

# Implement K-Means clustering/ hierarchical clustering

- ✓ on sales\_data\_sample.csv dataset. Determine the number of clusters using the elbow method.


```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
#Importing the required libraries.
```

```
from sklearn.cluster import KMeans, k_means #For clustering
from sklearn.decomposition import PCA #Linear Dimensionality reduction.
```

```
df = pd.read_csv("sales_data_sample.csv") #Loading the dataset.
```

## ✓ Preprocessing

```
df.head()
```




|   | ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES   | ORDERDATE  |
|---|-------------|-----------------|-----------|-----------------|---------|------------|
| 0 | 10107       | 30              | 95.70     | 2               | 2871.00 | 2/24/2001  |
| 1 | 10121       | 34              | 81.35     | 5               | 2765.90 | 5/7/2001   |
| 2 | 10134       | 41              | 94.74     | 2               | 3884.34 | 7/1/2001   |
| 3 | 10145       | 45              | 83.26     | 6               | 3746.70 | 8/25/2001  |
| 4 | 10159       | 49              | 100.00    | 14              | 5205.27 | 10/10/2001 |

5 rows x 25 columns

df.shape

 (2823, 25)

df.describe()



|       | ORDERNUMBER  | QUANTITYORDERED | PRICEEACH   | ORDERLINENUMBER | SALES        |
|-------|--------------|-----------------|-------------|-----------------|--------------|
| count | 2823.000000  | 2823.000000     | 2823.000000 | 2823.000000     | 2823.000000  |
| mean  | 10258.725115 | 35.092809       | 83.658544   | 6.466171        | 3553.889072  |
| std   | 92.085478    | 9.741443        | 20.174277   | 4.225841        | 1841.865106  |
| min   | 10100.000000 | 6.000000        | 26.880000   | 1.000000        | 482.130000   |
| 25%   | 10180.000000 | 27.000000       | 68.860000   | 3.000000        | 2203.430000  |
| 50%   | 10262.000000 | 35.000000       | 95.700000   | 6.000000        | 3184.800000  |
| 75%   | 10333.500000 | 43.000000       | 100.000000  | 9.000000        | 4508.000000  |
| max   | 10425.000000 | 97.000000       | 100.000000  | 18.000000       | 14082.800000 |

```
df.info()
```

```
↗ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 2823 entries, 0 to 2822
Data columns (total 25 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ORDERNUMBER           2823 non-null   int64
1   QUANTITYORDERED       2823 non-null   int64
2   PRICEEACH             2823 non-null   float64
3   ORDERLINENUMBER       2823 non-null   int64
4   SALES                 2823 non-null   float64
5   ORDERDATE             2823 non-null   object
6   STATUS                2823 non-null   object
7   QTR_ID                2823 non-null   int64
8   MONTH_ID              2823 non-null   int64
9   YEAR_ID               2823 non-null   int64
10  PRODUCTLINE           2823 non-null   object
11  MSRP                  2823 non-null   int64
12  PRODUCTCODE           2823 non-null   object
13  CUSTOMERNAME          2823 non-null   object
14  PHONE                 2823 non-null   object
15  ADDRESSLINE1          2823 non-null   object
16  ADDRESSLINE2          302 non-null    object
17  CITY                  2823 non-null   object
18  STATE                 1337 non-null   object
19  POSTALCODE            2747 non-null   object
20  COUNTRY               2823 non-null   object
21  TERRITORY             1749 non-null   object
22  CONTACTLASTNAME       2823 non-null   object
23  CONTACTFIRSTNAME      2823 non-null   object
24  DEALSIZE              2823 non-null   object
dtypes: float64(2), int64(7), object(16)
memory usage: 551.5+ KB
```

df.isnull().sum()

|   |                  |      |
|---|------------------|------|
| ⇌ | ORDERNUMBER      | 0    |
|   | QUANTITYORDERED  | 0    |
|   | PRICEEACH        | 0    |
|   | ORDERLINENUMBER  | 0    |
|   | SALES            | 0    |
|   | ORDERDATE        | 0    |
|   | STATUS           | 0    |
|   | QTR_ID           | 0    |
|   | MONTH_ID         | 0    |
|   | YEAR_ID          | 0    |
|   | PRODUCTLINE      | 0    |
|   | MSRP             | 0    |
|   | PRODUCTCODE      | 0    |
|   | CUSTOMERNAME     | 0    |
|   | PHONE            | 0    |
|   | ADDRESSLINE1     | 0    |
|   | ADDRESSLINE2     | 2521 |
|   | CITY             | 0    |
|   | STATE            | 1486 |
|   | POSTALCODE       | 76   |
|   | COUNTRY          | 0    |
|   | TERRITORY        | 1074 |
|   | CONTACTLASTNAME  | 0    |
|   | CONTACTFIRSTNAME | 0    |
|   | DEALSIZE         | 0    |
|   | dtype: int64     |      |

