



**UNIVERSITÉ
DE REIMS
CHAMPAGNE-ARDENNE**

CRTP1 IA et Données

INFO0503

07 novembre 2024

VERPOORTE William

Sommaire

1	Implémentation de KMeans	2
2	Analyse	6

1 Implémentation de KMeans

```
1 import matplotlib.pyplot as plt
2 from math import sqrt
3 from typing import List
4 import random
5 from sklearn.cluster import KMeans
6 from sklearn.metrics import confusion_matrix
7 import seaborn as sns
8 from sklearn.metrics import silhouette_score
9
10
11 class Point:
12     def __init__(self, x: float, y: float):
13         self.x = x
14         self.y = y
15
16     def __repr__(self):
17         return f"Point({self.x}, {self.y})"
18
19 def Kmeans(points: List[Point], centres: List[Point]):
20     max_iterations = 10
21     colors = ['red', 'blue', 'green', 'purple', 'orange',
22              'brown', 'pink', 'cyan', 'magenta', 'yellow']
23
24     plt.figure(figsize=(8, 6)) # Cr er la fen tre pour le
25     graphique
26
27     for iteration in range(max_iterations):
28         classe = {centre: [] for centre in centres}
29
30         for pt in points:
```

```
31         dist = sqrt((pt.x - centre.x)**2 + (pt.y -
32                 centre.y)**2)
33         if dist < dist_min:
34             dist_min = dist
35             centre_min = centre
36             classe[centre_min].append(pt)
37
38 # Effacer l'ancienne figure pour la mise à jour
39 plt.clf()
40
41 # Affichage des points associés à chaque centre
42 for idx, (centre, points_in_class) in
43     enumerate(classe.items()):
44         x_vals = [pt.x for pt in points_in_class]
45         y_vals = [pt.y for pt in points_in_class]
46         plt.scatter(x_vals, y_vals, c=colors[idx],
47                 label=f'Centre {centre} (Cluster {idx+1})',
48                 alpha=0.6)
49
50 # Affichage des centres de cluster
51 centre_x_vals = [centre.x for centre in centres]
52 centre_y_vals = [centre.y for centre in centres]
53 plt.scatter(centre_x_vals, centre_y_vals, c='black',
54         marker='X', label='Centres', s=200)
55
56 # Titre et légende
57 plt.title(f"K-means - Iteration {iteration + 1}")
58 plt.xlabel('X')
59 plt.ylabel('Y')
60 plt.legend(loc='upper right')
61 plt.grid(True)
62
63 new_centres = []
64 for centre in centres:
```

```
60         if classe[centre]:
61             moyenne_x = sum(pt.x for pt in classe[centre])
62                         / len(classe[centre])
63             moyenne_y = sum(pt.y for pt in classe[centre])
64                         / len(classe[centre])
65             new_centres.append(Point(moyenne_x, moyenne_y))
66         else:
67             new_centres.append(centre)
68
69     if all(abs(new.x - old.x) < 1e-4 and abs(new.y -
70         old.y) < 1e-4 for new, old in zip(new_centres,
71         centres)):
72         print("Convergence atteinte!")
73         break
74
75     centres = new_centres
76
77     plt.pause(3) # Pause de 3 secondes
78
79     plt.show()
80
81 # Exemple d'utilisation avec des points supplémentaires
82 points = [Point(random.uniform(0, 10), random.uniform(0, 10))
83           for _ in range(100)]
84 points2 = [[random.uniform(0, 10), random.uniform(0, 10)] for
85             _ in range(100)]
86
87 # Centres initiaux pour le k-means
88 centres = [
89     Point(2, 2), Point(5, 5), Point(8, 3), Point(10, 10)
90 ]
```

```

88 # Appel de la fonction kmeans avec la liste tendue de points
89 Kmeans(points, centres)
90
91 # Appliquer K-means avec 3 clusters
92 kmeans = KMeans(n_clusters=3, random_state=0)
93 kmeans.fit(points2)
94
95 # Afficher les r sultats
96 plt.scatter([p[0] for p in points2], [p[1] for p in points2],
97             c=kmeans.labels_)
98 plt.scatter(kmeans.cluster_centers_[0],
99             kmeans.cluster_centers_[1], s=200, c='red', marker='X')
100 plt.title('K-means Clustering with scikit-learn')
101 plt.show()

