



SIA 3611

Année universitaire 2022-2023

TP Machine Learning

TP 3 : Clustering

Vassilis.Christophides@ensea.fr
guillaume.renton@ensea.fr

Table des matières

1 Spatial datasets and first clusters	1
1.1 To do :	1
1.2 To do :	2
2 Spatial dataset normalization	3
2.1 To code :	3
2.2 To code :	3
2.3 To code :	3
2.4 Bonus :	3
3 Detection of caribbean island	4
3.1 To code :	4
3.2 To code :	4
4 Evaluating a cluster	4
4.1 To do :	4
4.2 To do :	5
4.3 To code :	6
4.4 Bonus :	6

Introduction

In machine learning, clustering is related to unsupervised learning approaches in which the algorithm fits from the distribution of given data. The main advantage of such methods is detection without a priori of patterns, sorting data and detecting outliers.

The given dataset is a combination of a dataset produced by the World Health Organization and the location of countries. It pooled the evolution of 23 features for 15 years among numerous countries. One of the goals of this TP3 is to detect continents and subcontinents.

Objectives :

- Visualize spatial datasets
- Apply K-Means and GMM on spatial datasets
- Normalize the datasets
- Propose a methodology to detect specific pools of data

1 Spatial datasets and first clusters

The first step is visualizing the mercator projection

1.1 To do :

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib.colors
5
6 df = pd.read_csv("data/Life_Expectancy_Data.csv")
7 df = df.dropna()
8 df.info()
9
10 df1 = df[(df.Year == 2013)]

```

```

1 df_X = df1[['Longitude', 'Latitude']]
2 df_Status = df1[['Continent']]
3
4 df_Y = df_Status.replace(['Africa', 'Asia', 'Europe', 'North America', 'South America', 'Seven
    seas', 'Oceania'], [0, 1, 2, 3, 4, 5, 6])
5 np_Y = df_Y.to_numpy()
6 np_Y = np_Y.reshape((np_Y.shape[0],))
7
8 np1 = df_X.to_numpy()
9 plt.scatter(np1[:,0], np1[:,1], c=np_Y, cmap=matplotlib.colors.ListedColormap(['red', 'green',
    'blue', 'purple', 'yellow', 'magenta', 'cyan']))
10 plt.show()

```

Question 1

Identify each class with the plot colors.

What do the coordinates correspond to?

1.2 To do :

Execute the following cell

```

1 from itertools import product
2 from sklearn import cluster
3
4 x_min, x_max = np1[:, 0].min() - 10, np1[:, 0].max() + 10
5 y_min, y_max = np1[:, 1].min() - 10, np1[:, 1].max() + 10
6 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
7                       np.arange(y_min, y_max, 0.1))
8
9 km2 = cluster.KMeans(n_clusters=2).fit(np1)
10 km3 = cluster.KMeans(n_clusters=3).fit(np1)
11 km4 = cluster.KMeans(n_clusters=4).fit(np1)
12 km6 = cluster.KMeans(n_clusters=6).fit(np1)
13
14 f, axarr = plt.subplots(2, 2, sharex='col', sharey='row', figsize=(10, 8))
15
16 for idx, km, tt in zip(product([0, 1], [0, 1]),
17                         [km2, km3, km4, km6],
18                         ["K=2, J=%.2f" % km2.inertia_,
19                         "K=3, J=%.2f" % km3.inertia_,
20                         "K=4, J=%.2f" % km4.inertia_,
21                         "K=6, J=%.2f" % km6.inertia_]):
22
23     Z = km.predict(np.c_[xx.ravel(), yy.ravel()])
24     Z = Z.reshape(xx.shape)
25
26     axarr[idx[0], idx[1]].contourf(xx, yy, Z, alpha=0.4)
27     axarr[idx[0], idx[1]].scatter(np1[:, 0], np1[:, 1], c=np_Y,
28                                   s=20, cmap=matplotlib.colors.ListedColormap(['red', 'green',
29                                   'blue', 'purple', 'yellow', 'magenta', 'cyan']))
30     axarr[idx[0], idx[1]].set_title(tt)
31 plt.show()

```

Question 2

How can inertia be used to compare cluster? Can you propose a better metric for spatial dataset?

Which is the principal problem with the mercator representation for continent detection?

2 Spatial dataset normalization

2.1 To code :

Displace the origin of longitude 30° east.

Question 3

What are the advantages of this normalization?

2.2 To code :

Apply KMeans on the new normalized dataset

Gaussian Mixture Model is a clustering method allowing soft boundaries.

This method can be used through a [sklearn function](#).

2.3 To code :

Apply GMM to the normalized dataset.

You have to test 3 conditions :

- GMM with default parameters for 2, 3, 4 and 6 components
 - GMM with diagonal covariance matrix for 2, 3, 4 and 6 components
 - GMM with random initialization for 2, 3, 4 and 6 components
-
-
-

Question 4

What is the default initialization for GMM?

Which is the best method to detect continents?

For this method, what would be the probability to find a country in the coordinate [-50,-40]? What methodology can you think of to detect unpopulated oceans?

2.4 Bonus

Rather than applying the previous longitudinal normalization in the mercator projected data and applying k-means, we could directly apply k-means on the unit-sphere representing the earth. The cosine similarity thus becomes a more suitable similarity measure than the euclidean distance.

Apply a KMeans with cosine similarity on the sphere instead of the previously tested euclidean distance in the plane.

3 Detection of caribbean island

3.1 To code

Isolate the North and South American continents

3.2 To code

Propose a methodology to create a cluster including caribbean independant nations.

Question 5

Which is the outlier of this clustering problem?

4 Evaluating a cluster

In this section, we want to evaluate the quality of the different clusters computed.

4.1 To do

Execute the following cells. For the first cell, you can add code to once again displace the longitude by 30° east.

