



Master Informatique EID2

# Traitement numérique des données

TP4 Clustering de données

```

1 import numpy as np
2 from matplotlib import pyplot as plt
3 import matplotlib.cm as cm
4 from sklearn.datasets.samples_generator import make_blobs
5 from sklearn.cluster import KMeans
6 import pandas as pd
7 from sklearn.metrics import silhouette_samples, silhouette_score
8 from sklearn.decomposition import PCA
9 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
10 from sklearn import preprocessing
11 from sklearn.cluster import AgglomerativeClustering
12 from sklearn import metrics

```

```

1 #=====
2 # Données générées #
3 #=====
4 X, y = make_blobs(n_samples=300, centers=4, cluster_std=0.60, random_state=0)

```

```

1 #####
2 # A) K-Moyennes #
3 #####
4
5 #-----#
6 ## 1- Écrivez en python l'algorithme des K-Moyennes sous la forme d'une fonction ##
7 #-----#
8
9 def K_means(X,k):
10     n = X.shape[0]
11     m=0
12     c = X.shape[1]
13     mean = np.mean(X, axis = 0)
14     std = np.std(X, axis = 0)
15     centre = np.random.randn(k,c)*std + mean
16     centre_old = np.zeros(centre.shape)
17     centre_new = np.copy(centre)
18     clusters = np.zeros(n)
19     distances = np.zeros((n,k))
20     error = np.linalg.norm(centre_new - centre_old)
21     while error !=0 and m != 120:
22         for i in range(k):
23             distances[:,i] = np.linalg.norm(X - centre[i], axis=1)
24             clusters = np.argmin(distances, axis = 1)
25             centre_old = np.copy(centre_new)
26             for i in range(k):
27                 centre_new[i] = np.mean(X[clusters == i], axis=0)
28             error = np.linalg.norm(centre_new - centre_old)
29             m = m+1
30     return clusters, centre_new
31
32 def afficher_Kmeans(X,c):
33     plt.scatter(X[:,0], X[:,1])
34     plt.scatter(c[:,0], c[:,1], s=300, c='red')
35     plt.show()

```

```

1 #-----#
2 # Comparez avec la fonction Kmeans de sklearn #
3 #-----#
4 def Kmeans(X,k):
5     kmeans = KMeans(n_clusters=4, init='k-means++', max_iter=len(X), n_init=10, random_state=0)
6     kmeans = kmeans.fit(X)
7     pred_y = kmeans.predict(X)
8     plt.scatter(X[:,0], X[:,1])
9     plt.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[1], s=300, c='red')
10    plt.show()

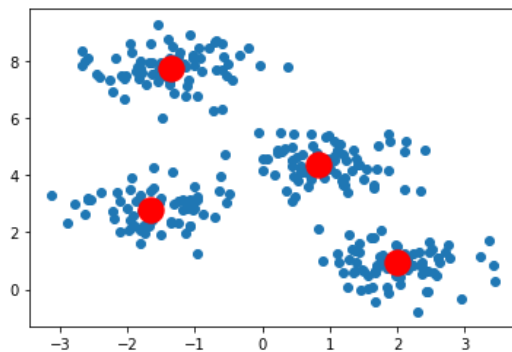
```

```

1 print("Mon Kmeans")
2 c,nc = K_means(X,4)
3 afficher_Kmeans(X,nc)

```

Mon Kmeans

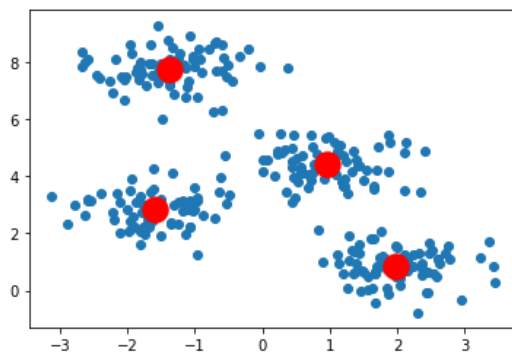


```

1 print("Sklearn Kmeans")
2 Kmeans(X,4)

