

## (1) ما هو التجمع؟ اشرح بمثال بطريقة هل يمكن أن تكون مفيدة؟

التجزئة (Clustering) هي عملية تنظيم البيانات المتشابهة معاً في مجموعات صغيرة تسمى (clusters)، حيث يتم تجميع البيانات التي تشترك في خصائص معينة معاً داخل نفس التجمع (cluster)، بينما تكون متميزة عن بيانات التجمعات الأخرى. بمعنى آخر، فإنها تسعى إلى فصل البيانات إلى مجموعات تتشابه داخلها وتختلف بينها.

على سبيل المثال، يمكن تجميع المستخدمين الذين يشاهدون ويفضلون مقاطع الفيديو ذات المحتوى التقني في مجموعة واحدة، بينما يمكن تجميع المستخدمين الذين يفضلون مقاطع الفيديو ذات المحتوى المنزلي في مجموعة أخرى.

مثل إذا كانت بيانات الاستخدام تشير إلى أن هناك مجموعة كبيرة من المستخدمين يشاهدون بشكل رئيسي مقاطع الفيديو التقنية والمرتبطة بالتكنولوجيا، فيمكن تجميعهم في تجمع واحد. بينما يمكن تجميع المستخدمين الذين يشاهدون بشكل رئيسي مقاطع الفيديو التي تتعلق بالطبخ والديكور في (cluster) آخر.

مثال آخر: لنفترض أن لدينا مجموعة من المشترين عبر الإنترنت، ونريد تقسيمهم إلى مجموعات استناداً إلى تفضيلاتهم فيما يتعلق بالمنتجات التقنية. إذا كانت بيانات المشترين تشير إلى أن هناك مجموعة يشتررون بشكل رئيسي الهواتف الذكية والأجهزة الإلكترونية، فيمكن تجميعهم في (cluster) واحد. بينما يمكن تجميع المشترين الذين يشتررون بشكل رئيسي الأجهزة المنزلية مثل الثلاجات والغسالات في (cluster) آخر.

هذا المثال يظهر كيف يمكن استخدام التجزئة لتحليل سلوك المشترين وتقسيمهم إلى مجموعات استناداً إلى تفضيلاتهم في المنتجات، مما يمكن للشركة من تحديد استراتيجيات التسويق والعروض الترويجية المستهدفة لكل مجموعة.

فيما يلي بعض الطرق التي يمكن أن تكون فيها التجزئة مفيدة:

1. تصنيف العملاء أو المستهلكين: من خلال تجزئة العملاء أو المستهلكين بناءً على عادات الشراء أو التفضيلات، يمكن للشركات توجيه استراتيجيات التسويق بشكل أفضل وتلبية احتياجات العملاء بشكل أكثر فعالية.
2. تجزئة الصور أو الفيديوهات: يمكن استخدام التجزئة في معالجة الصور والفيديوهات لتنظيمها وتصنيفها بناءً على الخصائص المشتركة مثل الألوان أو الأشكال أو النمط.
3. تحليل البيانات الجغرافية: يمكن استخدام التجزئة لتجميع البيانات الجغرافية مثل العناوين البريدية أو المواقع الجغرافية لتحديد مناطق محددة أو توزيعات جغرافية.
4. تحليل السلوك الاجتماعي عبر وسائل التواصل الاجتماعي: يمكن استخدام التجزئة لتقسيم المستخدمين على وسائل التواصل الاجتماعي إلى مجموعات استناداً إلى أنماط تفاعلهم ومشاركاتهم، مما يسمح للشركات والمعلنين بفهم أفضل لجمهورهم وتصميم الحملات التسويقية بشكل أكثر فعالية.
5. تقسيم المرضى في المجال الطبي: يمكن استخدام التجزئة لتقسيم المرضى إلى مجموعات استناداً إلى التشخيصات الطبية والعوامل الصحية المشتركة، مما يمكن الأطباء من تخصيص العلاجات والرعاية بشكل أكثر دقة وفعالية.

(2) ما هي أنواع مختلفة من التجمع؟ اشرح مع الأمثلة.

