



ÉCOLE NATIONALE DES SCIENCES APPLIQUÉES AGADIR

RAPPORT D'ANALYSE DES DONNÉES

Analyse de Classification

Élève :

Salma ELAISSARI

Enseignant :

Prof.BrahimEL ASRI

Filière:

FID

16 avril 2024

1

La classification est l'une des tâches fondamentales en analyse de données et en apprentissage automatique. Elle consiste à regrouper des éléments similaires dans des catégories ou des classes distinctes en fonction de leurs caractéristiques ou de leurs attributs. L'objectif principal de la classification est de prédire l'appartenance d'un nouvel élément à une classe donnée, en se basant sur les caractéristiques observées dans les données d'entraînement.

L'importance de la classification réside dans sa capacité à fournir des informations utiles et exploitables à partir de données brutes. Voici quelques-unes des raisons pour lesquelles la classification est largement utilisée en analyse de données :

- **Organisation des données** : La classification permet de structurer et d'organiser des ensembles de données complexes en regroupant des éléments similaires dans des catégories distinctes.
- **Prédiction et reconnaissance** : En identifiant des modèles dans les données, la classification permet de prédire ou de reconnaître des éléments inconnus et de prendre des décisions éclairées sur leur classification.
- **Segmentation de marché** : Dans le domaine du marketing, la classification est utilisée pour segmenter les clients en groupes homogènes afin de mieux cibler les offres et les campagnes publicitaires.

En résumé, la classification joue un rôle crucial dans l'analyse de données en permettant d'extraire des informations significatives à partir de données complexes, ce qui facilite la prise de décisions éclairées dans de nombreux domaines d'application.

Les méthodes de classification hiérarchique et de k-means sont deux techniques fondamentales largement utilisées dans divers domaines tels que la science des données, la finance, le marketing, etc., pour extraire des informations significatives à partir de données brutes.

La classification hiérarchique est une approche qui divise progressivement les données en clusters plus petits, formant ainsi une structure arborescente ou dendrogramme. Cette méthode peut être divisée en deux approches principales : la classification hiérarchique agglomérative et la classification hiérarchique divisive. La première commence par considérer chaque point de données comme un cluster séparé et fusionne ensuite progressivement les clusters les plus proches, tandis que la seconde commence par considérer tous les points de données comme un seul cluster et divise ensuite récursivement le cluster en clusters plus petits.

D'autre part, l'algorithme k-means est une méthode itérative de partitionnement de données en k clusters, où k est un nombre prédéfini. Cet algorithme commence par initialiser aléatoirement les centroïdes des clusters, puis il alterne entre deux étapes : attribution

des points de données au cluster dont le centroïde est le plus proche et mise à jour des centroïdes en calculant les moyennes des points de données dans chaque cluster. Ce processus est répété jusqu'à ce qu'une convergence soit atteinte, c'est-à-dire que les centroïdes ne changent pas de manière significative entre les itérations successives.

La classification hiérarchique est souvent préférée lorsque la structure hiérarchique des clusters est importante et que le nombre de clusters n'est pas prédéfini, tandis que k-means est plus adapté lorsque le nombre de clusters est connu à l'avance et que les clusters sont globulaires et bien séparés.

1.3 Objectifs du rapport

L'objectif principal de ce rapport est de présenter en détail le projet de développement d'un programme de classification utilisant les méthodes de classification hiérarchique (CAH) et de k-means en Python. Les objectifs spécifiques sont les suivants :

- Démontrer l'utilisation des fonctionnalités interactives de l'interface utilisateur pour faciliter l'analyse et la visualisation des données.
- Expliquer le choix des graphiques et des visualisations utilisés pour représenter les résultats de la classification, en mettant en évidence leur pertinence et leur utilité.

1.4 Description du projet

Le projet consiste à développer un programme de classification en Python, offrant aux utilisateurs une interface conviviale pour l'analyse de données à l'aide des méthodes de CAH et de k-means. Voici une description détaillée du projet :

- **Fonctionnalités du programme** : Le programme permet aux utilisateurs de charger des ensembles de données à partir de fichiers Excel, d'appliquer les algorithmes de classification (CAH et k-means) et d'explorer les résultats à travers des visualisations interactives.
- **Interface utilisateur intuitive** : L'interface utilisateur offre une expérience conviviale, permettant aux utilisateurs de paramétrer les algorithmes, d'interagir avec les données et de visualiser les résultats de manière intuitive.
- **Documentation du code** : Le rapport fournira une documentation complète du code source du programme, expliquant chaque fonction, classe et module utilisé, ainsi que leur contribution à la fonctionnalité globale du programme.
- **Exemples d'utilisation** : Le rapport présentera des exemples d'utilisation du programme avec différents ensembles de données, démontrant les capacités d'analyse et de visualisation offertes par les méthodes de CAH et de k-means.
- **Analyse des résultats** : En plus de présenter les fonctionnalités du programme, le rapport analysera les résultats obtenus avec différents ensembles de données, mettant en évidence les avantages et les limitations des méthodes de classification utilisées.

2 Méthodes de Classification

2.1 Explication détaillée de la CAH :

2.1.1 Principe de fonctionnement

La classification hiérarchique est une méthode de classification itérative qui permet de regrouper des individus en différentes classes selon un critère de ressemblance défini au préalable. Cette méthode est représentée par un dendrogramme, qui montre les différents niveaux de regroupement des individus.

La classification est ascendante car elle part des observations individuelles, et elle est hiérarchique car elle produit des classes ou groupes de plus en plus vastes, incluant des sous-groupes en leur sein. En découpant cet arbre à une certaine hauteur choisie, on produira la partition désirée.

Le critère de ressemblance, utilisé dans la classification hiérarchique, est calculé en se basant sur les mesures de distance ou de similarité entre les individus. Ces mesures peuvent varier en fonction du type de données et du contexte de l'analyse. Voici quelques exemples de calcul de critères de ressemblance couramment utilisés :

- **Distance euclidienne** : Pour des données numériques continues, la distance euclidienne entre deux individus i et j peut être calculée comme suit :

$$d_{ij} = \sqrt{\sum_{k=1}^n (x_{ik} - x_{jk})^2}$$

où x_{ik} et x_{jk} sont les valeurs des variables pour les individus i et j respectivement, et n est le nombre de variables.

Le **coefficient de corrélation** est une mesure statistique qui exprime à quel point deux ensembles de données sont liés linéairement les uns aux autres. Pour calculer le coefficient de corrélation, nous utilisons généralement le coefficient de corrélation de Pearson, qui est largement utilisé pour les données numériques continues. Voici comment il est calculé :

- **Calcul des moyennes** : Calculer les moyennes des ensembles de données X et Y , notées respectivement par \bar{x} et \bar{y} :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

- **Calcul des écarts-types** : Calculer les écarts-types des ensembles de données X et Y , notés respectivement par s_x et s_y :

$$s_x = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$s_y = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

- **Calcul de la covariance** : Calculer la covariance entre les ensembles de données X et Y , notée $\text{cov}(X, Y)$:

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- **Calcul du coefficient de corrélation de Pearson** : Utiliser la formule suivante pour calculer le coefficient de corrélation r :

$$r = \frac{\text{cov}(X, Y)}{s_x s_y}$$

Le coefficient de corrélation de Pearson varie entre -1 et 1. Un coefficient de 1 indique une corrélation linéaire positive parfaite, -1 indique une corrélation linéaire négative parfaite, et 0 indique l'absence de corrélation linéaire entre les ensembles de données.

Matrice de distance

Choix d'une mesure de distance : Sélectionnez une mesure de distance appropriée en fonction de la nature de vos données. Les mesures de distance couramment utilisées incluent :

- Distance euclidienne
- Distance de Manhattan (ou distance en "city block")

Calcul des distances entre les paires d'observations : Pour chaque paire d'observations i et j dans votre ensemble de données, calculez la distance selon la mesure choisie. Par exemple, si vous utilisez la distance euclidienne entre deux observations \mathbf{x}_i et \mathbf{x}_j dans un espace n -dimensionnel, la formule est :

Vous répétez le calcul de distance euclidienne pour toutes les paires d'observations pour construire la matrice de distance.