```

Sklearn Kmeans

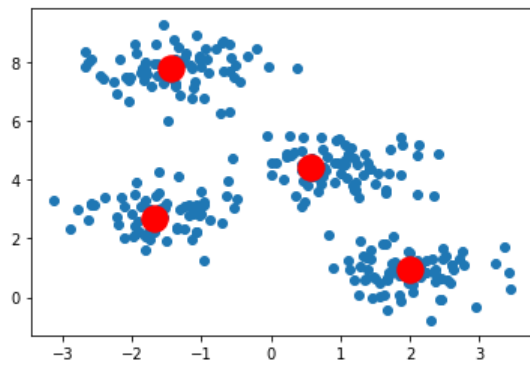


```

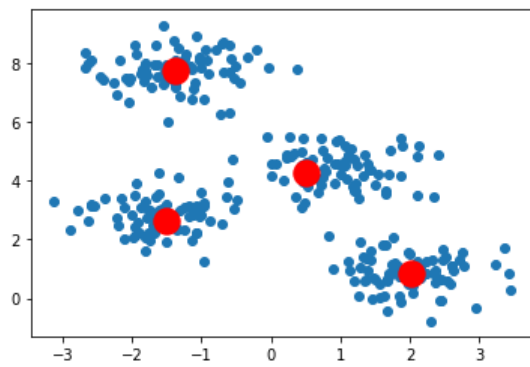
1 #-----#
2 ## 2. Expérimenter l'instabilité due à l'initialisation ##
3 #-----#
4 def instabilite(X,k,n):
5     for i in range(n):
6         c,nc = K_means(X,k)
7         print("Test n°",i+1)
8         afficher_Kmeans(X,nc)
9
10    instabilite(X,4,4)

```

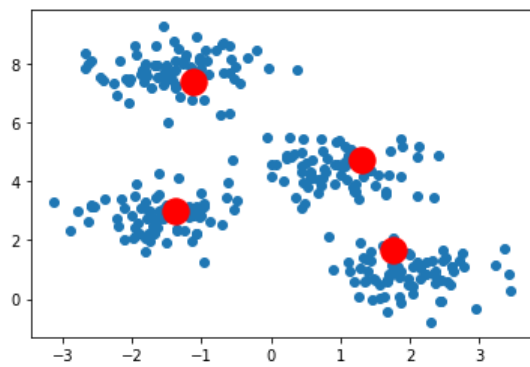
Test n° 1



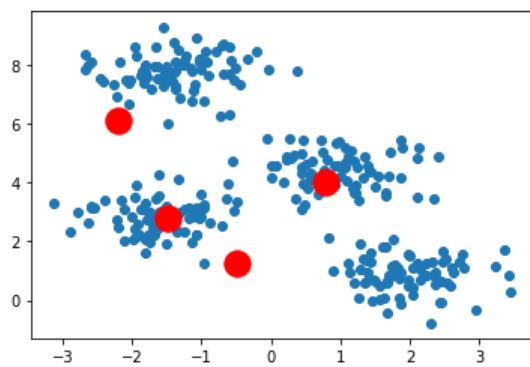
Test n° 2



Test n° 3



Test n° 4

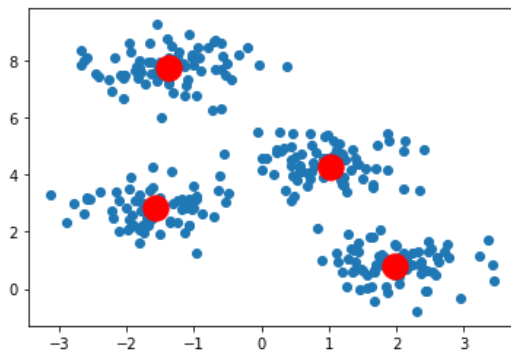


```

1 #-----#
2 ## 3. Utiliser l'indice de Silhouette ##
3 #-----#
4
5 def Silhouette(X,clusters):
6     return metrics.silhouette_score(X, clusters, metric='euclidean')
7
8 def protocole(X):
9     score = 0
10    for i in range(2,10):
11        for j in range(10):
12            c,nc = K_means(X,i)
13            b=Silhouette(X,c)
14            if score < b:
15                score = b
16                centre_new = nc
17    print("Meilleur indice de Silhoutte est : ",score)
18    afficher_Kmeans(X,centre_new)
19
20 protocole(X)

