

# Customer Segmentation using Clustering

- K-means clustering model

## Importing Dependencies

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

from sklearn.cluster import KMeans
```

## Data Collection and Analysis

```
In [2]: df = pd.read_csv("C://Users//ankit//anaconda3//jupyter_script//AI_pro//Mall_Customers.csv")
```

```
In [3]: df.head()
```

```
Out[3]:
```

	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 [4]: df.tail()
```

```
Out[4]:
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

```
In [5]: # shape of dataset
```

```
df.shape
```

```
Out[5]: (200, 5)
```

```
In [6]: # information about dataset
```

```
df.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

In [7]: `# datatypes of columns dataset`

`df.dtypes`

Out[7]:

CustomerID	int64
Gender	object
Age	int64
Annual Income (k\$)	int64
Spending Score (1-100)	int64

dtype: object

In [8]: `# statistical measures of dataset`  
`df.describe()`

Out[8]:

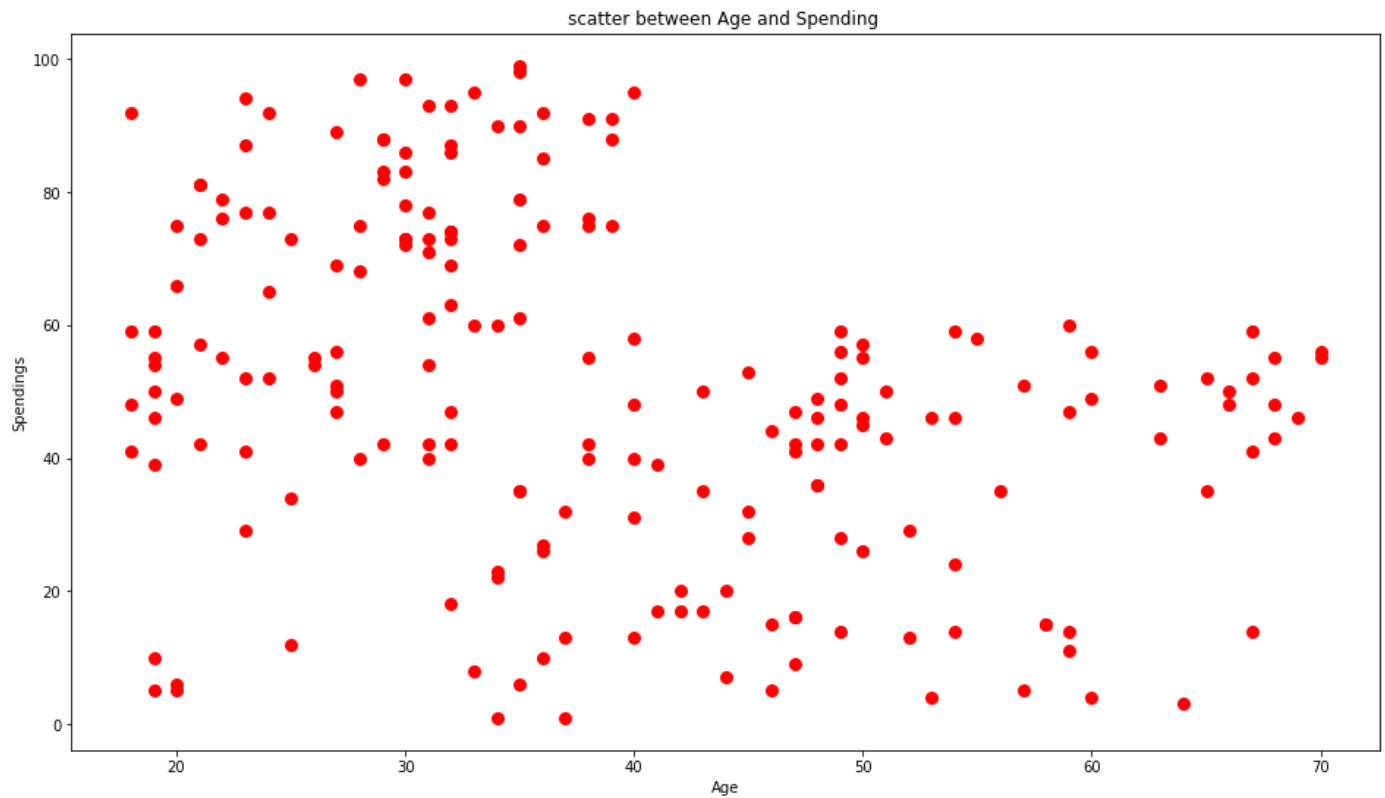
	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
<b>count</b>	200.000000	200.000000	200.000000	200.000000
<b>mean</b>	100.500000	38.850000	60.560000	50.200000
<b>std</b>	57.879185	13.969007	26.264721	25.823522
<b>min</b>	1.000000	18.000000	15.000000	1.000000
<b>25%</b>	50.750000	28.750000	41.500000	34.750000
<b>50%</b>	100.500000	36.000000	61.500000	50.000000
<b>75%</b>	150.250000	49.000000	78.000000	73.000000
<b>max</b>	200.000000	70.000000	137.000000	99.000000