Construction de la matrice de distance : Une fois que vous avez calculé les distances entre toutes les paires d'observations, vous les stockez dans une matrice. Cette matrice est appelée matrice de distance et a la forme d'une matrice carrée $n \times n$, où n est le nombre d'observations dans votre ensemble de données. Chaque élément de la matrice représente la distance entre deux observations.

Construction de l'arbre de clustering :

Ensuite, un algorithme de clustering est appliqué pour regrouper les individus en clusters en fonction de leurs similarités mesurées par la matrice de distance. L'algorithme de clustering le plus couramment utilisé pour la création de dendrogrammes est l'algorithme de clustering hiérarchique.

Construction du dendrogramme :

Le dendrogramme est construit en utilisant l'arbre de clustering résultant. Chaque individu est représenté par une feuille dans le dendrogramme, et les branches de l'arbre indiquent les regroupements successifs des individus en clusters. La longueur des branches représente la distance entre les individus ou les clusters. Plus la distance est grande, moins les individus sont similaires.

Interprétation du dendrogramme :

Le dendrogramme peut être interprété en coupant les branches à différentes hauteurs pour former des clusters. En coupant le dendrogramme à différentes hauteurs, vous

En résumé, le dendrogramme est une représentation graphique qui permet de visualiser les relations de regroupement entre les individus dans un ensemble de données, en utilisant une approche hiérarchique basée sur la distance entre les individus.

- **Flexibilité dans le choix du type de dissimilarité** : La CAH permet le choix d'un type de dissimilarité adapté au sujet d'étude et à la nature des données. Cette flexibilité permet d'appliquer la CAH à différents types de données et contextes.
- **Visualisation de la progression du regroupement** : La CAH produit un dendrogramme qui permet de visualiser la progression du regroupement des données. Cette visualisation facilite la détermination d'un nombre approprié de classes dans lesquelles les données peuvent être regroupées.
- **Adaptabilité à différents types de données** : La CAH peut être appliquée à différents types de données, telles que les données quantitatives ou qualitatives, et peut être adaptée à des contextes spécifiques, tels que le marketing, les médias ou l'économie.
- **Représentation hiérarchique** : La CAH produit une représentation hiérarchique des données, ce qui permet d'identifier différents niveaux de granularité dans la classification. Cette représentation hiérarchique peut être utile pour identifier des sous-groupes au sein de groupes plus larges ou pour identifier des clusters à différents niveaux de similarité.
- **Choix de la méthode d'agrégation** : La CAH permet le choix de la méthode d'agrégation, qui peut être adaptée au contexte spécifique des données et aux objectifs de l'analyse. Le choix de la méthode d'agrégation peut avoir un impact sur les résultats de l'analyse et l'interprétation des données.

- **Choix de la mesure de dissimilarité** : Bien que la CAH permette le choix d'une mesure de dissimilarité, ce choix peut être subjectif et peut affecter les résultats. D'autres méthodes de classification, telles que k-means, ne nécessitent pas le choix d'une mesure de dissimilarité.
- **Sensibilité aux valeurs aberrantes** : La CAH est sensible aux valeurs aberrantes, ce qui peut affecter significativement les résultats. D'autres méthodes de classification, telles que k-medoids, sont plus robustes aux valeurs aberrantes.

- **Complexité computationnelle** : La CAH peut être coûteuse en termes de calcul, surtout pour de grands ensembles de données. D'autres méthodes de classification, comme k-means, sont généralement plus rapides et plus évolutives.
- **Difficulté à choisir le nombre de clusters** : La CAH ne fournit pas de moyen clair de déterminer le nombre optimal de clusters. D'autres méthodes de classification, telles que la méthode du coude ou la méthode du silhouette, peuvent aider à déterminer le nombre optimal de clusters.
- **Absence d'interprétation probabiliste** : La CAH ne fournit pas d'interprétation probabiliste des clusters, ce qui peut limiter son utilité dans certaines applications. D'autres méthodes de classification, comme les modèles de mélange gaussien, fournissent une interprétation probabiliste des clusters.
- **Limitation aux structures hiérarchiques** : La CAH est limitée aux structures hiérarchiques, ce qui peut ne pas être approprié pour tous les ensembles de données. D'autres méthodes de classification, comme k-means, n'imposent pas une structure hiérarchique aux données.

2.1.4 Principes et Fonctionnement du code de CAH sur Python :

2.1.5 Packages utilisés :

```
1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from scipy.cluster.hierarchy import linkage, dendrogram
6 from sklearn.metrics import pairwise_distances
```

Listing 1 – python script1

Ce code Python utilise plusieurs packages pour effectuer différentes tâches. Voici le rôle de chaque package utilisé :

numpy (np) : NumPy est une bibliothèque Python utilisée pour effectuer des calculs numériques. Dans ce code, NumPy est utilisé pour manipuler des tableaux et effectuer des calculs matriciels, tels que le calcul des distances euclidiennes entre les points de données. **matplotlib.pyplot (plt)** : Matplotlib est une bibliothèque de visualisation en Python. Ici, le sous-module pyplot de Matplotlib est utilisé pour créer des graphiques, notamment pour afficher le dendrogramme généré par le clustering hiérarchique.

pandas (pd) : Pandas est une bibliothèque Python utilisée pour la manipulation et l'analyse des données. Dans ce code, Pandas est utilisé pour importer les données à partir d'un fichier Excel et les manipuler sous forme de DataFrame, ce qui facilite le travail avec les données tabulaires.

`from scipy.cluster.hierarchy import linkage, dendrogram` : Scipy est une bibliothèque Python qui contient des outils pour le calcul scientifique et l'analyse de données. Elle offre plusieurs sous-modules, dont `scipy.cluster.hierarchy`, qui contient des fonctions pour le clustering hiérarchique, les fonctions `linkage` et `dendrogram` de Scipy effectuent le clustering hiérarchique et visualisent les résultats sous forme de

2.1.6 Classes utilisés :

8


```

1 def dendrogram_visualization(data, labels):
2     X = data.values
3     Z = linkage(X, method='ward')
4     fig, ax = plt.subplots(figsize=(10, 5))
5     dn = dendrogram(Z, ax=ax)
6     ax.set_title("Dendrogram")
7     ax.set_xlabel("Data Points")
8     ax.set_ylabel("Distance")
9     st.pyplot(fig) # Display the plot in Streamlit

```

Listing 4 – python script1

dendrogram-visualization(data, labels) : La fonction est utilisée pour visualiser un dendrogramme à partir des données de clustering et des étiquettes de cluster associées. Voici un aperçu de ce que fait cette fonction :

Z = linkage(X, method='ward') : il calcule la matrice de liaison à partir des données X.

dn = dendrogram(Z, ax=ax) : il trace le dendrogramme à partir de la matrice de liaison Z sur les axes spécifiés ax.

2.2 Explication détaillée de la méthode kmeans :

2.2.1 Principe de fonctionnement :

K-means est un algorithme de clustering populaire qui vise à partitionner un ensemble de données en K clusters distincts et non chevauchants. L'algorithme affecte itérativement chaque point de données au centroïde le plus proche, puis met à jour les centroïdes en fonction de la moyenne des points de données assignés à chaque cluster. L'algorithme fonctionne comme suit :

— **1-Initialisation** :

- K centroïdes sont initialisés de manière aléatoire à partir de l'ensemble de données.

— **2-Affectation** :

- Chaque point de données est assigné au centroïde le plus proche en fonction d'une distance choisie, telle que la distance euclidienne et la distance de Manhattan. A noter : Soit un espace n-dimensionnel, la distance entre deux points p et q est comme suit :

$$\text{Distance euclidienne}(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

$$\text{Distance de Manhattan(City Block)}(p, q) = \sqrt{\sum_{i=1}^n |p_i - q_i|}$$

- Assignation de l'observation : L'observation est assignée au cluster dont le centroid est le plus proche, en fonction de la distance calculée.