df.dtypes

```
➡ ORDERNUMBER          int64
   QUANTITYORDERED      int64
   PRICEEACH            float64
   ORDERLINENUMBER      int64
   SALES                float64
   ORDERDATE            object
   STATUS               object
   QTR_ID               int64
   MONTH_ID             int64
   YEAR_ID              int64
   PRODUCTLINE          object
   MSRP                 int64
   PRODUCTCODE          object
   CUSTOMERNAME         object
   PHONE               object
   ADDRESSLINE1         object
   ADDRESSLINE2         object
   CITY                 object
   STATE               object
   POSTALCODE           object
   COUNTRY              object
   TERRITORY            object
   CONTACTLASTNAME      object
   CONTACTFIRSTNAME     object
   DEALSIZE             object
dtype: object
```

```
df_drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATUS', 'POSTALCODE', 'CITY', 'TEF
df = df.drop(df_drop, axis=1) #Dropping the categorical unnecessary columns al
```

df.isnull().sum()

```
➡ QUANTITYORDERED      0
   PRICEEACH            0
   ORDERLINENUMBER      0
   SALES                0
   ORDERDATE            0
   QTR_ID               0
   MONTH_ID             0
   YEAR_ID              0
   PRODUCTLINE          0
   MSRP                 0
   PRODUCTCODE          0
   COUNTRY              0
   DEALSIZE             0
dtype: int64
```

```
df.dtypes
```

```
➡ QUANTITYORDERED    int64
   PRICEEACH          float64
   ORDERLINENUMBER    int64
   SALES              float64
   ORDERDATE          object
   QTR_ID             int64
   MONTH_ID           int64
   YEAR_ID            int64
   PRODUCTLINE        object
   MSRP               int64
   PRODUCTCODE        object
   COUNTRY            object
   DEALSIZE           object
dtype: object
```

```
# Checking the categorical columns.
```

```
df['COUNTRY'].unique()
```

```
➡ array(['USA', 'France', 'Norway', 'Australia', 'Finland', 'Austria', 'UK',
        'Spain', 'Sweden', 'Singapore', 'Canada', 'Japan', 'Italy',
        'Denmark', 'Belgium', 'Philippines', 'Germany', 'Switzerland',
        'Ireland'], dtype=object)
```

```
df['PRODUCTLINE'].unique()
```

```
➡ array(['Motorcycles', 'Classic Cars', 'Trucks and Buses', 'Vintage Cars',
        'Planes', 'Ships', 'Trains'], dtype=object)
```

```
df['DEALSIZE'].unique()
```

```
➡ array(['Small', 'Medium', 'Large'], dtype=object)
```

```
productline = pd.get_dummies(df['PRODUCTLINE']) #Converting the categorical col
Dealsize = pd.get_dummies(df['DEALSIZE'])
```

```
df = pd.concat([df,productline,Dealsize], axis = 1)
```

```
df_drop = ['COUNTRY','PRODUCTLINE','DEALSIZE'] #Dropping Country too as there
df = df.drop(df_drop, axis=1)
```

```
df['PRODUCTCODE'] = pd.Categorical(df['PRODUCTCODE']).codes #Converting the dat
```

```
df.drop('ORDERDATE', axis=1, inplace=True) #Dropping the Orderdate as Month is
```

df.dtypes #All the datatypes are converted into numeric

```
➡ QUANTITYORDERED      int64
   PRICEEACH             float64
   ORDERLINENUMBER      int64
   SALES                 float64
   QTR_ID               int64
   MONTH_ID             int64
   YEAR_ID              int64
   MSRP                 int64
   PRODUCTCODE          int8
   Classic Cars          uint8
   Motorcycles           uint8
   Planes                uint8
   Ships                uint8
   Trains                uint8
   Trucks and Buses      uint8
   Vintage Cars          uint8
   Large                 uint8
   Medium                uint8
   Small                 uint8
dtype: object
```

## ✓ Plotting the Elbow Plot to determine the number of clusters.

```
distortions = [] # Within Cluster Sum of Squares from the centroid
K = range(1,10)
for k in K:
    kmeanModel = KMeans(n_clusters=k)
    kmeanModel.fit(df)
    distortions.append(kmeanModel.inertia_)    #Appeding the intertia to the Dis
```

```
➡ -----
-
NameError                                Traceback (most recent call
last)
<ipython-input-1-c91e41b47b20> in <module>
      2 K = range(1,10)
      3 for k in K:
----> 4     kmeanModel = KMeans(n_clusters=k)
      5     kmeanModel.fit(df)
      6     distortions.append(kmeanModel.inertia_)    #Appeding the
intertia to the Distortions
```

```
plt.figure(figsize=(16,8))
plt.plot(K, distortions, 'bx-')
plt.xlabel('k')
plt.ylabel('Distortion')
plt.title('The Elbow Method showing the optimal k')
plt.show()
```

-----

**NameError** Traceback (most recent call last)

<ipython-input-2-afc58ba4fec1> in <module>

----> 1 plt.figure(figsize=(16,8))

2 plt.plot(K, distortions, 'bx-')

3 plt.xlabel('k')

4 plt.ylabel('Distortion')

5 plt.title('The Elbow Method showing the optimal k')

**NameError**: name 'plt' is not defined

✓ As the number of k increases Inertia decreases.

