	Data Science and Business Analytics Internship Prediction using Unsupervised ML
	TASK-2: From the given 'Iris' dataset, predict the optimum number of clusters and represent it visually. Author - Vijaya Lakshmi Basireddy
In [1]:	<pre>importing the libraries # Import required libraries import numpy as np</pre>
	<pre>import pandas as pd import matplotlib.pyplot as plt import seaborn as sns print("libraries are imported successfully")</pre>
	libraries are imported successfully load the dataset
In [2]:	<pre>#loading the given dataset df = pd.read_csv('Iris.csv') print("data is imported successfully") data is imported successfully</pre>
In [3]:	<pre>#to get the type of data df.info() <class 'pandas.core.frame.dataframe'=""></class></pre>
	RangeIndex: 150 entries, 0 to 149 Data columns (total 6 columns): # Column Non-Null Count Dtype
	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
In [4]:	5 Species 150 non-null object dtypes: float64(4), int64(1), object(1) memory usage: 7.2+ KB #describing the given dataset
Out[4]:	df.describe() Id SepalLengthCm SepalWidthCm SepalWidthCm PetalLengthCm PetalWidthCm count 150.000000 150.000000 150.000000 150.000000 150.000000 150.000000
	mean 75.500000 5.843333 3.054000 3.758667 1.198667 std 43.445368 0.828066 0.433594 1.764420 0.763161
	min 1.000000 4.300000 2.000000 1.000000 0.100000 25% 38.250000 5.100000 2.800000 1.600000 0.300000 50% 75.500000 5.800000 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
In [5]: Out[5]:	
	145 146 6.7 3.0 5.2 2.3 Iris-virginica 146 147 6.3 2.5 5.0 1.9 Iris-virginica 147 148 6.5 3.0 5.2 2.0 Iris-virginica
	148 149 6.2 3.4 5.4 2.3 Iris-virginica 149 150 5.9 3.0 5.1 1.8 Iris-virginica
In [6]: In [7]:	<pre># Removal of Id column as it is not required for data analysis i.e clustering daf = df.iloc[: , 1:] #printing the head part of the given data</pre>
Out[7]:	df head() Id SepalLengthCm SepalWidthCm SepalWidthCm SepalWidthCm Species 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 3 4 4.6 3.1 1.5 0.2 Iris-setosa
	4 5 5.0 3.6 1.4 0.2 Iris-setosa
In [8]:	Separation of independent and dependent matrix # Independent matrix X = df.iloc[:,:-1].values # Dependent matrix ** The df.iloc[:,:-1] values
In [9]:	<pre>y = df.iloc[:,-1].values X[0] array([1 5 1 3 5 1 4 0 2])</pre>
Out[9]: In [10]:	
	Handling of categorical data
In [11]:	<pre>from sklearn.preprocessing import LabelEncoder y_labelencoder = LabelEncoder() y = y_labelencoder.fit_transform(y)</pre>
Out[11]:	y array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
	0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
In [13]:	Applying suitable ML algorithm from sklearn.cluster import KMeans
111 [10].	<pre>kmeans = KMeans(n_clusters=3, n_init=150) kmeans.fit(X)</pre>
Out[13]:	<pre>KMeans(n_clusters=3, n_init=150) y_kmeans = kmeans.predict(X)</pre>
Out[14]:	y_kmeans array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
	1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
In [15]:	Visualization of Iris flower dataset with optimum 3 clusters plt.scatter(X[:, 0], X[:, 2], c=y_kmeans)
	<pre>plt.xlabel('Sepal length') plt.ylabel('Petal length') plt.title('Sepal length V/S Petal length') pd.DataFrame({'Color':['purple', 'blue', 'yellow'], 'Iris species':['setosa', 'versicolor', 'virginica']})</pre>
Out[15]:	Color Iris species 0 purple setosa
	1 blue versicolor 2 yellow virginica Sepal length V/S Petal length
	4.5
	Hg 3.5 - The state of the state
	2.5 -
In [16]:	0 20 40 60 80 100 120 140 Sepal length plt.scatter(X[:, 3], X[:, 1], c=y_kmeans) plt.xlabel('Petal Width')
	plt.ylabel('Sepal Width') plt.title('Petal Width V/S Sepal Width') pd.DataFrame({'Color':['purple', 'blue', 'yellow'], 'Iris species':['setosa', 'versicolor', 'virginica']})
Out[16]:	Color Iris species 0 purple setosa 1 blue versicolor
	2 yellow virginica Petal Width V/S Sepal Width
	8.0 - 7.5 - 7.0 -
	## 6.5 - Fig. 6.0 - Fi
	5.0
	1 2 3 4 5 6 7 Petal Width Difference in predicted and actual clustering values
In [17]:	<pre>fig, axes = plt.subplots(1,2) axes[0].scatter(X[:, 0], X[:, 2], c=y_kmeans) axes[0].set_xlabel('Sepal length')</pre>
	<pre>axes[0].set_ylabel('Petal length') axes[0].set_title('Predicted') axes[1].scatter(X[:, 0], X[:, 2], c=y)</pre>
Out[17]:	<pre>axes[1].set_xlabel('Sepal length') axes[1].set_ylabel('Petal length') axes[1].set_title('Actual')</pre> Text(0.5, 1.0, 'Actual')
000[21].	Predicted Actual 4.5 Actual
	4.0 4.0
	2.5 - 2.5 -
	2.0 - 2.0 -
In [18]: Out[18]:	y_kmeans array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
	1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
In [19]: Out[19]:	pd.DataFrame({'Predicted':y_kmeans,'Actual':y}) Predicted Actual
	0 1 0 1 1 0 2 1 0
	3 1 0 4 1 0
	145 0 2 146 0 2 147 0 2
	148 0 2 149 0 2
	150 rows × 2 columns Direct cluster analysis using Seaborn
In [20]: Out[20]:	<pre>sns.swarmplot(x='Species', y='SepalWidthCm', data=df)</pre>

4.5

4.0

SepalWidthCm 0.8

2.5

2.0

PetalLengthCm w b G

2 -

Iris-setosa

lris-setosa

Iris-versicolor Species

In [21]: sns.stripplot(x='Species', y='PetalLengthCm', data=df)

Out[21]: <AxesSubplot:xlabel='Species', ylabel='PetalLengthCm'>

lris-versicolor Species lris-virginica

lris-virginica

GRIP - THE SPARK FOUNDATION