Data Science & Business Analytics Internship Task -2: Prediction using unsupervised ML By - Prasanta Mohanty problem statement: From the given "iris" Dataset, predict the optimum number of cluster and represent it visually. step1 - Importing nacessary libraries import numpy as np import pandas as pd import matplotlib.pyplot as plt importing the Dataset df = pd.read_csv("task_data") df.head(10) Id SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm **Species** Out[2]: 0 1 5.1 3.5 1.4 0.2 Iris-setosa **1** 2 4.9 3.0 1.4 0.2 Iris-setosa **2** 3 4.7 3.2 1.3 0.2 Iris-setosa 1.5 0.2 Iris-setosa 4.6 3.1 **4** 5 5.0 3.6 1.4 0.2 Iris-setosa 5 5.4 3.9 1.7 0.4 Iris-setosa **6** 7 4.6 3.4 1.4 0.3 Iris-setosa 5.0 3.4 1.5 0.2 Iris-setosa **8** 9 4.4 2.9 0.2 Iris-setosa 1.4 4.9 0.1 Iris-setosa **9** 10 3.1 1.5 check Any missing values df.isnull().sum() Ιd Out[3]: SepalLengthCm SepalWidthCm 0 PetalLengthCm 0 PetalWidthCm 0 Species dtype: int64 summary of Dataset In [4]: df.describe() Id SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm Out[4]: 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

<class 'pandas.core.frame.DataFrame'> RangeIndex: 150 entries, 0 to 149 Data columns (total 6 columns): Non-Null Count Dtype Column

memory usage: 7.2+ KB

metadata of the data

1.000000

38.250000

75.500000

75% 112.750000

max 150.000000

df.info()

Id

5 Species

plt.show()

700

ون 600 م

B 500

100 Suster

model = KMeans(

y_pred

plt.legend() plt.show()

4.5

4.0

3.5

3.0

2.8

In [10]:

visualizing Total

plt.figure(figsize=(10,6))

n_clusters=3, init='k-means++',

y_pred = model.fit_predict(x)

n_init=10, max_iter=300, random_state=0)

In [7]:

400

0

1

In [5]:

4.300000

5.100000

5.800000

6.400000

7.900000

150 non-null

150 non-null

SepalLengthCm 150 non-null

SepalWidthCm 150 non-null

PetalLengthCm 150 non-null

PetalWidthCm 150 non-null

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

wcss.append(model_1.inertia_)

plt.ylabel('witin_cluster_sum_of_squares')

The Elbow method

#plt.figure(figsize=(16,9)) plt.plot(range(1, 11), wcss) plt.title('The Elbow method') plt.xlabel('Number of clusters')

2.000000

2.800000

3.000000

3.300000

4.400000

float64

float64

float64

float64

object

1.000000

1.600000

4.350000

5.100000

6.900000

0.100000

0.300000

1.300000

1.800000

2.500000

min

25%

finding the number of optimum clusters using Elbow method In [6]: x = df.iloc[:,[1,2,3,4]].valuesfrom sklearn.cluster import KMeans wcss = [] #within_cluster_sum_of_squares **for** i **in** range(1, 11): $\#model_1 = KMeans(n_clusters=i, init = 'k-means++', max_iter= 300, n_init=10, random_state=0)$ $model_1 = KMeans(i)$ model_1.fit(x)

100 Number of clusters so from the Above graph we define that the optimum number is 3 Applying kmeans clustring algorithm

2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 0, 0, 0, 0, 0, 0, 0, 0,

visualization of the given Data

0, 2, 0, 0, 0, 0, 2, 0, 0, 0, 2, 0, 0, 0, 2, 0, 0, 2])

plt.figure(figsize=(10,6)) plt.scatter($x[y_pred == 0, 0], x[y_pred == 0, 1],$ s = 100, c='black', label="Iris-versicolor") $plt.scatter(x[y_pred == 1, 0], x[y_pred == 1, 1],$ s = 100, c='red', label="Iris-setosa") plt.scatter($x[y_pred == 2, 0], x[y_pred == 2, 1],$

> Iris-versicolor Iris-setosa Iris-virginica

7.0

s = 100, c='green', label="Iris-virginica")

0, 0, 0, 2, 2, 0, 0, 0, 0, 2, 0, 2, 0, 2, 0, 0, 2, 2, 0, 0, 0, 0,

2.5 2.0 4.5 5.5 6.5 plt.scatter(model.cluster_centers_[:,0], model.cluster_centers_[:,1], s = 100, c='blue', label="centroids") plt.legend() Out[9]: <matplotlib.legend.Legend at 0x1a74fab0700> centroids 3.3 3.2 3.1 3.0 2.9

5.00 5.25 5.50 5.75 6.00 6.25 6.50

 $plt.scatter(x[y_pred == 0, 0], x[y_pred == 0, 1],$

 $plt.scatter(x[y_pred == 1, 0], x[y_pred == 1, 1],$

 $plt.scatter(x[y_pred == 2, 0], x[y_pred == 2, 1],$

s = 100, c='blue', label="Iris-versicolor")

s = 100, c='green', label="Iris-virginica")

plt.scatter(model.cluster_centers_[:,0], model.cluster_centers_[:,1],

s = 100, c='red', label="Iris-setosa")

s = 100, c='yellow', label="centroids") plt.legend() plt.show() Iris-versicolor Iris-setosa Iris-virginica centroids 4.0 3.5

Thank You

2.5

2.0