## Implemention of Clustering Algoritm in Python : KMeans

```
In [3]:  import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         from sklearn.cluster import KMeans
         from sklearn.metrics import silhouette_score
```

```
In [4]:  customer_data = pd.read_csv ("C:/Users/Acer/Downloads/archive/Mall_Customers.csv", delimiter = ",")
```

```
In [7]:  customer_data
```

Out[7]:

| | 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 × 5 columns

```
In [9]:  customer_data.shape
```

Out[9]:  (200, 5)

```
In [19]:  #chossing the relevent feature
          X = customer_data [['Annual Income (k$)', 'Spending Score (1-100)']]
```
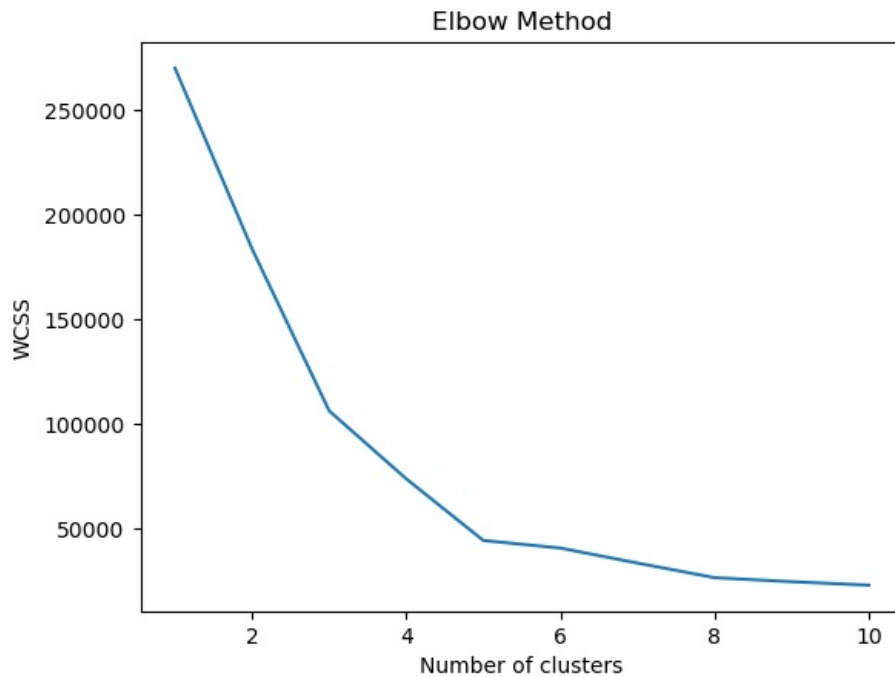
```
In [21]:  X
```

Out[21]:

| | Annual Income (k$) | Spending Score (1-100) |
|---|---|---|
| 0 | 15 | 39 |
| 1 | 15 | 81 |
| 2 | 16 | 6 |
| 3 | 16 | 77 |
| 4 | 17 | 40 |
| ... | ... | ... |
| 195 | 120 | 79 |
| 196 | 126 | 28 |
| 197 | 126 | 74 |
| 198 | 137 | 18 |
| 199 | 137 | 83 |