Observations: A Elbow can be observed at 3 and after that the curve decreases gradually.

X\_train = df.values #Returns a numpy array.

-----

**NameError** Traceback (most recent call last)

<ipython-input-3-5ccbd1c8ee61> in <module>

----> 1 X\_train = df.values #Returns a numpy array.

**NameError**: name 'df' is not defined



X\_train.shape

```
-----  
-  
NameError                                Traceback (most recent call  
last)  
<ipython-input-4-d2ba684acd0f> in <module>  
----> 1 X_train.shape  
  
NameError: name 'X train' is not defined
```

```
model = KMeans(n_clusters=3,random_state=2) #Number of cluster = 3  
model = model.fit(X_train) #Fitting the values to create a model.  
predictions = model.predict(X_train) #Predicting the cluster values (0,1,or 2)
```

```
unique,counts = np.unique(predictions,return_counts=True)
```

```
counts = counts.reshape(1,3)
```

```
counts_df = pd.DataFrame(counts,columns=['Cluster1','Cluster2','Cluster3'])
```

```
counts_df.head()
```

## ✓ Visualization

```
pca = PCA(n_components=2) #Converting all the features into 2 columns to make i
```

```
-----  
-  
NameError                                Traceback (most recent call  
last)  
<ipython-input-5-026957ba5f55> in <module>  
----> 1 pca = PCA(n_components=2) #Converting all the features into 2  
columns to make it easy to visualize using Principal COmponent Analysis.
```

```
reduced_X = pd.DataFrame(pca.fit_transform(X_train),columns=['PCA1','PCA2']) #C
```



```
-----  
-  
NameError                                Traceback (most recent call  
last)  
<ipython-input-6-af82ff426f69> in <module>  
----> 1 reduced_X = pd.DataFrame(pca.fit_transform(X_train),columns=  
['PCA1','PCA2']) #Creating a DataFrame.
```

```
reduced_X.head()
```

```
#Plotting the normal Scatter Plot  
plt.figure(figsize=(14,10))  
plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
```



```
-----  
-  
NameError                                Traceback (most recent call  
last)  
<ipython-input-7-8c1f248abdc8> in <module>  
      1 #Plotting the normal Scatter Plot  
----> 2 plt.figure(figsize=(14,10))  
      3 plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])  
  
NameError: name 'plt' is not defined
```

```
model.cluster_centers_ #Finding the centriods. (3 Centriods in total. Each Arra
```

```
reduced_centers = pca.transform(model.cluster_centers_) #Transforming the centr
```

```
reduced_centers
```

```
plt.figure(figsize=(14,10))
plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
plt.scatter(reduced_centers[:,0],reduced_centers[:,1],color='black',marker='x',
```

```

-----
-
NameError                                Traceback (most recent call
last)
<ipython-input-9-be5dfb45822d> in <module>
----> 1 plt.figure(figsize=(14,10))
      2 plt.scatter(reduced_X['PCA1'],reduced_X['PCA2'])
      3
plt.scatter(reduced_centers[:,0],reduced_centers[:,1],color='black',marker=
```

```
reduced_X['Clusters'] = predictions #Adding the Clusters to the reduced dataframe
```

```
reduced_X.head()
```

```
#Plotting the clusters
```

```
plt.figure(figsize=(14,10))
```

```
#                                taking the cluster number and first column                                taki
plt.scatter(reduced_X[reduced_X['Clusters'] == 0].loc[:, 'PCA1'],reduced_X[reduc
plt.scatter(reduced_X[reduced_X['Clusters'] == 1].loc[:, 'PCA1'],reduced_X[reduc
plt.scatter(reduced_X[reduced_X['Clusters'] == 2].loc[:, 'PCA1'],reduced_X[reduc
```

```
plt.scatter(reduced_centers[:,0],reduced_centers[:,1],color='black',marker='x',
```

```

-----
-
NameError                                Traceback (most recent call
last)
<ipython-input-10-7e26f16fd27d> in <module>
      1 #Plotting the clusters
----> 2 plt.figure(figsize=(14,10))
      3 #                                taking the cluster number and first column
taking the same cluster number and second column                                Assigning the color
      4 plt.scatter(reduced_X[reduced_X['Clusters'] ==
0].loc[:, 'PCA1'],reduced_X[reduced_X['Clusters'] ==
0].loc[:, 'PCA2'],color='slateblue')
```

Start coding or [generate](#) with AI.