,	Virtual Internship Program  Author:- Prateek Kumars Singh  Beginer Level Tasks  Task-1 Iris Flowers Classification MI Project
	This particular ML project is usually referred to as the "Hello World" of Machine Learning. The iris flowers dataset contains numeric attributes, and it is perfect for beginners to learn about supervised ML algorithms, mainly how to load and handle data. Also, since this is a small dataset, it can easily fit in memory without requiring special transformations or scaling capabilities.  Dataset link:-http://archive.ics.uci.edu/ml/datasets/lris  1. Importing Libraries
	<pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns from sklearn.model_selection import train_test_split from sklearn.linear_model import LogisticRegression from sklearn.metrics import accuracy_score import warnings warnings.filterwarnings('ignore') import os os.environ["OMP_NUM_THREADS"] = '1'</pre>
In [30]: In [22]: Out[22]: _	2. Importing The dataset  iris = pd.read_csv('Iris.csv')  iris.describe()  sepal_length
	mean         5.843333         3.054000         3.758667         1.198667           std         0.828066         0.433594         1.764420         0.763161           min         4.300000         2.000000         1.00000         0.100000           25%         5.100000         2.800000         1.600000         0.300000           5.80000         3.00000         4.350000         1.300000           75%         6.40000         3.300000         5.100000         1.800000           max         7.90000         4.40000         6.90000         2.500000
Out[23]: _	sepal_length         sepal_width         petal_width         sepcies           0         5.1         3.5         1.4         0.2         lris-setosa           1         4.9         3.0         1.4         0.2         lris-setosa           2         4.7         3.2         1.3         0.2         lris-setosa           3         4.6         3.1         1.5         0.2         lris-setosa
In [24]: Out[24]: _	4 5.0 3.6 1.4 0.2 lris-setosa  iris-tail()  sepal_length sepal_width petal_length petal_width species  145 6.7 3.0 5.2 2.3 lris-virginica  146 6.3 2.5 5.0 1.9 lris-virginica  147 6.5 3.0 5.2 2.0 lris-virginica
In [25]: Out[25]:	148         6.2         3.4         5.4         2.3         Iris-virginica           149         5.9         3.0         5.1         1.8         Iris-virginica    (150, 5)  iris.info()
	<pre><class 'pandas.core.frame.dataframe'=""> RangeIndex: 150 entries, 0 to 149 Data columns (total 5 columns): # Column Non-Null Count Dtype</class></pre>
In [35]: Out[35]:	iris['species'].value_counts()  Iris-setosa 50 Iris-versicolor 50 Iris-virginica 50 Name: species, dtype: int64  3. Preprocessing The Dataset
Out[38]:	# checking for missing values iris.isnull().sum()  sepal_length
Out[46]:	iris['sepal_length'].hist() <axessubplot:>  25  20  15  10  5  10  10  11  12  13  14  15  15  16  17  18  18  18  18  18  18  18  18  18</axessubplot:>
Out[49]:	1ris['sepal_width'].hist() <axessubplot:></axessubplot:>
In [51]:	25
In [53]:	35
	<pre><axessubplot:></axessubplot:></pre>
In [60]:	<pre># scatterplot colours = ['red', 'yellow', 'violet'] species = ['Iris-setosa', 'Iris-virginica', 'Iris-versicolor']  for i in range(3):     x = iris[iris['species'] == species[i]]     plt.scatter(x['sepal_length'], x['sepal_width'], c = colours[i], label=species[i]) species = ['Iris-setosa', 'Iris-virginica', 'Iris-versicolor']</pre>
Out[60]:	<pre>plt.scatter(x['sepal_length'], x['sepal_width'], c = colours[i], label=species[i]) plt.xlabel('sepal length') plt.ylabel('sepal width') plt.legend()  <matplotlib.legend.legend 0x214161d2bb0="" at="">  4.5 4.0 Iris-setosa Iris-virginica Iris-virginica Iris-virginica Iris-versicolor</matplotlib.legend.legend></pre>
	3.5 - 2.5 - 2.5 - 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0 sepal length
	<pre>for i in range(3):     x = iris[iris['species'] == species[i]]     plt.scatter(x['petal_length'], x['petal_width'], c = colours[i], label=species[i]) plt.xlabel('petal_length') plt.ylabel('petal_width') plt.legend()  <matplotlib.legend.legend 0x214159e8250="" at=""></matplotlib.legend.legend></pre> 25     iris-setosa     iris-virginica     iris-virginica     iris-versicolor
	## 15 -
	<pre>for i in range(3):     x = iris[iris['species'] == species[i]]     plt.scatter(x['sepal_length'], x['petal_length'], c = colours[i], label=species[i]) plt.xlabel('sepal_length') plt.ylabel('petal_length') plt.legend() </pre> <pre> </pre> <pre> </pre> <pre> <pre></pre></pre>
	6 Iris-versicolor  5 3 2 1
