[3]:	TASK 2 : Prediction using unsupervised ML Problem Statement : Predict the optimum number of clusters and represent it visualy NAME: Divyashree K Importing Libraries import numpy as np import pandas as pd import matplotlib.pyplot as plt
4]: 5]:	<pre>import seaborn as sns import warnings warnings.filterwarnings("ignore") import io %cd "E:\Internship\Spark Foundation" E:\Internship\Spark Foundation Read data data = pd.read_csv("iris.csv")</pre>
5]:	data
6]:	148
3]:	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 float64 5 Species 150 non-null object dtypes: float64(4), int64(1), object(1) memory usage: 7.2+ KB data.columns Index(['Id', 'SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm', 'Species'],
)]:)]:	# droping id columns data= data.drop(['Id'], axis=1) SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm Species 0 5.1 3.5 1.4 0.2 Iris-setosa 1 4.9 3.0 1.4 0.2 Iris-setosa 2 4.7 3.2 1.3 0.2 Iris-setosa
	3 4.6 3.1 1.5 0.2 Iris-setosa 4 5.0 3.6 1.4 0.2 Iris-setosa
)]:	<pre>'Species'], dtype='object') # check for null values data.isnull().sum().sort_values(ascending = False) Species 0</pre>
2]:	PetallwidthCm 0 PetallengthCm 0 SepallwidthCm 0 SepallengthCm 0 dtype: int64 # check for duplicate values print("number of duplicate rows: ",data.duplicated().sum()) number of duplicate rows: 3 # drop duplicate rows data.drop_duplicates(inplace=True)
3]: 4]: 4]:	data.shape[0] #shape[0] shows number of rows 147 data.shape[1] # shape[1] shows number columns 5
	sis - Story () ((at a [4])) 11 slow()) 45
3]:	<pre>q1,q3 = np.percentile(data['SepalWidthCm'],[25,75]) iqr = q3-q1 lower_fence = q1-(1.5*iqr) upper_fence = q3+(1.5*iqr)</pre>
7]:	data['SepalWidthCm'] = data['SepalWidthCm'].apply(lambda x: upper_fence if x>upper_fence else lower_fence if x <lower_fence <axessubplot:xlabel="SepalWidthCm" else="" sns.boxplot(data['sepalwidthcm'])="" x)=""> </lower_fence>
3]:	Understanding the data # Target class data.Species.value_counts() sns.countplot(data.Species) <axessubplot:xlabel='species', ylabel="count"> 50-</axessubplot:xlabel='species',>
	data.describe()
)))	
)]:)]: L]:	<pre>max 7.90000 4.05000 6.90000 2.500000 data.Species.unique() array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object) #Distribution of features by Species for i in data.columns[: -1]: sns.kdeplot(data = data.loc[data.Species == 'Iris-setosa'][i],label = 'Iris-setosa',shade = True) sns.kdeplot(data = data.loc[data.Species == 'Iris-versicolor'][i], label = 'Iris-versicolor', shade = True) sns.kdeplot(data = data.loc[data.Species == 'Iris-virginica'][i], label = 'Iris-virginica', shade = True) plt.title(i)</pre>
	Sepalwidthon Sepalwidthon Sepalwidthon Fitall engthon Petallengthon Petalwiddhon Fitallengthon Fitalleng
2]: 2]:	SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm SepalLengthCm 1.000000 -0.110155 0.871305 0.817058 SepalWidthCm -0.110155 1.000000 -0.420140 -0.355139 PetalLengthCm 0.871305 -0.420140 1.000000 0.961883 PetalWidthCm 0.817058 -0.355139 0.961883 1.000000 plt.figure(figsize=(10,5)) sns.heatmap(abs(data.corr()), cmap = 'GnBu', annot=True) annot=True)
