

In [1]:

```

1 import numpy as np # Linear algebra
2 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
3 import matplotlib.pyplot as plt #Data Visualization
4 import seaborn as sns #Python Library for Vidualization

```

In [2]:

```

1 df = pd.read_csv('Mall_Customers.csv')
2 df.head(10)

```

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
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 [3]:

```

1 #total rows and colums in the dataset
2 df.shape

```

Out[3]:

(200, 5)

In [4]:

```

1 df.info() # there are no missing values as all the columns has 200 entries properly

```

```

<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 [5]:

```
1 #Missing values computation
2 df.isnull().sum()
```

Out[5]:

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

In [6]:

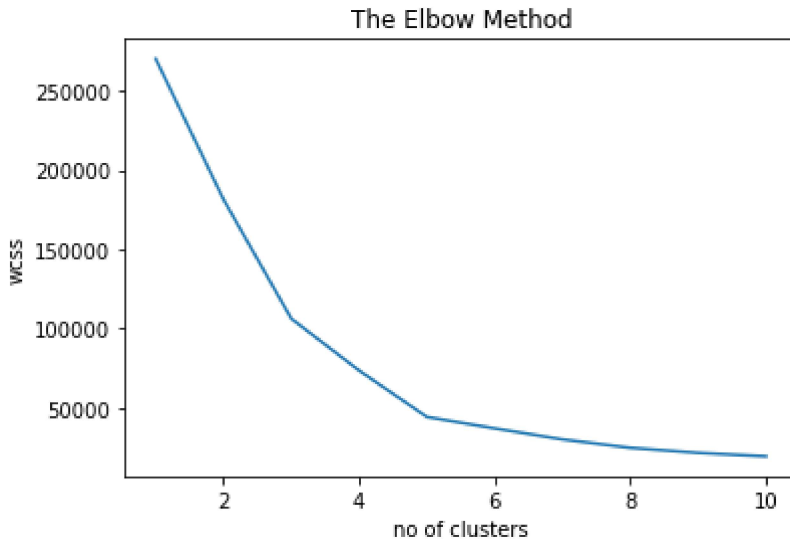
```
1 ### Feature sleection for the model
2 #Considering only 2 features (Annual income and Spending Score) and no Label available
3 X= df.iloc[:, [3,4]].values
```

In [7]:

```
1 #Building the Model
2 #KMeans Algorithm to decide the optimum cluster number , KMeans++ using Elbow Mmethod
3 #to figure out K for KMeans, I will use ELBOW Method on KMEANS++ Calculation
4 from sklearn.cluster import KMeans
5 wcss=[]
6
7 #we always assume the max number of cluster would be 10
8 #you can judge the number of clusters by doing averaging
9 ###Static code to get max no of clusters
10
11 for i in range(1,11):
12     kmeans = KMeans(n_clusters= i, init='k-means++', random_state=0)
13     kmeans.fit(X)
14     wcss.append(kmeans.inertia_)
15
16     #inertia_ is the formula used to segregate the data points into clusters
```

In [8]:

```
1 #Visualizing the ELBOW method to get the optimal value of K
2 plt.plot(range(1,11), wcss)
3 plt.title('The Elbow Method')
4 plt.xlabel('no of clusters')
5 plt.ylabel('wcss')
6 plt.show()
```



In [9]:

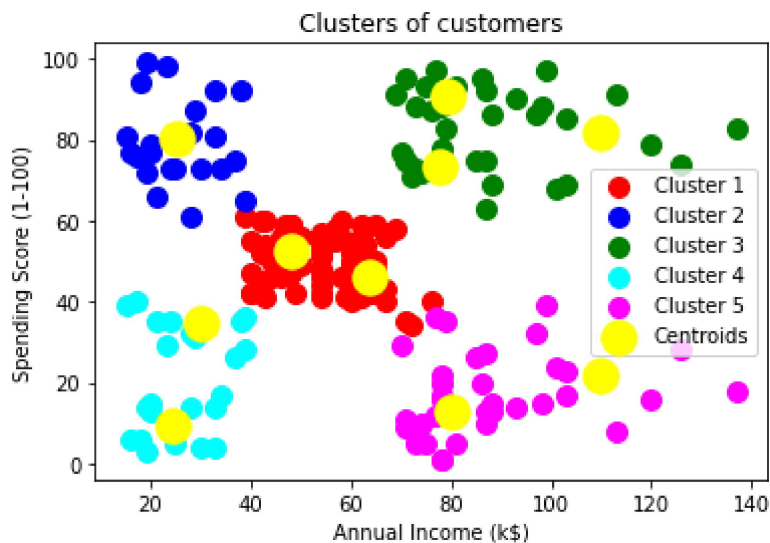
```
1 #If you zoom out this curve then you will see that last elbow comes at k=5
2 #no matter what range we select ex- (1,21) also i will see the same behaviour but if we
3 #that is why we usually prefer range (1,11)
4 ##Finally we got that k=5
5
6 #Model Build
7 kmeansmodel = KMeans(n_clusters= 5, init='k-means++', random_state=0)
8 y_kmeans= kmeansmodel.fit_predict(X)
9
10 #For unsupervised learning we use "fit_predict()" wherein for supervised learning we use
11 #y_kmeans is the final model . Now how and where we will deploy this model in production
12 #This use case is very common and it is used in BFS industry(credit card) and retail for
```

In [10]:

```

1 #Visualizing all the clusters
2
3 plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 100, c = 'red', label = 'Clus
4 plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 100, c = 'blue', label = 'Clu
5 plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 100, c = 'green', label = 'Cl
6 plt.scatter(X[y_kmeans == 3, 0], X[y_kmeans == 3, 1], s = 100, c = 'cyan', label = 'Clu
7 plt.scatter(X[y_kmeans == 4, 0], X[y_kmeans == 4, 1], s = 100, c = 'magenta', label = '
8 plt.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[1], s = 300, c =
9 plt.title('Clusters of customers')
10 plt.xlabel('Annual Income (k$)')
11 plt.ylabel('Spending Score (1-100)')
12 plt.legend()
13 plt.show()

```



In [11]:

```

1 ###Model Interpretation
2 #Cluster 1 (Red Color) -> earning high but spending Less
3 #cluster 2 (Blue Color) -> average in terms of earning and spending
4 #cluster 3 (Green Color) -> earning high and also spending high [TARGET SET]
5 #cluster 4 (cyan Color) -> earning Less but spending more
6 #Cluster 5 (magenta Color) -> Earning Less , spending Less
7
8
9 #####We can put Cluster 3 into some alerting system where email can be send to them or
10 #wherein others we can set like once in a week or once in a month

```