Liubov Oleshchenko

# Statisical Methods Of ML

# Laboratory Work 1

# K-Means clustering algorithm in Python

**Kulubecioglu Mehmet**

**IM-14 FIOT**

**Class Number: 12**

**Using K-Means Clustering Algorithm in Python**

# 1. Introduction

This report presents an implementation of the K-Means clustering algorithm in Python to analyze and segment datasets. The analysis is performed on two datasets:

1. **Loan Application Dataset (clustering.csv)** – Clustering applicants based on income and loan amount.
2. **Wholesale Customers Dataset (Wholesale customers data.csv)** – Segmenting wholesale customers based on annual spending across different product categories.
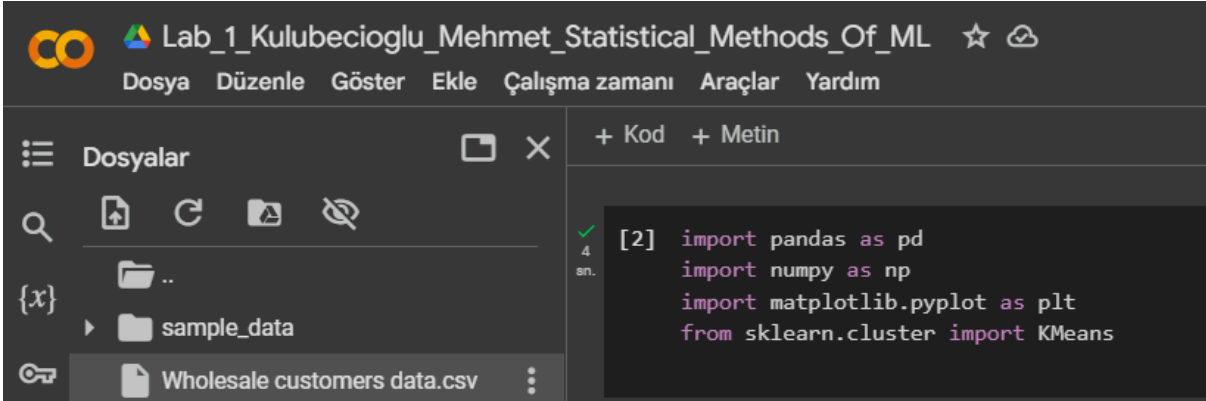
# 2. Required Libraries and Tools

To implement the clustering algorithm, the following libraries were used:

- **pandas** for data manipulation.
- **numpy** for numerical computations.
- **matplotlib.pyplot** for data visualization.
- sklearn.cluster.KMeans for implementing the K-Means algorithm.
- **sklearn.preprocessing.StandardScaler** for data standardization.