3]:	SepalLengthCm 1 0.11 0.87 0.82 -0.9 SepalWidthCm 0.11 1 0.42 0.36 -0.7 PetalLengthCm 0.87 0.42 1 0.96 -0.5 PetalWidthCm 0.82 0.36 0.96 1 -0.3 SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm
4]: 5]:	<pre>for i in range(1,10): kmeans = KMeans(n_jobs = -1,n_clusters=i,init = 'k-means++') kmeans.fit(data.iloc[: ,[0,1,2,3]]) SSE.append(kmeans.inertia_) df = pd.DataFrame({'Cluster':range(1,10), 'SSE':SSE}) plt.figure(figsize=(12,6))</pre>
6]:	plt.plot(df['Cluster'],df['SSE'],marker = '0') plt.xlabel('Number of clusters') plt.ylabel('Intertia') plt.title("ELBOW MEATHOD TO DETERMINE OPTIONAL VALUE OF 'K' \n ") Text(0.5, 1.0, "ELBOW MEATHOD TO DETERMINE OPTIONAL VALUE OF 'K' \n ") ELBOW MEATHOD TO DETERMINE OPTIONAL VALUE OF 'K'
	500 - 400 - 200 - 100 -
7]: 7]:	[5.01041667, 3.41979167, 1.4625 , 0.25], [6.85 , 3.07368421, 5.74210526, 2.07105263]])
3]: 3]: 9]:	Array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
) *	0 5.1 3.5 1.4 0.2 Iris-setosa 1 1 4.9 3.0 1.4 0.2 Iris-setosa 1 2 4.7 3.2 1.3 0.2 Iris-setosa 1 3 4.6 3.1 1.5 0.2 Iris-setosa 1 4 5.0 3.6 1.4 0.2 Iris-setosa 1 145 6.7 3.0 5.2 2.3 Iris-virginica 2 146 6.3 2.5 5.0 1.9 Iris-virginica 0
)]:	147 6.5 3.0 5.2 2.0 Iris-virginica 2 148 6.2 3.4 5.4 2.3 Iris-virginica 2 149 5.9 3.0 5.1 1.8 Iris-virginica 0 147 rows × 6 columns display(data['Cluster'].value_counts(), data['Species'].value_counts()) 0 61 1 48 2 38
l]: l]:	Name: Cluster, dtype: int64 Iris-versicolor 50 Iris-virginica 49 Iris-setosa 48 Name: Species, dtype: int64 plt.figure(figsize = (10,5)) plt.scatter(data['SepalLengthCm'], data['SepalWidthCm'], c=data.Cluster) plt.title('predicted cluster \n') plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1],s = 200,c='red', label ='centroid') plt.show <function block="None)" matplotlib.pyplot.show(close="None,"></function>
	4.00 - 3.75 - 3.50 - 3.25 - 3.00 - 2.75 -
	2.50 - 2.25 - 2.00 - 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0 data.loc[data['Species']== 'Iris-setosa']['Cluster'].value_counts()
2]:	<pre>data.loc[data['Species']== 'Iris-virginica']['Cluster'].value_counts() 2 36 0 13 Name: Cluster, dtype: int64 data.loc[data['Species'] == 'Iris-versicolor']['Cluster'].value_counts() 0 48 2 2 Name: Cluster, dtype: int64 data['Species_encoded']=data['Species'].apply(lambda x: 1 if x=='Iris-setosa' else 2 if x== 'Iris-virginica' else 0)</pre>
2]: 3]: 3]: 4]:	data
2]: 3]: 4]: 4]:	1 4.9 3.0 1.4 0.2 Iris-setosa 1 1 2 4.7 3.2 1.3 0.2 Iris-setosa 1 1 3 4.6 3.1 1.5 0.2 Iris-setosa 1 1 4 5.0 3.6 1.4 0.2 Iris-setosa 1 1
2]: 3]: 4]: 5]:	2 4.7 3.2 1.3 0.2 Iris-setosa 1 1 3 4.6 3.1 1.5 0.2 Iris-setosa 1 1 4 5.0 3.6 1.4 0.2 Iris-setosa 1 1
2]: 2]: 3]: 3]: 4]: 5]: 6]:	2 4.7 3.2 1.3 0.2 Iris-setosa 1 1 3 4.6 3.1 1.5 0.2 Iris-setosa 1 1 4 5.0 3.6 1.4 0.2 Iris-setosa 1 1
2]: 3]: 4]: 5]:	2 4.7 32 1.3 0.2 fris-settosa 1 1 3 4.6 3.1 1.5 0.2 fris-settosa 1 1 4 5.0 3.6 1.4 0.2 fris-settosa 1 1 1