3-Mise à jour :

- Les centroïdes sont mis à jour en calculant la moyenne de tous les points de données assignés à chaque cluster.

4-Itération :

- Les étapes 2 et 3 sont répétées jusqu'à convergence, c'est-à-dire lorsque les centroïdes ne changent plus significativement ou qu'un nombre maximal d'itérations est atteint.

L'algorithme K-Means est sensible au choix initial des centroids et peut converger vers un optimum local. Pour atténuer cet effet, des techniques telles que l'initialisation K-Means++ peuvent être utilisées pour initialiser les centroids de manière plus intelligente. De plus, le choix du nombre de clusters k est souvent déterminé en utilisant des méthodes.

2.2.2 Avantages de K-means :

K-means est un algorithme de clustering qui présente plusieurs avantages qui en font un choix populaire pour l'analyse de données et les tâches d'apprentissage automatique.

- **Facilité d'utilisation** : K-means est un algorithme simple et facile à mettre en œuvre qui ne nécessite pas de calculs mathématiques complexes ou de compétences en programmation avancées.

2.2.6 Classes utilisés :

```

1 class KMeansClustering:
2     def __init__(self, k=3, distance_metric='euclidean', random_state=
None):
3         self.k = k
4         self.distance_metric = distance_metric
5         self.centroids = None
6         self.random_state = random_state
7
8     def calculate_distance(self, data_point, centroids):
9         if self.distance_metric == 'euclidean':
10             return np.sqrt(np.sum((centroids - data_point) ** 2, axis=1)
)
11         elif self.distance_metric == 'city_block':
12             return np.sum(np.abs(centroids - data_point), axis=1)
13         else:
14             raise ValueError("Invalid distance metric. Supported metrics
are 'euclidean' and 'city_block'.")
15
16     def fit(self, X, max_iterations=200):
17         np.random.seed(self.random_state) # Set the random seed
18
19         # Initialize centroids randomly
20         self.centroids = X[np.random.choice(X.shape[0], self.k, replace=
False)]
21
22         inertia_values = []
23
24         for _ in range(max_iterations):
25             # Assign each data point to the nearest centroid
26             distances = np.vstack([self.calculate_distance(data_point,
self.centroids) for data_point in X])
27             cluster_assignment = np.argmin(distances, axis=1)
28
29             # Update centroids
30             new_centroids = np.array([X[cluster_assignment == i].mean(
axis=0) for i in range(self.k)])
31
32             # Calculate inertia
33             inertia = np.sum((X - new_centroids[cluster_assignment]) **
2)
34             inertia_values.append(inertia)
35
36             # Check convergence
37             if np.allclose(self.centroids, new_centroids):
38                 break
39
40             self.centroids = new_centroids
41
42             # Calculate silhouette score
43             labels = np.argmin(distances, axis=1)
44             silhouette_avg = silhouette_score(X, labels)
45
46             return cluster_assignment, self.centroids, inertia_values,
silhouette_avg

```

Listing 6 – python script1

- `__init__` : Initialise la classe avec les paramètres tels que le nombre de clusters, la métrique de distance, etc.
- `calculate_distance` : Calcule la distance entre un point de données donné et tous les centroids actuels.
- `fit` : Cette méthode effectue l'algorithme de clustering K-Means. Elle attribue itérativement chaque point de données au centroïde le plus proche, met à jour les centroïdes en fonction de la moyenne des points de données attribués, et répète le processus jusqu'à ce qu'il y ait convergence ou jusqu'à ce que le nombre maximal d'itérations (max-iterations) soit atteint. Elle calcule également les valeurs d'inertie (somme des distances au carré des échantillons jusqu'au centre de leur cluster le plus proche), retourne les affectations de cluster, les centroïdes finaux, les valeurs d'inertie au fil des itérations, et le score de silhouette comme mesure de qualité du clustering.

```

1 def run_clustering(X, k, distance_metric, use_pca):
2     if use_pca:
3         pca = PCA(n_components=2)
4         X = pca.fit_transform(X)
5
6     kmeans = KMeansClustering(k=k, distance_metric=distance_metric,
7                               random_state=42)
8     labels, centroids, inertia_values, silhouette_avg = kmeans.fit(X)
9     return labels, centroids, inertia_values, silhouette_avg
10
11 def plot_result(X_reduced, labels, centroids):
12     if X_reduced.shape[1] == 2:
13         df = pd.DataFrame(X_reduced, columns=['Component 1', 'Component
14         2'])
15     else:
16         df = pd.DataFrame(X_reduced[:, :2], columns=['Component 1', '
17         Component 2'])
18     df['Cluster'] = labels.astype(str) # Convert labels to string for
19     coloring purposes
20
21     fig = px.scatter(df, x='Component 1', y='Component 2', color='
22     Cluster',
23                     title='Clustering Result', hover_name=df.index)
24     fig.add_trace(go.Scatter(
25         x=centroids[:, 0],
26         y=centroids[:, 1],
27         mode='markers',
28         marker=dict(size=10, color='black', symbol='star'),
29         name='Centroids'
30     ))
31     return fig
32
33 def plot_convergence(inertia_values):
34     fig = go.Figure()
35     fig.add_trace(go.Scatter(

```

```

31     x=list(range(1, len(inertia_values) + 1)),
32     y=inertia_values,
33     mode='lines+markers',
34     marker=dict(color='blue'),
35     name='Inertia'
36 ))
37 fig.update_layout(
38     title='Convergence Plot',
39     xaxis=dict(title='Iteration'),
40     yaxis=dict(title='Inertia'),
41     hovermode='closest'
42 )
43 return fig
44
45 def plot_silhouette_score(silhouette_values):
46     fig = go.Figure()
47     fig.add_trace(go.Scatter(
48         x=list(range(2, 12)),
49         y=silhouette_values,
50         mode='lines+markers',
51         marker=dict(color='green'),
52         name='Silhouette Score'
53     ))
54     fig.update_layout(
55         title='Silhouette Score',
56         xaxis=dict(title='Number of Clusters (K)'),
57         yaxis=dict(title='Silhouette Score'),
58         hovermode='closest'
59     )
60     return fig

```

Listing 7 – python script1

Fonction plot_result : La fonction `plot_result` génère un graphique interactif de la visualisation des résultats de clustering. Elle prend en entrée les données réduites en dimensionnalité, les affectations de cluster et les centroids finaux.

Fonction run_clustering : La fonction `run_clustering` exécute le processus complet de clustering K-means sur les données fournies. Elle prend en entrée les données, le nombre de clusters, la métrique de distance à utiliser et un indicateur pour indiquer si une PCA doit être utilisée.

Fonction plot_convergence : La fonction `plot_convergence` génère un graphique interactif montrant la convergence de l'algorithme K-means au fil des itérations. Elle prend en entrée les valeurs d'inertie à chaque itération.