```

Meilleur indice de Silhoutte est : 0.6636114905089738



```

1 #-----#
2 ## 4. Utiliser une ACP ##
3 #-----#
4 PCA = PCA(n_components=2).fit(X).transform(X)
5 print('PCA (10 premières lignes):\n',PCA[:10])
6 #-----#
7 ## 4. Utiliser une LDA ##
8 #-----#
9 LDA = LinearDiscriminantAnalysis(n_components=2).fit(X, y).transform(X)
10 print('LDA (10 premières lignes):\n',LDA[:10])

```

PCA (10 premières lignes):

```

[[-2.01616694 -0.01679795]
 [ 3.71613677 -0.13003844]
 [ 0.55503775 -1.52383607]
 [ 3.9234751  -0.6566597 ]
 [-2.41689163 -0.31285899]
 [-4.79345065 -1.61841309]
 [-1.421835   2.61647075]
 [-0.2061925  -1.46454628]
 [ 4.15211641 -1.65200822]
 [ 4.56816034 -1.28424075]]

```

LDA (10 premières lignes):

```

[[-3.38223087  0.04572696]
 [ 6.24076844  0.18367997]
 [ 0.99513795  2.53132643]
 [ 6.61072394  1.05834219]
 [-4.04220558  0.54201203]
 [-7.97515004  2.73587476]
 [-2.49531416 -4.34222911]
 [-0.28461715  2.43935557]
 [ 7.03606014  2.7129495 ]
 [ 7.71873507  2.0971838 ]]

```

```

1 #####
2 # B) Analyse des données « choix projet » #
3 #####
4
5 #-----#
6 ## Importation de fichier csv ##
7 #-----#
8
9 CP = pd.read_csv("choixprojetstab.csv" , sep=';')
10 C = np.array(CP['étudiant-e'])
11 M = np.array(CP.drop(['étudiant-e'], axis=1))
12
13 print('Code étudiant-e :\n',C,'\n')
14 print('Data :\n',M,'\n')
15 print('Test bien numerique : M[0,0]+M[0,0] = ', M[0,0] + M[0,0])

```

```

Code étudiant-e :
['b1/.vSDYCGrSs' 'b1/1NiMubceBs' 'b1/dvgMTLVsvk' 'b11NWhKcNADF2'
'b11ao5B7htJfQ' 'b11u./AF8TEp6' 'b12LwDaKpKT/c' 'b12jhb3v2qUhc'
'b12rqWgKava4o' 'b14wbw9k3zw/c' 'b15.Z2M26c/5o' 'b15/1tjKhKrAE'
'b155W0uSMrpz6' 'b165TErI.VjAA' 'b19BdxOPgV0XU' 'b19HR9LuSmrMo'
'b1A5K4G1KIEN2' 'b1B60Jg1i1NV2' 'b1C888Fe1.osw' 'b1DPCjzWCczr6'
'b1FD1aQDfmpH.' 'b1Fq6YzuUpuG6' 'b1GDihFnA6.LE' 'b1HIVuF4ftibA'
'b1HxzJKy1kWeU' 'b1If9FrZBwJyk' 'b1K.9jBCD/SSw' 'b1KXgoRO0CXOE'
'b1KsQZnr3tnvw' 'b1MKF1GVYbnSA' 'b1NDVy9IQ.pyc' 'b1QdUJch/nfdg'
'b1SI8b6JyYqHY' 'b1UiqJyD6t5eo' 'b1Xfmkaw817T6' 'b1XtkgxAdNzb.'
'b1Yb2ABO7EU.M' 'b1Zxs/nnNLCos' 'b1ZzwRPNDZpRY' 'b1aADEMnN043A'
'b1aIQ61DLW6uc' 'b1bKjW2k4.8xo' 'b1bfV11TjpkQY' 'b1bk19iy6NSzQ'
'b1ddve1gjnW96' 'b1fIGd5mN48z2' 'b1g10tiwNMDAs' 'b1g8PMUnKJxes'
'b1ginkLBh05f.' 'b1h91EioXhmgc' 'b1hai4aeyq8rU' 'b1hnnMwGPnuNc'
'b1iYbAmgGkzSI' 'b1jdvFfVzFYBk' 'b1lOE3T7MU9nE' 'b1leUFdn4IybY'
'b1mLNMnzoOUT6' 'b1n85uWpKDNjM' 'b1nahdfBfYod.' 'b1pJfyKbetUsc'
'b1pNI0oxLIG.' 'b1tZdophgVNQ2' 'b1tppVKeGreVc' 'b1tsUkerTn/gA'
'b1luruT/s0vVm6' 'b1uwV0aELRqFQ' 'b1vqKhCDhQab.' 'b1wyn40XjgvEs'
'b1y3Qqsvf6WlU' 'b1y69tf4z1FiM' 'b1zon470EHguA']

