithmes de clustering. Testez les différents algorithmes du package et proposez le meilleur clustering possible des données selon

with : Kmeans

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\cluster\ affinity propagation.py:152: FutureWarning: 'ran dom state' has been introduced in 0.23. It will be set to None starting from 0.25 which means that results w ill differ at every function call. Set 'random state' to None to silence this warning, or to 0 to keep the b

with: Birch

with : MiniBatchKMeans

with : AffinityPropagation

with : GaussianMixture

3

0

1

1

0

0

à réaliser en binômes. Chaque étudiant·e a fait des voeux sur les differents projets en leur donnant une mention très bien, bien ou

1. Utilisez le package csv (ou l'importation de variable de Spyder) pour lire le fichier et

et la matrice « M » des données (tout sauf la première ligne et la première colonne). La

(L'ACP a tendance à entraîner de meilleurs résultats de classification dans une tâche de reconnaissance d'image si le nombre d'échantillons pour une classe donnée était relativement faible.) Contrairement à PCA, LDA tente de trouver un sous-espace de fonctionnalités qui maximise la séparabilité des classes. LDA fait des hypothèses sur les classes normalement distribuées et les

4. Utiliser une ACP (fonction PCA de sklearn) pour vérifier visuellement la cohérence des groupes obtenus. Vérifier aussi visuellement la séparabilité et la compacité de ces groupes à

| 9/9 [00:11<00:00,

LDA visualisation

<matplotlib.collections.PathCollection at 0x1ad04f7c348>

le meilleur score de Silhouette.

from sklearn.decomposition import PCA

fig, (ax1, ax2) = plt.subplots(1, 2)

cluster labels = K Moyennes (X, 3)

plt.title('PCA visualisation')

plt.title('LDA visualisation')

fig.set size inches(15, 5)

PCA = PCA (n components=2)dataPCA = PCA.fit transform(X)

LDA = LDA (n components=2)

plt.subplot(1, 2, 1)

plt.subplot(1, 2, 2)

covariances de classe égales.

import csv

(71, 80)

df.head()

import numpy as np

# print(C)

[[1 1 1 ... 1 1 3] [1 1 1 ... 1 1 1] [0 0 0 ... 1 0 0]

[1 1 1 ... 1 0 2] [1 1 0 ... 1 1 1] [1 1 1 ... 0 1 1]]

import pandas as pd

bl/.vSDYCGrSs

bl/1NiMubceBs

bl/dvgMTLVSvk

bl1ao5B7htJfQ

l'indice Silhouette.

import numpy as np

from sklearn import cluster

# Create cluster objects

clustering\_algorithms = (

ehavior of versions <0.23.

FutureWarning)

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js

Links

from sklearn import metrics, mixture

birch = cluster.Birch(n\_clusters= 4)

kmeans = cluster.KMeans(n\_clusters=4)

('Kmeans' , kmeans),

('GaussianMixture', gmm)

> silhouette\_score = 0.29958429496509337 > silhouette\_score = 0.29958429496509337

> silhouette\_score = 0.2905635020378965

> silhouette score = 0.1840557144219895

• E-mail: zakaria.abbou199434@gmail.com GitHub: github.com/ZakariaAABBOU

Linkedin: linkedin.com/in/zakaria-aabbou/

> silhouette\_score = 0.24603465972457622

('Birch', birch),

# 'très bien' = 3, 'bien' = 2, 'moyens' = 1, 'jamais' = 0

two\_means = cluster.MiniBatchKMeans(n\_clusters= 4)

('MiniBatchKMeans', two means),

for name, algorithm in clustering\_algorithms:

cluster\_labels = algorithm.fit\_predict(M)

3 bl1NWhKcNADF2

5 rows × 81 columns

C = np.array(rows[1:,0])

print(M.shape,"\n", M)

plt.show()

1.5

1.0

0.5

0.0

-0.5

-1.0

cluster labels = K Moyennes(X, k)

scores.append( metrics.silhouette score(X, cluster labels) )

l'aide d'une ADL (fonction LinearDiscriminantAnalysis de sklearn).

from sklearn.discriminant analysis import LinearDiscriminantAnalysis as LDA

plt.scatter(dataLDA[:, 0], dataLDA[:, 1], c=cluster labels, cmap='viridis')

plt.scatter(dataPCA[:,0], dataPCA[:,1], c=cluster labels)

dataLDA = LDA.fit(X, cluster labels).transform(X)

PCA visualisation

Quelle est la différence entre les deux méthodes ?

B) Analyse des données « choix projet » :

méthode astype de numpy en cas de besoin.

M = np.array(rows[1:,1:]).astype(int)

csv\_reader = csv.reader(csv\_file, delimiter=';') rows = np.asarray([row for row in csv reader])

jamais = 0. L'objectif est de former des clusters d'étudiant·e·s ayant faits des choix similaires. </P>

with open('choixprojetstab.csv', newline='', encoding='utf-8') as csv\_file:

df = pd.read\_csv('choixprojetstab.csv', delimiter=';', encoding='utf-8')

1

affinity\_propagation = cluster.AffinityPropagation(damping=.9, preference=-200)

gmm = mixture.GaussianMixture( n\_components= 2, covariance\_type='full')

('AffinityPropagation', affinity\_propagation),

score = metrics.silhouette\_score(M, cluster\_labels)

print("> silhouette\_score = ", score, "\t with : "+name)

Le meilleur score de Silhouette pour 90 iteration est : 0.6810461692117459 Avec k = 2

from tqdm import tqdm

for k in tqdm(range(2,11)):

for i in range(10):

scores = []

100%|

In [44]:

In [47]:

In [48]:

