



ÉCOLE NATIONALE DES SCIENCES APPLIQUÉES AGADIR

RAPPORT D'ANALYSE DES DONNÉES

Analyse de Classification

Élève :

Salma ELAISSARI

Enseignant :

Prof.BrahimEL ASRI

Filière:

FID

14 avril 2024

1	Introduction	2
1.1	Importance de la classification	2
1.2	Introduction aux méthodes de classification hiérarchique et de k-means . .	2
1.3	Objectifs du rapport	3
1.4	Description du projet	3
2	Méthodes de Classification	4
2.1	Explication détaillée de la CAH :	4
2.1.1	Principe de fonctionnement	4
2.1.2	Avantages de la CAH :	6
2.1.3	Inconvénients de la CAH :	6
2.1.4	Réalisation de CAH sur Python :	7
2.1.5	Principes et Fonctionnement du code :	9
2.2	Explication détaillée de la méthode kmeans :	10
2.2.1	Principe de fonctionnement :	10
2.2.2	Avantages de K-means :	11
2.2.3	Inconvénients de K-means :	12
2.2.4	Réalisation de la méthode kmeans sur Python :	12
2.2.5	Principes et Fonctionnement du code :	14
3	La classification selon kmeans et classification hiérarchique :	16
3.1	Code total sur python :	16
3.2	Description du programme :	21
3.3	Principes et Fonctionnement du code :	21
3.4	Interface et utilisation du programme :	22
4	Conclusion :	28

1 Introduction

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.

1.1 Importance de la classification

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.

1.2 Introduction aux méthodes de classification hiérarchique et de k-means

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 :

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

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

- $$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$
- **Calcul du coefficient de corrélation de Pearson** : Utiliser la formule suivante pour calculer le coefficient de corrélation r :

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

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

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

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

- **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 Réalisation de CAH sur Python :

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 from scipy.cluster import hierarchy
5
6 class HierarchicalClustering:
7     def __init__(self, k):
8         self.k = k
9
10    def fit(self, data):
11        self.data = data
12        n = len(data)
13        distances = np.zeros((n, n))
14
15        # Calculate pairwise distances
16        for i in range(n):
17            for j in range(n):
18                distances[i, j] = self.euclidean_distance(data[i],
19 data[j])
20
21        # Initialize clusters
22        clusters = [[i] for i in range(n)]
23
24        # Hierarchical clustering algorithm
25        while len(clusters) > self.k:
26            min_dist = np.inf
27            for i in range(len(clusters)):
28                for j in range(i + 1, len(clusters)):
29                    cluster1 = clusters[i]
30                    cluster2 = clusters[j]
31                    avg_dist = 0
32                    count = 0
33                    for idx1 in cluster1:
34                        for idx2 in cluster2:

```



```

34         if not np.isnan(distances[idx1, idx2]):
35             avg_dist += distances[idx1, idx2]
36             count += 1
37         if count > 0:
38             avg_dist /= count
39             if avg_dist < min_dist:
40                 min_dist = avg_dist
41                 merge_index = (i, j)
42     i, j = merge_index
43     new_cluster = clusters[i] + clusters[j]
44     clusters[i] = new_cluster
45     clusters.pop(j)
46
47     self.labels_ = np.zeros(n)
48     for i, cluster in enumerate(clusters):
49         for idx in cluster:
50             self.labels_[idx] = i
51
52     def euclidean_distance(self, x1, x2):
53         x1 = np.array(x1, dtype=np.float64) # Convert x1 to numpy
54         x2 = np.array(x2, dtype=np.float64) # Convert x2 to numpy
55         if np.any(np.isnan(x1)) or np.any(np.isnan(x2)):
56             return np.nan
57         return np.sqrt(np.sum((x1 - x2) ** 2))
58
59 # Fonction pour importer les données depuis un fichier Excel
60 def import_data_from_excel(file_path):
61     data = pd.read_excel(file_path)
62     # Vérifier si la première ligne ou la première colonne
63     # contient des valeurs non numériques
64     first_row_is_str = any(isinstance(val, str) for val in data.
65                             iloc[0])
66     first_column_is_str = isinstance(data.iloc[:, 0].values[0], str)
67     if first_row_is_str or first_column_is_str:
68         data = data.iloc[1:, 1:] # Exclure la première ligne et
69         # la première colonne si elles contiennent des valeurs non
70         # numériques
71     return data
72
73 # Fonction améliorée pour visualiser le dendrogramme
74 def plot_dendrogram(data):
75     Z = hierarchy.linkage(data, method='ward')
76     plt.figure(figsize=(10, 5))
77     dn = hierarchy.dendrogram(Z)
78     plt.title("Dendrogram")
79     plt.xlabel("Data Points")
80     plt.ylabel("Distance")
81     plt.show()
82     file_path = "data.xlsx" # Path to the Excel file
83     data = import_data_from_excel(file_path)
84     clustering_algorithm = HierarchicalClustering(k=3) # Initialize
85     # the hierarchical clustering algorithm
86     clustering_algorithm.fit(data.values) # Fit the data
87     labels = clustering_algorithm.labels_ # Get the cluster labels
88     print(labels) # Print the cluster labels

```

