

Analyse de données - TP 4

ISEP – 10 Novembre 2020

Instructions : Préparez un rapport incluant le code source et vos résultats, et déposez-le sur Moodle. Pas plus de 2 personnes par groupe. N'oubliez pas de mettre les 2 noms sur le rapport, ou de faire 2 rendus.

Librairies

Ce TP nécessite les librairies pandas, matplotlib, numpy, scipy et sklearn.

A Analyse des Iris de Fisher avec l'algorithme K-Means

1. Ouvrez les données iris du TP2, ou récupérez les données iris directement dans le module datasets de sklearn.
2. La dernière colonne de vos données représente le label correspondant à l'espèce d'Iris associée. Stockez ces labels dans un vecteur séparé et enlevez-le de votre jeu de données.
3. Utilisez la commande **PCA(.)** de sklearn.decomposition pour réaliser une analyse en composante principale de vos données. Puis, afin de récupérer vos données projetées dans 2 composantes, utilisez les commandes suivantes :

```
1  pca = PCA(n_components=2)
2  PCA_val = pca.fit_transform(X)
3  df_iris_PCA = pd.DataFrame(data = PCA_val, columns = ['PC1', 'PC2'])
4  df_iris_PCA_class = pd.concat([df_iris_PCA, Y], axis = 1)
```

Vous devriez maintenant avoir 2 jeux de données, le jeux de données original stocké dans une première variable ("*data*" par exemple), et le même jeu de données projeté sur 2 composantes stockée dans la variable *df_iris_PCA*.

4. En utilisant les 2 premières composantes, projetez les labels des 3 espèces d'iris en 2 dimensions. Que pouvez-vous dire sur la séparabilité des classes ? Quelle peuvent être les implications si on utilise un algorithme de clustering sur ces données ?
5. Utilisez l'algorithme K-Means de sklearn.cluster sur vos données originales afin d'obtenir 3 clusters et visualisez les résultats. Pour ce faire, vous pouvez utiliser le code ci-dessous. Expliquez les paramètres de la fonction KMeans.

```
1  #Kmeans code, change X by your dataset
2  kmeans = KMeans(n_clusters=3, n_init=5, max_iter=300, random_state=0).fit(X)
3  kmeans.score(X)
4  prediction = kmeans.predict(X)
```

6. Répétez la question 4) plusieurs fois avec des valeurs différentes de `random_state`. Que constatez-vous ? Expliquez.
7. Affichez la matrice de confusion comparant vos résultats aux labels théoriques.
8. Choisissez le clustering donnant la meilleure matrice de confusion et évaluez le avec l'indice de Silhouette associé à vos données (`silhouette_score()` dans `sklearn.metrics`). Commentez.

B Clustering hiérarchique

Cet exercice est un tutoriel sur l'utilisation du clustering hiérarchique avec Scipy.

1. Commencez par importer les packages nécessaires.

```
1 # needed imports
2 from matplotlib import pyplot as plt
3 from scipy.cluster.hierarchy import dendrogram, linkage, fcluster
4 from scipy.cluster.hierarchy import cophenet, inconsistent, maxRstat
5 from scipy.spatial.distance import pdist
6 import numpy as np
```

2. Nous allons maintenant générer un jeu de données aléatoire avec 2 classes à partir du code ci-dessous. Essayez de modifier les paramètres des fonctions pour créer *a* et *b* afin de voir à quoi ils correspondent.

```
1 # np.random.seed(0) # uncomment for repeatability
2 a = np.random.multivariate_normal([10, 0], [[3, 1], [1, 4]], size=
    =[100,])
3 b = np.random.multivariate_normal([0, 20], [[3, 1], [1, 4]], size=
    =[50,])
4 X = np.concatenate((a, b),)
5 plt.scatter(X[:,0], X[:,1])
6 plt.title('My data distribution')
7 plt.show()
```