Fonction plot_silhouette_score : La fonction `plot_silhouette_score` génère un graphique interactif montrant l'évolution du score de silhouette en fonction du nombre de clusters. Elle prend en entrée les valeurs de score de silhouette pour différents nombres de clusters.

```

1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 from sklearn.decomposition import PCA
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.metrics import silhouette_score
7 import plotly.graph_objs as go
8 import plotly.express as px
9 import matplotlib.pyplot as plt
10 from scipy.cluster.hierarchy import linkage, dendrogram # Modified
11     import statement
12 import base64
13 import requests
14
15 # Set background image
16 background_image_url = "https://img.freepik.com/free-vector/ai-
17     technology-brain-background-vector-digital-transformation-
18     concept_53876-117820.jpg?w=900&t=st=1713012632~exp=1713013232~hmac=05
19     da65e2d9e6da77202ded9006cfecb86c0c11fbc70346fb0532307b18d8ac3f"
20
21 def get_base64_of_bin_file(url):
22     data = requests.get(url).content
23     return base64.b64encode(data).decode()
24
25 def set_png_as_page_bg(png_file):
26     bin_str = get_base64_of_bin_file(png_file)
27     st.markdown(
28         f'<style>
29         f'.stApp {{
30         f'background-image: url("data:image/png;base64,{bin_str}")';
31         f'background-size: cover;
32         f'}}'
33         f'</style>',
34         unsafe_allow_html=True
35     )
36
37 set_png_as_page_bg(background_image_url)
38
39 class KMeansClustering:
40     def __init__(self, k=3, distance_metric='euclidean', random_state=
41         None):
42         self.k = k
43         self.distance_metric = distance_metric
44         self.centroids = None
45         self.random_state = random_state
46
47     def calculate_distance(self, data_point, centroids):
48         if self.distance_metric == 'euclidean':
49             return np.sqrt(np.sum((centroids - data_point) ** 2, axis=1))
50         elif self.distance_metric == 'city_block':

```

95

```

96         # Initialize clusters
97         clusters = [[i] for i in range(n)]
98
99         # Hierarchical clustering algorithm
100        while len(clusters) > self.k:
101            min_dist = np.inf
102            for i in range(len(clusters)):
103                for j in range(i + 1, len(clusters)):
104                    cluster1 = clusters[i]
105                    cluster2 = clusters[j]
106                    avg_dist = 0
107                    count = 0
108                    for idx1 in cluster1:
109                        for idx2 in cluster2:
110                            if not np.isnan(distances[idx1, idx2]):
111                                avg_dist += distances[idx1, idx2]
112                                count += 1
113                    if count > 0:
114                        avg_dist /= count
115                        if avg_dist < min_dist:
116                            min_dist = avg_dist
117                            merge_index = (i, j)
118            i, j = merge_index
119            new_cluster = clusters[i] + clusters[j]
120            clusters[i] = new_cluster
121            clusters.pop(j)
122
123        self.labels_ = np.zeros(n)
124        for i, cluster in enumerate(clusters):
125            for idx in cluster:
126                self.labels_[idx] = i
127
128        def euclidean_distance(self, x1, x2):
129            x1 = np.array(x1, dtype=np.float64) # Convert x1 to numpy array
130            x2 = np.array(x2, dtype=np.float64) # Convert x2 to numpy array
131            if np.any(np.isnan(x1)) or np.any(np.isnan(x2)):
132                return np.nan
133            return np.sqrt(np.sum((x1 - x2) ** 2))
134
135        def plot_dendrogram(data):
136            Z = linkage(data, method='ward') # Use linkage from scipy.cluster.
137            hierarchy directly
138            plt.figure(figsize=(10, 5))
139            dn = dendrogram(Z)
140            plt.title("Dendrogram")
141            plt.xlabel("Data Points")
142            plt.ylabel("Distance")
143            st.pyplot(plt.gcf()) # Pass the current figure as an argument to st
144            .pyplot()
145
146        def plot_result(X_reduced, labels, centroids):
147            if X_reduced.shape[1] == 2:
148                df = pd.DataFrame(X_reduced, columns=['Component 1', 'Component
149                2'])
150            else:

```

200

```

201     if selected == "KMeans Clustering":
202         st.title("KMeans Clustering")
203         uploaded_file = st.file_uploader("Upload Excel file", type=["
xlsx", "xls"])
204         if uploaded_file is not None:
205             data = pd.read_excel(uploaded_file, index_col=0)
206             st.write("Data imported successfully!")
207
208             X = data.values
209             X = X[:, 2:] # Adjust this according to your data
210
211             scaler = StandardScaler()
212             X_scaled = scaler.fit_transform(X)
213
214             k = st.sidebar.slider("Number of Clusters (K)", min_value=2,
max_value=10, value=3)
215             distance_metric = st.sidebar.radio("Distance Metric", ["
Euclidean", "City Block"], index=0)
216             distance_metric = 'euclidean' if distance_metric == '
Euclidean' else 'city_block'
217
218             use_pca = st.sidebar.checkbox("Use PCA")
219
220             labels, centroids, inertia_values, silhouette_avg = None,
None, None, None # Initialize variables
221
222             if st.sidebar.button("Cluster"):
223                 kmeans = KMeansClustering(k=k, distance_metric=
distance_metric, random_state=42)
224                 labels, centroids, inertia_values, silhouette_avg =
kmeans.fit(X_scaled)
225
226                 st.success(f"Convergence achieved in {len(inertia_values
)} iterations.")
227                 st.write(f"Number of clusters: {k}")
228                 st.write(f"Number of data points: {len(X)}")
229                 st.write(f"Number of features: {X.shape[1]}")
230                 st.write(f"Silhouette Score: {silhouette_avg:.4f}")
231
232
233                 plot_choice = st.sidebar.selectbox("Select Plot", ["
Clustering Result", "Convergence Plot", "Silhouette Score"])
234
235                 if plot_choice == "Clustering Result" and labels is not None
:
236                     fig = plot_result(X_scaled, labels, centroids)
237                     st.plotly_chart(fig, use_container_width=True)
238                 elif plot_choice == "Convergence Plot" and inertia_values is
not None:
239                     fig = plot_convergence(inertia_values)
240                     st.plotly_chart(fig, use_container_width=True)
241                 elif plot_choice == "Silhouette Score" and silhouette_avg is
not None:
242                     silhouette_values = []
243                     for k in range(2, 12):
244                         kmeans = KMeansClustering(k=k, random_state=42)
245                         _, _, _, silhouette_avg = kmeans.fit(X_scaled)

```


- Ces packages sont utilisés pour différentes fonctionnalités de l'application de clustering, telles que le chargement des données, la manipulation et l'analyse des données, la création de visualisations, le clustering et l'évaluation des modèles.

Sélection de l'option : Lorsque vous ouvrez l'application, vous êtes invité à sélectionner une option parmi celles disponibles, telles que "KMeans Clustering" ou "Hierarchical Clustering".

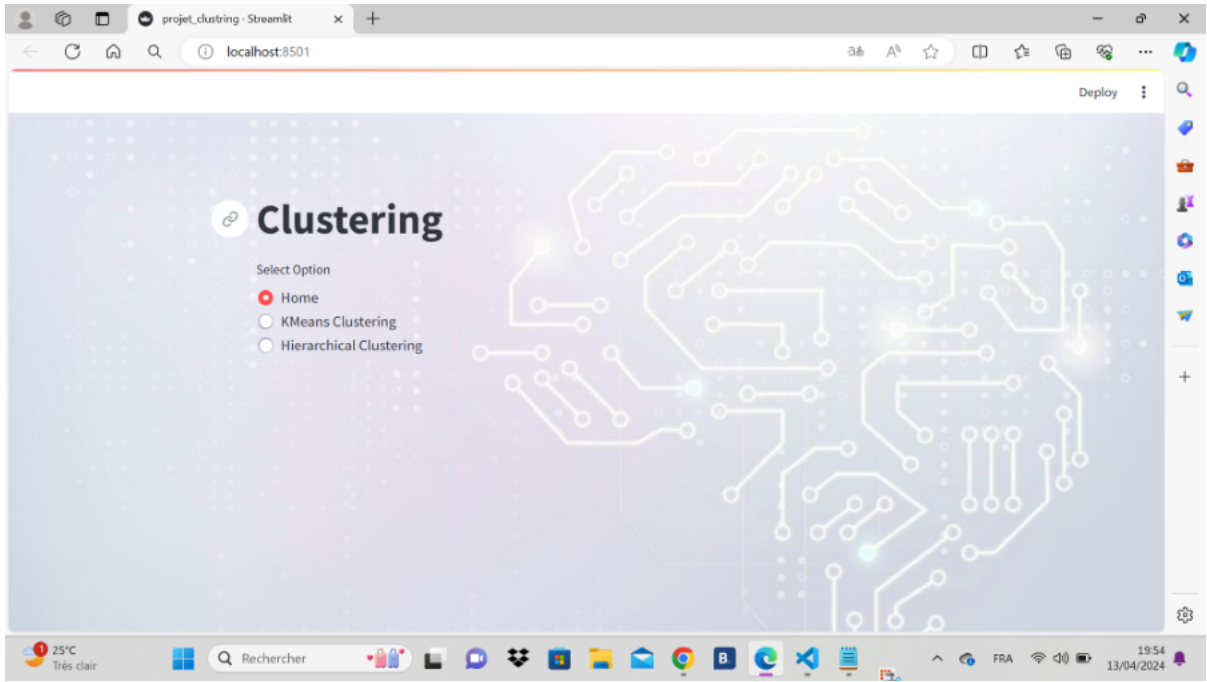


FIGURE 1 – étape 1

Chargement des données : Si vous choisissez l'option "KMeans Clustering", vous pouvez charger vos données en utilisant le bouton "Upload Excel file".

