

Name:: SHIVAM VILAS GAIKWAD

GMAIL:: shivamgaikwad39@gmail.com

Task 2:: From the given ‘Iris’ dataset, predict the optimum number of clusters and represent it visually.

Dataset:: <https://drive.google.com/file/d/11lq7YvbWZbt8VXjfm06brx66b10YiwK/view>

In [1]:

```
# Importing the required libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from sklearn import datasets
```

Loading the Iris Dataset into the notebook

In [2]:

```
# Loading the iris dataset
iris = datasets.load_iris()
iris_df = pd.DataFrame(iris.data, columns = iris.feature_names)
iris_df.head() # The first 5 rows
```

Out[2]:

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

Finding the optimal number of clusters for K-Means and determining the value of K

In [3]:

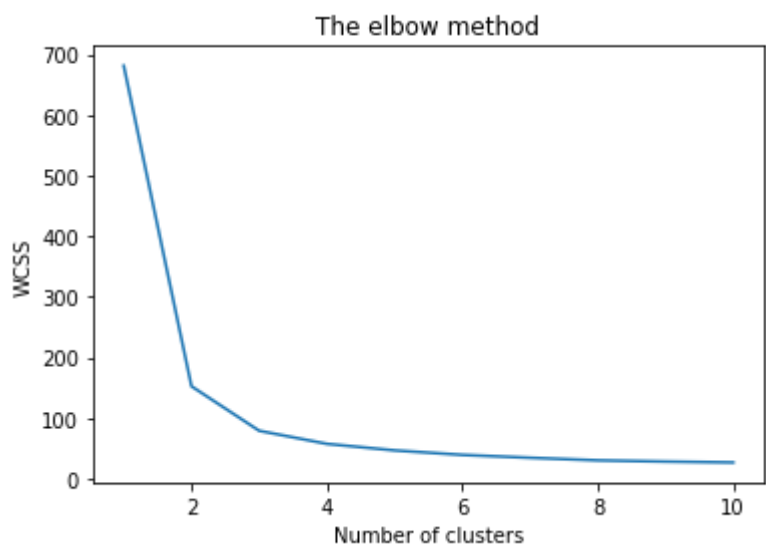
```
# Finding the optimum number of clusters for k-means classification
x = iris_df.iloc[:, [0, 1, 2, 3]].values

from sklearn.cluster import KMeans
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters = i, init = 'k-means++',
                    max_iter = 300, n_init = 10, random_state = 0)
    kmeans.fit(x)
    wcss.append(kmeans.inertia_)
```

Plotting the graph onto a line graph to observe the pattern

In [5]:

```
# Plotting the results onto a line graph,
# `allowing us to observe 'The elbow'
plt.plot(range(1, 11), wcss)
plt.title('The elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS') # Within cluster sum of squares
plt.show()
```



"The elbow method" got its name from the elbow pattern forming something like above. The optimal clusters are formed where the elbow occurs. This is when the WCSS(Within Cluster Sum of Squares) doesn't decrease with every iteration significantly.

Here we choose the number of clusters as '3'.

Creating K-Means Classifier

In [6]:

```
# Applying kmeans to the dataset
# Creating the kmeans classifier

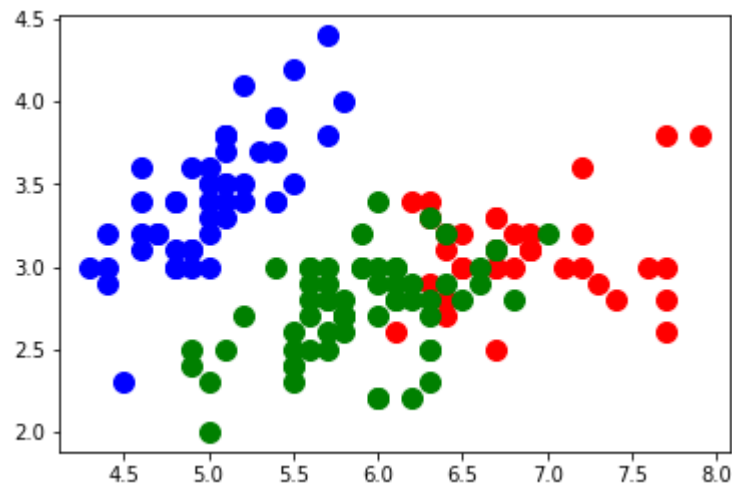
kmeans = KMeans(n_clusters = 3, init = 'k-means++',
                max_iter = 300, n_init = 10, random_state = 0)
y_kmeans = kmeans.fit_predict(x)
```

Visualizing the cluster data

In [7]:

```
# Visualising the clusters
# Preferably on the first two columns
plt.scatter(x[y_kmeans == 0, 0], x[y_kmeans == 0, 1],
            s = 100, c = 'red', label = 'Iris-setosa')
plt.scatter(x[y_kmeans == 1, 0], x[y_kmeans == 1, 1],
            s = 100, c = 'blue', label = 'Iris-versicolour')
plt.scatter(x[y_kmeans == 2, 0], x[y_kmeans == 2, 1],
            s = 100, c = 'green', label = 'Iris-virginica')
```

Out[7]: <matplotlib.collections.PathCollection at 0x5689be0>

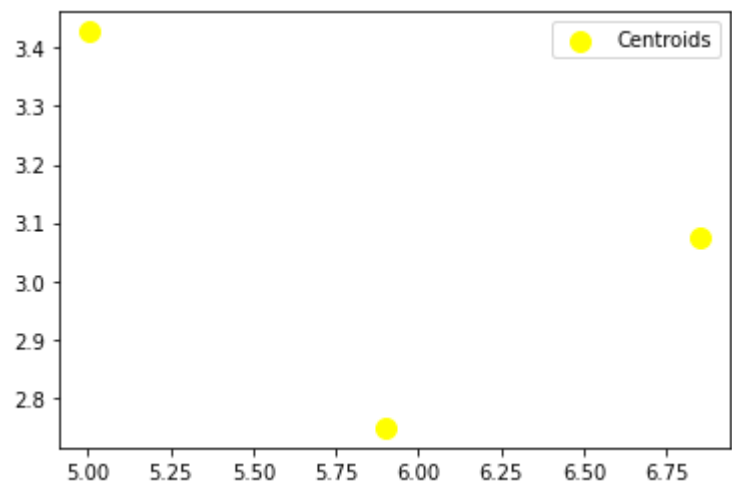


In [8]:

```
# Plotting the centroids of the clusters
plt.scatter(kmeans.cluster_centers[:, 0], kmeans.cluter_centers[:,1],
            s = 100, c = 'yellow', label = 'Centroids')

plt.legend()
```

Out[8]: <matplotlib.legend.Legend at 0x56a1f88>



Now Combining both the above graphs together

In [9]:

```
# Visualising the clusters
plt.scatter(x[y_kmeans == 0, 0], x[y_kmeans == 0, 1],
            s = 100, c = 'red', label = 'Iris-setosa')
plt.scatter(x[y_kmeans == 1, 0], x[y_kmeans == 1, 1],
            s = 100, c = 'blue', label = 'Iris-versicolour')
plt.scatter(x[y_kmeans == 2, 0], x[y_kmeans == 2, 1],
            s = 100, c = 'green', label = 'Iris-virginica')

# Plotting centroids of the clusters
plt.scatter(kmeans.cluster_centers[:, 0], kmeans.cluter_centers[:,1],
            s = 100, c = 'yellow', label = 'Centroids')

plt.legend()
```

Out[9]: <matplotlib.legend.Legend at 0x570c430>

