

K-Means Clustering

```
In [1]: # Importing the Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

```
In [2]: # Importing the dataset
dataset = pd.read_csv('Mall_Customers.csv')
dataset.head()
```

```
Out[2]:
```

	CustomerID	Genre	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 [3]: dataset.shape
```

```
Out[3]: (200, 5)
```

```
In [4]: dataset.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]: dataset.isnull().sum()
```

```
Out[5]: CustomerID            0
Genre              0
Age                0
Annual Income (k$)  0
Spending Score (1-100) 0
dtype: int64
```

```
In [6]: X = dataset.iloc[:, [3, 4]].values
```

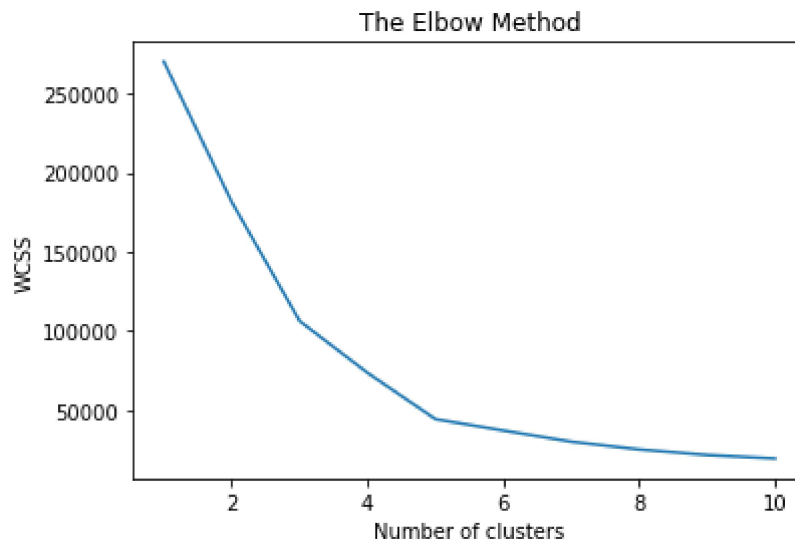
```
In [7]: # Kmeans Model
from sklearn.cluster import KMeans
```

```
In [8]: # Using the elbow method to find the optimal number of clusters
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 [9]: wcss
```

```
Out[9]: [269981.28000000014,
181363.59595959607,
106348.37306211119,
73679.78903948837,
44448.45544793369,
37265.86520484345,
30241.34361793659,
25336.94686147186,
21850.16528258562,
19634.554629349972]
```

```
In [10]: #Elbow Curve
plt.plot(range(1, 11), wcss)
plt.title('The Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()
```



```
In [11]: # fit the K-Means model on the data
kmeans = KMeans(n_clusters = 5, init = 'k-means++', random_state = 42)
```

```
In [12]: # predict
y_kmeans = kmeans.fit_predict(X)
```

In [13]: *# Visualising the clusters*

```
plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1],
            s = 100, c = 'red', label = 'Cluster 1')

plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1],
            s = 100, c = 'blue', label = 'Cluster 2')

plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1],
            s = 100, c = 'green', label = 'Cluster 3')

plt.scatter(X[y_kmeans == 3, 0], X[y_kmeans == 3, 1],
            s = 100, c = 'cyan', label = 'Cluster 4')

plt.scatter(X[y_kmeans == 4, 0], X[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 = 'Centroids')

plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```

