

K Means Clustering Algorithm

Dataset- Mall Customer dataset

importing the libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
```

Loading the data

```
In [2]: customer_data = pd.read_csv('C:/Users/lekshmi/Downloads/Mall_Customers.csv')
```

Understanding the data and PreProcessing

```
In [12]: customer_data.head()
```

Out[12]:

	CustomerID	Gender	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 [13]: customer_data.shape
```

Out[13]: (200, 5)

```
In [14]: customer_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column                      Non-Null Count  Dtype
---  -
0   CustomerID                  200 non-null   int64
1   Gender                      200 non-null   object
2   Age                         200 non-null   int64
3   Annual Income (k$)          200 non-null   int64
4   Spending Score (1-100)      200 non-null   int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
```

checking missing values

```
In [3]: customer_data.isnull().sum()
```

```
Out[3]: CustomerID          0
        Gender             0
        Age                0
        Annual Income (k$)  0
        Spending Score (1-100)  0
        dtype: int64
```

choosing attributes Annual income and Spending score

```
In [5]: X = customer_data.iloc[:,[3,4]].values
```

```
X
```

Scaling the variables

```
In [7]: scaler = StandardScaler()
        X = scaler.fit_transform(X)
```

choosing number of clusters

wcss- Within Clusters Sum of Squares

In [8]: *# finding wcss value for different number of clusters(upto 10)*

```
wcss = []
```

```
for i in range(1,11):
```

```
    kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42)
```

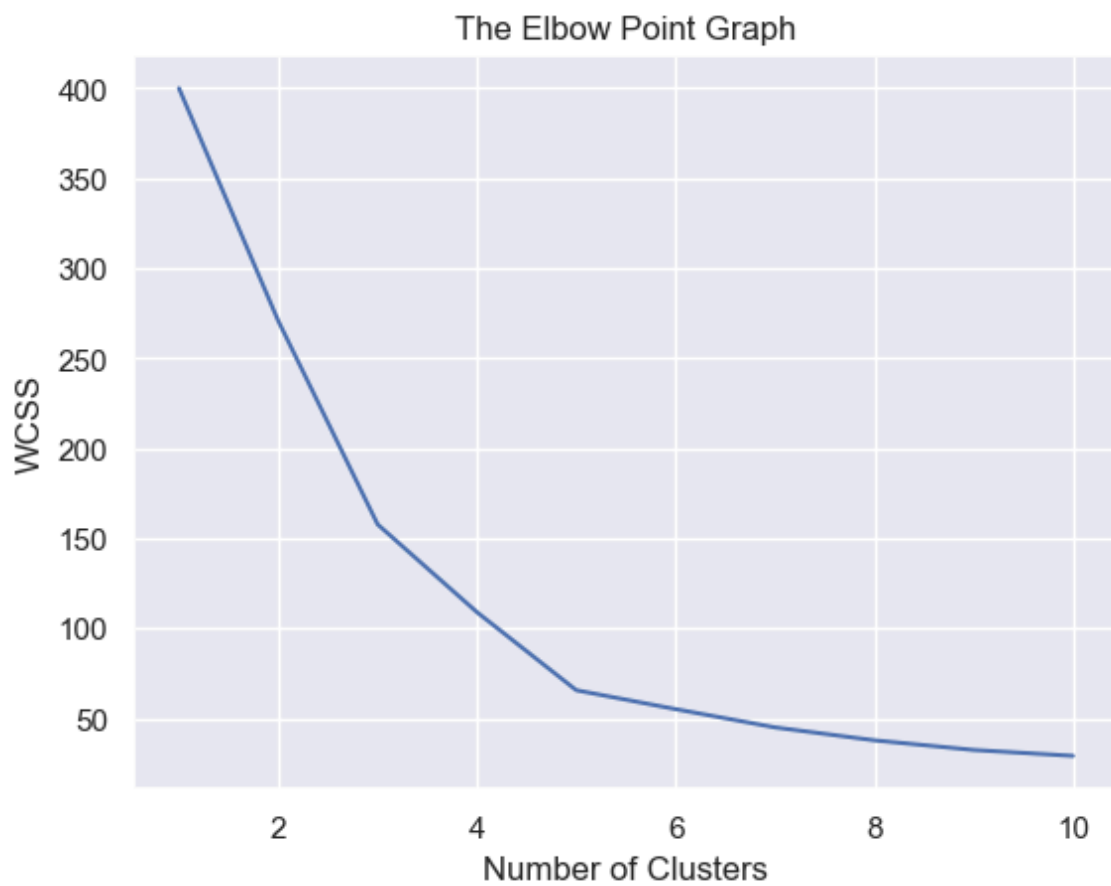
```
    kmeans.fit(X)
```

```
    wcss.append(kmeans.inertia_)
```

```
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
```

```
setting the environment variable OMP_NUM_THREADS=1.
warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
warnings.warn(
```

```
In [9]: sns.set()
plt.plot(range(1,11), wcss)
plt.title('The Elbow Point Graph')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS')
plt.show()
```



