	Prediction using Unsupervised ML Objective: Predict the optimum number of clusters and represent visually. Method Used: K-Means Clustering Dataset: https://bit.ly/3kXTdox By: MAYURI ARUN PATHAK
	Iris Versicolor Iris Setosa Iris Virginica
	Steps that we follow: 1) Reading the dataset 2) Checking and cleaning the data 3) EDA 4) Preparing Data for Modelling 5) Clustering Using K-means 6) Visualize the cluster.
[2]:	<pre>import os os.chdir("H:\\Data Science\\Internship\\Spark") Importing Libraries: from sklearn.cluster import KMeans from sklearn import datasets import pandas as pd import seaborn as sns import numpy as np import matplotlib.pyplot as plt Mmatplotlib inline # Supress warnings import warnings warnings.filterwarnings("ignore") #For K-Mean Cluster import KMeans from sklearn.cluster import KMeans from sklearn.cluster import KMeans from sklearn.cluster import KMeans from sklearn.cluster import silhouette_score</pre>
[4]: [5]:	Understanding Dataset # Import Data data_iris = pd.read_csv('iris.csv') print("data load successfully") data load successfully # View some sample records data_iris.head() to SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm Species to SepalLengthCm SepalWidthCm PetalLengthCm Species to SepalLengthCm SepalWidthCm PetalLengthCm Species to SepalLengthCm SepalWidthCm PetalWidthCm Species to SepalLengthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalLengthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalWidthCm SepalLengthCm SepalWidthCm SepalLengthCm SepalWidthCm SepalW
[6]: [6]: [7]:	<pre>#shape of the Dataframe data_iris.shape (150, 6) #Checking data type of variables data_iris.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 150 entries, 0 to 149 Data columns (total 6 columns): # Column Non-Null Count Dtype</class></pre>
[8]: [8]:	0 Id 150 non-null int64 1 SepalLengthCm 150 non-null float64 2 SepalWidthCm 150 non-null float64 3 PetalLengthCm 150 non-null float64 4 PetalWidthCm 150 non-null object dtypes: float64(4), int64(1), object(1) memory usage: 7.2+ KB ##Column wise Null value data_iris.isnull().sum() Id 0 SepalLengthCm 0 SepalWidthCm 0 PetalLengthCm 0 SepalWidthCm 0 Species 0 dtype: int64 ##ROW wise Null value ##ROW wise Null value ##ROW wise Null value ##ROW wise Null value
[9]: 10]:	data_iris.isna().sum(axis=1) 0
10]: 11]:	0 False 1 False 2 False 3 False 4 False 145 False 146 False 147 False 148 False 149 False 149 False Length: 150, dtype: bool #Counting Species data_iris['Species'].value_counts() Iris-setosa 50
12]:	<pre>Tris-versicolor 50 Iris-virginica 50 Name: Species, dtype: int64 • All species are in same quantity EDA # Distribution plot of all the continuous variables in the dataset plt.figure(figsize = (15,8)) plt.subplot(2,2,1) sns.distplot(data_iris["SepalLengthCm"]) plt.subplot(2,2,2)</pre>
12]:	sns.distplot(data_iris["PetalLengthCm"]) plt.subplot(2,2,3) sns.distplot(data_iris["PetalLengthCm"]) plt.subplot(2,2,4) sns.distplot(data_iris["PetalWidthCm"]) AxesSubplot:xlabel='PetalWidthCm", ylabel='Density'> 14 12 10 20 20 25 20 25 30 35 40 45 50 00 05 66 05 05 00 05 06 05 06 05 06 05 06 06 05 06 06 06 06 06 06 07 06 06 07 06 06 07 07 06 07 07 07 06 07 07 07 06 07 07 07 07 07 07 07 07 07 07 07 07 07
13]:	Observation • Variables PetalLengthCm and PetalWidthCm have chance of internal grouping.Hence, For these variables we can do our further analysis. #Pairplot of iris dataset sns.pairplot(data_iris , hue = "Species" , markers = "x") plt.show()
	150 100 120 120 120 120 120 120 12
	After graphing the features in a pair plot, it is clear that the relationship beatween pairs of features of iris-setosa (in blue) is distinctly different from those of othe two species. There is some overlap in the pairwise relationships of the othe two species, iris-versicolor(orange) and iris-virginica(green).
14]:	######################################
15]:	<pre># Checking Outliers plt.figure(figsize =(16,8)) plt.subplot(3,3,1) sns.boxplot(y = data_iris["SepalLengthCm"]) plt.subplot(3,3,2) sns.boxplot(y = data_iris["PetalLengthCm"]) plt.subplot(3,3,3) sns.boxplot(y = data_iris["PetalLengthCm"]) plt.subplot(3,3,4) sns.boxplot(y = data_iris["PetalWidthCm"]) plt.subplot(3,3,4) sns.boxplot(y = data_iris["PetalWidthCm"]) plt.show()</pre>
16].	• we can see that sepal width has some outliers.