```

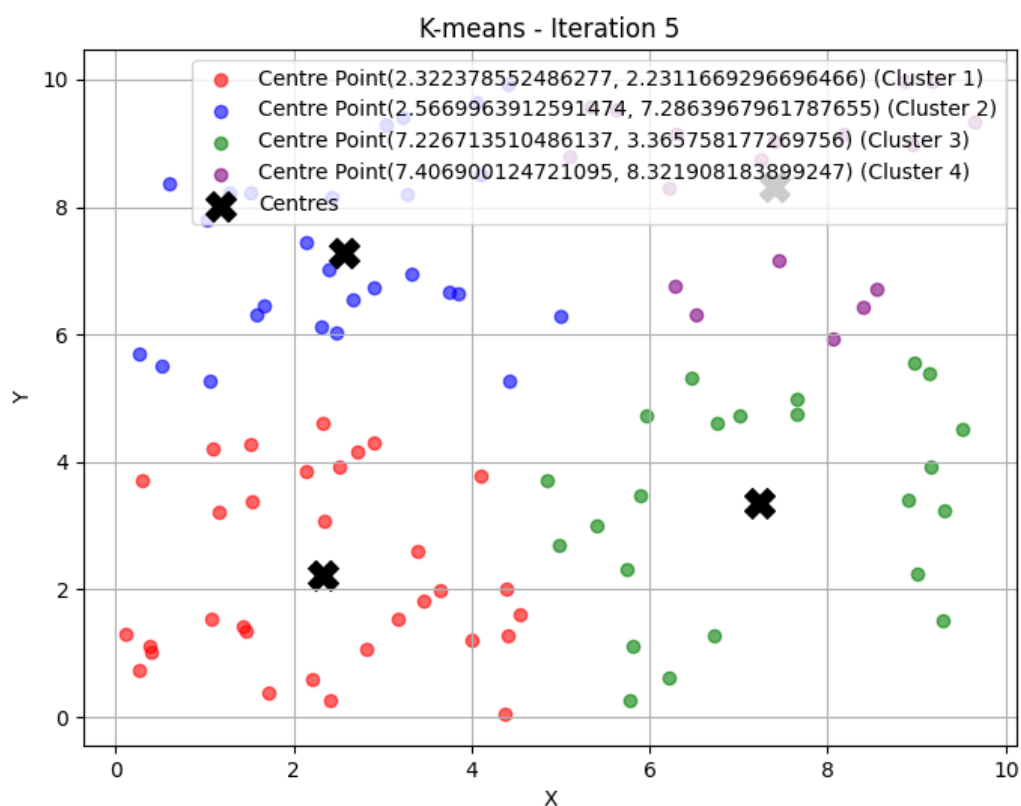


Figure 1: KMeans à la main

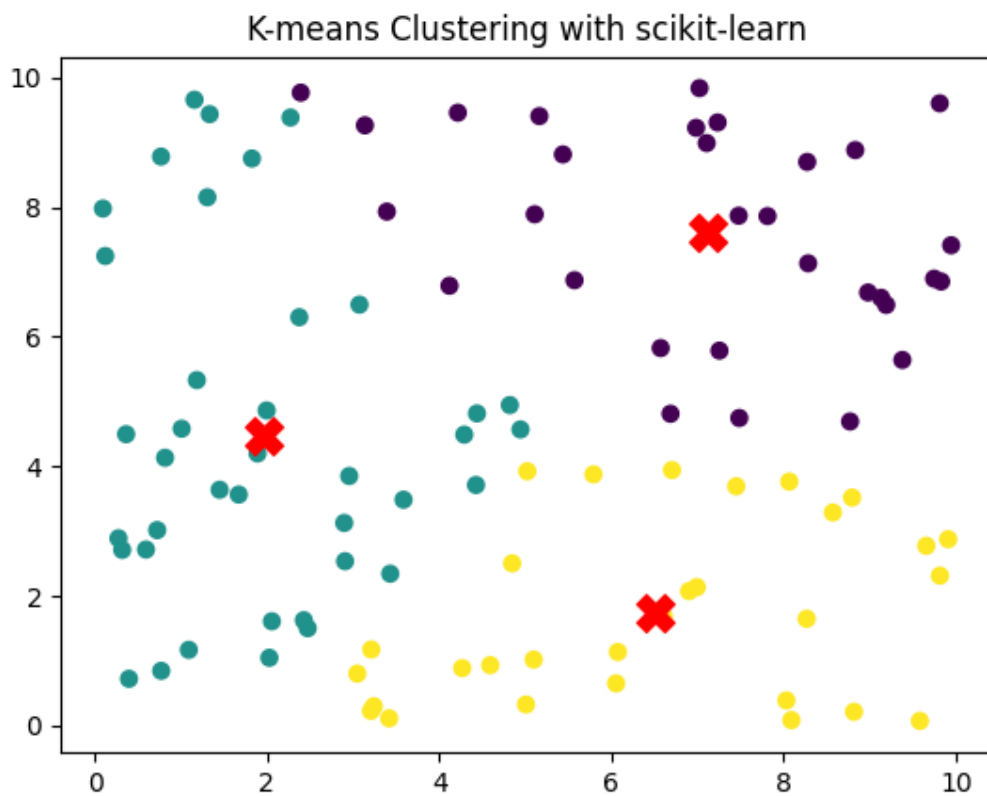


Figure 2: KMeans de sklearn

2 Analyse

En utilisant des codes tests obtenues sur le site de sklearn avec de plus en plus de points et changement des centres de départs nous en conclure plusieurs choses.