#### 1. التجزئة النمطية: (Partitioning Clustering)

في هذا النوع من التجزئة، يتم تقسيم مجموعة البيانات إلى عدد محدد من التجمعات، وتعتمد هذه التجميعات على معايير محددة مثل المسافات بين النقاط. واحدة من أشهر أنواع التجزئة النمطية هي خوارزمية K-means.

مثال: تقسيم مجموعة من العملاء في متجر إلكتروني إلى مجموعات استنادًا إلى عادات شرائهم، حيث يمكن تقسيمهم إلى مجموعة من يشترون ملابس رياضية ومجموعة أخرى من يشترون ملابس أنيقة.

#### 2. التجزئة بوسائل القياس: (K-Means Clustering)

في هذا النوع من التجزئة، يتم تقسيم مجموعة البيانات إلى عدد محدد من التجمعات، حيث يعتمد تكوين هذه التجمعات على المسافات بين النقاط. واحدة من أشهر أنواع التجزئة بوسائل القياس هي خوارزمية K-Means. مثال: يمكن تقسيم مجموعة من العملاء في متجر إلكتروني إلى مجموعات استنادًا إلى عادات شرائهم. على سبيل المثال، يمكن تقسيمهم إلى مجموعة من يشترون ملابس رياضية ومجموعة أخرى من يشترون ملابس أنيقة.

#### 3. التجمعات الهرمية: (Hierarchical Clustering)

هذا النوع من التجزئة يقوم بتقسيم المجموعة إلى تجمعات فرعية تدريجيًا بناءً على مستويات مختلفة من التفاعل بين البيانات. يمكن أن يكون التجزئة هرمية متقدمة لتكون تجزئة مفصلة تصل إلى مستويات صغيرة جدًا (وتسمى "التجزئة النمطية"). مثال: تقسيم مجموعة من السكان في منطقة معينة إلى مناطق فرعية بناءً على الموصفات الجغرافية، مثل البلدان أو المدن.

#### 4. التجمعات المركزية: (Centroid Clustering)

في هذا النوع من التجمعات، يتم تقسيم البيانات استنادًا إلى المسافة من نقاط مركزية محددة مسبقًا (مراكز). يتم اختيار هذه المراكز بحيث تمثل مواقع متوسطة لمجموعات البيانات. مثال: تجميع السكان في المدن بناءً على المسافة من وسط المدينة، يتم تجزئة السكان في المدن استنادًا إلى مسافتهم من وسط المدينة. يتم اختيار وسط المدينة كمركز للتجمع، ومن ثم يتم تقسيم السكان إلى مجموعات استنادًا إلى المسافة من هذا المركز.

#### 5. التجمعات الكثافية: (Density-based Clustering)

في هذا النوع من التجمعات، يتم تحديد التجمعات بناءً على كثافة البيانات في المساحة، حيث تكون المناطق ذات كثافة عالية من المعطيات تمثل تجمعات. مثال: تجزئة المناطق الحضرية إلى أحياء، يتم تحديد التجمعات بناءً على كثافة السكان في المناطق الحضرية. يتم تحديد التجمعات عن طريق تحديد المناطق التي تحتوي على كثافة عالية من السكان وتعتبر نقاط انطلاق لتجمعات جديدة.

#### 6. التجمعات الطبيعية: (Natural Clustering)

يشير هذا النوع من التجمعات إلى تجمعات تظهر بشكل طبيعي في البيانات دون وجود تدخل خارجي في عملية التجميع، مثل تجمعات المدن في البيانات الجغرافية. مثال: تجزئة البيانات الجغرافية للمدن والقرى، تظهر تجمعات طبيعية في البيانات الجغرافية للمدن والقرى بناءً على المواقع الجغرافية والعوامل البيئية والاقتصادية المحيطة. على سبيل المثال، يمكن تجزئة البيانات الجغرافية لمنطقة إلى تجمعات طبيعية مثل المدن والقرى والمناطق الريفية.