```
83 plot_dendrogram(data.values) # Visualize the dendrogram
```

Listing 1 – python script1

2.1.5 Principes et Fonctionnement du code :

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.

scipy.cluster.hierarchy : Ce sous-module de SciPy contient des fonctions pour effectuer le clustering hiérarchique. Il est utilisé ici pour effectuer le clustering hiérarchique et générer le dendrogramme.

Ces packages sont essentiels pour effectuer différentes étapes du processus de clustering hiérarchique, notamment la manipulation des données, le calcul des distances et la visualisation des résultats.

La classe **HierarchicalClustering** implémente l'algorithme de clustering hiérarchique. Voici le rôle et le fonctionnement détaillé de chaque méthode de cette classe : **init(self, k)** : Cette méthode est le constructeur de la classe. Elle initialise un objet de la classe HierarchicalClustering avec un paramètre k qui représente le nombre de clusters souhaités.

fit(self, data) : Cette méthode prend en entrée les données à clusteriser et applique l'algorithme de clustering hiérarchique. Voici son fonctionnement :

self.data = data : Stocke les données dans l'objet.

n = len(data) : Calcule le nombre de points de données.

distances = np.zeros((n, n)) : Initialise une matrice de distances pour stocker les distances euclidiennes entre chaque paire de points de données.

Calcule les distances euclidiennes entre chaque paire de points de données à l'aide de la méthode euclidean distance.

Initialise les clusters en tant que listes de points individuels.

Applique l'algorithme de clustering hiérarchique : Tant que le nombre de clusters est supérieur à k, fusionne les clusters les plus proches en termes de distance moyenne. Met à jour les distances entre les clusters fusionnés. Répète le processus jusqu'à ce que le nombre de clusters atteigne k.

Affecte des étiquettes de cluster à chaque point de données en fonction des clusters

euclidean distance(self, x1, x2) : Cette méthode calcule la distance euclidienne entre deux points de données. Voici son fonctionnement :

import-data-from-excel(file-path) : Cette fonction importe les données à partir d'un fichier Excel. Voici son fonctionnement :

Check if the first row or the first column contains non-numeric values :
Vérifie si la première ligne ou la première colonne contient des valeurs non numériques. Cela est fait en vérifiant si au moins un des éléments de la première ligne est de type str et si le premier élément de la première colonne est de type str.

Return data : Renvoie le DataFrame contenant les données, avec les ajustements effectués si nécessaire.

2.2 Explication détaillée de la méthode kmeans :

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 :

- ## Rapport - Analyse de Classification

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

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

11

2.2.3 Inconvénients de K-means :

Les inconvénients de l'algorithme de clustering K-means comprennent :

- **Sensibilité aux Conditions Initiales** : K-means est sensible au placement initial des centroids, ce qui peut conduire à différentes affectations de clusters et résultats en fonction de la sélection aléatoire initiale des centroids.
- **Nécessité de Pré-définir le Nombre de Clusters (K)** : Une des limitations de K-means est que l'utilisateur doit spécifier le nombre de clusters au préalable, ce qui peut être difficile à déterminer avec précision en pratique.
- **Absence de Développement d'un Ensemble de Clusters Optimal** : K-means ne fournit pas intrinsèquement un ensemble optimal de clusters et nécessite que l'utilisateur pré-définisse le nombre de clusters, ce qui peut ne pas toujours correspondre à la structure sous-jacente des données.
- **Incohérence des Résultats** : K-means peut produire des résultats variables entre différentes exécutions de l'algorithme en raison de l'initialisation aléatoire des centroids, ce qui entraîne une incohérence dans les résultats du clustering.
- **Hypothèse de Clusters Sphériques** : K-means suppose que les clusters sont sphériques et isotropes, ce qui peut ne pas être vrai pour tous les types de distributions de données, limitant son applicabilité dans certaines situations.
- **Difficulté dans le Choix du Nombre de Clusters** : Sélectionner le nombre optimal de clusters (K) dans K-means peut être difficile, car cela nécessite une réflexion attentive et peut ne pas toujours avoir de solution simple.

Ces inconvénients mettent en lumière certaines des limitations et des défis associés à l’algorithme de clustering K-means, soulignant l’importance de comprendre ses contraintes et de considérer des méthodes de clustering alternatives lorsque c’est nécessaire.

2.2.4 Réalisation de la méthode kmeans sur Python :

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.decomposition import PCA
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.metrics import silhouette_score
6 import plotly.graph_objs as go
7 import plotly.express as px
8
9 class KMeansClustering:
10     def __init__(self, k=3, distance_metric='euclidean', random_state=
None):
11         self.k = k
12         self.distance_metric = distance_metric
13         self.centroids = None
14         self.random_state = random_state
15
16     def calculate_distance(self, data_point, centroids):
17         if self.distance_metric == 'euclidean':
18             return np.sqrt(np.sum((centroids - data_point) ** 2, axis=1))
19
20         elif self.distance_metric == 'city_block':

```

```

20         return np.sum(np.abs(centroids - data_point), axis=1)
21     else:
22         raise ValueError("Invalid distance metric. Supported metrics
23                             are 'euclidean' and 'city_block'.")
24
25 def fit(self, X, max_iterations=200):
26     np.random.seed(self.random_state) # Set the random seed
27
28     # Initialize centroids randomly
29     self.centroids = X[np.random.choice(X.shape[0], self.k, replace=
30 False)]
31
32     inertia_values = []
33
34     for _ in range(max_iterations):
35         # Assign each data point to the nearest centroid
36         distances = np.vstack([self.calculate_distance(data_point,
37 self.centroids) for data_point in X])
38         cluster_assignment = np.argmin(distances, axis=1)
39
40         # Update centroids
41         new_centroids = np.array([X[cluster_assignment == i].mean(
42 axis=0) for i in range(self.k)])
43
44         # Calculate inertia
45         inertia = np.sum((X - new_centroids[cluster_assignment])**
46 2)
47
48         inertia_values.append(inertia)
49
50         # Check convergence
51         if np.allclose(self.centroids, new_centroids):
52             break
53
54         self.centroids = new_centroids
55
56         # Calculate silhouette score
57         labels = np.argmin(distances, axis=1)
58         silhouette_avg = silhouette_score(X, labels)
59
60     return cluster_assignment, self.centroids, inertia_values,
61 silhouette_avg
62
63 def main():
64     # Load data from Excel file
65     file_path = "your_excel_file.xlsx"
66     data = pd.read_excel(file_path)
67
68     X = data.values
69     X = X[:, 2:] # Adjust this according to your data
70
71     scaler = StandardScaler()
72     X_scaled = scaler.fit_transform(X)
73
74     k = 3 # You can set the number of clusters here
75
76     kmeans = KMeansClustering(k=k, distance_metric='euclidean',
77 random_state=42)

```



```

70 labels, centroids, _, _ = kmeans.fit(X_scaled)
71
72 # Use PCA for visualization if needed
73 pca = PCA(n_components=2)
74 X_reduced = pca.fit_transform(X_scaled)
75
76 # Plot clustering result
77 df = pd.DataFrame(X_reduced, columns=['Component 1', 'Component 2'])
78 df['Cluster'] = labels.astype(str) # Convert labels to string for
coloring purposes
79
80 fig = px.scatter(df, x='Component 1', y='Component 2', color='
Cluster',
81                  title='Clustering Result', hover_name=df.index)
82 fig.add_trace(go.Scatter(
83     x=centroids[:, 0],
84     y=centroids[:, 1],
85     mode='markers',
86     marker=dict(size=10, color='black', symbol='star'),
87     name='Centroids'
88 ))
89 fig.show()
90
91 if __name__ == "__main__":
92     main()