## Data Visulalization

In [9]:

```
plt.figure(figsize = (16,9))
plt.scatter(x = df.Age, y = df['Spending Score (1-100)'], c = "red", linewidth = 3)
plt.title('scatter between Age and Spending ')
plt.ylabel("Spending")
plt.xlabel("Age")

plt.show()
```

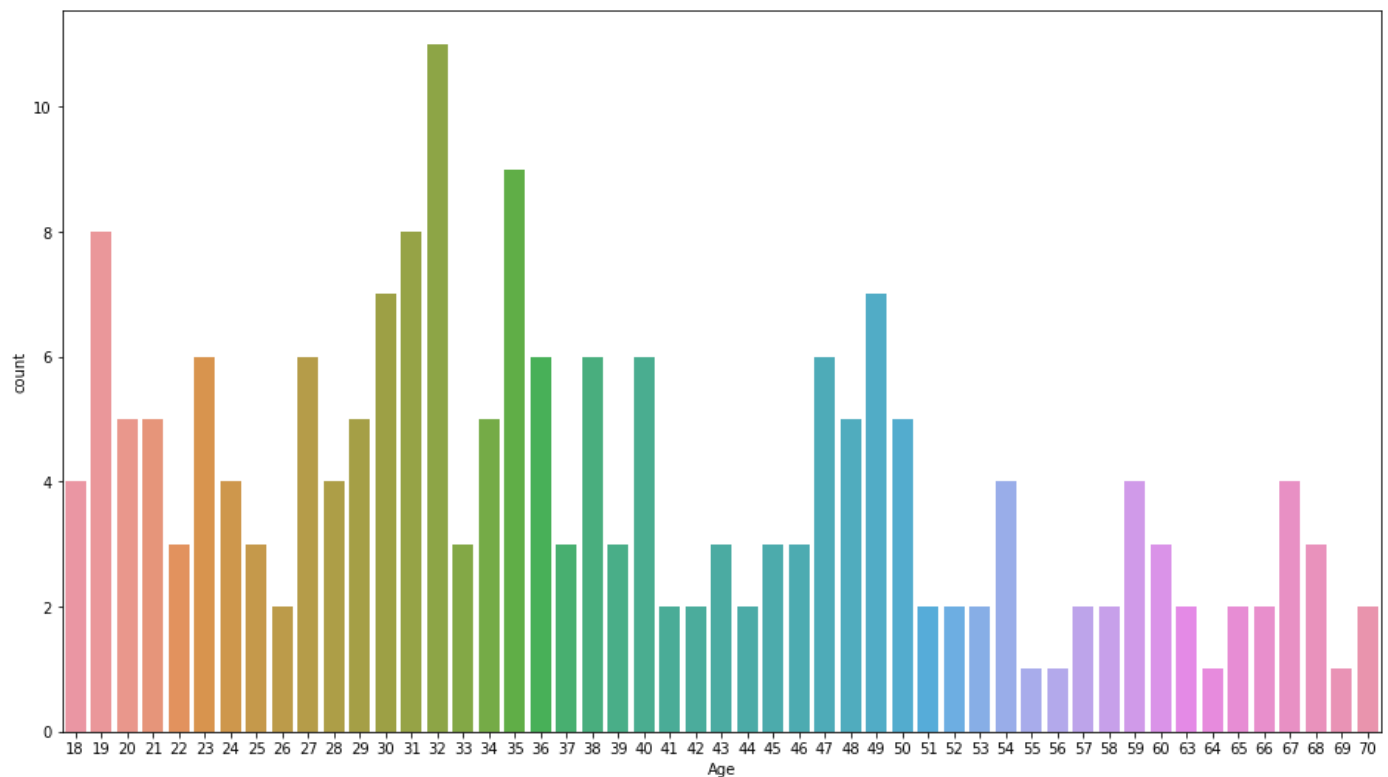


```
In [10]: plt.figure(figsize = (16,9))
sns.countplot(df.Age)
# range of age
```