### 7. التجمعات المناسبة للغرض: (Task-specific Clustering)

في هذا النوع من التجمعات، يتم تكييف عملية التجميع بشكل خاص لحل مشكلة محددة أو تحقيق هدف محدد، مما يتطلب استخدام معايير خاصة ومتخصصة في عملية التجميع.  
مثال: تجزئة العملاء في متجر التجزئة لإعداد حملات تسويقية مستهدفة، يتم تكييف عملية التجميع بشكل خاص لتحليل بيانات العملاء في متجر التجزئة بهدف إعداد حملات تسويقية مستهدفة. يمكن استخدام معايير مثل نمط الشراء والميزانية لتقسيم العملاء إلى مجموعات مستهدفة.

### (3) وضح مشكلتين على الأقل مع خوارزمية K-Means Clustering واضف الحلول المحتملة ؟

#### المشكلة الأولى: النقاط الأولية: (Initial Centroids)

- المشكلة: يعتمد أداء خوارزمية K-Means بشكل كبير على اختيار النقاط المركزية الأولية. إذا تم اختيار النقاط الأولية بشكل عشوائي، فقد يؤدي ذلك إلى وصول الخوارزمية إلى حلول محلية غير مثلى، خاصة في الحالات التي تحتوي على تجمعات مختلفة بأحجام متفاوتة أو شكل غير منتظم.
- الحل: من أجل التغلب على هذه المشكلة، يمكن استخدام تقنية ++K-means لاختيار النقاط المركزية الأولية بشكل ذكي. تقوم هذه التقنية بتحديد النقاط المركزية الأولية بحيث تكون متباعدة بشكل جيد، مما يزيد من فرص الوصول إلى تجزئة أكثر دقة واستقراراً. بالإضافة إلى ذلك، يمكن أيضاً استخدام أساليب أخرى مثل الاستنتاج من بيانات سابقة أو استخدام تقنيات تحليل البيانات لتحديد النقاط المركزية الأولية بشكل أفضل.

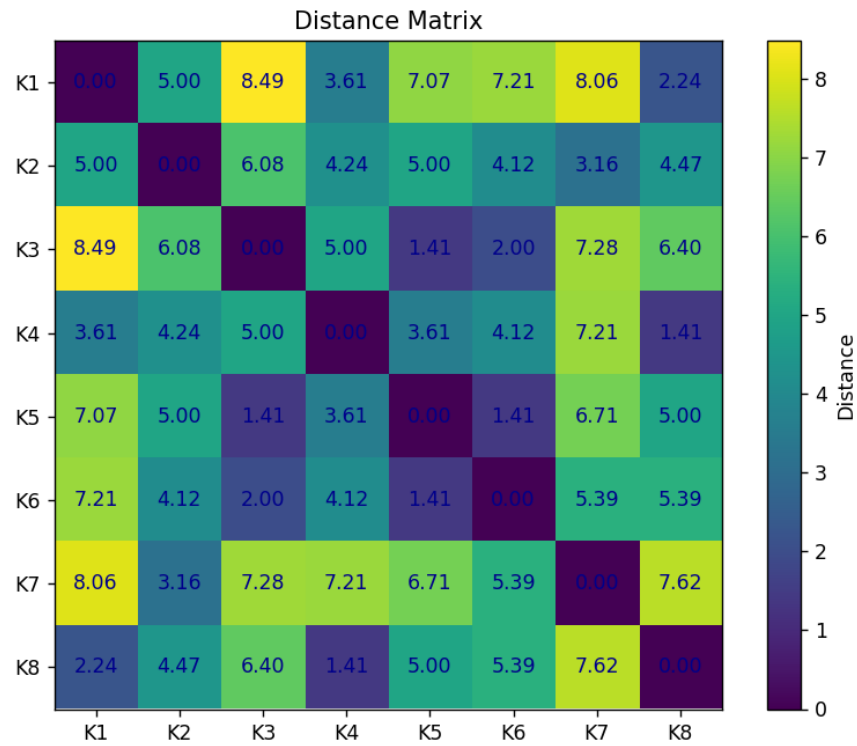
#### المشكلة الثانية: عدد الـ K في K-Means