```

Listing 2 – python script1

2.2.5 Principes et Fonctionnement du code :

pandas (pd) :Pandas est une bibliothèque Python qui fournit des structures de données et des outils d'analyse de données faciles à utiliser. La structure de données principale de pandas est le DataFrame, qui est une structure de données tabulaire bidimensionnelle avec des étiquettes d'axe (lignes et colonnes). Les fonctionnalités de pandas incluent le chargement de données depuis diverses sources, la manipulation et la transformation des données, ainsi que l'analyse exploratoire des données.

numpy (np) : NumPy est une bibliothèque Python pour le calcul numérique qui prend en charge les tableaux multidimensionnels et les fonctions mathématiques pour les manipuler. NumPy fournit des outils efficaces pour effectuer des opérations mathématiques et statistiques sur les tableaux, ce qui en fait une bibliothèque fondamentale pour la science des données et l'apprentissage automatique.

sklearn.decomposition.PCA :PCA (Principal Component Analysis) est une technique de réduction de dimensionnalité largement utilisée pour compresser les données en conservant les principales caractéristiques. L'objet PCA de scikit-learn permet d'effectuer une analyse en composantes principales sur les données pour projeter les observations dans un espace de dimension inférieure.

sklearn.preprocessing.StandardScaler : StandardScaler de scikit-learn est une méthode de prétraitement qui standardise les fonctionnalités en supprimant la moyenne et en mettant à l'échelle jusqu'à la variance unitaire. Cela permet de normaliser les caractéristiques des données, ce qui est souvent nécessaire pour de nombreux algorithmes d'apprentissage automatique qui supposent que les données sont centrées et ont des variances similaires.

sklearn.metrics.silhouette-score : La silhouette est une mesure de la cohésion et de la séparation des clusters dans un ensemble de données. silhouette-score de scikit-learn calcule la silhouette moyenne sur toutes les observations, fournissant ainsi une indication de la qualité du clustering.

plotly.graph_objs (go) : Plotly est une bibliothèque graphique interactive qui permet de créer des graphiques et des visualisations de données hautement personnalisables. plotly.graph_objs fournit une interface basée sur des objets pour créer des graphiques, permettant un contrôle précis sur les paramètres de visualisation tels que les axes, les couleurs et les annotations.

plotly.express (px) : Plotly Express est une interface de haut niveau pour créer des visualisations de données avec Plotly. Il offre une syntaxe simple pour créer rapidement des graphiques interactifs tels que des diagrammes à barres, des diagrammes circulaires, des nuages de points et bien d'autres, avec une prise en charge automatique de la mise en forme et de la configuration des graphiques.

Classe KMeansClustering : La classe **KMeansClustering** représente l'algorithme de clustering K-means. Elle possède les méthodes suivantes :

- **__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** : Exécute l'algorithme K-means pour ajuster les centroids aux données fournies.

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 # Added import for matplotlib
10 from scipy.cluster import hierarchy
11 import base64
12 import requests
13
14 # Set background image
15 background_image_url = "https://img.freepik.com/free-vector/ai-
    technology-brain-background-vector-digital-transformation-
    concept_53876-117820.jpg?w=900&t=st=1713012632~exp=1713013232~hmac=05
    da65e2d9e6da77202ded9006cfecb86c0c11fbc70346fb0532307b18d8ac3f"
16
17 def get_base64_of_bin_file(url):
18     data = requests.get(url).content
19     return base64.b64encode(data).decode()
20
21 def set_png_as_page_bg(png_file):
22     bin_str = get_base64_of_bin_file(png_file)
23     st.markdown(
24         f'<style>
25         f'.stApp {{
26         f'background-image: url("data:image/png;base64,{bin_str}")';
27         f'background-size: cover;
28         f'}}'
29         f'</style>',
30         unsafe_allow_html=True
31     )
32
33 set_png_as_page_bg(background_image_url)
34
35 class KMeansClustering:
36     def __init__(self, k=3, distance_metric='euclidean', random_state=
    None):
37         self.k = k
38         self.distance_metric = distance_metric
39         self.centroids = None
40         self.random_state = random_state
41
42     def calculate_distance(self, data_point, centroids):
43         if self.distance_metric == 'euclidean':
44             return np.sqrt(np.sum((centroids - data_point) ** 2, axis=1))
45
46         elif self.distance_metric == 'city_block':
47             return np.sum(np.abs(centroids - data_point), axis=1)

```

```

47         else:
48             raise ValueError("Invalid distance metric. Supported metrics
are 'euclidean' and 'city_block'.")
49
50     def fit(self, X, max_iterations=200):
51         np.random.seed(self.random_state) # Set the random seed
52
53         # Initialize centroids randomly
54         self.centroids = X[np.random.choice(X.shape[0], self.k, replace=
False)]
55
56         inertia_values = []
57
58         for _ in range(max_iterations):
59             # Assign each data point to the nearest centroid
60             distances = np.vstack([self.calculate_distance(data_point,
self.centroids) for data_point in X])
61             cluster_assignment = np.argmin(distances, axis=1)
62
63             # Update centroids
64             new_centroids = np.array([X[cluster_assignment == i].mean(
axis=0) for i in range(self.k)])
65
66             # Calculate inertia
67             inertia = np.sum((X - new_centroids[cluster_assignment])**
2)
68             inertia_values.append(inertia)
69
70             # Check convergence
71             if np.allclose(self.centroids, new_centroids):
72                 break
73
74             self.centroids = new_centroids
75
76             # Calculate silhouette score
77             labels = np.argmin(distances, axis=1)
78             silhouette_avg = silhouette_score(X, labels)
79
80             return cluster_assignment, self.centroids, inertia_values,
silhouette_avg
81
82 class HierarchicalClustering:
83     def __init__(self, k):
84         self.k = k
85
86     def fit(self, data):
87         self.data = data
88         n = len(data)
89         distances = np.zeros((n, n))
90
91         # Calculate pairwise distances
92         for i in range(n):
93             for j in range(n):
94                 distances[i, j] = self.euclidean_distance(data[i], data[
j])
95
96         # Initialize clusters

```

