

The Sparks Foundation Task 2 - Prediction using Unsupervised ML

Author : ATANU DAS

From the given 'Iris' dataset, predict the optimum number of clusters and represent it visually

Iris Dataset

Iris dataset contains five columns such as Petal Length, Petal Width, Sepal Length, Sepal Width and Species Type. Iris is a flowering plant, the researchers have measured various features of the different iris flowers and recorded digitally

```
In [1]: #importing all the required libraries
import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
%matplotlib inline
```

```
In [2]: iris=datasets.load_iris()
df=pd.DataFrame(iris.data, columns=iris.feature_names)
df.head(150)
```

```
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

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
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
...
145	6.7	3.0	5.2	2.3
146	6.3	2.5	5.0	1.9
147	6.5	3.0	5.2	2.0
148	6.2	3.4	5.4	2.3
149	5.9	3.0	5.1	1.8

150 rows × 4 columns

K-Means Clustering

K-means clustering aims to partition data into k clusters in a way that data points in the same cluster are similar and data points in the different clusters are farther apart. The K-means algorithm identifies k number of centroids, and then allocates every data point to the nearest cluster, while keeping the centroids as small as possible.

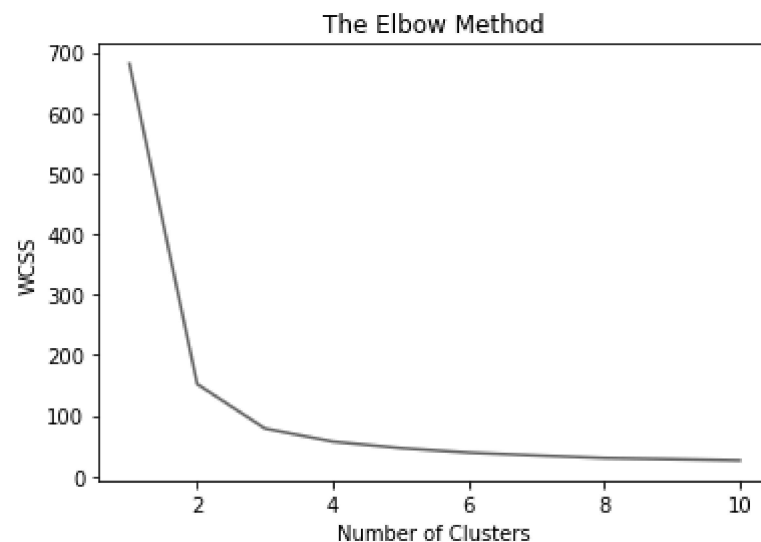
Inertia measures how well a dataset was clustered by K-Means. It is calculated by measuring the distance between each data point and its centroid, squaring this distance, and summing these squares across one cluster. A good model is one with low inertia AND a low number of clusters(K)

In [3]:

```
x=df.iloc[:,[0,1,2,3]].values
from sklearn.cluster import KMeans
list=[]

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)
    list.append(kmeans.inertia_)
```

```
mt.plot(range(1,11),list)
mt.title("The Elbow Method")
mt.xlabel("Number of Clusters")
mt.ylabel("WCSS")
mt.show()
```



```
In [4]: #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)
y_kmeans
```

[illegible]

```
In [5]: #visualising th clusters
mt.scatter(x[y_kmeans==0,0], x[y_kmeans==0,1], s=100, c='steelblue', label='Iris-setosa')
mt.scatter(x[y_kmeans==1,0], x[y_kmeans==1,1], s=100, c='coral', label='Iris-versicolor')
mt.scatter(x[y_kmeans==2,0], x[y_kmeans==2,1], s=100, c='olive', label='Iris-verginica')

#plotting centriods of clusters
```

```
mt.scatter(kmeans.cluster_centers_[ :,0], kmeans.cluster_centers_[ :,1], s=100, c='black', label='Centriods')  
mt.legend()
```

Out[5]: <matplotlib.legend.Legend at 0x1aa2b463ca0>

