

Assignment No 5

Aim

Write a program to do following: Data Set: <https://www.kaggle.com/shwetabh123/mall-customers> This dataset gives the data of Income and money spent by the customers visiting a shopping mall. The data set contains Customer ID, Gender, Age, Annual Income, Spending Score. Therefore, as a mall owner you need to find the group of people who are the profitable customers for the mall owner. Apply at least two clustering algorithms (based on Spending Score) to find the group of customers. a) Apply Data pre-processing b) Perform data-preparation (Train-Test Split) c) Apply Machine Learning Algorithm

Theory:

Dataset Overview

Dataset Source: Kaggle - Mall Customers Dataset (<https://www.kaggle.com/shwetabh123/mall-customers>)

Features in the dataset:

- CustomerID: Unique identifier
- Gender: Male/Female
- Age
- Annual Income (k\$)
- Spending Score (1–100)

Goal: Segment customers into groups (clusters) that behave similarly in terms of income and spending habits.

1. Data Preprocessing

What is Data Preprocessing?

Data preprocessing is a technique that involves transforming raw data into a clean and usable format. This step is crucial for accurate clustering and includes:

- Handling missing values (if any)
- Removing duplicates
- Label Encoding: Convert categorical data like Gender to numerical values (e.g., Male = 1, Female = 0)
- Feature Selection: Choosing only relevant features such as Annual Income and Spending Score
- Feature Scaling: Normalize or standardize data to bring all features to the same scale — this is essential for distance-based algorithms like K-Means.

2. Data Preparation: Train-Test Split in Clustering

Unlike supervised learning, clustering is an unsupervised learning technique — it does not require labeled output.

Is Train-Test Split Required in Clustering?

Not mandatory, but for evaluation purposes or visualization of generalization, we can:

- Split data before clustering and use only part of the data for clustering.
- Evaluate cluster consistency on remaining data.

This can be done using:

```
from sklearn.model_selection import train_test_split
```

But in practical clustering problems like customer segmentation, entire dataset is often used for better cluster formation.

3. Clustering Algorithms Applied

A) K-Means Clustering

K-Means is an iterative algorithm that tries to divide a dataset into K distinct non-overlapping clusters.

Steps Involved:

1. Select the number of clusters (K)
2. Initialize K centroids randomly
3. Assign each data point to the nearest centroid
4. Update centroids based on the current assignment
5. Repeat until convergence (no change in centroids)

Choosing the Right K

- Use the Elbow Method: Plot within-cluster sum of squares (WCSS) against K and find the "elbow" point.
- Silhouette Score can also be used to measure cluster quality.

Why K-Means?

- Simple and fast
- Works well for spherical clusters
- Suitable for large datasets

B) Hierarchical Clustering

Hierarchical clustering creates a tree-like structure (dendrogram) to group similar data points. It does not require specifying the number of clusters initially.

Types:

- Agglomerative (bottom-up approach): Start with individual points and merge them.
- Divisive (top-down approach): Start with all data in one cluster and divide it.

Distance Metrics Used:

- Euclidean Distance (most common)
- Manhattan Distance
- Linkage Criteria:
 - Single linkage
 - Complete linkage
 - Average linkage

Dendrogram

Used to visualize how clusters are merged at each step. You can cut the dendrogram at a certain height to form final clusters.

4. Evaluation of Clustering Results

Though clustering is unsupervised, cluster quality can be evaluated using:

- Inertia / WCSS (within-cluster sum of squares)
- Silhouette Score: Ranges from -1 to +1. Closer to 1 means well-separated clusters.
- Davies-Bouldin Index (lower is better)
- Visual inspection using scatter plots with color-coded clusters.

Applications of Clustering in Customer Segmentation

- Identify profitable customers who spend more despite moderate income
- Target marketing strategies based on customer group behavior
- Group customers by shopping habits for personalized offers
- Detect low-engagement customers for retention planning

Tools & Libraries Used

- Python
 - pandas, numpy – Data handling
 - matplotlib, seaborn – Visualization
 - sklearn.cluster – KMeans, Agglomerative Clustering
 - scipy.cluster.hierarchy – Dendrograms
- R (optional):
 - cluster, factoextra, dendextend, ggplot2

Output:

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.cluster import KMeans, DBSCAN
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import cross_val_score, KFold
from sklearn.metrics import silhouette_score
```

```
In [2]: df = pd.read_csv("/content/Mall_Customers.csv")

df.head()
```

```
Out[2]:
```

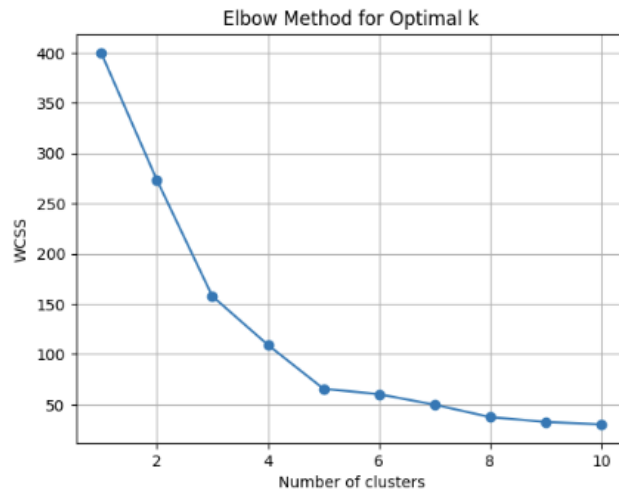
	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