```

```

Data :
[[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]]

```

Test bien numerique : M[0,0]+M[0,0] = 2

```

1 #-----#
2 ## Testez les différents algorithmes du package et proposez le meilleur ##
3 ## ----- clustering possible des données selon l'indice Silhouette ----- ##
4 #-----#
5
6 def protocole2(X,Y,f):
7     range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
8     for n_clusters in range_n_clusters:
9         fig, (ax1, ax2) = plt.subplots(1, 2)
10        fig.set_size_inches(12, 4)
11        ax1.set_xlim([-0.1, 1])
12        ax1.set_ylim([0, len(X) + (n_clusters + 1) * 10])
13        if f=='K':
14            clusterer = KMeans(n_clusters=n_clusters, random_state=10)
15            cluster_labels = clusterer.fit_predict(X)
16        if f=='A':
17            clusterer = AgglomerativeClustering(n_clusters=n_clusters, linkage='ward').fit(X)
18            cluster_labels = clusterer.labels_
19        silhouette_avg = silhouette_score(X, cluster_labels)
20        print("Pour n_clusters =", n_clusters,
21              "La moyenne silhouette_score est :", silhouette_avg)
22        sample_silhouette_values = silhouette_samples(X, cluster_labels)
23        y_lower = 10

```

```

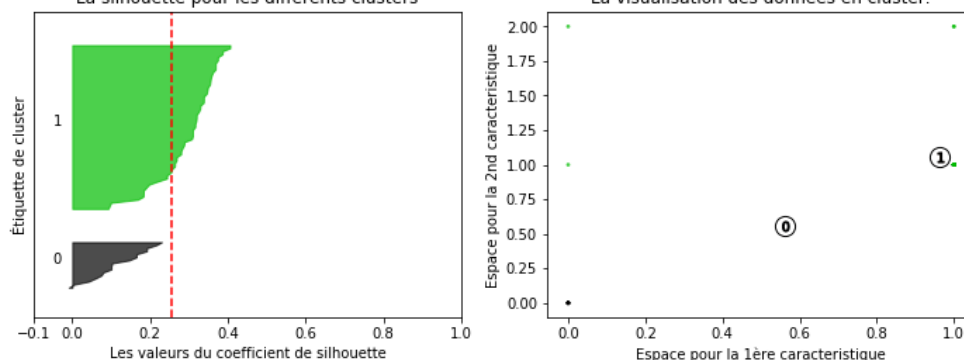
24
25     for i in range(n_clusters):
26         ith_cluster_silhouette_values = \
27             sample_silhouette_values[cluster_labels == i]
28
29         ith_cluster_silhouette_values.sort()
30
31         size_cluster_i = ith_cluster_silhouette_values.shape[0]
32         y_upper = y_lower + size_cluster_i
33
34         color = cm.nipy_spectral(float(i) / n_clusters)
35         ax1.fill_betweenx(np.arange(y_lower, y_upper),
36                          0, ith_cluster_silhouette_values,
37                          facecolor=color, edgecolor=color, alpha=0.7)
38         ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))
39         y_lower = y_upper + 10
40     ax1.set_title("La silhouette pour les différents clusters")
41     ax1.set_xlabel("Les valeurs du coefficient de silhouette")
42     ax1.set_ylabel("Étiquette de cluster")
43     ax1.axvline(x=silhouette_avg, color="red", linestyle="--")
44     ax1.set_yticks([])
45     ax1.set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])
46     colors = cm.nipy_spectral(cluster_labels.astype(float) / n_clusters)
47     ax2.scatter(X[:, 0], X[:, 1], marker='.', s=30, lw=0, alpha=0.7,
48                c=colors, edgecolor='k')
49     if f == 'K':
50         centers = clusterer.cluster_centers_
51         ax2.scatter(centers[:, 0], centers[:, 1], marker='o',
52                    c="white", alpha=1, s=200, edgecolor='k')
53         for i, c in enumerate(centers):
54             ax2.scatter(c[0], c[1], marker='%d$' % i, alpha=1,
55                        s=50, edgecolor='k')
56         ax2.set_title("La visualisation des données en cluster.")
57         ax2.set_xlabel("Espace pour la 1ère caractéristique")
58         ax2.set_ylabel("Espace pour la 2nd caractéristique")
59     if f == 'K':
60         plt.suptitle(("Analyse de silhouette pour le regroupement de KMeans sur des exemples de données "
61                      "avec n_clusters = %d" % n_clusters),
62                      fontsize=14, fontweight='bold')
63     if f == 'A':
64         plt.suptitle(("Analyse des silhouettes pour le regroupement WARD sur des exemples de données "
65                      "avec n_clusters = %d" % n_clusters),
66                      fontsize=14, fontweight='bold')
67     plt.show()

