

Customer Segmentation

Importing the dependencies

In [2]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
import warnings
warnings.filterwarnings('ignore')
```

Data Collection

In [3]:

```
data=pd.read_csv("C:\\Users\\DELL\\Downloads\\Mall_Customers.csv")
```

In [4]:

```
data.head(10)
```

Out[4]:

	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
5	6	Female	22	17	76
6	7	Female	35	18	6
7	8	Female	23	18	94
8	9	Male	64	19	3
9	10	Female	30	19	72

In [5]:

```
data.tail(10)
```

Out[5]:

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
190	191	Female	34	103	23
191	192	Female	32	103	69
192	193	Male	33	113	8
193	194	Female	38	113	91
194	195	Female	47	120	16
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 [6]:

```
data.shape
```

Out[6]:

(200, 5)

In [7]:

```
data.describe()
```

Out[7]:

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

In [8]:

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
```

In [9]:

data.isnull().sum()

Out[9]:

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

In [10]:

data.drop(['CustomerID'],axis=1,inplace=True)

In [11]:

data.head()

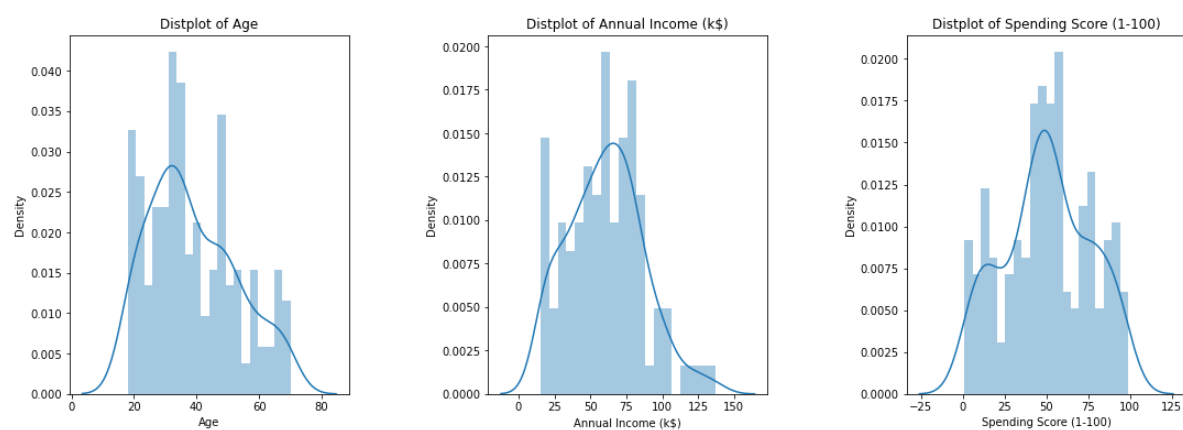
Out[11]:

	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

Exploratory Data Analysis

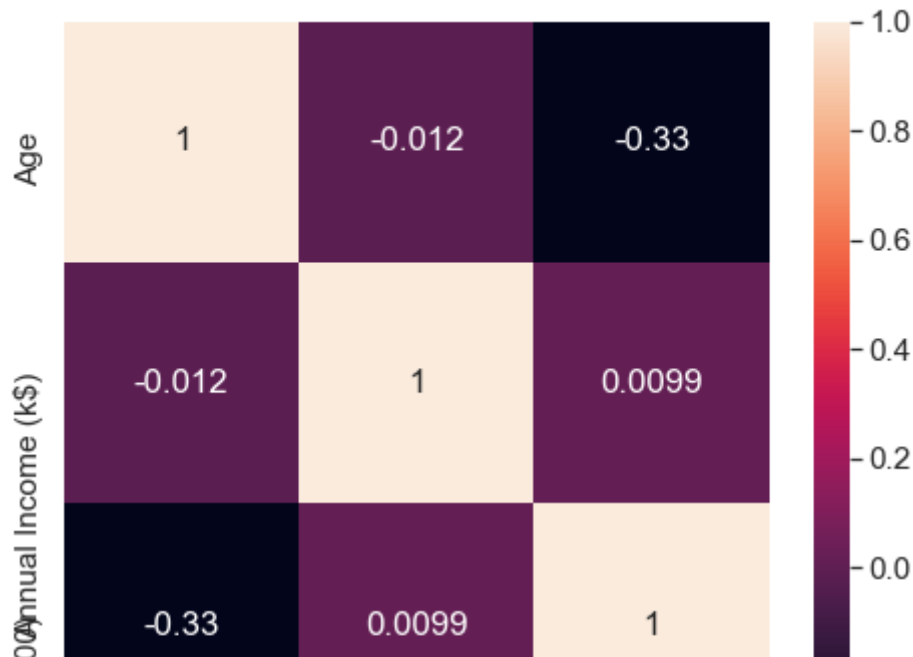
In [12]:

```
plt.figure(1, figsize=(25,6))
n=0
for x in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
    n+=1
    plt.subplot(1,4,n )
    plt.subplots_adjust(hspace=0.5,wspace=0.5)
    sns.distplot(data[x],bins=20)
    plt.title('Distplot of {}'.format(x))
plt.show()
```