200 rows × 2 columns

```
In [ ]:  #Determine optimal number of cluster using the elbow method
         wcss = []# Within-cluster sum of square(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 [43]:  plt.plot(range(1,11),wcss)
          plt.title('Elbow Method')
```

```python
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()
```
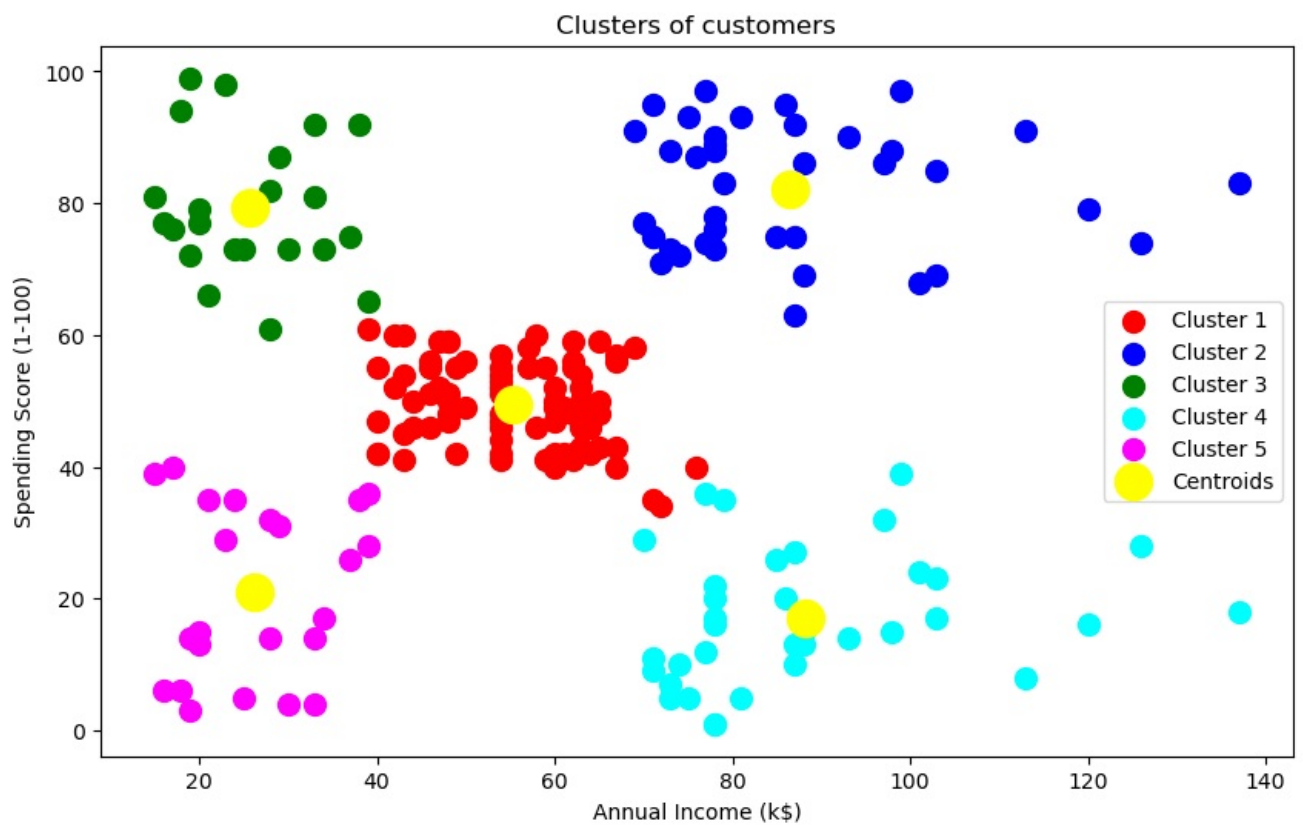


Elbow Method

In [ ]:
```python
#Apply K-Means CLustering
kmeans= KMeans(n_clusters=5, init='k-means++' , random_state= 42)
y_kmeans = kmeans.fit_predict(X)
```

In [51]:
```python
y_kmeans
```

Out[51]:
```
array([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, 0,
       4, 2, 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, 0, 0, 0, 0, 0, 0, 1, 3, 1, 0, 1, 3, 1, 3, 1,
       0, 1, 3, 1, 3, 1, 3, 1, 3, 1, 0, 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, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1,
       3, 1])
```

In [55]:
```python
#Visualize the Clusters
plt.figure(figsize=(10,6))
plt.scatter(X.iloc[y_kmeans == 0, 0], X.iloc[y_kmeans == 0, 1], s = 100, c = 'red', label = 'Cluster 1')
plt.scatter(X.iloc[y_kmeans == 1, 0], X.iloc[y_kmeans == 1, 1], s = 100, c = 'blue', label = 'Cluster 2')
plt.scatter(X.iloc[y_kmeans == 2, 0], X.iloc[y_kmeans == 2, 1], s = 100, c = 'green', label = 'Cluster 3')
plt.scatter(X.iloc[y_kmeans == 3, 0], X.iloc[y_kmeans == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4')
plt.scatter(X.iloc[y_kmeans == 4, 0], X.iloc[y_kmeans == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5')
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s = 300, c = 'yellow', label = 'Centr
plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```

Clusters of customers

In [57]:
```python
sil_score = silhouette_score(X, kmeans.labels_)
print('Silhouette Score: %.3f' % sil_score)
```

Silhouette Score: 0.554

In [ ]: