

About our Dataset

The Mall customers dataset contains information about people visiting the mall. The dataset has gender, customer id, age, annual income, and spending score. It collects insights from the data and group customers based on their behaviors.

Our Target

Segment the customers based on the age, gender, interest. Customer segmentation is an important practise of dividing customers base into individual groups that are similar. It is useful in customised marketing.

Importing all the important library

```
In [1]: import pandas as pd  
import numpy as np  
import seaborn as sns  
import matplotlib.pyplot as plt  
import plotly.express as px  
%matplotlib inline
```

Now reading our.csv file using Pandas Library

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

```
Out[2]:   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
...
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

200 rows x 5 columns

Checking for Null Values or missing Values in our Dataset

```
In [3]: df.isna().sum()
```

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

Data Exploration

```
In [4]: 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   Genre            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 [5]: df.describe()
```

```
Out[5]:   CustomerID  Age  Annual Income (k$)  Spending Score (1-100)
```

count	200.000000	200.000000	200.000000
mean	100.500000	38.850000	60.560000
std	57.871985	13.969007	26.647211
min	1.000000	18.000000	15.000000
25%	50.750000	28.750000	41.500000
50%	100.500000	36.000000	61.500000
75%	150.250000	49.000000	78.000000
max	200.000000	70.000000	137.000000

Plotting all the features against all other features.

We do this to get the main attributes that should be used for classifications

```
In [6]: sns.pairplot(df , diag_kind = 'kde' , height = 3)
```

```
Out[6]: <seaborn.axisgrid.PairGrid at 0x234095fd0c0>
```



Removing the unneeded data in our data. Here that is Customer ID

```
In [9]: df = df.iloc[:,1:]  
df
```

```
Out[9]:   Gender  Age  Annual Income (k$)  Spending Score (1-100)
```

0	Male	19	15	39
1	Male	21	15	81
2	Female	20	16	6
3	Female	23	16	77
4	Female	31	17	40
...
195	Female	35	120	79
196	Female	45	126	28
197	Male	32	126	74
198	Male	32	137	18
199	Male	30	137	83

200 rows x 4 columns

Transforming the Gender or ['Genre'] column into numeric data.

We are using LabelEncoding for this

```
In [10]: from sklearn.preprocessing import LabelEncoder  
encoder = LabelEncoder()  
df['Genre'] = encoder.fit_transform(df['Genre'])  
df
```

```
Out[10]:   Gender  Age  Annual Income (k$)  Spending Score (1-100)
```

0	1	19	15	39
1	1	21	15	81
2	0	20	16	6
3	0	23	16	77
4	0	31	17	40
...
195	0	35	120	79
196	0	45	126	28
197	1	32	126	74
198	1	32	137	18
199	1	30	137	83

200 rows x 4 columns

Creating a New DataFrame for the Centers of Clusters as it helps in Plotting

```
In [21]: df2 = pd.DataFrame(kmeans.cluster_centers_)
```

```
Out[21]:   0   1   2
```

0	32.692308	86.538462	82.12820513
1	45.217391	26.30434783	20.91304348
2	40.6666667	87.75	17.58333333
3	25.52173913	26.30434783	78.56521739
4	43.08860759	55.29113924	49.56962025

The Predicted Info about which data goes to which cluster

```
In [17]: label
```

```
Out[17]:   CustomerID  Age  Annual Income (k$)  Spending Score (1-100)
```

0	19	15	39
1	21	15	81
2	20	16	6
3	23	16	77
4	31	17	40
...
195	35	120	79
196	45	126	28
197	32	126	74
198	32	137	18
199	30	137	83

200 rows x 4 columns

Plotting the Correlation using Heatmap

```
In [8]: plt.figure(figsize=(15,8))  
sns.heatmap(df.corr(), cmap = 'YlGnBu', annot = True , linewidths = 2, vmax = 1, vmin = -0.5 , square= True)
```

```
Out[8]: <AxesSubplot: ax[0,0]>
```


Removing the unneeded data in our data. Here that is Customer ID

```
In [9]: df = df.iloc[:,1:]  
df
```

```
Out[9]:   Gender  Age  Annual Income (k$)  Spending Score (1-100)
```

0	Male	19	15	39
1	Male	21	15	81
2	Female	20	16	6
3	Female	23	16	77
4	Female	31	17	40
...
195	Female	35	120	79
196	Female	45	126	28
197	Male	32	126	74
198	Male	32	137	18
199	Male	30	137	83

200 rows x 4 columns

Getting the Optimum Value of K using 'Elbow Method'

```
In [13]: from sklearn.cluster import KMeans  
temp = []
```

```
for i in range(1,11):  
    km = KMeans(n_clusters = i , init='k-means++')  
    km.fit(X)  
    temp.append(km.inertia_)
```

```
Out[14]: [ 32692308, 2630434783, 2091304348, 1758333333, 14308860759]
```


We get our K to be equal to '5'.

Now fitting our dataset to the Model using the best/optimal K value

```
In [15]: kmeans = KMeans(n_clusters = 5 , init='k-means++')  
label = kmeans.fit_predict(X)
```

```
Out[15]: [ 0, 1, 2, 3, 4 ]
```

These are the Centers of our 5 Clusters

```
In [16]: kmeans.cluster_centers_
```

```
Out[16]: array([[32
```