La rapidité et l'efficacité de l'algorithme changent beaucoup selon les centres de départs:

```
1 import numpy as np
2
3 from sklearn.datasets import load_digits
4
5 data, labels = load_digits(return_X_y=True)
6 (n_samples, n_features), n_digits = data.shape,
   np.unique(labels).size
7
8 print(f"# digits: {n_digits}; # samples: {n_samples}; #
```

```
features {n_features}")

9
10 from time import time
11
12 from sklearn import metrics
13 from sklearn.pipeline import make_pipeline
14 from sklearn.preprocessing import StandardScaler
15
16
17 def bench_k_means(kmeans, name, data, labels):
18     """Benchmark to evaluate the KMeans initialization methods.
19
20     Parameters
21     -----
22     kmeans : KMeans instance
23         A :class:`~sklearn.cluster.KMeans` instance with the
24         initialization
25         already set.
26     name : str
27         Name given to the strategy. It will be used to show
28         the results in a
29         table.
30     data : ndarray of shape (n_samples, n_features)
31         The data to cluster.
32     labels : ndarray of shape (n_samples,)
33         The labels used to compute the clustering metrics
34         which requires some
35         supervision.
36
37     """
38     t0 = time()
39     estimator = make_pipeline(StandardScaler(),
40                               kmeans).fit(data)
41     fit_time = time() - t0
42     results = [name, fit_time, estimator[-1].inertia_]
```



```
38
39     # Define the metrics which require only the true labels
        and estimator
40     # labels
41     clustering_metrics = [
42         metrics.homogeneity_score,
43         metrics.completeness_score,
44         metrics.v_measure_score,
45         metrics.adjusted_rand_score,
46         metrics.adjusted_mutual_info_score,
47     ]
48     results += [m(labels, estimator[-1].labels_) for m in
        clustering_metrics]
49
50     # The silhouette score requires the full dataset
51     results += [
52         metrics.silhouette_score(
53             data,
54             estimator[-1].labels_,
55             metric="euclidean",
56             sample_size=300,
57         )
58     ]
59
60     # Show the results
61     formatter_result = (
62         "{:9s}\t{:.3f}s\t{:.0f}\t{:.3f}\t{:.3f}\t{:.3f}\t{:.3f}\t{:.3f}\t{:.3f}"
63     )
64
65     print(formatter_result.format(*results))
66
67 from sklearn.cluster import KMeans
68 from sklearn.decomposition import PCA
```

```

70 print(82 * "_")
71 print("init\t\ttime\tinertia\tthomo
72 \tcompl\tv-meas\tARI\tAMI\tsilhouette")
73
74 kmeans = KMeans(init="k-means++", n_clusters=n_digits,
75                 n_init=4, random_state=0)
76
77 bench_k_means(kmeans=kmeans, name="k-means++", data=data,
78               labels=labels)
79
80 kmeans = KMeans(init="random", n_clusters=n_digits, n_init=4,
81                 random_state=0)
82 bench_k_means(kmeans=kmeans, name="random", data=data,
83               labels=labels)
84
85 kmeans = KMeans(init=pca.components_, n_clusters=n_digits,
86                 n_init=1)
87 bench_k_means(kmeans=kmeans, name="PCA-based", data=data,
88               labels=labels)
89
90 print(82 * "_")

```

Ici on teste 3 méthodes pour initialiser les centres : la méthode de kmeans++, une version aléatoire et la version PCA, le code ci-dessus nous donne la sorties suivantes :

```

PS C:\Users\willi\OneDrive\Documents\INF08503\TP> & C:/Users/willl/AppData/Local/Programs/Python/Python313/python.exe c:/Users/willl/OneDrive/Docu
ments/INF08503/TP/testkmeans.py
# digits: 10; # samples: 1797; # features 64

```

init	time	inertia	homo	compl	v-meas	ARI	AMI	silhouette
k-means++	0.108s	69545	0.598	0.645	0.621	0.469	0.617	0.146
random	0.030s	69735	0.681	0.723	0.701	0.574	0.698	0.177
PCA-based	0.010s	69513	0.600	0.647	0.622	0.468	0.618	0.146

Figure 3: Sortie de code

on remarque que le temps d'exécution varie plus ou moins selon la méthode utilisé et donc incite plus ou moins à changer selon les besoins.

```

1 import matplotlib.pyplot as plt
2
3 reduced_data = PCA(n_components=2).fit_transform(data)