```

clusters = [[i] for i in range(n)]

# Hierarchical clustering algorithm
while len(clusters) > self.k:
    min_dist = np.inf
    for i in range(len(clusters)):
        for j in range(i + 1, len(clusters)):
            cluster1 = clusters[i]
            cluster2 = clusters[j]
            avg_dist = 0
            count = 0
            for idx1 in cluster1:
                for idx2 in cluster2:
                    if not np.isnan(distances[idx1, idx2]):
                        avg_dist += distances[idx1, idx2]
                        count += 1
            if count > 0:
                avg_dist /= count
                if avg_dist < min_dist:
                    min_dist = avg_dist
                    merge_index = (i, j)
    i, j = merge_index
    new_cluster = clusters[i] + clusters[j]
    clusters[i] = new_cluster
    clusters.pop(j)

self.labels_ = np.zeros(n)
for i, cluster in enumerate(clusters):
    for idx in cluster:
        self.labels_[idx] = i

def euclidean_distance(self, x1, x2):
    x1 = np.array(x1, dtype=np.float64) # Convert x1 to numpy array
of numbers
    x2 = np.array(x2, dtype=np.float64) # Convert x2 to numpy array
of numbers
    if np.any(np.isnan(x1)) or np.any(np.isnan(x2)):
        return np.nan
    return np.sqrt(np.sum((x1 - x2) ** 2))

def plot_dendrogram(data):
    Z = hierarchy.linkage(data, method='ward')
    plt.figure(figsize=(10, 5))
    dn = hierarchy.dendrogram(Z)
    plt.title("Dendrogram")
    plt.xlabel("Data Points")
    plt.ylabel("Distance")
    st.pyplot(plt.gcf()) # Pass the current figure as an argument to st
.pyplot()

def plot_result(X_reduced, labels, centroids):
    if X_reduced.shape[1] == 2:
        df = pd.DataFrame(X_reduced, columns=['Component 1', 'Component
2'])
    else:
        df = pd.DataFrame(X_reduced[:, :2], columns=['Component 1', '
Component 2'])

```

```

149 df['Cluster'] = labels.astype(str) # Convert labels to string for
coloring purposes
150
151 fig = px.scatter(df, x='Component 1', y='Component 2', color='
Cluster',
152                  title='Clustering Result', hover_name=df.index)
153 fig.add_trace(go.Scatter(
154     x=centroids[:, 0],
155     y=centroids[:, 1],
156     mode='markers',
157     marker=dict(size=10, color='black', symbol='star'),
158     name='Centroids'
159 ))
160 return fig
161
162 def plot_convergence(inertia_values):
163     fig = go.Figure()
164     fig.add_trace(go.Scatter(
165         x=list(range(1, len(inertia_values) + 1)),
166         y=inertia_values,
167         mode='lines+markers',
168         marker=dict(color='blue'),
169         name='Inertia'
170     ))
171     fig.update_layout(
172         title='Convergence Plot',
173         xaxis=dict(title='Iteration'),
174         yaxis=dict(title='Inertia'),
175         hovermode='closest'
176     )
177     return fig
178
179 def plot_silhouette_score(silhouette_values):
180     fig = go.Figure()
181     fig.add_trace(go.Scatter(
182         x=list(range(2, 12)),
183         y=silhouette_values,
184         mode='lines+markers',
185         marker=dict(color='green'),
186         name='Silhouette Score'
187     ))
188     fig.update_layout(
189         title='Silhouette Score',
190         xaxis=dict(title='Number of Clusters (K)'),
191         yaxis=dict(title='Silhouette Score'),
192         hovermode='closest'
193     )
194     return fig
195
196 def main():
197     st.title("Clustering")
198
199     selected = st.radio("Select Option", ["Home", "KMeans Clustering", "
Hierarchical Clustering"])
200
201     if selected == "KMeans Clustering":
202         st.title("KMeans Clustering")

```