C:\Users\ankit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(

```
Out[10]: <AxesSubplot:xlabel='Age', ylabel='count'>
```

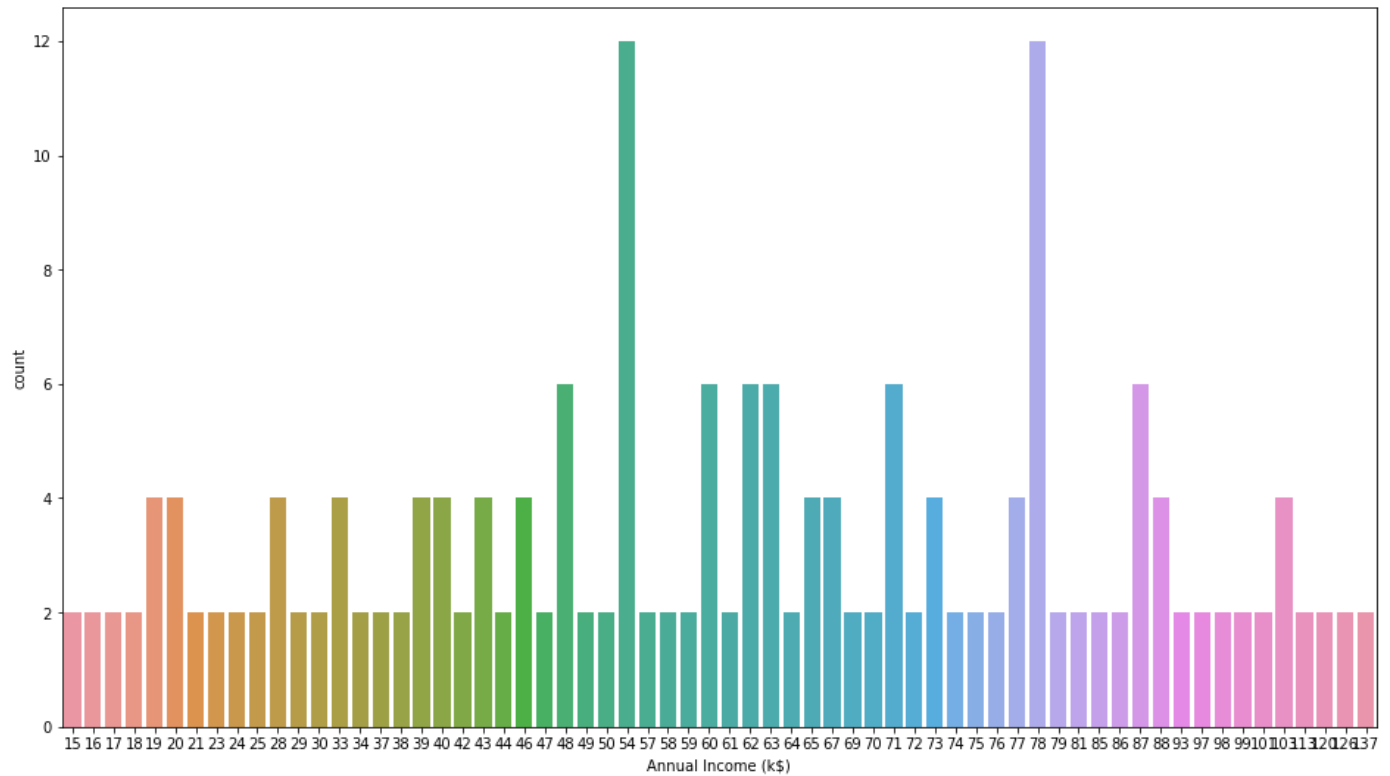


```
In [11]: plt.figure(figsize = (16,9))
sns.countplot(df['Annual Income (k$)'])
```