```
1 df_X = df1[['Longitude', 'Latitude']]
2 df_Status = df1[['Continent']]
3
4 df_Y = df_Status.replace(['Africa', 'Asia', 'Europe', 'North America', 'South America', 'Seven
5 seas', 'Oceania'], [0, 1, 2, 3, 4, 5, 6])
6 np_Y = df_Y.to_numpy()
7 np_Y = np_Y.reshape((np_Y.shape[0],))
8
9 np1 = df_X.to_numpy()
10 plt.scatter(np1[:,0], np1[:,1], c=np_Y, cmap=matplotlib.colors.ListedColormap(['red', 'green',
    'blue', 'purple', 'yellow', 'magenta', 'cyan']))
11 plt.show()
```

```
1 x_min, x_max = np1[:, 0].min() - 10, np1[:, 0].max() + 10
2 y_min, y_max = np1[:, 1].min() - 10, np1[:, 1].max() + 10
3 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
4                       np.arange(y_min, y_max, 0.1))
```

```
1 from sklearn.metrics import silhouette_samples, silhouette_score
2
3 K = 3
4
5 fig, (ax1, ax2) = plt.subplots(1, 2)
6 ax1.set_xlim([-0.1, 1])
7 ax1.set_ylim([0, len(np1) + (K + 1) * 10])
8
9 km = cluster.KMeans(n_clusters=K, random_state=10)
10 cluster_labels = km.fit_predict(np1)
11
12 silhouette_avg = silhouette_score(np1, cluster_labels)
13 print("For n_clusters =", K, "The average silhouette_score is :", silhouette_avg)
14 sample_silhouette_values = silhouette_samples(np1, cluster_labels)
15
16 y_lower = 10
17 for i in range(K):
18     ith_cluster_silhouette_values = sample_silhouette_values[cluster_labels == i]
19     ith_cluster_silhouette_values.sort()
20
21     size_cluster_i = ith_cluster_silhouette_values.shape[0]
22     y_upper = y_lower + size_cluster_i
23
```

```

24     ax1.fill_betweenx(
25         np.arange(y_lower, y_upper),
26         0,
27         ith_cluster_silhouette_values,
28         alpha=0.7,
29     )
30
31     ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))
32
33     y_lower = y_upper + 10 # 10 for the 0 samples
34
35     ax1.set_title("The silhouette plot for the various clusters.")
36     ax1.set_xlabel("The silhouette coefficient values")
37     ax1.set_ylabel("Cluster label")
38
39     # The vertical line for average silhouette score of all the values
40     ax1.axvline(x=silhouette_avg, color="red", linestyle="--")
41
42     ax1.set_yticks([]) # Clear the yaxis labels / ticks
43     ax1.set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])
44
45     # 2nd Plot showing the actual clusters formed
46     Z = km.predict(np.c_[xx.ravel(), yy.ravel()])
47     Z = Z.reshape(xx.shape)
48
49     ax2.contourf(xx, yy, Z, alpha=0.4)
50     ax2.scatter(np1[:, 0], np1[:, 1], c=np_Y,
51                 s=20, cmap=matplotlib.colors.ListedColormap(['red', 'green', 'blue', 'purple',
52                 'yellow', 'magenta', 'cyan']))
53
54     # Labeling the clusters
55     centers = km.cluster_centers_
56     # Draw white circles at cluster centers
57     ax2.scatter(
58         centers[:, 0],
59         centers[:, 1],
60         marker="o",
61         c="white",
62         alpha=1,
63         s=200,
64         edgecolor="k",
65     )
66
67     for i, c in enumerate(centers):
68         ax2.scatter(c[0], c[1], marker="%d$" % i, alpha=1, s=50, edgecolor="k")
69
70     ax2.set_title("The visualization of the clustered data.")
71     ax2.set_xlabel("Feature space for the 1st feature")
72     ax2.set_ylabel("Feature space for the 2nd feature")
73
74     plt.suptitle(
75         "Silhouette analysis for KMeans clustering on sample data with n_clusters = %d"
76         % K,
77         fontsize=14,
78         fontweight="bold",
79     )
80
81     plt.show()

```

4.2 To do

Apply the silhouette computation for different values of K (at least 2,3,4 and 6).

Question 6

According to the silhouette score and the silhouette analysis, which is the most relevant value of K? Justify your response.

Bonus

Compute the silhouette score and analyze it for different parameters of Gaussian mixture.

4.3 To code

Another way to evaluate the quality of the clustering is through the homogeneity. This metric requires a ground truth, so it can't be computed for every clustering problem. Luckily, we do have a ground truth here. Compute the [homogeneity score](#) for different values of K.

Question 7

According to the homogeneity score, which is the most relevant value of K?

4.4 Bonus

Compute the homogeneity score and analyze it for different number of components and parameters of Gaussian Mixture.

Bonus : Alcoholism

This step is entirely optional and combines all the methods you used for this 3 TPs course.

The main goal is to develop a complete methodology to answer general questions.

All questions have to be justified by your homemade methodology and your methodology has also to be justified.

Question 1

In the year 2000, which countries are heavily concerned by an Alcohol issue?

Question 2

In these countries and in 2000, which are the parameters linked with Alcoholism? How do you explain these links?

Question 3

Which is the evolution trend in these countries between 2000 and 2015? Try to separate these different trends.

Question 4

By selecting a specific country, can you explain a decrease or an increase through specific policies?