```

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

```

```

248         st.plotly_chart(fig, use_container_width=True)
249
250     elif selected == "Hierarchical Clustering":
251         st.title("Hierarchical Clustering")
252         uploaded_file = st.file_uploader("Upload Excel File", type=["
xls", "xls"])
253         if uploaded_file is not None:
254             data = pd.read_excel(uploaded_file, index_col=0)
255             clustering_algorithm = HierarchicalClustering(k=3)  #
Initialize the clustering algorithm
256             clustering_algorithm.fit(data.values)  # Fit the data
257             labels = clustering_algorithm.labels_  # Get the cluster
labels
258             st.write(labels)  # Display the cluster labels
259             plot_dendrogram(data.values)  # Visualize the dendrogram
260
261 if __name__ == "__main__":
262     st.set_option('deprecation.showPyplotGlobalUse', False)  # Disable
the warning about st.pyplot() usage
263     main()

```

Listing 3 – python script1

3.2 Description du programme :

Mon programme est une application de clustering avec une interface utilisateur conviviale, construite avec Streamlit. Voici une description de ses principales fonctionnalités :

1. **Chargement des données** : Les utilisateurs peuvent charger leurs données à partir de fichiers Excel en utilisant l'interface de téléchargement de fichiers.
2. **K-means Clustering** : L'application permet aux utilisateurs d'effectuer le clustering K-means sur leurs données chargées. Ils peuvent spécifier le nombre de clusters (K) et la distance métrique (euclidienne ou de bloc de ville). Une fois les paramètres sélectionnés, les utilisateurs peuvent lancer le clustering et visualiser les résultats.
3. **Visualisation des résultats** : Les utilisateurs peuvent choisir parmi plusieurs types de visualisations pour explorer les résultats du clustering. Les options incluent un graphique de dispersion des données avec