- المشكلة: قد يكون من الصعب تحديد عدد الـ K الأمثل، أي عدد التجمعات التي يجب تقسيم البيانات إليها. اختيار قيمة غير مناسبة لـ K قد يؤدي إلى تجزئة غير دقيقة للبيانات، حيث يمكن أن تكون التجمعات مفرطة في التجزئة (عدد كبير من التجمعات) أو غير مكتملة (عدد قليل من التجمعات).
- الحل: هناك عدة طرق لتحديد عدد الـ K بشكل مناسب، منها:
  - أسلوب الكوع (Elbow Method): يقوم هذا الأسلوب بتجريب قيم مختلفة لـ K وقياس متوسط مربعات الانحراف (SSE) لكل قيمة من الـ K. يتم اختيار قيمة الـ K التي تظهر فيها منحنى SSE انحناءً حاداً، مما يشير إلى أن تجزئة البيانات بهذه القيمة من الـ K هي الأكثر دقة. ويتم توضيحها بالرسم.
  - أسلوب Silhouette: يستخدم هذا الأسلوب قياس الانفصال بين التجمعات لتحديد الـ K الأمثل. يتم اختيار القيمة التي تحقق أعلى قيمة لمقياس Silhouette، حيث يشير ذلك إلى أن البيانات تكون مجمعة بشكل جيد في التجمعات المحددة.
  - الاستنتاج البصري: قد يكون من المفيد أيضاً فحص البيانات بشكل بصري وتقدير عدد التجمعات بناءً على فهمك للبيانات والهدف من التجزئة.

(4) اوجد المسافة ال Euclidean  $d(K1, K2)$  أ-

$$\text{Sqrt}((2-2)**2+(10-5)**2)=5$$

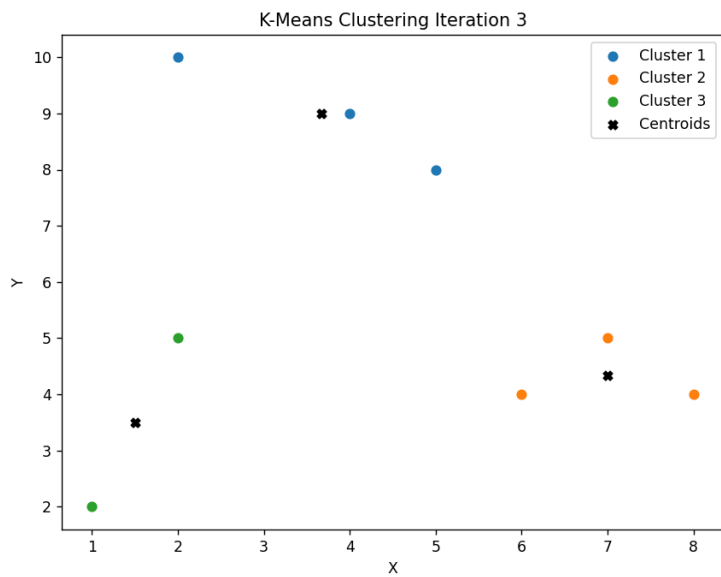
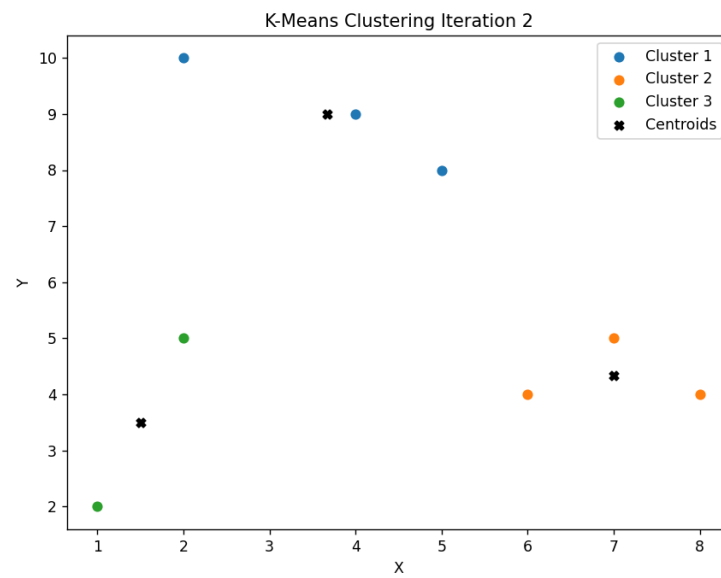
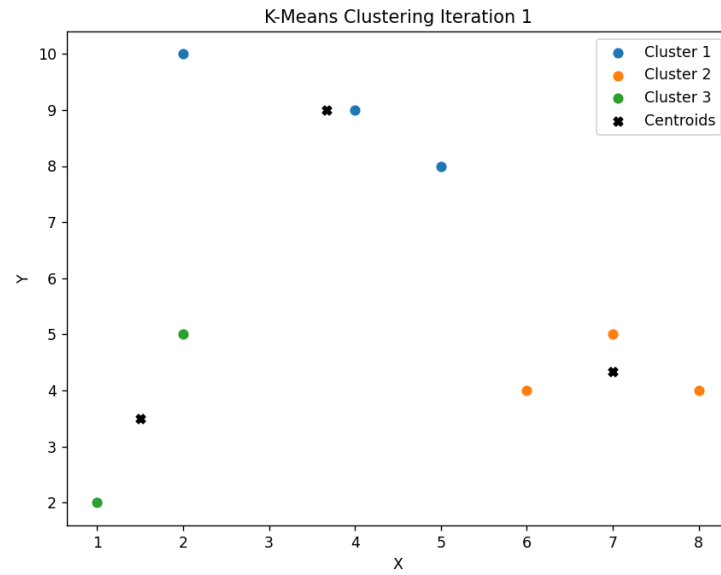
بعد الحل تبين ان المسافة بين (K1, K2) تساوي 5