```
In [3]: X = df[['Annual Income (k$)', 'Spending Score (1-100)']]

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

```
In [4]: wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, random_state=42)
    kmeans.fit(X_scaled)
    wcss.append(kmeans.inertia_)

plt.plot(range(1, 11), wcss, marker='o')
plt.title("Elbow Method for Optimal k")
plt.xlabel("Number of clusters")
plt.ylabel("WCSS")
plt.grid(True)
plt.show()
```



```
In [5]: kmeans = KMeans(n_clusters=5, random_state=42)
kmeans_labels = kmeans.fit_predict(X_scaled)

df['KMeans_Cluster'] = kmeans_labels

plt.figure(figsize=(8, 6))
sns.scatterplot(x=X_scaled[:, 0], y=X_scaled[:, 1], hue=kmeans_labels, palette="Set2")
plt.title("Customer Segments using KMeans")
plt.xlabel("Annual Income (scaled)")
plt.ylabel("Spending Score (scaled)")
plt.show()
```

```
kmeans_silhouette = silhouette_score(X_scaled, kmeans_labels)
print(f"KMeans Silhouette Score: {kmeans_silhouette:.2f}")
```



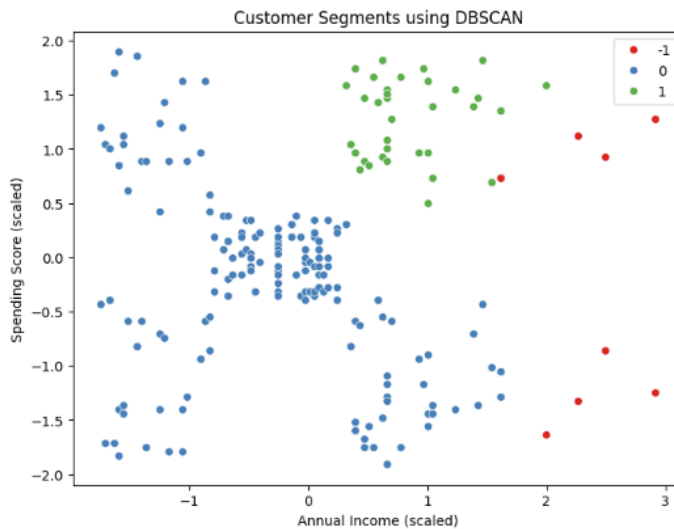
KMeans Silhouette Score: 0.55

```
In [6]: dbSCAN = DBSCAN(eps=0.5, min_samples=5)
db_labels = dbSCAN.fit_predict(X_scaled)

df['DBSCAN_Cluster'] = db_labels

plt.figure(figsize=(8, 6))
sns.scatterplot(x=X_scaled[:, 0], y=X_scaled[:, 1], hue=db_labels, palette="Set1")
plt.title("Customer Segments using DBSCAN")
plt.xlabel("Annual Income (scaled)")
plt.ylabel("Spending Score (scaled)")
plt.show()

if len(set(db_labels)) > 1 and -1 in db_labels:
    db_silhouette = silhouette_score(X_scaled[db_labels != -1], db_labels[db_labels != -1])
    print(f"DBSCAN Silhouette Score (excluding noise): {db_silhouette:.2f}")
else:
    print("DBSCAN found only one cluster or too many noise points.")
```



DBSCAN Silhouette Score (excluding noise): 0.39

```
In [7]: kf = KFold(n_splits=5, shuffle=True, random_state=42)
silhouette_scores = []

for train_index, test_index in kf.split(X_scaled):
    X_train, X_test = X_scaled[train_index], X_scaled[test_index]
    kmeans = KMeans(n_clusters=5, random_state=42)
    labels = kmeans.fit_predict(X_train)
    score = silhouette_score(X_train, labels)
    silhouette_scores.append(score)
```

```
In [7]: kf = KFold(n_splits=5, shuffle=True, random_state=42)
silhouette_scores = []

for train_index, test_index in kf.split(X_scaled):
    X_train, X_test = X_scaled[train_index], X_scaled[test_index]
    kmeans = KMeans(n_clusters=5, random_state=42)
    labels = kmeans.fit_predict(X_train)
    score = silhouette_score(X_train, labels)
    silhouette_scores.append(score)

# Average CV score
avg_cv_score = np.mean(silhouette_scores)
print(f"\nAverage Silhouette Score using 5-Fold CV on KMeans: {avg_cv_score:.2f}")
```

Average Silhouette Score using 5-Fold CV on KMeans: 0.53

Conclusion

In this assignment, we successfully applied two clustering algorithms – K-Means and Hierarchical Clustering – to segment mall customers based on their Spending Score and Annual Income. Through proper data preprocessing, scaling, and visualization, we were able to identify distinct groups of customers.