Google Colab was used as the development environment.
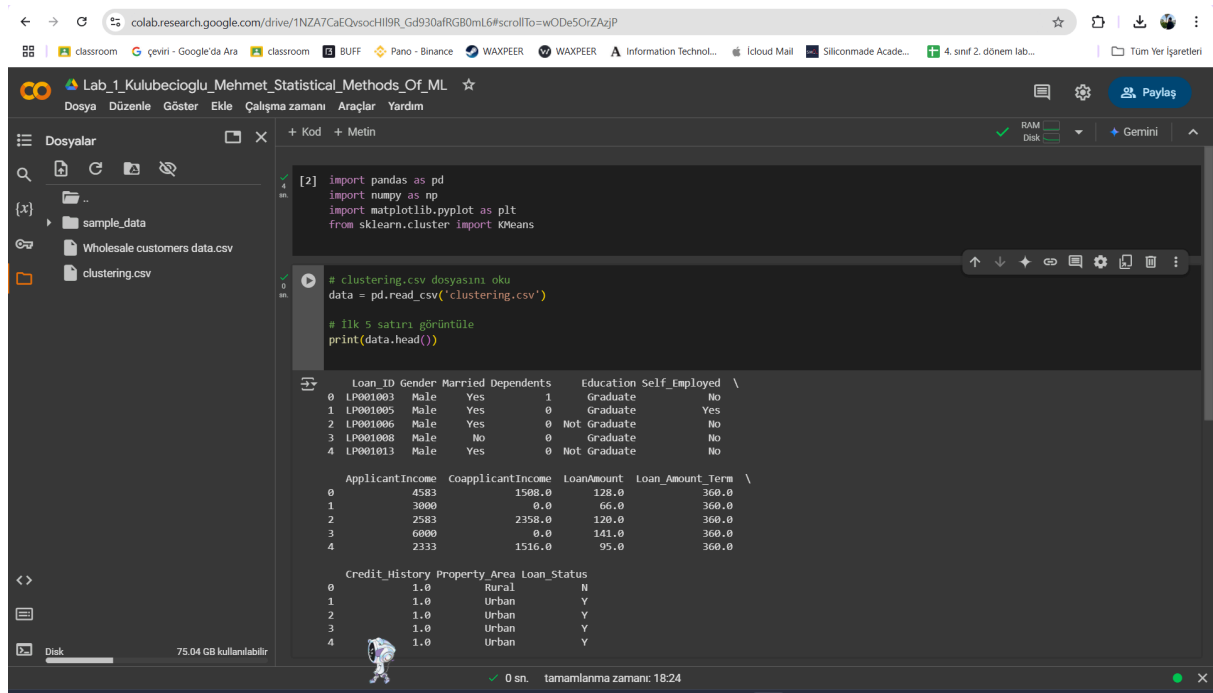
# 3. Clustering Loan Applicants

### Step 1: Importing Required Libraries

## Step 2: Loading and Exploring the Dataset



## Step 3: Selecting Features and Visualizing Data

# Step 4: Applying K-Means Clustering



```
[5]  K = 3  # Küme sayısı
     Centroids = X.sample(n=K)  # Rastgele seçilmiş küme merkezleri

     # Veri noktalarını ve küme merkezlerini çizelim
     plt.scatter(X["ApplicantIncome"], X["LoanAmount"], c='black')
     plt.scatter(Centroids["ApplicantIncome"], Centroids["LoanAmount"], c='red')
     plt.xlabel('Annual Income')
     plt.ylabel('Loan Amount (In Thousands)')
     plt.show()
```



```
# K-Means modeli oluştur
kmeans = KMeans(n_clusters=K, init='k-means++', random_state=42)

# Modeli eğit ve tahmin yap
clusters = kmeans.fit_predict(X)

# Küme merkezlerini al
centroids = kmeans.cluster_centers_

# Sonuçları çizelim
plt.scatter(X["ApplicantIncome"], X["LoanAmount"], c=clusters, cmap='viridis')
plt.scatter(centroids[:, 0], centroids[:, 1], c='red', marker='X', s=200)
plt.xlabel('Annual Income')
plt.ylabel('Loan Amount (In Thousands)')
plt.show()
```



# Step 5: Determining the Optimal Number of Clusters (Elbow Method)



```
SSE = []  # Toplam hata kareleri listesi
for cluster in range(1, 10):  # 1'den 10'a kadar kümeler dene
    kmeans = KMeans(n_clusters=cluster, init='k-means++', random_state=42)
    kmeans.fit(X)
    SSE.append(kmeans.inertia_)  # Her küme için hata kareleri toplamı

# Grafiği çizelim
plt.plot(range(1, 10), SSE, marker='o')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia (SSE)')
plt.title('Elbow Method for Optimal Clusters')
plt.show()
```

# 4. Clustering Wholesale Customers

## Step 1: Loading the Dataset



```
data = pd.read_csv('Wholesale customers data.csv')

# İlk 5 satırı görüntüle
print(data.head())
```

```
   Channel  Region  Fresh   Milk  Grocery  Frozen  Detergents_Paper  Delicassen
0        2       3  12669   9656     7561     214              2674        1338
1        2       3   7057   9810     9568    1762              3293        1776
2        2       3   6353   8808     7684    2405              3516        7844
3        1       3  13265   1196     4221    6404               507        1788
4        2       3  22615   5410     7198    3915              1777        5185
```