```

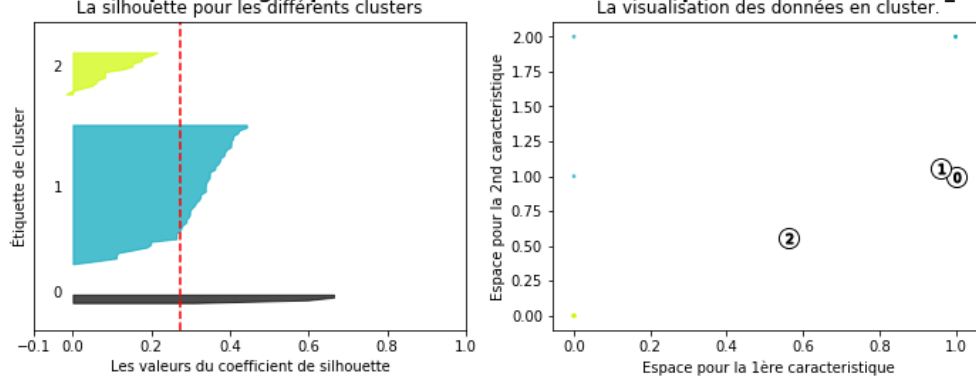
```
1 protocole2(M,C,'K')
```

Pour n\_clusters = 2 La moyenne silhouette\_score est : 0.2537954711744975  
 Pour n\_clusters = 3 La moyenne silhouette\_score est : 0.27395879223997877  
 Pour n\_clusters = 4 La moyenne silhouette\_score est : 0.26800412978792404  
 Pour n\_clusters = 5 La moyenne silhouette\_score est : 0.3176827275171515  
 Pour n\_clusters = 6 La moyenne silhouette\_score est : 0.336781722416005

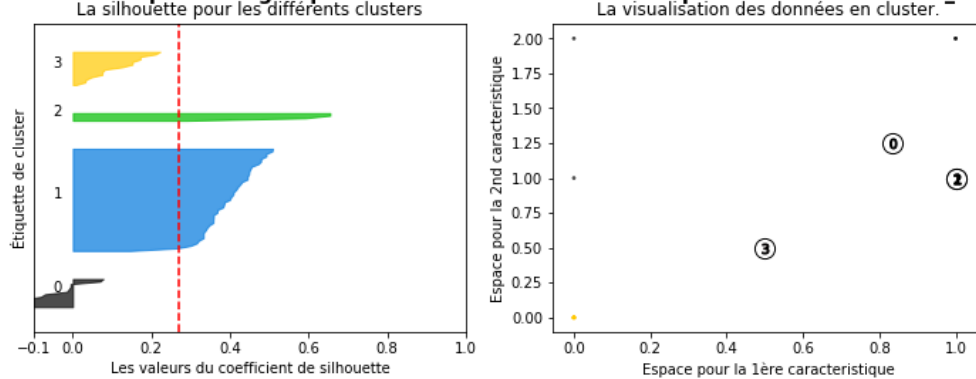
### **Analyse de silhouette pour le regroupement de KMeans sur des exemples de données avec n\_clusters = 2**