16]:	plt.subplot(2,2,i[0]+1) sns.barplot(x=data_iris['Species'],y= data_iris[i[1]]) #################################
17]: 17]:	## Preparing data for Modelling data_iris.describe()
18]: 19]:	std 43.445388 0.828066 0.433594 1.764420 0.763161 min 1.000000 4.300000 2.000000 1.000000 0.300000 25% 38.250000 5.100000 2.800000 1.600000 0.300000 50% 75.500000 5.800000 3.00000 1.300000 75% 112.750000 6.400000 3.300000 5.100000 1.800000 max 150.000000 7.900000 4.400000 6.900000 2.500000 MOdelling import sklearn from sklearn.preprocessing import Standardscaler # Rescalling data for numeric columns scaler = StandardScaler() column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_iris.drop(['Id', 'Species'], axis = 1)) **The column_scaled = scaler.fit_transform(data_i
	Clustering - [Using K-mean Clustering] **Rchecking whether clustering is possible or not using Hopkins from sklearn.neighbors import NearestNeighbors from random import sample from numpy.random import uniform import unmpy.random import uniform import unmpy as np from math import isnan def hopkins(X): d = X. shape(1) n = len(X)
20]: 19]:	<pre>if isnan(H): print(ujd,wjd) H = 0 return H data_check_clustering = data_iris.drop(["Id", "Species"], axis = 1) hopkins(data_check_clustering) 0.8506584915489421 • we knew it Hopkins Statistics is a way of measuring the clustering tendency of a # Finding optimal number of clusters using elbow ssd = [] range_n = [2,3,4,5,6,7,8,9] for i in range_n: kmeans = KMeans(n_clusters=i , max_iter = 500 , random_state=50) kmeans = KMeans(n_clusters=i , max_iter = 500 , random_state=50) kmeans.fit(data) ssd.append([i, kmeans.inertia_])</pre>
	ssd_dataframe = pd.Dataframe(ssd) plt.plot(ssd_dataframe[0], ssd_dataframe[1], 'go' , color = 'purple') plt.title("The Elbow Method") plt.ylabel("Wumber of clusters") plt.grid() plt.show() The Elbow Method The Elbow Method 40 40 40 40 40 40 40 40 40 4
20]:	#Finding optimal clusters using sheloitte score #Silhouette score measures that how a data point is similar to it's own cluster as compared to other cluster range_n = [2,3,4,5,6,7] ss = [] for i in range_n:
32]: 33]:	for n cluster 5 Silhouette score is 0.4885175508886279 for n cluster 6 Silhouette score is 0.3695162748599971 for n cluster 7 Silhouette score is 0.3695162748599971 • So after checking the Elbow Curve & Silhouette Score we can choose our final k as 3.Because we have seen that in as sum of squared distance curve the 1st knee like bend has come at the level3, and in silhouette score, the score = 3 is also considerable for clustering. Though in the silhouette score for 2 is quite ipressive from score of 3, but intuitively we can say from the point of our dataset that chossing the point of cluster at may not give us a optimal result that we can expect for. So for better clustering and better output and better output of our analysis. We will choose number of cluster as 3. Hierarchical Clustering # Reconfirming the number of cluster from scipy.cluster.hierarchy import linkage
33]:	<pre>from scipy.cluster.hierarchy import linkage from scipy.cluster.hierarchy import dendrogram from scipy.cluster.hierarchy import cut_tree mergin = linkage(column_scaled,method="complete",metric="euclidean") dendrogram(mergin) plt.show()</pre>
24]:	• We can clearly see, This dendrogram clearly indicate k=3. If we cut the tree between 4-5 and draw a straight line we can get 3 distinct set of clusters. We will move ahead for k-means cluster ##Assuming k = 3 kmeans = KMeans(n_clusters = 3,max_iter= 500 , random_state = 50) Cluster_Id= kmeans.fit_predict(data) Cluster_Id array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
39]:	#Visualizing Cluster using petal length and petal width column plt.scatter(data[Cluster_Id == 0 , 2], data[Cluster_Id == 0 , 3] ,s = 100 ,marker = '*',c = "red", label = "Iris-Versiscolor") plt.scatter(data[Cluster_Id == 1 , 2], data[Cluster_Id == 1 , 3] ,s = 100 ,marker = '*',c = "indigo", label = "Iris-Setosa") plt.scatter(data[Cluster_Id == 2 , 2], data[Cluster_Id == 2 , 3] ,s = 100 ,marker = '*',c = "green", label = "Iris-Virginica") #Centroids of the clusters plt.scatter(kmeans.cluster_centers_[:,2], kmeans.cluster_centers_[:,3], s = 50 ,marker = '*' ,c = "blue", label = "Centroids") plt.legend() plt.show() 25
	Conclusion: • Therefore, the optimum number of distinct clusters are 3.