	<pre>for i in range(3):     x = iris[iris['species'] == species[i]]     plt.scatter(x['sepal_width'], x['petal_width'], c = colours[i], label=species[i]) plt.ylabel('sepal width') plt.ylabel('pepal width') plt.legend()</pre> <pre><matplotlib.legend.legend 0x21415fe0700="" at=""></matplotlib.legend.legend></pre>
	2.5 - 2.0 - Fig. 1.5 - Iris-setosa Iris-virginica Iris-versicolor
	<pre>sns.violinplot(y='species', x='sepal_length', data=iris, inner='quartile') plt.show() sns.violinplot(y='species', x='sepal_width', data=iris, inner='quartile') plt.show() sns.violinplot(y='species', x='sepal_width', data=iris, inner='quartile') plt.show() sns.violinplot(y='species', x='petal_length', data=iris, inner='quartile') plt.show()</pre>
	sns.violinplot(y='species', x='petal_width', data=iris, inner='quartile') plt.show()    Iris-setosa
	Iris-virginica - 4 5 6 7 8 sepal_length  Iris-setosa
	Iris-virginica -
	Iris-versicolor -  Iris-virginica -
	lris-setosa -
	Iris-virginica  Output  Output  Description  Description
In [65]: In [67]:	5. Level Encoding  Label Encoding is a popular encoding technique for handling categorical variables. In this technique, each label is assigned a unique integer based on alphabetical ordering. Let's see how to implement label encoding in Python using the scikit-learn library and also understand the challenges with label encoding  from sklearn.preprocessing import LabelEncoder le = LabelEncoder()  iris['species'] = le.fit_transform(iris['species']) iris.head()
	sepal_length         sepal_width         petal_length         species           0         5.1         3.5         1.4         0.2         0           1         4.9         3.0         1.4         0.2         0           2         4.7         3.2         1.3         0.2         0           3         4.6         3.1         1.5         0.2         0           4         5.0         3.6         1.4         0.2         0
Out[68]: _	
Out[70]: '	iris.species  0
	147 2 148 2 149 2 Name: species, Length: 150, dtype: int32  6. Model Training  Model training is the phase in the data science development lifecycle where practitioners try to fit the best combination of weights and bias to a machine learning algorithm to minimize a loss function over the prediction range
In [72]: In [73]: In [83]:	<pre>X = iris.drop(columns='species', axis=1) Y = iris['species']  X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.20, stratify=Y, random_state=2)  model = LogisticRegression()  model.fit(X_train, Y_train) LogisticRegression()</pre>
In [76]: In [77]:	<pre># accuracy on training data X_train_prediction = model.predict(X_train) training_data_accuracy = accuracy_score(X_train_prediction, Y_train)  print('Accuracy on Training data : ', training_data_accuracy*100) Accuracy on Training data : 96.6666666666666666666666666666666666</pre>
In [78]:	test_data_accuracy = accuracy_score(X_test_prediction, Y_test)  print('Accuracy on Test data : ', test_data_accuracy*100)  Accuracy on Test data : 100.0  7. Optimal Number Of Clustring For K-Means Classification Algoritham  # Finding the optimal number of clusters for k-means classification Algoritham  x = iris.iloc[:, [0, 1, 2, 3, 4]].values
	<pre># we import KMeans algorithm using sklearn library from sklearn.cluster import KMeans # we use the very first method is Elbow Method wcss = [] # WCSS means Within Cluster Sum of Squares  for i in range(1, 11):</pre>
	# plotting above result in line graph format plt.plot(range(1, 11), wcss) plt.title('The Elbow method') plt.xlabel('Number of clusters') plt.ylabel('Wcss') plt.show()  The Elbow method
	600 - 500 - 200 - 100 - 2 4 6 8 10
In [85]:	In cluster analysis, the elbow method is a heuristic used in determining the number of clusters in a data set. The method consists of plotting the explained variation as a function of the number of clusters and picking the elbow of the curve as the number of clusters to use.  8. Applying K-MEANS Method On Given Dataset  kmeans = KMeans(n_clusters=4,init='k-means++',max_iter=300,
In [86]:	print(y_kmeans)  [1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	<pre>[5.9</pre>
Out[89]:	plt.scatter(x[y_kmeans == 2,0],x[y_kmeans == 2,1],s = 100,
In [90]:	# plotting Centroids of the clusters in above graph plt.figure(figsize = (16,9)) style.use('ggplot') plt.scatter(x[y_kmeans == 0,0],x[y_kmeans == 0,1],s = 100,
Out[90]:	<pre>c = 'blue', label = 'Iris-versicolour', marker='*') plt.scatter(x[y_kmeans == 2,0], x[y_kmeans == 2,1],s = 100,</pre>
	3.0 - * * * * * * * * * * * * * * * * * *
	2.0 - 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0  Thank You