```

```
4 kmeans = KMeans(init="k-means++", n_clusters=n_digits,
    n_init=4)
5 kmeans.fit(reduced_data)
6
7 # Step size of the mesh. Decrease to increase the quality of
    the VQ.
8 h = 0.02 # point in the mesh [x_min, x_max]x[y_min, y_max].
9
10 # Plot the decision boundary. For that, we will assign a color
    to each
11 x_min, x_max = reduced_data[:, 0].min() - 1, reduced_data[:,
    0].max() + 1
12 y_min, y_max = reduced_data[:, 1].min() - 1, reduced_data[:,
    1].max() + 1
13 xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
    np.arange(y_min, y_max, h))
14
15 # Obtain labels for each point in mesh. Use last trained model.
16 Z = kmeans.predict(np.c_[xx.ravel(), yy.ravel()])
17
18 # Put the result into a color plot
19 Z = Z.reshape(xx.shape)
20 plt.figure(1)
21 plt.clf()
22 plt.imshow(
23     Z,
24     interpolation="nearest",
25     extent=(xx.min(), xx.max(), yy.min(), yy.max()),
26     cmap=plt.cm.Paired,
27     aspect="auto",
28     origin="lower",
29 )
30
31 plt.plot(reduced_data[:, 0], reduced_data[:, 1], "k.",
```

```
        markersize=2)
32 # Plot the centroids as a white X
33 centroids = kmeans.cluster_centers_
34 plt.scatter(
35     centroids[:, 0],
36     centroids[:, 1],
37     marker="x",
38     s=169,
39     linewidths=3,
40     color="w",
41     zorder=10,
42 )
43 plt.title(
44     "K-means clustering on the digits dataset (PCA-reduced
45     data)\n"
46     "Centroids are marked with white cross"
47 )
48 plt.xlim(x_min, x_max)
49 plt.ylim(y_min, y_max)
50 plt.xticks(())
51 plt.yticks(())
52 plt.show()
```

En ajoutant cette partie de code on peut donc obtenir l’affichage suivant :

nous montrant les différents clusters de données (chaque couleur) et leur centre (croix blanche).

On remarque aussi une tendance des clusters à ne pas se mélanger, chacun se retrouve de son côté, on le voit sur la figure en dessous mais aussi dans différent test que l’on peut retrouver sur le site de sklearn et leur différent tests.

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross

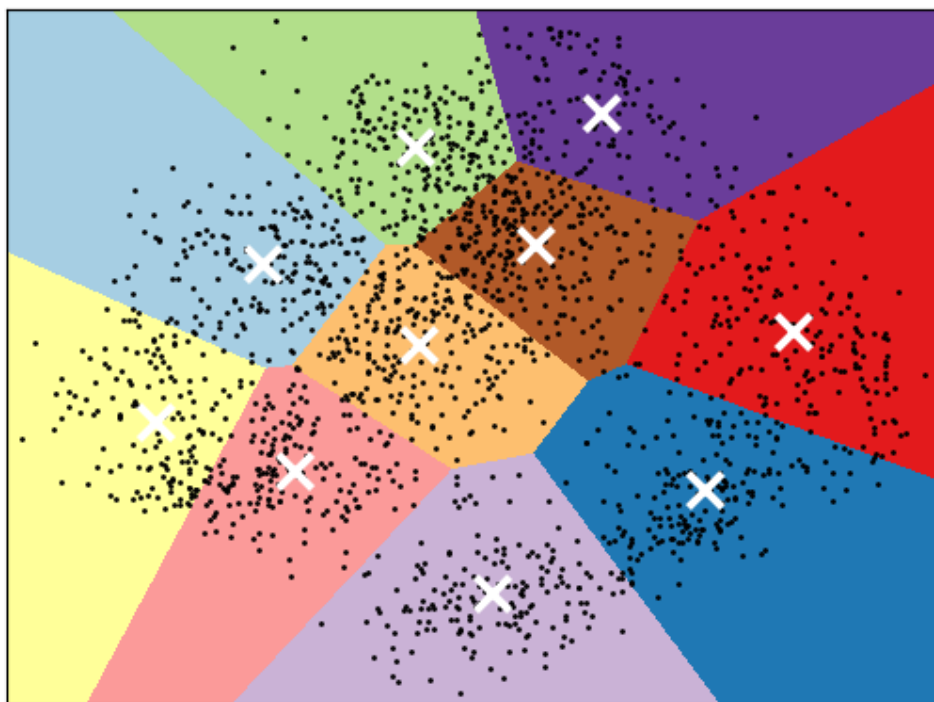


Figure 4: Affichage du code