ب-

في البداية تبدو ان التكرارات هي نفسها ولكن عند التدقيق في الاحداثيات يتبين انها مختلفة قليلا، ومع ذلك فإن التباين بين النقطة الوسطى في التكرارات هو تغير ادنى.

نلاحظ انه من الرغم انه تتشابه بصريا (يصعب ملاحظتها بالعين) الا انها تتغير وتتلاقى في التكرار 2.



## Clustering Analysis Report

### Introduction:

This report provides an analysis of clustering techniques applied to a dataset using the K-Means and hierarchical clustering algorithms. The dataset contains points in a 2-dimensional space represented by the features X and Y.

### Dataset Overview:

The dataset consists of points distributed in a 2-dimensional space. Each point is represented by its coordinates (X, Y). Before clustering, the dataset is visualized to understand its distribution and characteristics.

### K-Means Clustering:

K-Means clustering is applied to the dataset with a predetermined number of clusters ( $n\_clusters = 3$ ). The algorithm iteratively assigns points to the nearest centroid and updates the centroids until convergence. Three iterations are performed to ensure stability in the results. The final centroids and cluster labels are visualized along with the original dataset.

### Hierarchical Clustering:

Hierarchical clustering is applied using three different linkage methods: single, complete, and average. Each method defines the distance between clusters differently, leading to varying cluster structures. The resulting dendrograms are visualized to illustrate the hierarchical clustering process for each method.

### Method Explanations:

- Single Linkage:
  - Defines the distance between two clusters as the shortest distance between any two points in the clusters.
  - Tends to produce elongated clusters sensitive to outliers and noise.
- Complete Linkage:
  - Defines the distance between two clusters as the maximum distance between any two points in the clusters.
  - Tends to produce compact, spherical clusters less sensitive to outliers.
- Average Linkage:
  - Defines the distance between two clusters as the average distance between all pairs of points in the clusters.
  - Balances sensitivity to outliers and cluster compactness, producing moderate compactness.