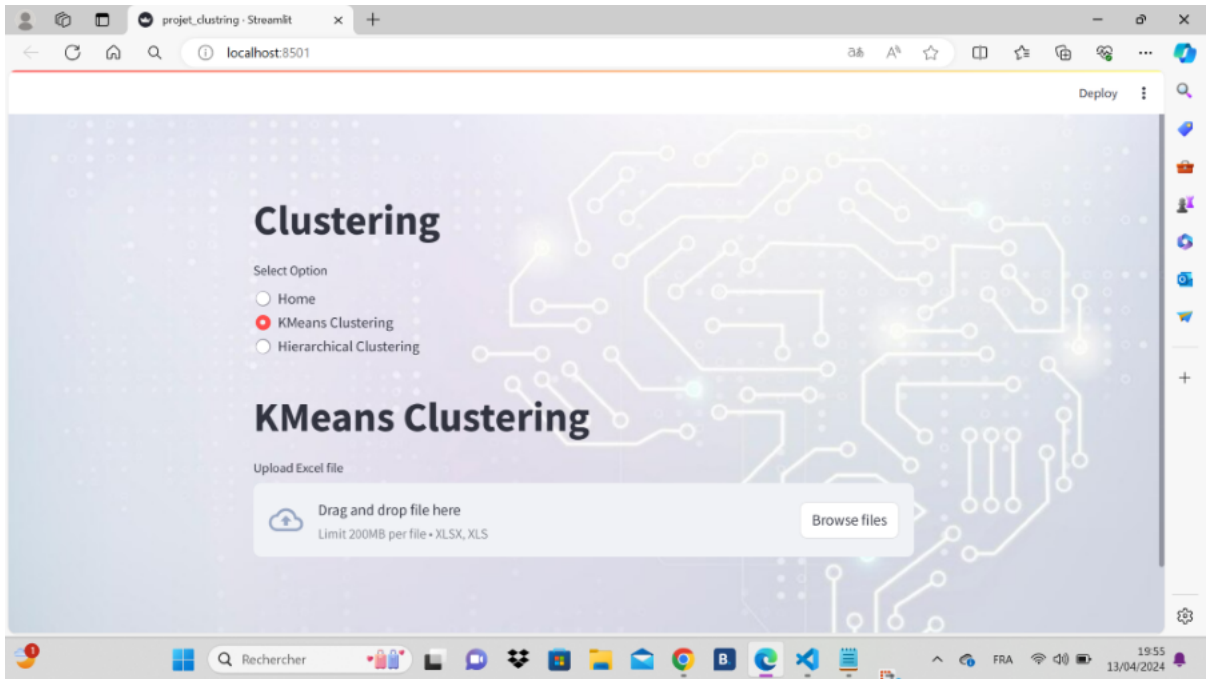


FIGURE 2 – étape 2

Vous pouvez télécharger un fichier Excel contenant vos données pour l'analyse

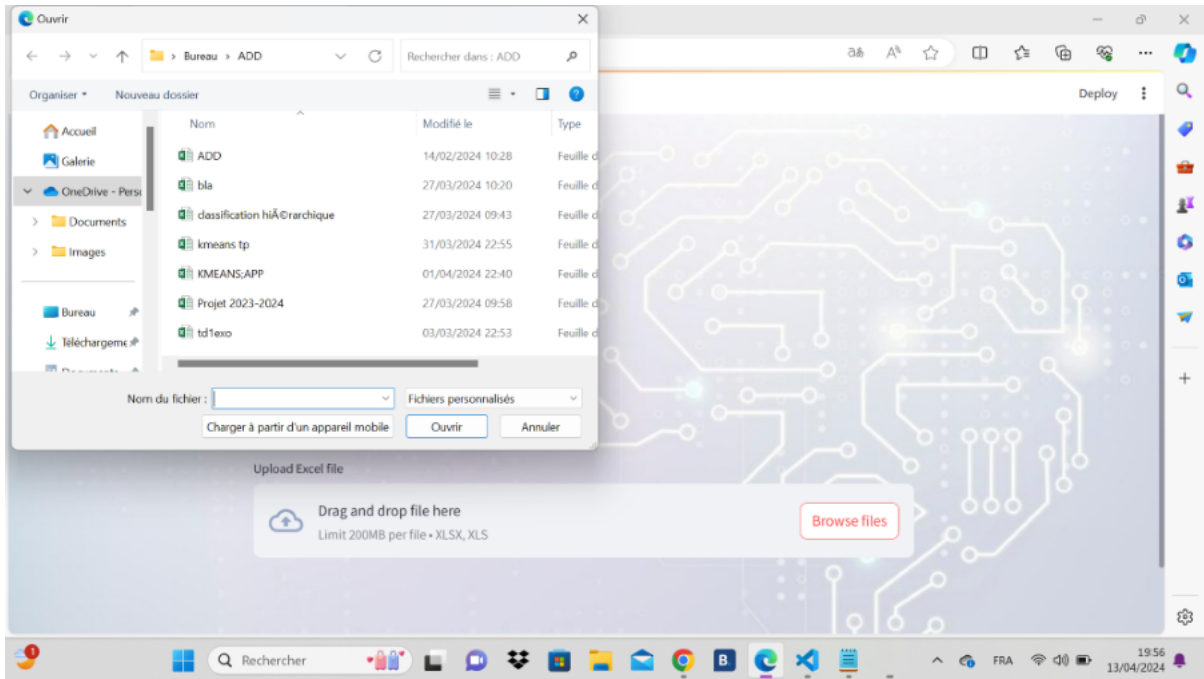


FIGURE 3 – étape 3

Une fois les données chargées, vous pouvez effectuer : Choix de la méthode de prétraitement : Vous avez la possibilité de spécifier la méthode de prétraitement des données. Par exemple, vous pouvez choisir d'appliquer une Analyse en Composantes Principales (ACP) pour réduire la dimensionnalité des données avant d'appliquer les méthodes de clustering.

Spécification des distances : Pour l'option "KMeans Clustering", vous pouvez spécifier la distance métrique à utiliser pour mesurer la similarité entre les points de données. Les distances les plus couramment utilisées sont la distance euclidienne et la distance de Manhattan.

Choix du nombre de clusters (k) : Vous pouvez spécifier le nombre de clusters (k) que vous souhaitez obtenir à partir de l'ensemble de données. Ce choix peut être basé sur des connaissances préalables du domaine, des analyses exploratoires des données ou des méthodes d'évaluation automatique du nombre optimal de clusters.