C:\Users\ankit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

Out[11]: <AxesSubplot:xlabel='Annual Income (k\$)', ylabel='count'>

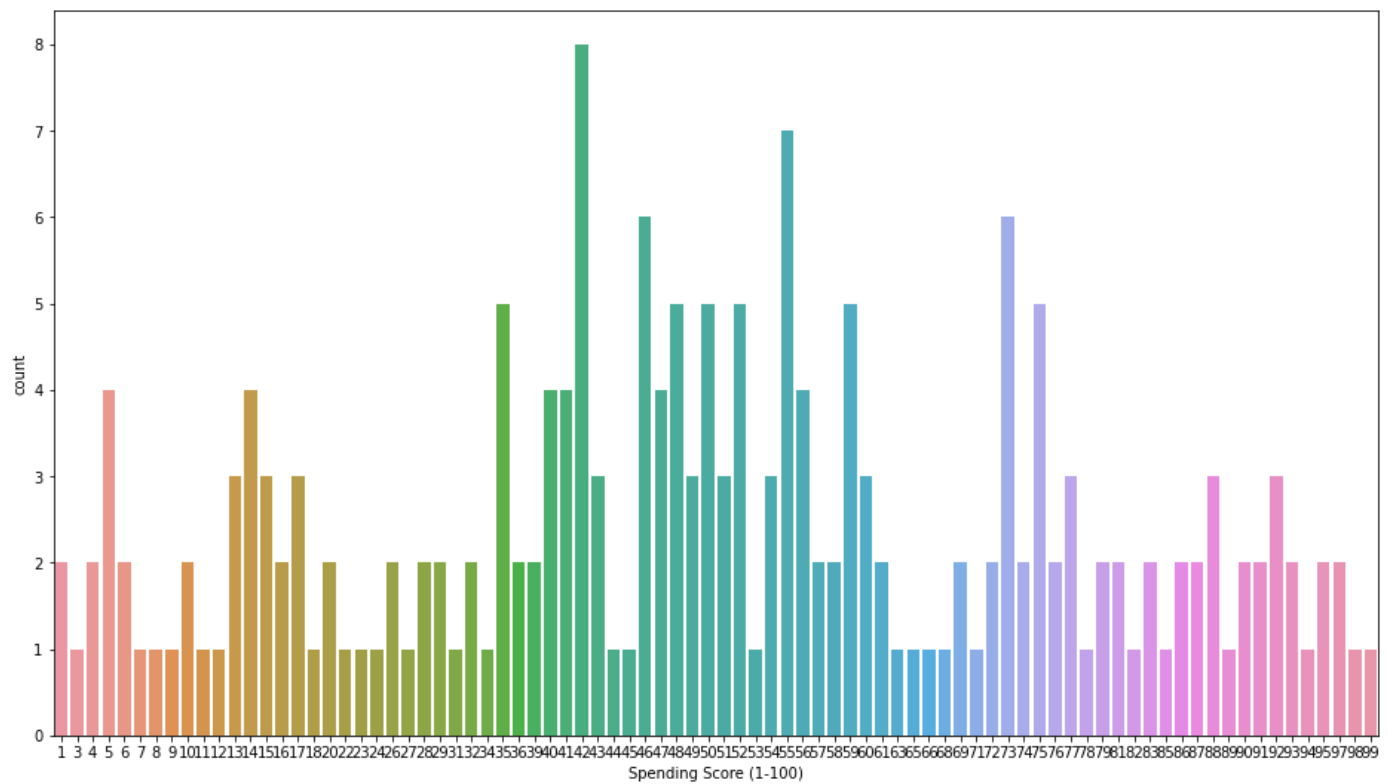


In [12]: `plt.figure(figsize = (16,9))`  
`sns.countplot(df['Spending Score (1-100)'])`