**Elbow Method:**

The elbow method is employed to determine the optimal number of clusters for K-Means clustering. It calculates the within-cluster sum of squares (inertia) for different numbers of clusters and identifies the point where the inertia starts to decrease at a slower rate, suggesting the optimal number of clusters.

**Conclusion:**

- K-Means clustering and hierarchical clustering with different linkage methods provide insights into the structure of the dataset.
- Each clustering technique has its advantages and limitations, which should be considered based on the specific characteristics of the dataset and the analysis goals.
- Further analysis and interpretation can be conducted based on the clustering results to extract meaningful patterns or insights from the data.

class K-Means:

// تهيئة عدد التجمعات والحد الأقصى لعدد التكرارات

def \_\_init\_\_(self, n\_clusters, max\_iters=1000):

self.n\_clusters = n\_clusters

self.max\_iters = max\_iters

// تهيئة المراكز

def fit(self, data, centroids):

self.centroids = centroids

// إنشاء مصفوفة فارغة لتخزين تعيينات التجمع لكل نقطة بيانات

self.labels = np.zeros(len(data))

// الدوران حتى التقارب أو الوصول إلى max\_iters

for \_ in range(self.max\_iters):

// تعيين كل نقطة بيانات إلى المركز الأقرب لها

for i in range(len(data)):

distances = [self.euclidean\_distance(data[i], centroid) for centroid in self.centroids]

self.labels[i] = np.argmin(distances)

// تحديث المراكز استنادًا إلى المتوسط الحسابي لنقاط البيانات في كل تجمع

for cluster in range(self.n\_clusters):

cluster\_points = data[self.labels == cluster]

if len(cluster\_points) > 0:

self.centroids[cluster] = np.mean(cluster\_points, axis=0)

return self.labels, self.centroids

// حساب المسافة اليوروكليدية بين نقطتين

def euclidean\_distance(self, p1, p2):

return np.sqrt((p1[0] - p2[0])\*\*2 + (p1[1] - p2[1])\*\*2)



```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.cluster import KMeans
5  from scipy.cluster.hierarchy import linkage, dendrogram
6
7  #تضمن وظائف لتنفيذ عمليات التجميع ورسم البيانات والنتائج
8  class Clustering:
9
10     def __init__(self, df):
11         self.df = df
12
13     # رسم مجموعة البيانات قبل التجميع
14     def plot_dataset_before_clustering(self):
15         plt.figure(figsize=(8, 6))
16         plt.scatter(self.df['X'], self.df['Y'], color='red')
17         plt.title('Dataset Before Clustering')
18         plt.xlabel('X')
19         plt.ylabel('Y')
20         plt.show()
21
22     # رسم مصفوفة المسافات
23     def plot_distance_matrix(self, distance_matrix):
24         plt.figure(figsize=(8, 6))
25         plt.imshow(distance_matrix, cmap='viridis')
26         plt.colorbar(label='Distance')
27         plt.title('Distance Matrix')
28         plt.xticks(np.arange(len(self.df)), self.df['Point'])
29         plt.yticks(np.arange(len(self.df)), self.df['Point'])
30
31         for i in range(len(self.df)):
32             for j in range(len(self.df)):
33                 plt.text(j, i, f'{distance_matrix[i, j]:.2f}', ha='center', va='center', color='darkblue')
34         plt.show()
35
36     # رسم التجميعات بعد تطبيق تجميع k-means
37     def plot_clusters(self, centroids, labels, n_clusters):
38         for iteration in range(3):
39             plt.figure(figsize=(8, 6))
40             for cluster in range(n_clusters):
41                 cluster_points = self.df.iloc[labels == cluster]
42                 plt.scatter(cluster_points['X'], cluster_points['Y'], label=f'Cluster {cluster + 1}')
43             plt.scatter(centroids[:, 0], centroids[:, 1], color='black', marker='X', label='Centroids')
44             plt.title(f'K-Means Clustering Iteration {iteration + 1}')
45             plt.xlabel('X')
46             plt.ylabel('Y')
47             plt.legend()
48             plt.show()
49
50     # رسم شجرة التجميع الهرمي
51     def plot_hierarchical_clustering(self, Z, method):
52         plt.figure(figsize=(8, 6))
53         dendrogram(Z)
54         plt.title(f'Hierarchical Clustering ({method.capitalize()} Linkage)')
55         plt.xlabel('Sample Index')