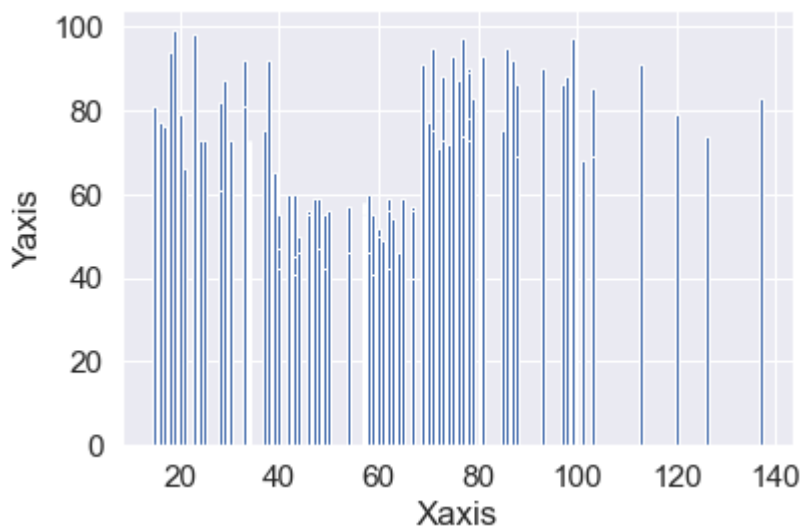
In [40]:

```
cor = data.corr()  
sns.set(font_scale=1.4)  
plt.figure(figsize=(7,8))  
sns.heatmap(cor, annot=True)  
plt.tight_layout()  
plt.show()
```



In [14]:

```
plt.bar(data['Annual Income (k$)'], data['Spending Score (1-100)'])  
plt.ylabel('Yaxis')  
plt.xlabel('Xaxis')  
plt.show()
```

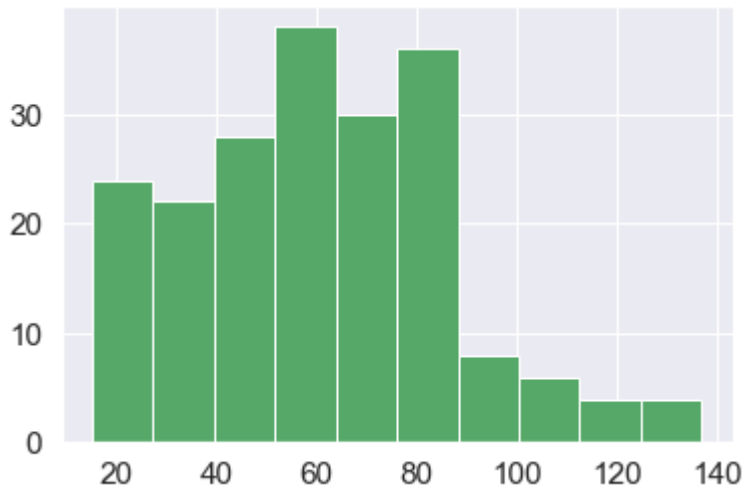


In [15]:

```
plt.hist(data['Annual Income (k$)'],color='g')
```

Out[15]:

```
(array([24., 22., 28., 38., 30., 36., 8., 6., 4., 4.]),  
 array([ 15. , 27.2, 39.4, 51.6, 63.8, 76. , 88.2, 100.4, 112.6,  
        124.8, 137. ]),  
<BarContainer object of 10 artists>)
```