C:\Users\ankit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

Out[12]: <AxesSubplot:xlabel='Spending Score (1-100)', ylabel='count'>

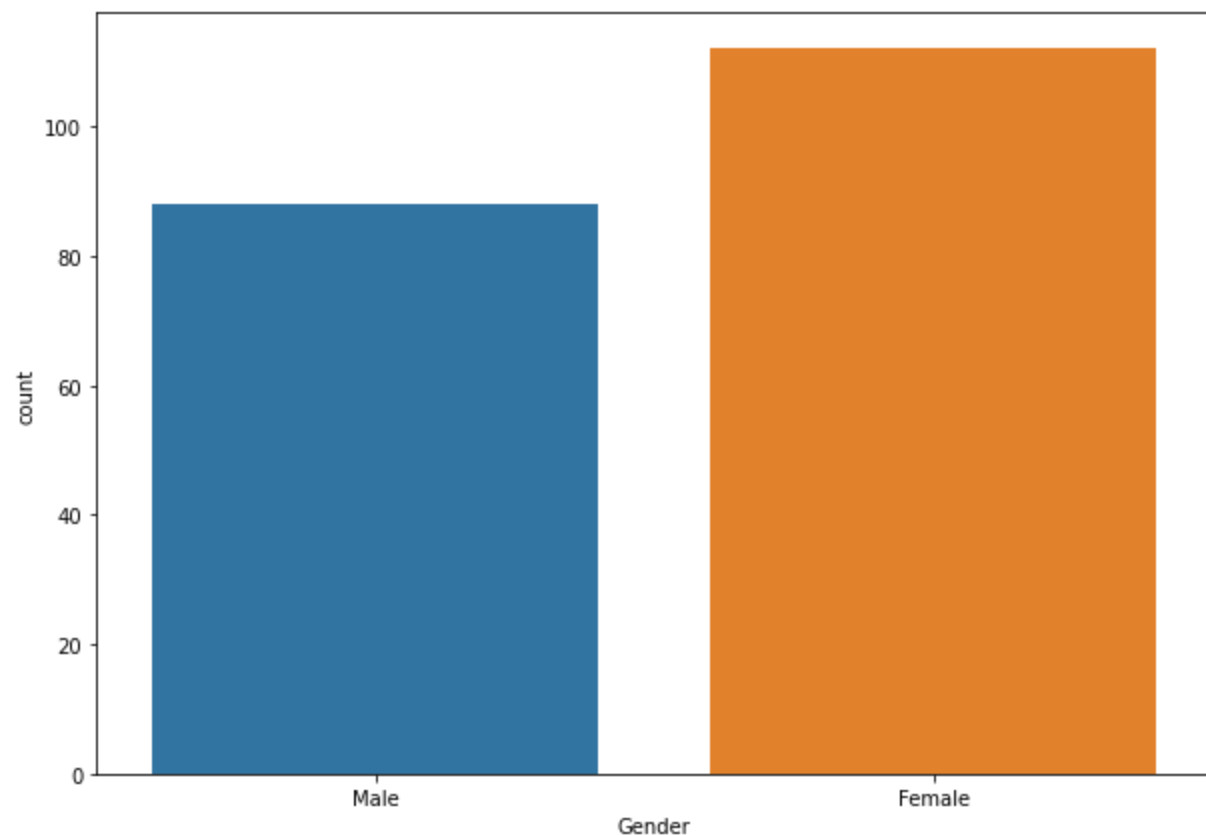


```
In [13]: plt.figure(figsize = (10,7))
sns.countplot(df['Gender'])
```

C:\Users\ankit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(

```
Out[13]: <AxesSubplot:xlabel='Gender', ylabel='count'>
```



**Choosing the Annual Income Column and Spending Columns**

```
In [14]: df.head()
```

```
Out[14]:
```

	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 [15]: df.columns
```

```
Out[15]: Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',  
              'Spending Score (1-100)'],  
              dtype='object')
```

```
In [16]: # using iloc for picking 3rd and 4th columns  
X = df.iloc[:, [3,4]].values
```

```
In [17]: print(X[:10,:10])  
# printing first 10 values of 3rd and 4th columns
```