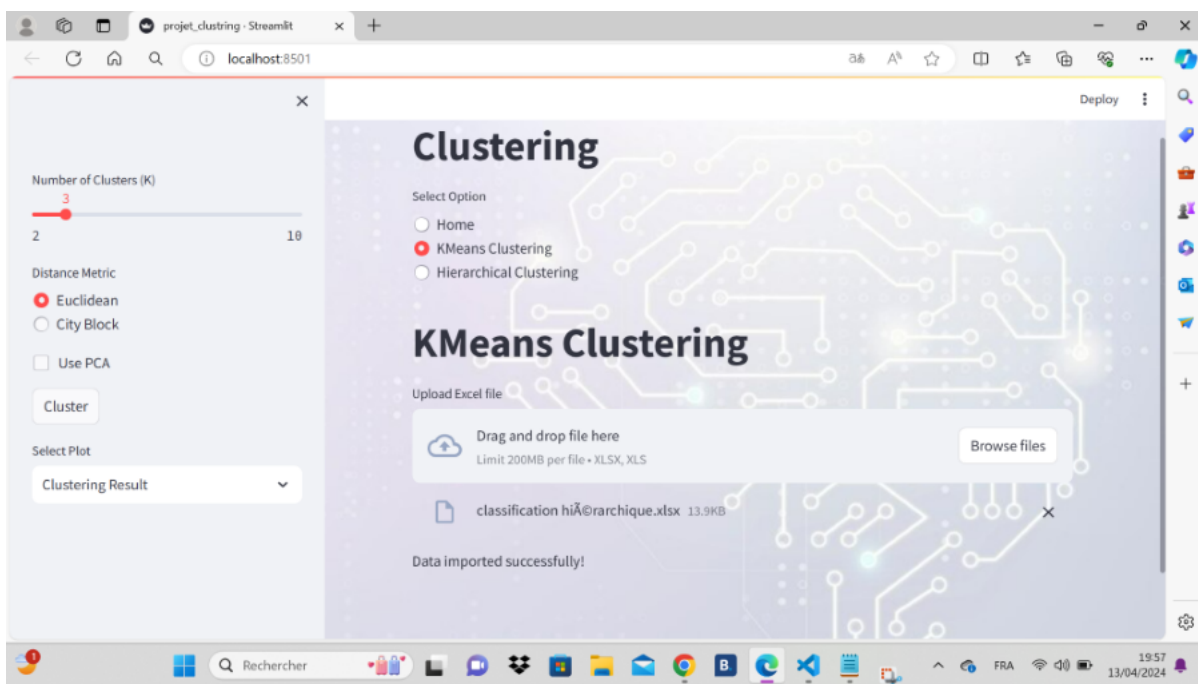


FIGURE 4 – étape 4

Sélection du type de plot : Une fois que le clustering est terminé et que les résultats sont prêts à être visualisés, vous pouvez choisir le type de plot que vous souhaitez afficher.

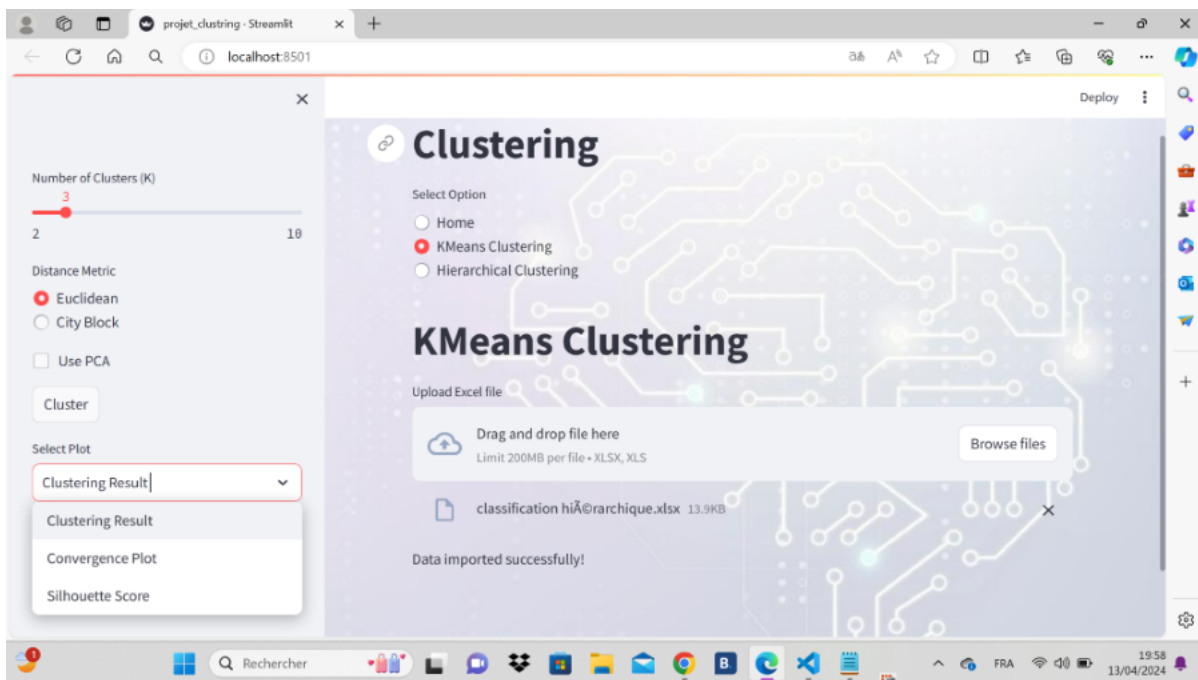


FIGURE 5 – étape 5

Si vous choisissez l'option "Clustering Result", vous obtiendrez un plot qui représente les résultats du clustering.

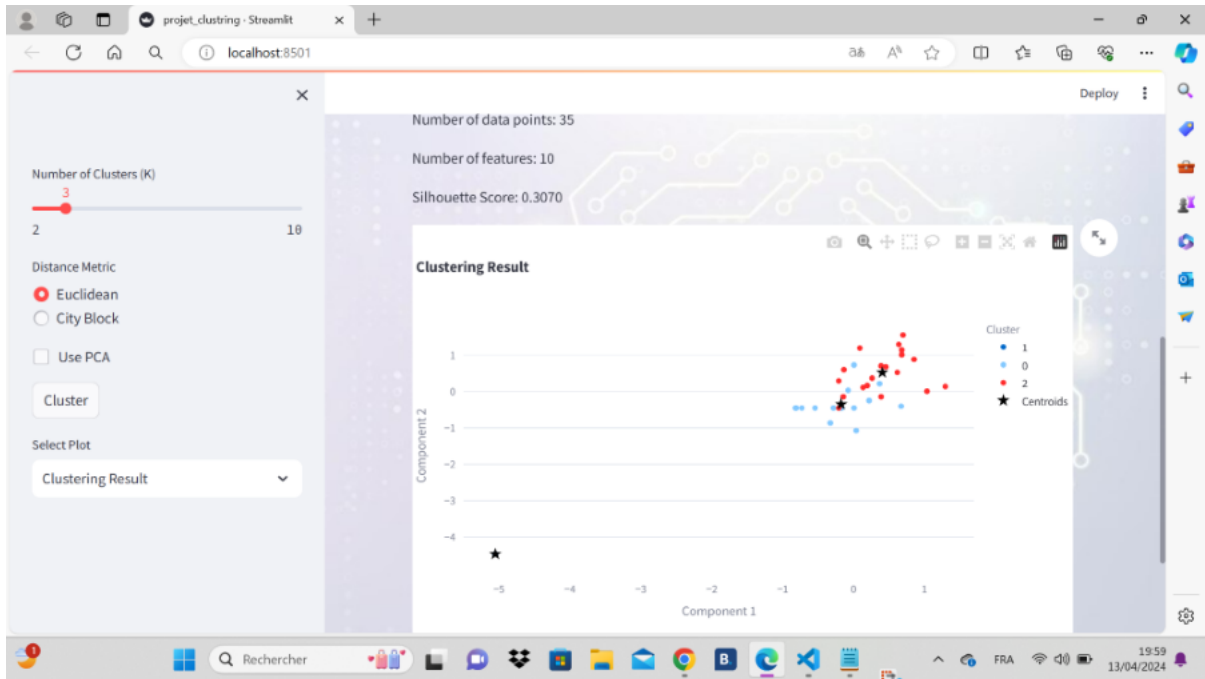


FIGURE 6 – étape 6

Si vous choisissez l'option "Convergence Plot", vous obtiendrez un graphique qui représente la convergence de clustering, dans laquelle l'inertie évolue au fil des itérations.