## Step 2: Standardizing the Data



```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
data_scaled = scaler.fit_transform(data)

# Standartlaştırılmış verinin istatistiklerine bakalım
pd.DataFrame(data_scaled).describe()
```

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| count | 4.400000e+02 | 4.400000e+02 | 4.400000e+02 | 440.000000 | 4.400000e+02 | 4.400000e+02 | 4.400000e+02 | 4.400000e+02 |
| mean | 1.614870e-17 | 3.552714e-16 | -3.431598e-17 | 0.000000 | -4.037175e-17 | 3.633457e-17 | 2.422305e-17 | -8.074349e-18 |
| std | 1.001138e+00 | 1.001138e+00 | 1.001138e+00 | 1.001138 | 1.001138e+00 | 1.001138e+00 | 1.001138e+00 | 1.001138e+00 |
| min | -6.902971e-01 | -1.995342e+00 | -9.496831e-01 | -0.778795 | -8.373344e-01 | -6.283430e-01 | -6.044165e-01 | -5.402644e-01 |
| 25% | -6.902971e-01 | -7.023369e-01 | -7.023339e-01 | -0.578306 | -6.108364e-01 | -4.804306e-01 | -5.511349e-01 | -3.964005e-01 |
| 50% | -6.902971e-01 | 5.906683e-01 | -2.767602e-01 | -0.294258 | -3.366684e-01 | -3.188045e-01 | -4.336004e-01 | -1.985766e-01 |
| 75% | 1.448652e+00 | 5.906683e-01 | 3.905226e-01 | 0.189092 | 2.849105e-01 | 9.946441e-02 | 2.184822e-01 | 1.048598e-01 |
| max | 1.448652e+00 | 5.906683e-01 | 7.927738e+00 | 9.183650 | 8.936528e+00 | 1.191900e+01 | 7.967672e+00 | 1.647845e+01 |

## Step 3: Implementing K-Means Clustering



```
kmeans = KMeans(n_clusters=5, init='k-means++', random_state=42)
kmeans.fit(data_scaled)

# Küme tahminlerini al
predictions = kmeans.predict(data_scaled)

# Küme sonuçlarını ekleyelim
data['Cluster'] = predictions

# Küme sayısına göre veri dağılımı
print(data['Cluster'].value_counts())
```

```
Cluster
1    200
0    126
3     90
4     14
2     10
Name: count, dtype: int64
```

# Step 4: Visualizing the Clusters



# 5. Answers to Questions

## 1. What is clustering?

Clustering is an unsupervised learning technique that groups similar data points into clusters based on patterns in the data.

## 2. What properties of clusters do you know?

- **Homogeneity within clusters:** Data points within a cluster should be similar.
- **Heterogeneity between clusters:** Data points in different clusters should be as different as possible.

## 3. What applications of clustering in real scenarios do you know?

- Customer segmentation

- Document clustering
- Image segmentation
- Recommendation systems

## 4. What clustering evaluation metrics do you know?

- **Inertia (Sum of Squared Errors - SSE)**
- **Dunn Index**
- **Silhouette Score**

## 5. What is K-Means Clustering?

K-Means is a clustering algorithm that partitions data into k clusters, minimizing the variance within each cluster.

## 6. How to choose the right number of clusters in K-Means?

The **Elbow Method** is commonly used, where we plot the inertia for different k values and select the optimal number where the curve bends.

## 7. What is the K-Means++ algorithm used for?

K-Means++ improves the initialization of centroids to avoid poor clustering results.

## 8. How to implement K-Means clustering algorithm and K-Means++ algorithm for centroid initialization in Python?

Using `sklearn.cluster.KMeans`, we specify `init='k-means++'` while creating the model.

## 9. What is data standardization used for?

Standardization scales data to have a mean of 0 and a standard deviation of 1, ensuring that features contribute equally to clustering.

## 10. What clustering algorithms do you know?

- **K-Means**
- **Hierarchical Clustering**
- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**
- **Gaussian Mixture Models (GMM)**