

Task 2: PREDICT THE OPTIMUM NUMBER OF CLUSTERS AND REPRESENT IT VISUALLY

dataset: <https://bit.ly/3kXTdox>

```
In [1]: # Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from sklearn import datasets
```

```
In [2]: #importing the dataset
iris=pd.read_csv(r"C:\Users\akank\Downloads\Iris.csv")
```

```
In [3]: iris = datasets.load_iris()
iris_df = pd.DataFrame(iris.data, columns = iris.feature_names)
iris_df.head() # See the first 5 rows
```

```
Out[3]:
```

	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

```
In [5]: x = iris_df.iloc[:, [0, 1, 2, 3]].values

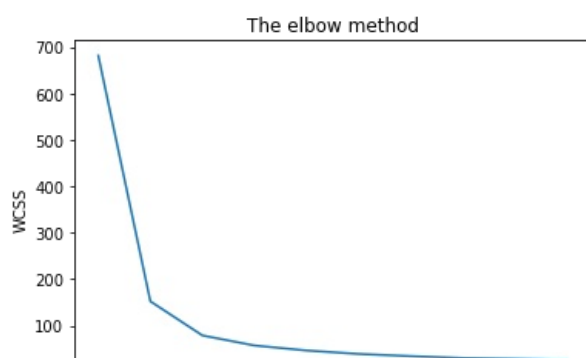
from sklearn.cluster import KMeans
wcss = []

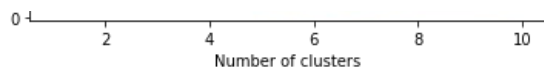
for i in range(1, 11):
    kmeans = KMeans(n_clusters = i, init = 'k-means++',
                    max_iter = 300, n_init = 10, random_state = 0)
    kmeans.fit(x)
    wcss.append(kmeans.inertia_)
print(wcss)
```

C:\Users\akank\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:881: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
warnings.warn(

[681.3705999999996, 152.34795176035797, 78.851441426146, 57.22847321428572, 46.47223015873018, 39.03998724608725, 34.299712121212146, 30.063110617452732, 28.27172172856384, 26.094324740540422]

```
In [6]: # 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()
```





In [7]: *#Number of clusters here can be seen to be 3 ,as after 3 there is no significant decrease in wccs*

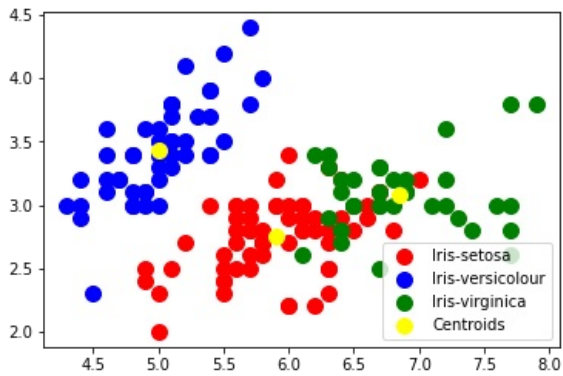
In [22]: *# Applying kmeans to the dataset / Creating the kmeans classifier*
 kmeans = KMeans(n_clusters = 3, init = 'k-means++',
 max_iter = 300, n_init = 10, random_state = 0)
 y_kmeans = kmeans.fit_predict(x)

In [23]: *# Visualising the clusters - 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')

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

 plt.legend()

Out[23]: <matplotlib.legend.Legend at 0x1b86474bf70>



In []:

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js