FIGURE 7 – étape 7

Si vous choisissez l'option "Silhouette Score", vous obtiendrez un graphique qui représente l'évaluation de la qualité du clustering à l'aide du score de silhouette.

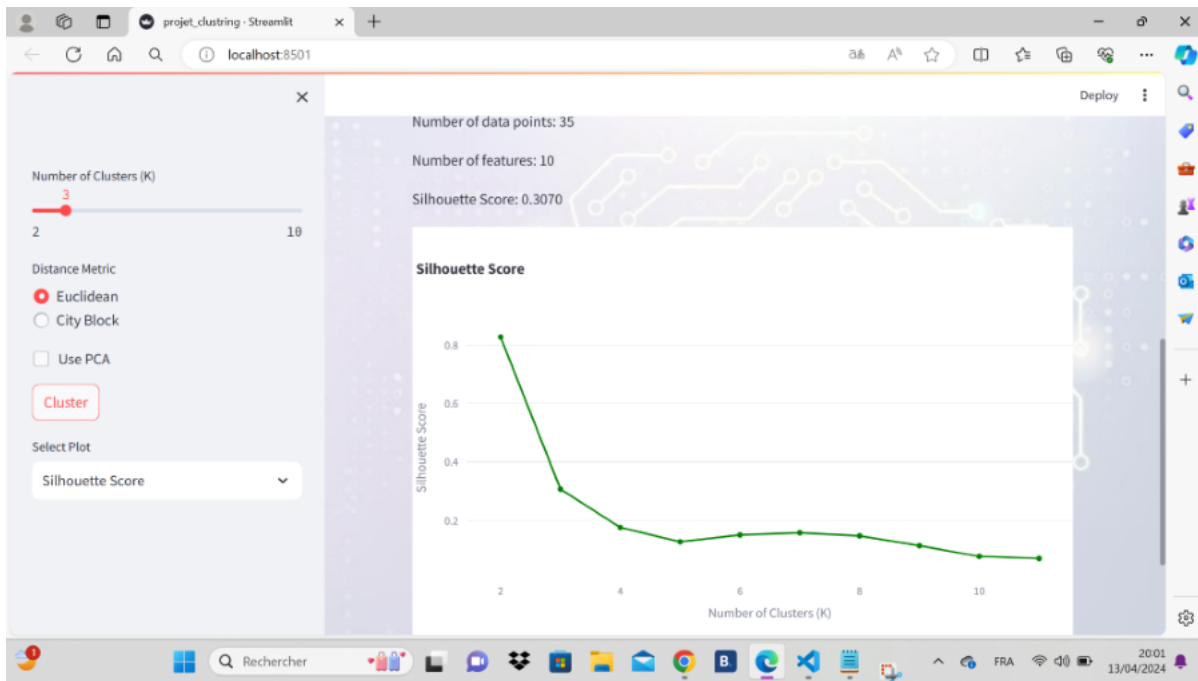


FIGURE 8 – étape 8

Si vous choisissez l'option "Hierarchical Clustering", vous pouvez charger vos données en utilisant le bouton "Upload Excel file"

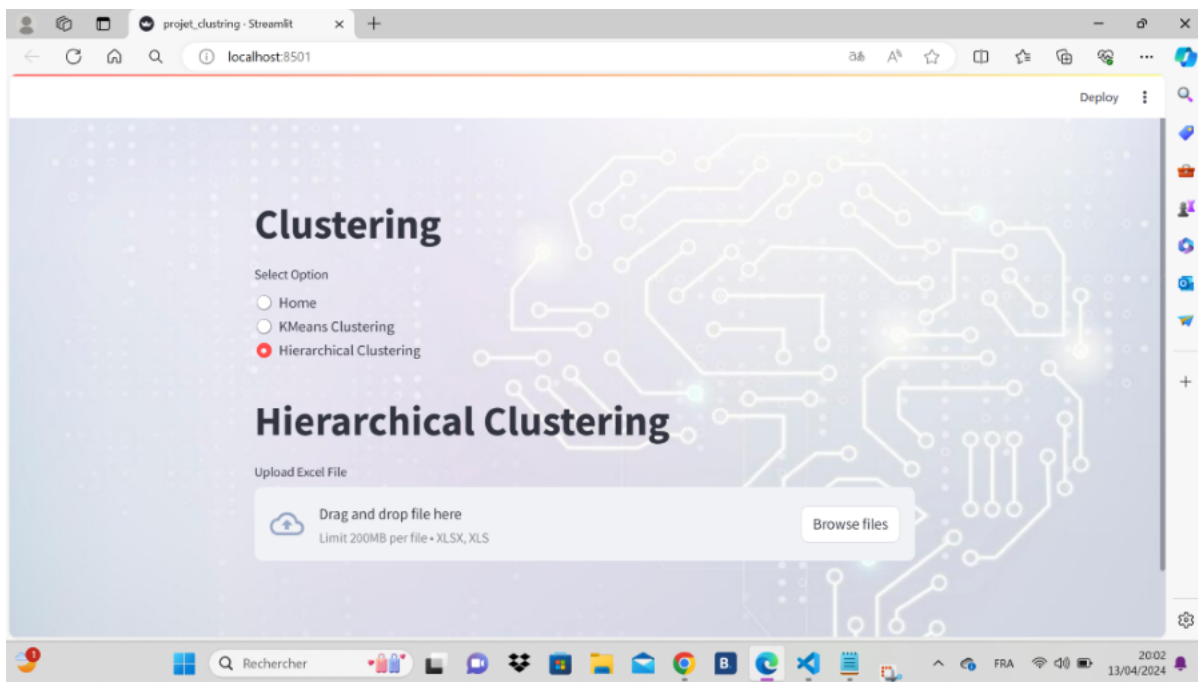


FIGURE 9 – étape 9

Après l'exécution du clustering, vous pouvez visualiser le dendrogramme résultant pour explorer la structure hiérarchique des clusters. Vous pouvez couper le dendrogramme à différentes hauteurs pour obtenir différents niveaux de granularité dans les clusters.



FIGURE 10 – étape 10

4 Conclusion :

Dans ce rapport, nous avons exploré les concepts de base de l'apprentissage non supervisé, en mettant l'accent sur la classification hiérarchique agglomérative (CAH) et l'algorithme K-Means. Nous avons abordé le prétraitement des données, les étapes de ces algorithmes, les mesures d'évaluation, ainsi que la visualisation des résultats à l'aide de dendrogrammes, ou les plots tel que convergence plot, clustering result. Cette exploration nous a permis de mieux comprendre l'application pratique de ces techniques dans l'analyse de données. Voici le lien vers mon référentiel GitHub où vous trouverez mon projet : **GitHub** :<https://github.com/sasoosaa>

```

class KMeansClustering:
    def __init__(self, k=3, distance_metric='euclidean', random_state=
None):
        self.k = k
        self.distance_metric = distance_metric
        self.centroids = None
        self.random_state = random_state

    def calculate_distance(self, data_point, centroids):
        if self.distance_metric == 'euclidean':
            return np.sqrt(np.sum((centroids - data_point) ** 2, axis=1))
        elif self.distance_metric == 'city_block':
            return np.sum(np.abs(centroids - data_point), axis=1)
        else:
            raise ValueError("Invalid distance metric. Supported metrics
are 'euclidean' and 'city_block'.")

    def fit(self, X, max_iterations=200):
        np.random.seed(self.random_state) # Set the random seed

        # Initialize centroids randomly
        self.centroids = X[np.random.choice(X.shape[0], self.k, replace=
False)]

        inertia_values = []

        for _ in range(max_iterations):
            # Assign each data point to the nearest centroid
            distances = np.vstack([self.calculate_distance(data_point,
self.centroids) for data_point in X])
            cluster_assignment = np.argmin(distances, axis=1)

            # Update centroids
            new_centroids = np.array([X[cluster_assignment == i].mean(
axis=0) for i in range(self.k)])

            # Calculate inertia
            inertia = np.sum((X - new_centroids[cluster_assignment]) **
2)

            inertia_values.append(inertia)

            # Check convergence
            if np.allclose(self.centroids, new_centroids):
                break

            self.centroids = new_centroids

        # Calculate silhouette score
        labels = np.argmin(distances, axis=1)
        silhouette_avg = silhouette_score(X, labels)

