GRIP @ The Sparks Foundation

Task 2: Prediction using Unsupervised Machine Learning In this K-means clustering task I tried to predict the optimum number of clusters and represent it visually from the given Iris dataset

Author: VENEPALLY THARUN

In [2]:

Out[2]:

In [6]:

1.Importing all the required libraries from sklearn import datasets

import matplotlib.pyplot as plt import pandas as pd

Technical Stack: Scikit Learn, Numpy Array, Scipy, Pandas, Matplotlib

import numpy as np from sklearn.cluster import KMeans import matplotlib.patches as mpatches import sklearn.metrics as sm

from mpl_toolkits.mplot3d import Axes3D from scipy.cluster.hierarchy import linkage,dendrogram from sklearn.cluster import DBSCAN

from sklearn.decomposition import PCA

2.Loading the Dataset iris = datasets.load_iris()

iris_df = pd.DataFrame(iris.data, columns = iris.feature_names) iris_df.head(10) #displaying first 10 rows sepal length (cm) sepal width (cm) petal length (cm) petal width (cm)

0 0.2 5.1 3.5 1.4 0.2 1 4.9 3.0 1.4 2 0.2 4.7 3.2 1.3 3 0.2 4.6 3.1 1.5 4 5.0 3.6 1.4 0.2 5 5.4 3.9 0.4 1.7 6 4.6 3.4 1.4 0.3

7 5.0 3.4 1.5 0.2 8 4.4 2.9 1.4 0.2 4.9 3.1 1.5 0.1

print(iris.target_names) In [3]: ['setosa' 'versicolor' 'virginica'] print(iris.target)

In [4]: 2 2] In [5]: x = iris.datay = iris.target 3. Visualizing the input and its Hierarchy

plt.title("Iris Clustering K Means=3", fontsize=14)

plt.title('Iris Hierarchical Clustering Dendrogram')

#Plotting fig = plt.figure(1, figsize=(7,5)) ax = Axes3D(fig, rect=[0, 0, 0.95, 1], elev=48, azim=134)ax.scatter(x[:, 3], x[:, 0], x[:, 2], edgecolor="k", s=50)ax.set_xlabel("Petal width")

ax.set_ylabel("Sepal length") ax.set_zlabel("Petal length")

plt.show() #Hierachy Clustering hier=linkage(x, "ward") $max_d=7.08$ plt.figure(figsize=(15,8))

truncate_mode='lastp',

plt.xlabel('Species') plt.ylabel('distance')

dendrogram(hier,

p=50,

2.5

2.0

Petal Width 10

leaf_rotation=90., leaf_font_size=8., plt.axhline(y=max_d, c='k') plt.show() Iris Clustering K Means=3

7.5

0.0 8.0

5.0

Iris Hierarchical Clustering Dendrogram 30 25 20 15 10 Species 4.Data Preprocessing x = pd.DataFrame(iris.data, columns=['Sepal Length', 'Sepal Width', 'Petal Length', 'Petal Width']) y = pd.DataFrame(iris.target, columns=['Target']) In [8]: x.head()

4.6 4 5.0

0

0

5.Model training

5.1

4.9

4.7

Out[8]:

0

1

2

3

4

In [9]: y.head() **Target** Out[9]: 0 1 0

Sepal Length Sepal Width Petal Length Petal Width

3.5

3.0

3.2

3.1

3.6

1.4

1.4

1.3

1.5

1.4

0.2

0.2

0.2

0.2

0.2

Out[10]: KMeans(n_clusters=3) In [11]: print(iris_k_mean_model.labels_)

[[5.006

[6.85

In [13]:

3.428

In [10]: iris_k_mean_model = KMeans(n_clusters=3)

iris_k_mean_model.fit(x)

2 1] In [12]: print(iris_k_mean_model.cluster_centers_)

> 6. Visualising the model cluster plt.figure(figsize=(14,6)) colors = np.array(['red', 'green', 'blue']) predictedY = np.choose(iris_k_mean_model.labels_, [1, 0, 2]).astype(np.int64)

plt.title('Before classification')

plt.title("Model's classification")

1.462

3.07368421 5.74210526 2.07105263]]

[5.9016129 2.7483871 4.39354839 1.43387097]

0.246

plt.scatter(x['Petal Length'], x['Petal Width'], c=colors[y['Target']])

plt.scatter(x['Petal Length'], x['Petal Width'], c=colors[predictedY])

red_patch = mpatches.Patch(color='red', label='The red data')

plt.legend(handles=[red_patch, green_patch, blue_patch])

green_patch = mpatches.Patch(color='green', label='The green data') blue_patch = mpatches.Patch(color='blue', label='The blue data') plt.legend(handles=[red_patch, green_patch, blue_patch])

plt.subplot(1, 2, 2)

1.0

0.5

THANK YOU

plt.subplot(1, 2, 1)

Out[13]: <matplotlib.legend.Legend at 0x216e7c99460> Model's classification Before classification The red data The red data 2.5 The green data The green data The blue data The blue data 2.0 2.0 1.5 1.5

1.0

0.5

7. Calculating Accuracy and Confusion Matrx In [14]: sm.accuracy_score(predictedY, y['Target']) Out[14]: 0.24 sm.confusion_matrix(predictedY, y['Target']) In [15]: Out[15]: array([[0, 48, 14], [50, 0, 0],

> In a confusion matrix, the predicted class labels (0, 1, 2) are written along the top (column names). The true class labels (Iris-setosa, etc.) are written along the right side. Each cell in the matrix is a count of how many instances of a true class where classified as each of the predicted classes.

[0, 2, 36]], dtype=int64)