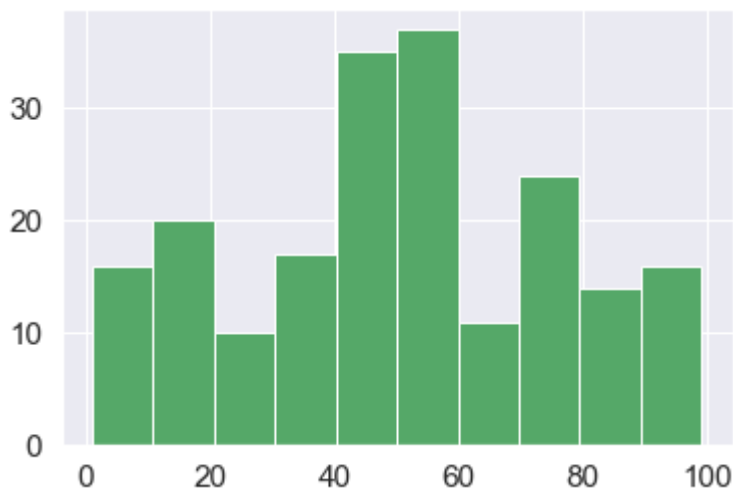


In [16]:

```
plt.hist(data['Spending Score (1-100)'],color='g')
```

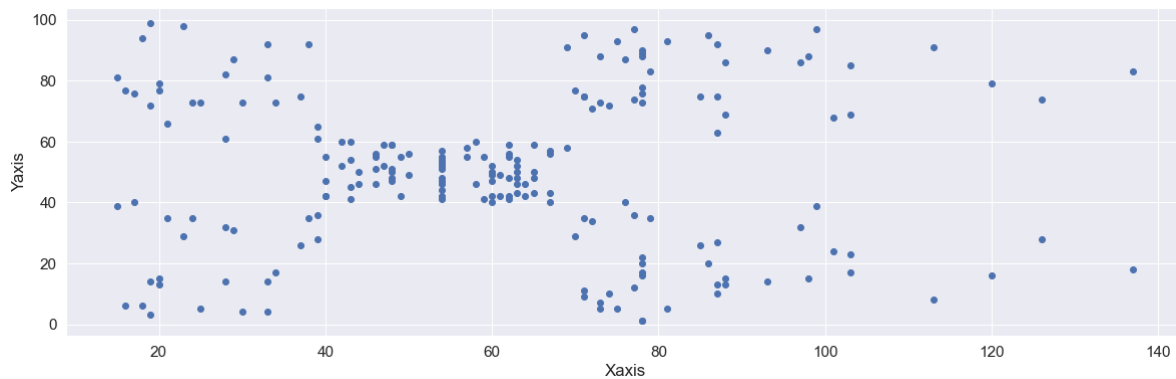
Out[16]:

```
(array([16., 20., 10., 17., 35., 37., 11., 24., 14., 16.]),  
 array([ 1. , 10.8, 20.6, 30.4, 40.2, 50. , 59.8, 69.6, 79.4, 89.2, 99. ]),  
<BarContainer object of 10 artists>)
```



In [25]:

```
plt.figure(figsize=(20,6))
plt.scatter(data['Annual Income (k$)'],data['Spending Score (1-100)'])
plt.ylabel('Yaxis')
plt.xlabel('Xaxis')
plt.show()
```

