

Graduate Rotational Internship Program (GRIP)

The Sparks Foundation

Task 2 - Prediction using Unsupervised ML

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

Programming Language - Python

IDE - Jupyter Notebook

AUTHOR - Kalyankar Sreeja

#Import all the required libraries import pandas as pd

In [19]:

STEP 1 - Import Libraries

```
import numpy as np
 import matplotlib.pyplot as plt
 import seaborn as sns
 %matplotlib inline
STEP 2 - Import Dataset
 #Displaying the whole dataset
```

Species

Iris-setosa

0.2

PetalWidthCm

0.100000

0.300000

Id SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm

1

Out[9]:

1.000000

38.250000

min

25%

5.1

df = pd.read csv('iris.csv')

```
4.9
                                                       3.0
                                                                         1.4
                                                                                          0.2
                                                                                                 Iris-setosa
              2
                    3
                                     4.7
                                                      3.2
                                                                                          0.2
                                                                        1.3
                                                                                                 Iris-setosa
              3
                                                      3.1
                                                                         1.5
                                                                                          0.2
                                     4.6
                                                                                                 Iris-setosa
              4
                    5
                                     5.0
                                                      3.6
                                                                        1.4
                                                                                          0.2
                                                                                                 Iris-setosa
            145
                 146
                                     6.7
                                                       3.0
                                                                                          2.3 Iris-virginica
                                                                        5.2
            146
                 147
                                                       2.5
                                                                        5.0
                                                                                             Iris-virginica
                                     6.3
            147 148
                                     6.5
                                                      3.0
                                                                        5.2
                                                                                          2.0 Iris-virginica
            148
                 149
                                     6.2
                                                       3.4
                                                                        5.4
                                                                                             Iris-virginica
            149 150
                                     5.9
                                                      3.0
                                                                        5.1
                                                                                          1.8 Iris-virginica
           150 rows × 6 columns
In [9]:
             #Viewing a series of numeric values
             df.describe()
```

3.5

```
count 150.000000
                         150.000000
                                          150.000000
                                                           150.000000
                                                                           150.000000
        75.500000
                           5.843333
                                            3.054000
                                                             3.758667
                                                                             1.198667
mean
        43.445368
                           0.828066
                                            0.433594
                                                             1.764420
                                                                             0.763161
  std
```

1.000000

1.600000

2.000000

2.800000

SepalLengthCm SepalWidthCm PetalLengthCm

50%	75.500000	5.800000	3.000000	4.350000	1.300000
75%	112.750000	6.400000	3.300000	5.100000	1.800000
max	150.000000	7.900000	4.400000	6.900000	2.500000
<pre>#Printing a concise summary of a DataFrame df.info()</pre>					
Range Data	s 'pandas.core Index: 150 entr columns (total Column	ries, 0 to 14	9		

4.300000

5.100000

SepalLengthCm 150 non-null

PetalLengthCm 150 non-null

PetalWidthCm 150 non-null

x = df.iloc[:, [0, 1, 2, 3]].values

from sklearn.cluster import KMeans

dtypes: float64(4), int64(1), object(1)

SepalWidthCm

Species

Method"

memory usage: 7.2+ KB

for i **in** range(1, 11):

150 non-null

150 non-null

wcss = [] #wcss stands for within-cluster sum of squares

kmeans = KMeans(n_clusters = i, init = 'k-means++', max iter = 300, n init = 10, n wcss.append(kmeans.inertia)

float64

float64

float64

STEP 3 - Find the optimum number of clusters using "Elbow

#Finding the optimum number of clusters for k-means classification and also showing h

```
#Plotting the results onto a line graph,
#Allows us to observe "The Elbow"
plt.figure(figsize = [11, 5])
plt.plot(range(1, 11), wcss)
plt.title('The Elbow Method')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS')
plt.grid(True)
plt.show()
                                         The Elbow Method
 250000
 200000
 150000
 100000
  50000
```

#Creating the kmeans classifier kmeans = KMeans(n clusters = 3, init = 'k-means++', max iter = 300, n init = 10, rando y kmeans = kmeans.fit predict(x) print(y kmeans)

plt.scatter($x[y_kmeans == 0, 0]$, $x[y_kmeans == 0, 1]$, s = 100, color = 'red', label = plt.scatter($x[y_kmeans == 1, 0]$, $x[y_kmeans == 1, 1]$, s = 100, color = 'blue', label = plt.scatter($x[y_kmeans == 2, 0]$, $x[y_kmeans == 2, 1]$, s = 100, color = 'green', label

The optimum clusters are the place the elbow happens. This is the point at which the WCSS doesn't reduce essentially with each iteration. From

Number of Clusters

10

plt.figure(figsize = [12, 7])

STEP 5 - Visualize Model

1 1]

this, we choose the number of clusters as "3".

STEP 4 - Train KMeans Model

```
#Plotting the centroids of the clusters
           plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s = 100, c =
           plt.legend()
Out[11]: <matplotlib.legend.Legend at 0x23eee81aeb8>
           8.0
                   Iris-setosa
                   Iris-versicolour
                   Iris-virginica
           7.5
                   Centroids
```

7.0 6.5 6.0 5.5 5.0 4.5 60 100 140 The plotted centroids of the cluster in the graph shows the predicted

optimum number of clusters from the iris dataset

CONCLUSION - The optimum number of clusters in the iris dataset are "3"