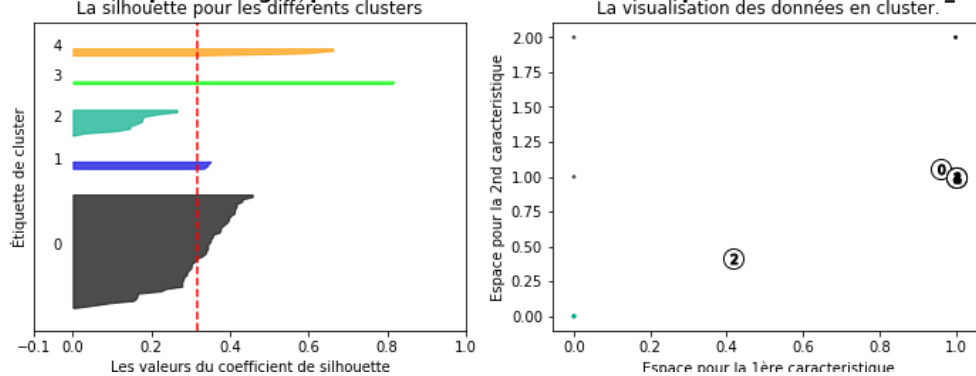
### Analyse de silhouette pour le regroupement de KMeans sur des exemples de données avec $n\_clusters = 3$



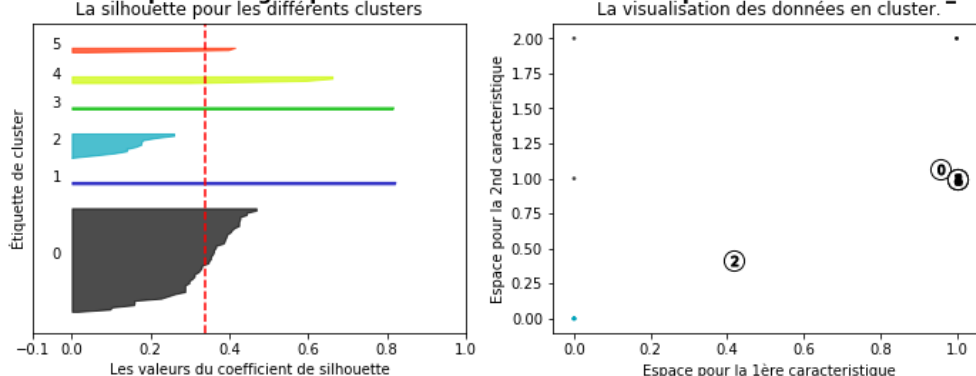
### Analyse de silhouette pour le regroupement de KMeans sur des exemples de données avec $n\_clusters = 4$



### Analyse de silhouette pour le regroupement de KMeans sur des exemples de données avec $n\_clusters = 5$



### Analyse de silhouette pour le regroupement de KMeans sur des exemples de données avec $n\_clusters = 6$



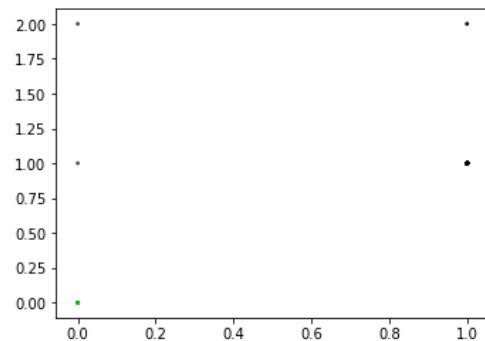
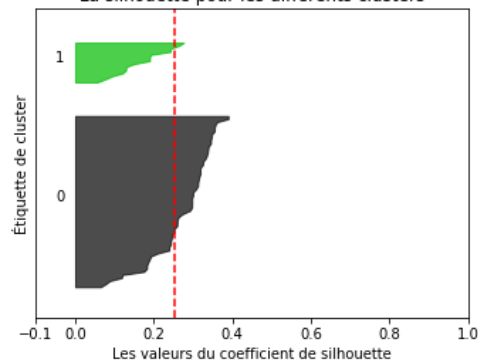


```
1 protocole2(M,C,'A')
```

Pour  $n\_clusters = 2$  La moyenne silhouette\_score est : 0.25270656690459653  
 Pour  $n\_clusters = 3$  La moyenne silhouette\_score est : 0.25719312860651694  
 Pour  $n\_clusters = 4$  La moyenne silhouette\_score est : 0.2884593801353125  
 Pour  $n\_clusters = 5$  La moyenne silhouette\_score est : 0.302588527642423  
 Pour  $n\_clusters = 6$  La moyenne silhouette\_score est : 0.3253001235762009