```
[[15 39]  
 [15 81]  
 [16  6]  
 [16 77]  
 [17 40]  
 [17 76]  
 [18  6]  
 [18 94]  
 [19  3]  
 [19 72]]
```

## choosing the number of clusters

WCSS : within Clusters Sum of Squares

```
In [18]: # finding WCSS values for different number of clusters  
  
wcscs = []  
  
for i in range(1,11):  
    # 1 and 11 will be excluded  
    kmeans = KMeans(n_clusters=i, init = 'k-means++', random_state = 42)  
    kmeans.fit(X)  
    wcscs.append(kmeans.inertia_)
```

```
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarning:  
The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of  
'n_init' explicitly to suppress the warning  
    warnings.warn(  

```

```
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:  
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th  
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA  
DS=1.  
    warnings.warn(  

```

```
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarning:  
The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of  
'n_init' explicitly to suppress the warning
```

```
warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
    warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
```

```

an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
warnings.warn(
C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning:
KMeans is known to have a memory leak on Windows with MKL, when there are less chunks th
an an available threads. You can avoid it by setting the environment variable OMP_NUM_THREA
DS=1.
warnings.warn(

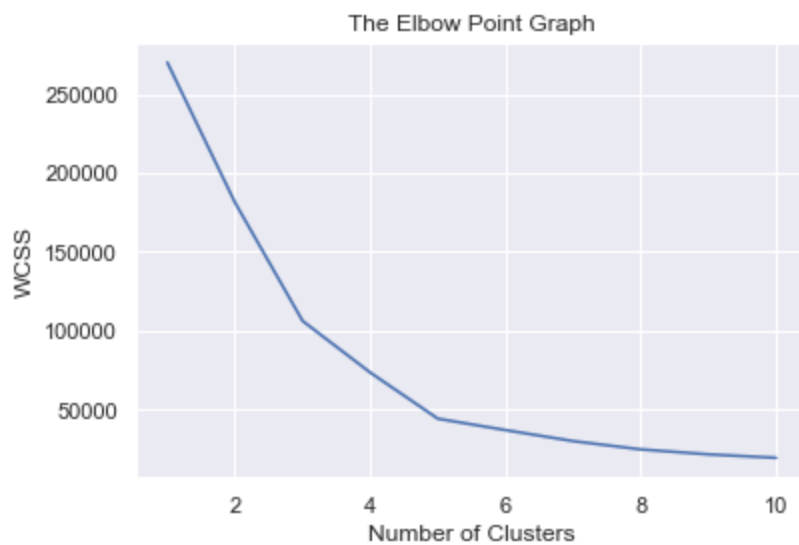
```

In [19]: *# plot an elbow graph*

```

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 K-Means Clustering model

```

In [20]: kmeans = KMeans(n_clusters =5, init = 'k-means++', random_state =0)

# return a label for each data point based on their clusters

```

```

In [21]: Y = kmeans.fit_predict(X)

print(Y)

```

```

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

```

```

C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: FutureWarnin
g: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of
`n_init` explicitly to suppress the warning
warnings.warn(

```



C:\Users\ankit\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.

```
warnings.warn(
```

## Plotting data in graph

```
In [22]: # plotting all the clusters and their centroids

plt.figure(figsize = (16,9))
plt.title("Customer Segmentation", fontsize = 15)
plt.scatter(X[Y==0,0], X[Y==0,1], s = 50, c = 'orange', label = 'Cluster 1')
plt.scatter(X[Y==1,0], X[Y==1,1], s = 50, c = 'green', label = 'Cluster 1')
plt.scatter(X[Y==2,0], X[Y==2,1], s = 50, c = 'red', label = 'Cluster 2')
plt.scatter(X[Y==3,0], X[Y==3,1], s = 50, c = 'blue', label = 'Cluster 3')
plt.scatter(X[Y==4,0], X[Y==4,1], s = 50, c = 'yellow', label = 'Cluster 4')

# plot the centroids
plt.xlabel("Annual Income")
plt.ylabel("Spending Score")
plt.scatter(kmeans.cluster_centers[:,0], kmeans.cluster_centers[:,1], s = 100, c = 'black')
```

```
Out[22]: <matplotlib.collections.PathCollection at 0x1ec3b5ef850>
```