```

```

8 class Clustering:
51 def plot_hierarchical_clustering(self, Z, method):
52     plt.figure(figsize=(8, 6))
53     dendrogram(Z)
54     plt.title(f'Hierarchical Clustering ({method.capitalize()} Linkage)')
55     plt.xlabel('Sample Index')
56     plt.ylabel('Distance')
57     plt.show()
58
59 # حساب مصفوفة المسافات
60 def calculate_distance_matrix(self):
61     distance_matrix = np.zeros((len(self.df), len(self.df)))
62
63     # حساب المسافة اليوروكليدية بين نقطتين p1 و p2 في الفضاء
64     def euclidean_distance(p1, p2):
65         return np.sqrt((p1[0] - p2[0])**2 + (p1[1] - p2[1])**2)
66
67     for i, row1 in self.df.iterrows():
68         for j, row2 in self.df.iterrows():
69             distance_matrix[i, j] = euclidean_distance((row1['X'], row1['Y']), (row2['X'], row2['Y']))
70
71     return distance_matrix
72
73 # تطبيق تجميع k-means
74 def kmeans_clustering(self, n_clusters, centroids, n_iterations=3, convergence_threshold=1e-4):
75
76     for iteration in range(n_iterations):
77         kmeans = KMeans(n_clusters=n_clusters, init=centroids, n_init=1)
78         labels = kmeans.fit_predict(self.df[['X', 'Y']])
79         new_centroids = kmeans.cluster_centers_
80
81         if np.allclose(centroids, new_centroids, atol=convergence_threshold):
82             print(f"Converged after {iteration + 1} iterations.")
83             break
84
85         centroids = new_centroids
86
87     return centroids, labels
88
89 # تطبيق تجميع هرمي
90 def hierarchical_clustering(self, method):
91     X = self.df[['X', 'Y']].values
92     Z = linkage(X, method=method)
93     return Z
94
95 # حساب elbow curve للعثور على العدد الأمثل للمجموعات باستخدام طريقة
96 def calculate_elbow(self, max_clusters=10):
97     # قائمة لتخزين التكاليف لعدد مختلف من
98     costs = []
99
100     for num_clusters in range(1, max_clusters + 1):
101         kmeans = KMeans(n_clusters=num_clusters)
102         kmeans.fit(self.df[['X', 'Y']])
103         # Inertia هو مجموع المسافات المربعة للعينات إلى أقرب مركز cluster لها
104         costs.append(kmeans.inertia_)

```

```

Assignment6-SulimanAlgaramanli.py
C: > Users > sulim > (method) def calculate_elbow(
                        self: Self@Clustering,
                        max_clusters: int = 10
                        ) -> None
121
122 clustering = Clustering(df)
123
124 # رسم مجموعة البيانات قبل التجميع
125 clustering.plot_dataset_before_clustering()
126
127 # حساب ورسم مصفوفة المسافات
128 distance_matrix = clustering.calculate_distance_matrix()
129 clustering.plot_distance_matrix(distance_matrix)
130
131 # بعلامات محددة k-means تطبيق تجميع
132 n_clusters = 3
133 initial_centroids = df.loc[df['Point'].isin(['K1', 'K4', 'K7']), ['X', 'Y']].values
134 centroids, labels = clustering.kmeans_clustering(n_clusters, initial_centroids)
135
136 # رسم التجميعات
137 clustering.plot_clusters(centroids, labels, n_clusters)
138
139 # تطبيق تجميع هرمي بطرق مختلفة
140 methods = ['single', 'complete', 'average']
141 for method in methods:
142     Z = clustering.hierarchical_clustering(method)
143     clustering.plot_hierarchical_clustering(Z, method)
144
145 # حساب ورسم منحنى elbow
146 clustering.calculate_elbow(max_clusters=8)
147

```