3. Le code ci-après permet de choisir un linkage (ward dans cet exemple), de générer la matrice du dendrogramme et d'afficher celui-ci. Après avoir fixé la seed aléatoire dans le morceau de code précédent, essayez plusieurs linkages. Commentez.

```
1 # alternative linkage methods: 'single', 'complete', 'average', '
    euclidean' (default), 'cityblock' aka Manhattan, 'hamming', '
    cosine' ...
2
3 Z = linkage(X, 'ward', optimal_ordering=True)
4
5
6 c, coph_dists = cophenet(Z, pdist(X))
7 print('Cophenetic Correlation: %1.2f' % c)
8
9
10 plt.figure(figsize=(25, 10))
11 plt.title('Hierarchical Clustering Dendrogram (full)')
12 plt.xlabel('sample clusters')
13 plt.ylabel('distance')
14
15 dendrogram(Z, leaf_rotation=90., leaf_font_size=8.,)
16
17 plt.show()
18
19 # truncate the dendrogram for better visibility
20
21 R = inconsistent(Z)
```

4. Dans le code ci-dessous, que contient la variable Z . Affichez son contenu et commentez.
5. Cherchez la définition du coefficient cophénétique décrit dans le code précédent. A quoi sert-il et comment l'interpréter ?
6. Comme vous le voyez, le dendrogramme original est massif et difficile à interpréter. Il est possible d'obtenir une version réduite en utilisant le code ci-après. Dans ce code, expliquez le rôle des paramètres "truncate_mode" et "p".

```

1  #display truncated dendrogram
2  plt.title('Hierarchical Clustering Dendrogram (truncated)')
3  plt.xlabel('sample index')
4  plt.ylabel('distance')
5  dendrogram(
6      Z,
7      truncate_mode='lastp', #to explain
8      p=12, #to explain
9      show_leaf_counts=False,
10     leaf_rotation=90.,
11     leaf_font_size=12.,
12     show_contracted=True,
13 )
14 plt.show()

```

7. Enfin, on souhaite maintenant couper le dendrogramme afin d'obtenir les clusters. Deux méthodes possibles sont proposées dans le code ci-après. Expliquez en quoi elles diffèrent, les avantages de l'une ou de l'autre, et ce que font les différents paramètres. Essayez de modifier les paramètres pour voir ce qui change.

```

1  max_d = 14
2  clusters = fcluster(Z, max_d, criterion='distance')
3
4  plt.figure(figsize=(10, 8))
5  plt.scatter(X[:,0], X[:,1], c=clusters, cmap='prism')
6  plt.show()

```

OR

```

1  k=4
2  clusters = fcluster(Z,k,criterion='maxclust')
3
4  plt.figure(figsize=(10, 8))
5  plt.scatter(X[:,0], X[:,1], c=clusters, cmap='prism')
6  plt.show()

```

8. Refaite les questions 3 à 7 avec d'autres linkages et commentez sur les différences de résultats.

C Données Wisconsin Data Breast Cancer

En vous inspirant des exercices précédents, ouvrez les données WDBC, analysez les données et la séparabilité des classes. Faites tourner une ou plusieurs méthodes de clustering, et analysez la qualité de vos résultats. N'hésitez pas à faire les analyses classiques des attributs de ces données, et à faire des projections pour visualiser.

D Nombre optimal de clusters sur les données atmosphère d'exoplanète

Dans cet exercice, on va évaluer le nombre de clusters optimal dans un jeu de données artificiel décrivant des compositions atmosphérique d'exoplanètes.

1. Ouvrez le fichier *exo4_atm_extr.csv* et retirez la dernière colonne qui contient les labels.
2. Rappelez les propriétés de l'indice de Calinski-Harabasz et celles de l'indice de Davies-Bouldin.
3. Utilisez `sklearn.metrics` pour déterminer le nombre optimal de clusters pour l'algorithme des K-Means à partir des indices de Davies-Bouldin et Calinski-Harabasz. Expliquez avec vos mots ce que fait cette fonction.
4. Commentez vos résultats. Vous pourrez utiliser des graphes, des projections ACP et les labels pour appuyer vos explications.