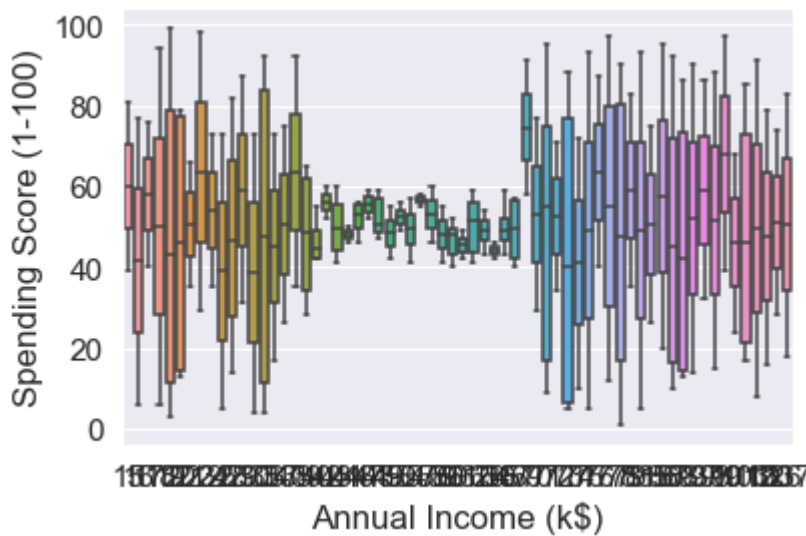


In [26]:

```
sns.boxplot(x='Annual Income (k$)',y='Spending Score (1-100)',data=data)
```

Out[26]:

```
<AxesSubplot:xlabel='Annual Income (k$)', ylabel='Spending Score (1-100)'\>
```



Choose Annual Income and Spendind score

In [27]:

```
X=data.iloc[:,[2,3]].values
```

In [28]:

```
print(X)
```

```
[[ 15  39]
 [ 15  81]
 [ 16   6]
 [ 16  77]
 [ 17  40]
 [ 17  76]
 [ 18   6]
 [ 18  94]
 [ 19   3]
 [ 19  72]
 [ 19  14]
 [ 19  99]
 [ 20  15]
 [ 20  77]
 [ 20  13]
 [ 20  79]
 [ 21  35]
 [ 21  66]
 [ 23  29]
 [ 23  88]]
```

Choosing the no of clusters

In [29]:

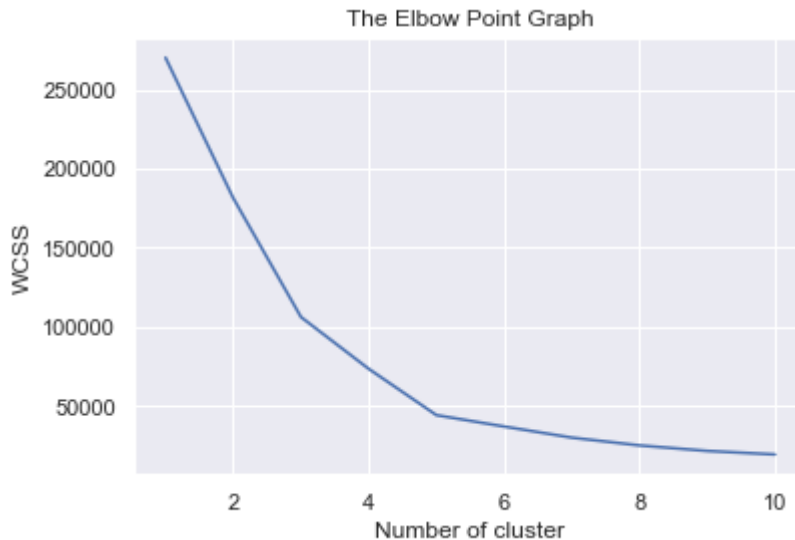
```
wcss=[]

for i in range(1,11):
    kmeans=KMeans(n_clusters=i,init='k-means++',random_state=42)
    kmeans.fit(X)

    wcss.append(kmeans.inertia_)
```


In [30]:

```
sns.set()
plt.plot(range(1,11),wcss)
plt.title("The Elbow Point Graph")
plt.xlabel("Number of cluster")
plt.ylabel('WCSS')
plt.show()
```



Optimum Number Of Cluster=5

Training the model

In [31]:

```
kmeans=KMeans(n_clusters=5,init='k-means++',random_state=0)
Y=kmeans.fit_predict(X)
print(Y)
```

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

Visualizing all clusters

In [39]:

```
plt.figure(figsize=(7,7))
plt.scatter(X[Y==0,0],X[Y==0,1],s=50,c='red',label='Cluster 1')
plt.scatter(X[Y==1,0],X[Y==1,1],s=50,c='blue',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='green',label='Cluster 5')

plt.scatter(kmeans.cluster_centers_[ :,0],kmeans.cluster_centers_[ :,1],s=70,c='black',label=
plt.title('Customer Cluster')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.legend()
plt.show()
```



Conclusion

Thus we have divided group of customers into 5 clusters based on their annual income and spending score.