```

```

47         return cluster_assignment, self.centroids, inertia_values,
           silhouette_avg
48
49 def run_clustering(X, k, distance_metric, use_pca):
50     if use_pca:
51         pca = PCA(n_components=2)
52         X = pca.fit_transform(X)
53
54         kmeans = KMeansClustering(k=k, distance_metric=distance_metric,
           random_state=42)
55         labels, centroids, inertia_values, silhouette_avg = kmeans.fit(X)
56         return labels, centroids, inertia_values, silhouette_avg
57
58 def plot_result(X_reduced, labels, centroids):
59     if X_reduced.shape[1] == 2:
60         df = pd.DataFrame(X_reduced, columns=['Component 1', 'Component
           2'])
61     else:
62         df = pd.DataFrame(X_reduced[:, :2], columns=['Component 1', '
           Component 2'])
63     df['Cluster'] = labels.astype(str) # Convert labels to string for
           coloring purposes
64
65     fig = px.scatter(df, x='Component 1', y='Component 2', color='
           Cluster',
66                      title='Clustering Result', hover_name=df.index)
67     fig.add_trace(go.Scatter(
68         x=centroids[:, 0],
69         y=centroids[:, 1],
70         mode='markers',
71         marker=dict(size=10, color='black', symbol='star'),
72         name='Centroids'
73     ))
74     return fig
75
76 def plot_convergence(inertia_values):
77     fig = go.Figure()
78     fig.add_trace(go.Scatter(
79         x=list(range(1, len(inertia_values) + 1)),
80         y=inertia_values,
81         mode='lines+markers',
82         marker=dict(color='blue'),
83         name='Inertia'
84     ))
85     fig.update_layout(
86         title='Convergence Plot',
87         xaxis=dict(title='Iteration'),
88         yaxis=dict(title='Inertia'),
89         hovermode='closest'
90     )
91     return fig
92
93 def plot_silhouette_score(silhouette_values):
94     fig = go.Figure()
95     fig.add_trace(go.Scatter(
96         x=list(range(2, 12)),
97         y=silhouette_values,

```

```

98     mode='lines+markers',
99     marker=dict(color='green'),
100     name='Silhouette Score'
101 ))
102 fig.update_layout(
103     title='Silhouette Score',
104     xaxis=dict(title='Number of Clusters (K)'),
105     yaxis=dict(title='Silhouette Score'),
106     hovermode='closest'
107 )
108 return fig
109
110 def main():
111     st.title("KMeans Clustering")
112
113     uploaded_file = st.file_uploader("Upload Excel file", type=["xlsx",
114 "xls"])
115     if uploaded_file is not None:
116         data = pd.read_excel(uploaded_file, index_col=0)
117         st.write("Data imported successfully!")
118
119         X = data.values
120         X = X[:, 2:] # Adjust this according to your data
121
122         scaler = StandardScaler()
123         X_scaled = scaler.fit_transform(X)
124
125         k = st.sidebar.slider("Number of Clusters (K)", min_value=2,
126 max_value=10, value=3)
127         distance_metric = st.sidebar.radio("Distance Metric", ["
128 Euclidean", "City Block"], index=0)
129         distance_metric = 'euclidean' if distance_metric == 'Euclidean'
130 else 'city_block'
131
132         use_pca = st.sidebar.checkbox("Use PCA")
133
134         labels, centroids, inertia_values, silhouette_avg = None, None,
135 None, None # Initialize variables
136
137         if st.sidebar.button("Cluster"):
138             labels, centroids, inertia_values, silhouette_avg =
139 run_clustering(X_scaled, k, distance_metric, use_pca)
140
141             st.success(f"Convergence achieved in {len(inertia_values)}
142 iterations.")
143             st.write(f"Number of clusters: {k}")
144             st.write(f"Number of data points: {len(X)}")
145             st.write(f"Number of features: {X.shape[1]}")
146             st.write(f"Silhouette Score: {silhouette_avg:.4f}")
147
148             plot_choice = st.sidebar.selectbox("Select Plot", ["Clustering
149 Result", "Convergence Plot", "Silhouette Score", "Elbow Method"])
150
151             if plot_choice == "Clustering Result" and labels is not None:
152                 fig = plot_result(X_scaled, labels, centroids)
153                 st.plotly_chart(fig, use_container_width=True)

```

```

146         elif plot_choice == "Convergence Plot" and inertia_values is not
None:
147             fig = plot_convergence(inertia_values)
148             st.plotly_chart(fig, use_container_width=True)
149         elif plot_choice == "Silhouette Score" and silhouette_avg is not
None:
150             silhouette_values = []
151             for k in range(2, 12):
152                 kmeans = KMeansClustering(k=k, random_state=42)
153                 _, _, _, silhouette_avg = kmeans.fit(X_scaled)
154                 silhouette_values.append(silhouette_avg)
155             fig = plot_silhouette_score(silhouette_values)
156             st.plotly_chart(fig, use_container_width=True)
157
158 if __name__ == "__main__":
159     main()

```

Listing 9 – python script1

Code de la méthode CAH :

```

1
2 class HierarchicalClustering:
3     def __init__(self, k):
4         self.k = k
5
6     def fit(self, data, distances):
7         n = len(data)
8         clusters = [[i] for i in range(n)]
9
10        while len(clusters) > self.k:
11            min_dist = np.inf
12            for i in range(len(clusters)):
13                for j in range(i + 1, len(clusters)):
14                    cluster1 = clusters[i]
15                    cluster2 = clusters[j]
16                    avg_dist = 0
17                    count = 0
18                    for idx1 in cluster1:
19                        for idx2 in cluster2:
20                            if not np.isnan(distances[idx1, idx2]):
21                                avg_dist += distances[idx1, idx2]
22                                count += 1
23                    if count > 0:
24                        avg_dist /= count
25                        if avg_dist < min_dist:
26                            min_dist = avg_dist
27                            merge_index = (i, j)
28            i, j = merge_index
29            new_cluster = clusters[i] + clusters[j]
30            clusters[i] = new_cluster
31            clusters.pop(j)
32
33        self.labels_ = np.zeros(n)
34        for i, cluster in enumerate(clusters):
35            for idx in cluster:
36                self.labels_[idx] = i
37

```

```

38 def import_data_from_excel(file_path):
39     data = pd.read_excel(file_path)
40     # Check if the first row or first column contains non-numeric values
41     first_row_is_str = any(isinstance(val, str) for val in data.iloc[0])
42     first_column_is_str = isinstance(data.iloc[:, 0].values[0], str)
43     if first_row_is_str or first_column_is_str:
44         data = data.iloc[1:, 1:] # Exclude first row and first column
45         if they contain non-numeric values
46         return data
47
48 def dendrogram_visualization(data, labels):
49     X = data.values
50     Z = linkage(X, method='ward')
51     fig, ax = plt.subplots(figsize=(10, 5))
52     dn = dendrogram(Z, ax=ax)
53     ax.set_title("Dendrogram")
54     ax.set_xlabel("Data Points")
55     ax.set_ylabel("Distance")
56     st.pyplot(fig) # Display the plot in Streamlit
57
58 # Example usage:
59 uploaded_file = st.file_uploader("Upload Excel File", type=["xlsx", "xls
60 ])
61 if uploaded_file is not None:
62     data = import_data_from_excel(uploaded_file)
63     k = 3 # Number of clusters
64     distances = pairwise_distances(data.values, metric='euclidean')
65     clustering_algorithm = HierarchicalClustering(k=k)
66     clustering_algorithm.fit(data, distances)
67     labels = clustering_algorithm.labels_
68     dendrogram_visualization(data, labels)

```

Listing 10 – python script1