Optimum Number of Clusters = 5

Training the Kmeans Clustering model

```
In [10]: kmeans = KMeans(n_clusters=5, init='k-means++', random_state=42)

# return a label for each data point based on their cluster
Y = kmeans.fit_predict(X)

print(Y)
```

```
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:87
0: FutureWarning: The default value of `n_init` will change from 10 to 'au
to' in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\lekshmi\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:13
82: UserWarning: KMeans is known to have a memory leak on Windows with MK
L, when there are less chunks than available threads. You can avoid it by
setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(

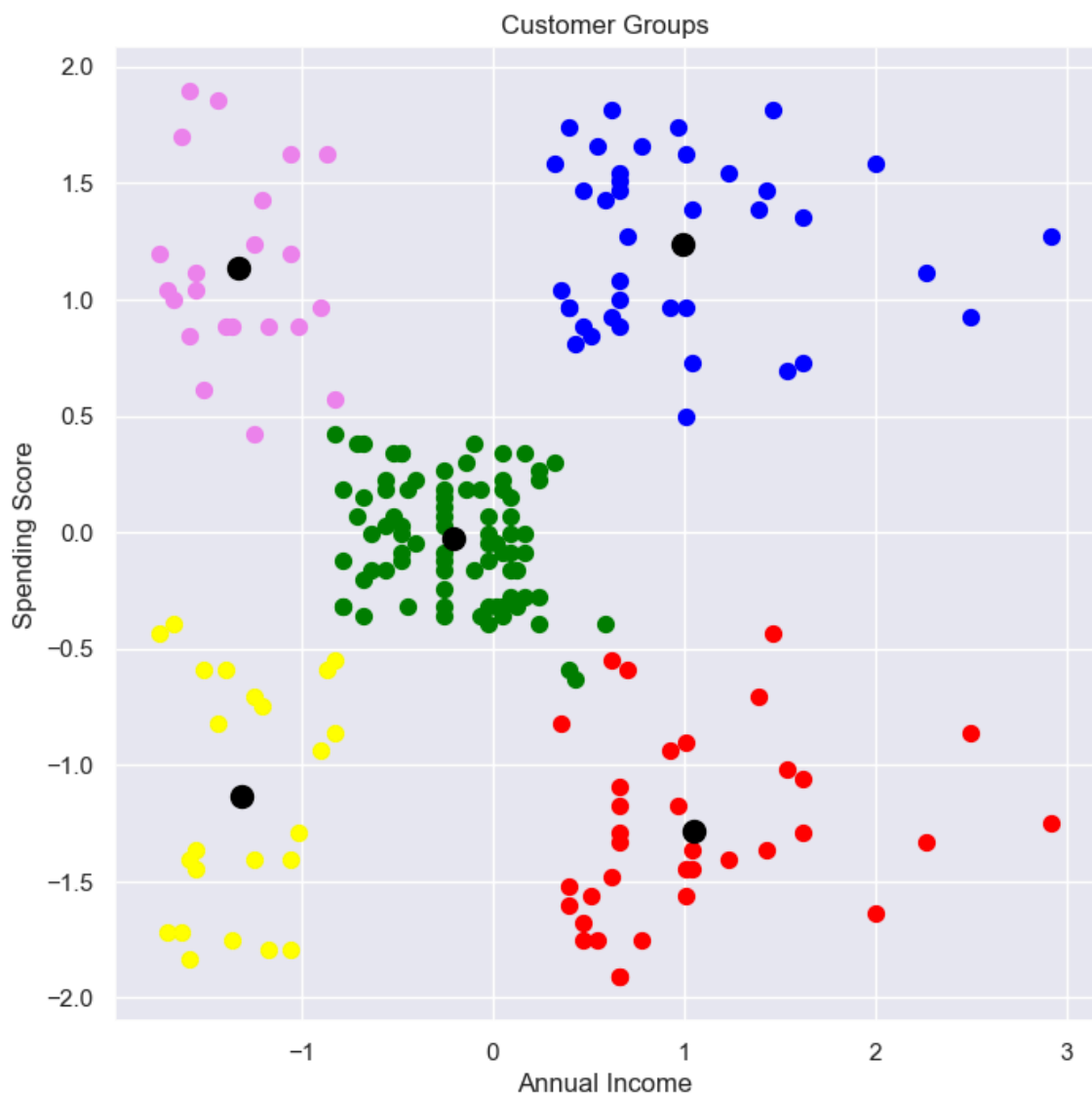
[2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2
 3 2 3 2 3 2 0 2 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 4 1 4 0 4 1 4 1 4 0 4 1 4 1 4 1 4
1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1
4 1 4 1 4 1 4 1 4 1 4 1 4]
```

Visualising clusters

```
In [11]: plt.figure(figsize=(8,8))
plt.scatter(X[Y==0,0], X[Y==0,1], s=50, c='green', label='Cluster 1')
plt.scatter(X[Y==1,0], X[Y==1,1], s=50, c='red', label='Cluster 2')
plt.scatter(X[Y==2,0], X[Y==2,1], s=50, c='yellow', label='Cluster 3')
plt.scatter(X[Y==3,0], X[Y==3,1], s=50, c='violet', label='Cluster 4')
plt.scatter(X[Y==4,0], X[Y==4,1], s=50, c='blue', label='Cluster 5')

plt.scatter(kmeans.cluster_centers_[0,0], kmeans.cluster_centers_[0,1], s=100, c='black')
plt.scatter(kmeans.cluster_centers_[1,0], kmeans.cluster_centers_[1,1], s=100, c='black')
plt.scatter(kmeans.cluster_centers_[2,0], kmeans.cluster_centers_[2,1], s=100, c='black')
plt.scatter(kmeans.cluster_centers_[3,0], kmeans.cluster_centers_[3,1], s=100, c='black')
plt.scatter(kmeans.cluster_centers_[4,0], kmeans.cluster_centers_[4,1], s=100, c='black')

plt.title('Customer Groups')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.show()
```



cluster 4 has high spending score but low annual income

cluster 2 has high annual income but low spending score

cluster 5 has both high annual income and spending score

K-Means is a popular unsupervised machine learning algorithm used for clustering, where it partitions a dataset into a predetermined number of clusters by iteratively

assigning data points to the cluster with the nearest centroid and updating the centroids based on the mean of the assigned points. The algorithm aims to minimize the within-cluster sum of squared distances, providing a simple and efficient method for grouping data into distinct clusters based on similarity.