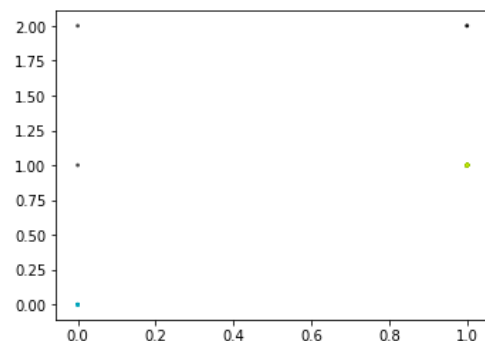
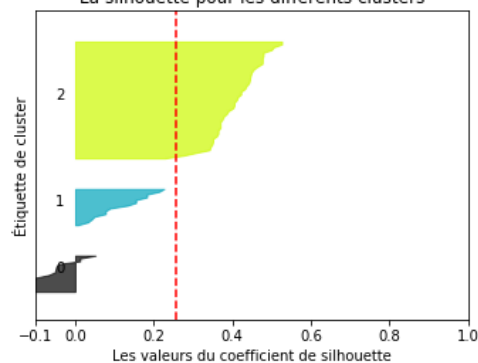
#### Analyse des silhouettes pour le regroupement WARD sur des exemples de données avec $n\_clusters = 2$

La silhouette pour les différents clusters



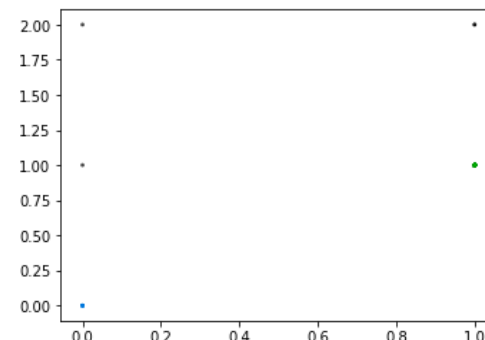
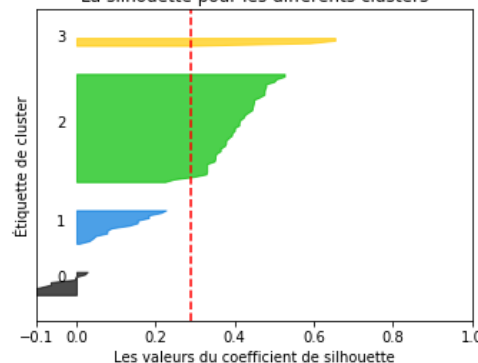
#### Analyse des silhouettes pour le regroupement WARD sur des exemples de données avec $n\_clusters = 3$

La silhouette pour les différents clusters



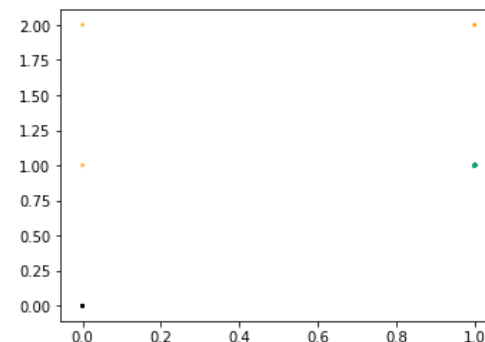
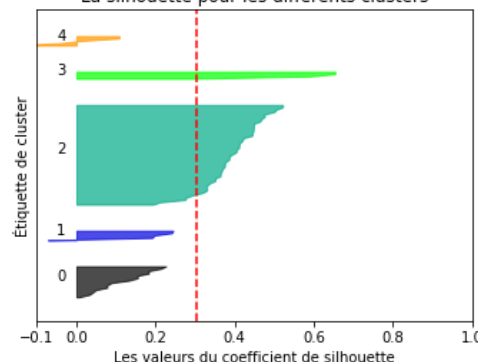
#### Analyse des silhouettes pour le regroupement WARD sur des exemples de données avec $n\_clusters = 4$

La silhouette pour les différents clusters



#### Analyse des silhouettes pour le regroupement WARD sur des exemples de données avec $n\_clusters = 5$

La silhouette pour les différents clusters



**Analyse des silhouettes pour le regroupement WARD sur des exemples de données avec  $n\_clusters = 6$**