les clusters colorés et les centroids, un graphique de convergence montrant l'évolution de l'inertie au fil des itérations, et un graphique de score de silhouette pour évaluer la qualité du clustering.
4. **Interprétation des résultats** : L'application fournit des informations sur le nombre de clusters, la convergence du clustering, et le score de silhouette pour évaluer la qualité des clusters obtenus.

En utilisant cette application, les utilisateurs peuvent facilement explorer leurs données, identifier des structures de clustering, et prendre des décisions basées sur les insights générés par les méthodes de clustering K-means et CAH.

3.3 Principes et Fonctionnement du code :

Les packages utilisés dans le code Python pour l'application de clustering sont les suivants :

3.4 Interface et utilisation du programme :

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".

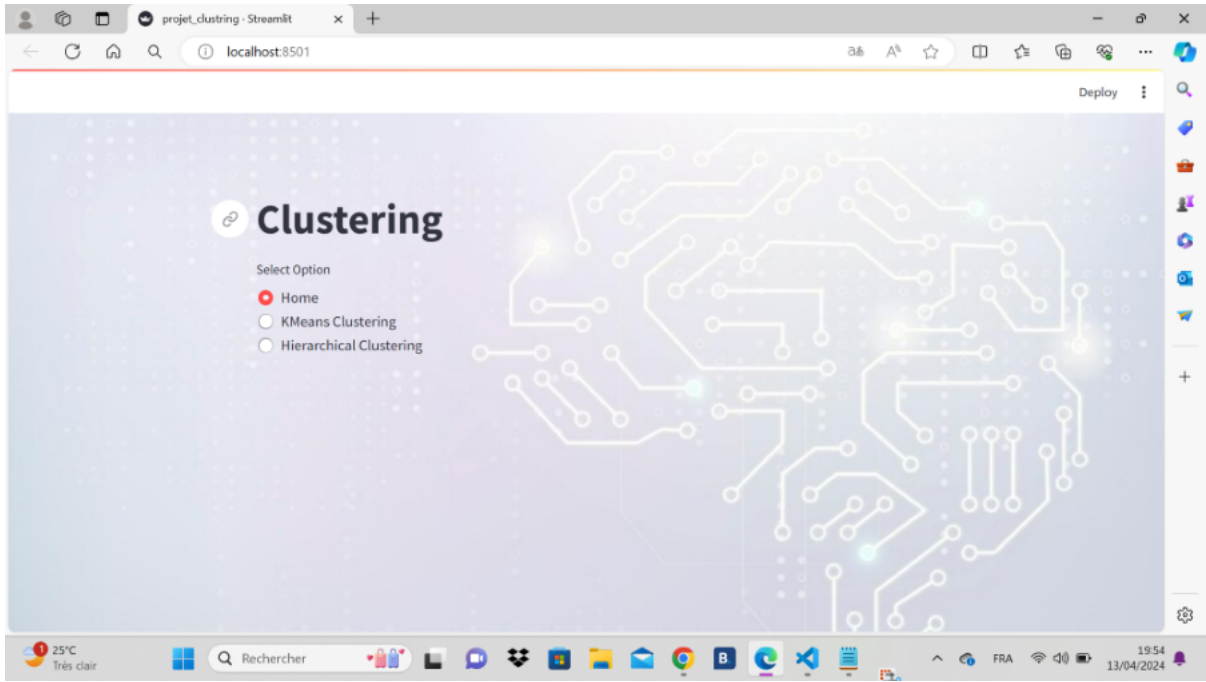


FIGURE 1 – étape 1

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

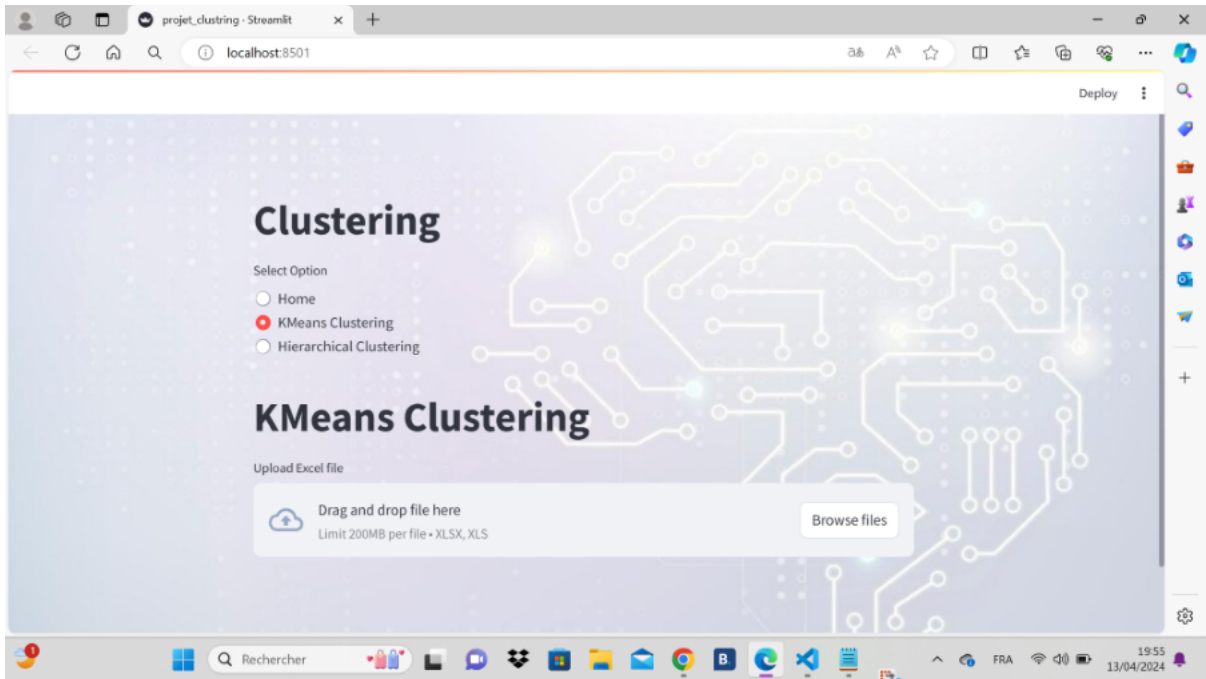


FIGURE 2 – étape 2

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

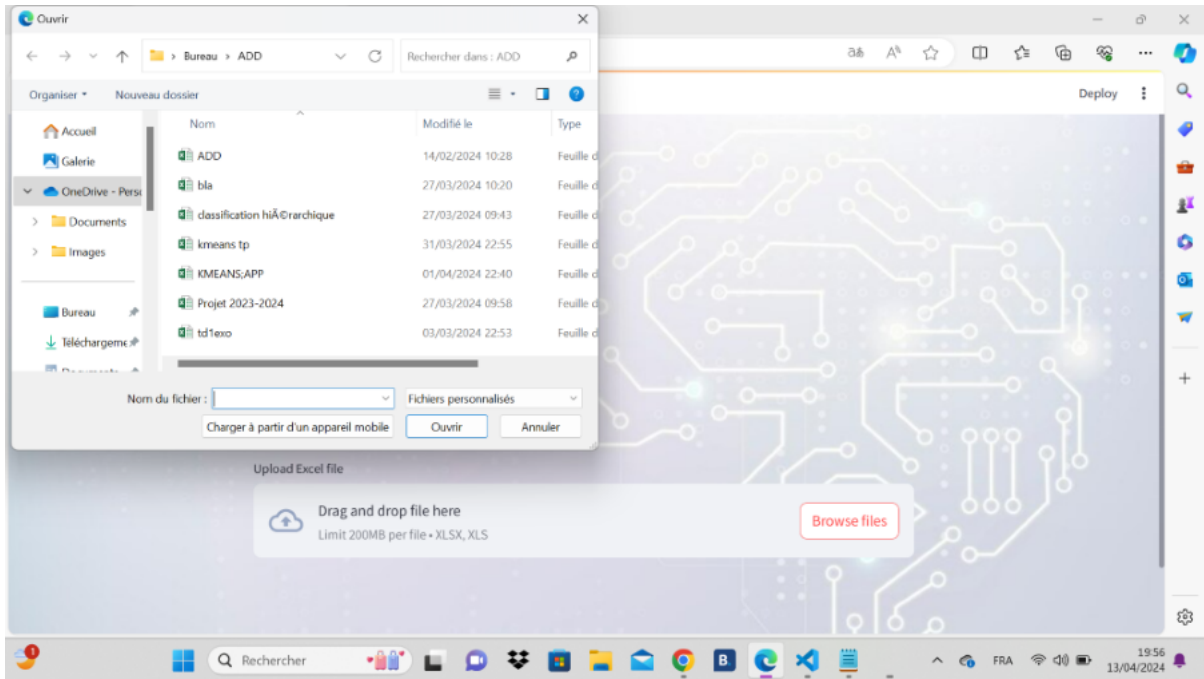


FIGURE 3 – étape 3

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.

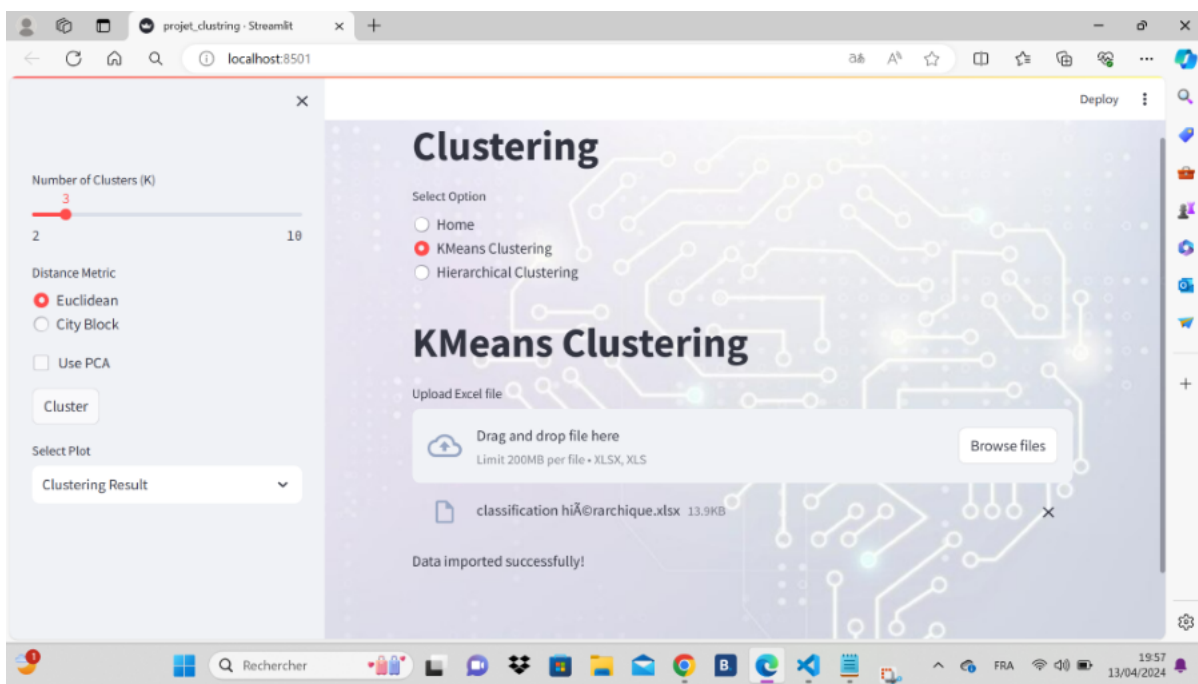


FIGURE 4 – étape 4

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.

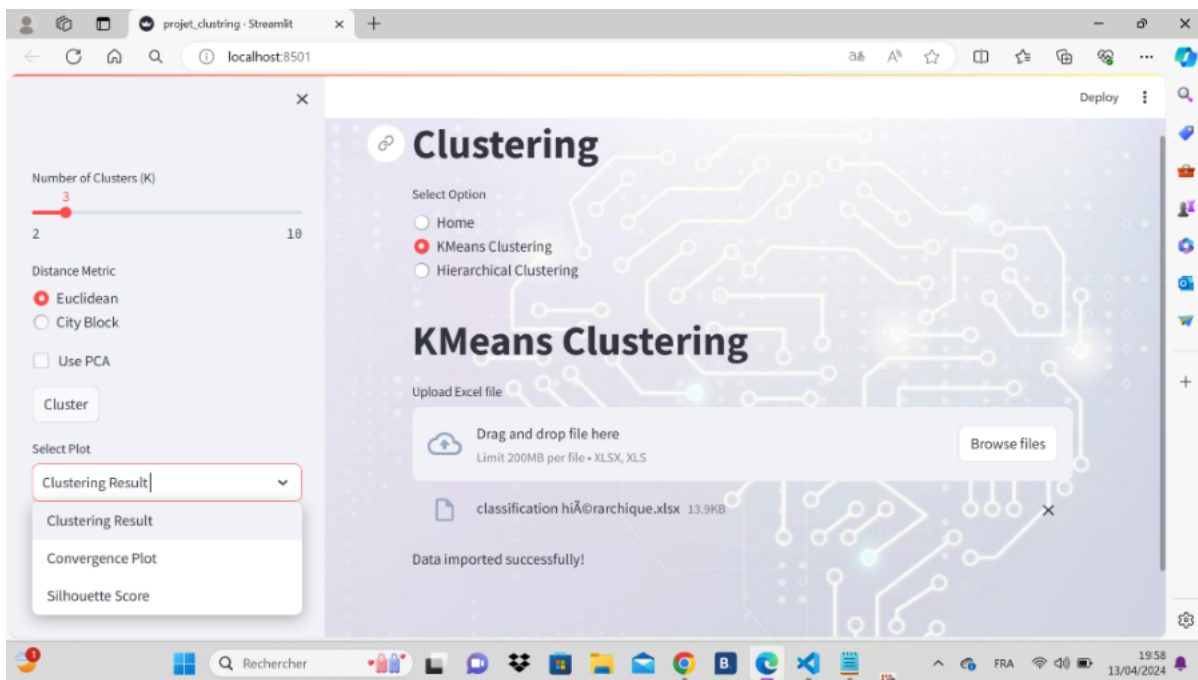


FIGURE 5 – étape 5

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

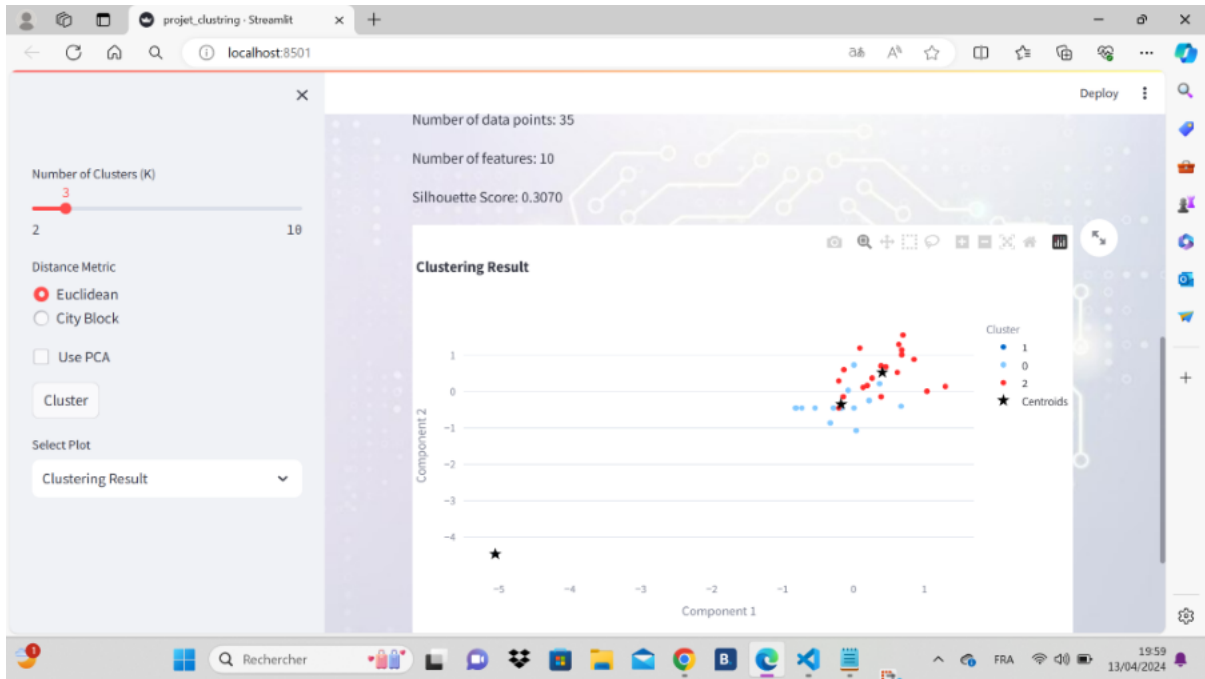


FIGURE 6 – étape 6

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.

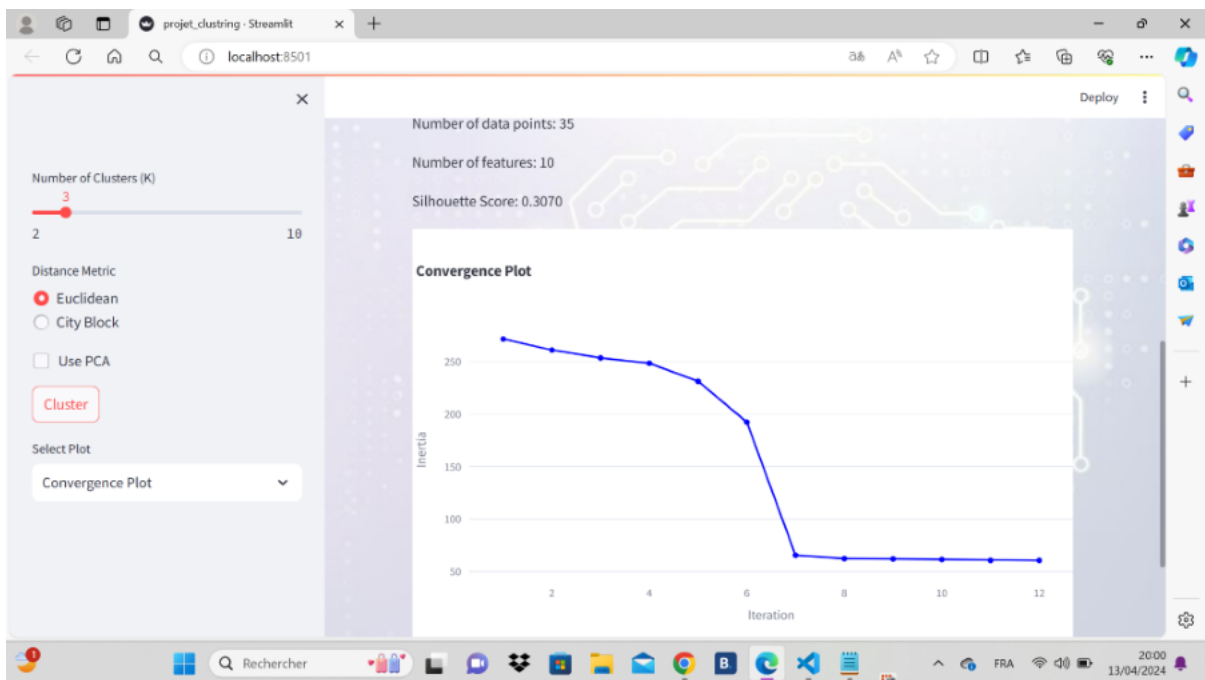


FIGURE 7 – étape 7

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.

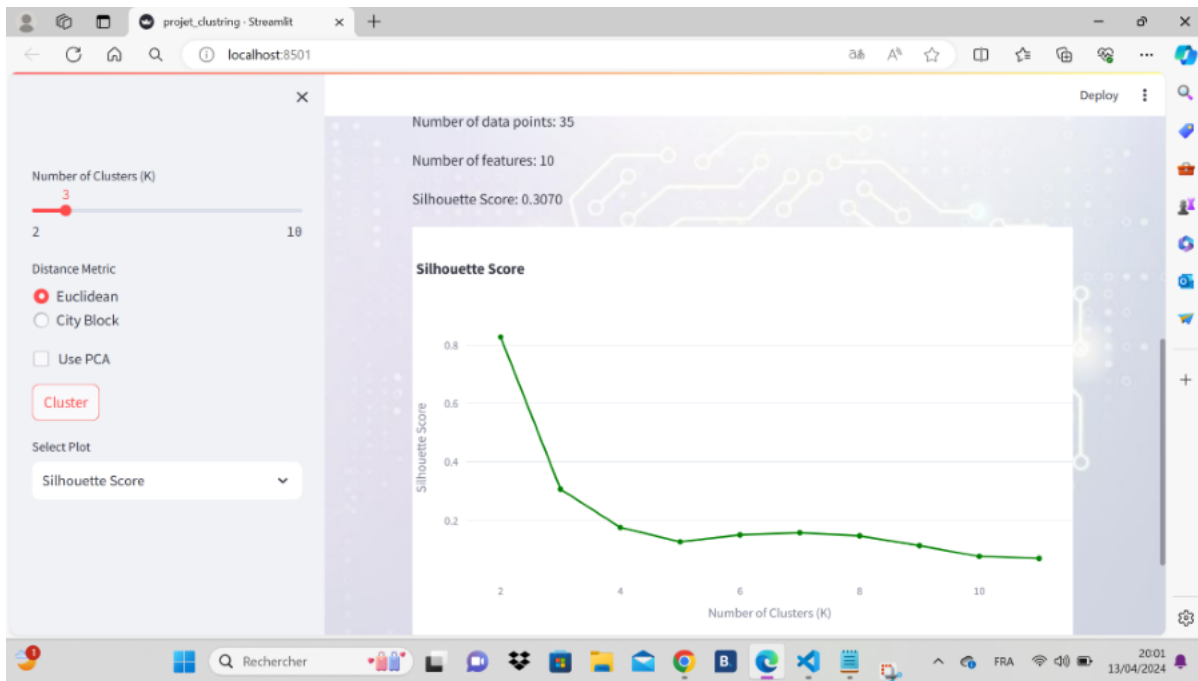


FIGURE 8 – étape 8

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

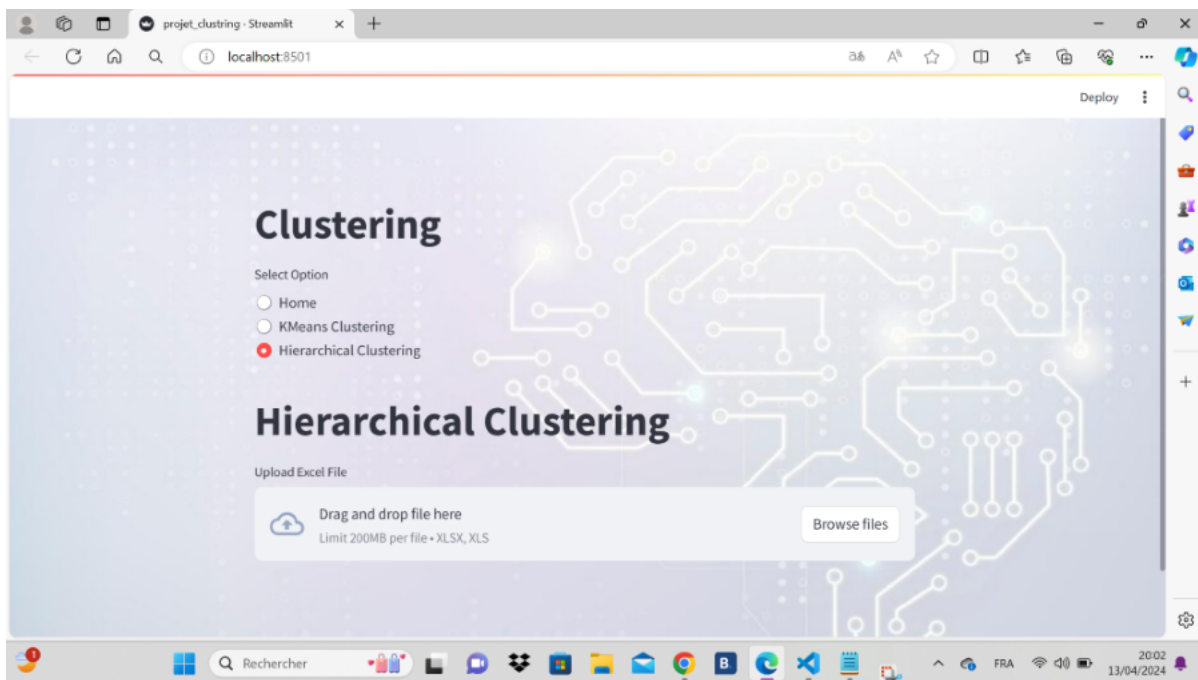


FIGURE 9 – étape 9

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>