**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/Iris 1. Importing Libraries In [82]: 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' 2. Importing The dataset In [30]: iris = pd.read\_csv('Iris.csv') iris.describe() Out[22]: sepal\_length sepal\_width petal\_length petal\_width **count** 150.000000 150.000000 150.000000 150.000000 mean 5.843333 3.054000 3.758667 1.198667 0.763161 std 0.828066 0.433594 1.764420 4.300000 2.000000 1.000000 0.100000 min 0.300000 5.100000 1.600000 25% 2.800000 **50**% 5.800000 3.000000 4.350000 1.300000 6.400000 **75**% 3.300000 5.100000 1.800000 max 7.900000 4.400000 6.900000 2.500000 iris.head() In [23]: sepal\_length sepal\_width petal\_length petal\_width species Out[23]: 5.1 3.5 1.4 0.2 Iris-setosa 4.9 3.0 1.4 0.2 Iris-setosa 2 4.7 3.2 1.3 0.2 Iris-setosa 3.1 4.6 1.5 0.2 Iris-setosa 5.0 3.6 1.4 0.2 Iris-setosa sepal\_length sepal\_width petal\_length petal\_width species 145 6.7 3.0 5.2 2.3 Iris-virginica 146 6.3 2.5 5.0 1.9 Iris-virginica 147 6.5 3.0 5.2 2.0 Iris-virginica 148 6.2 3.4 5.4 2.3 Iris-virginica 149 5.9 3.0 5.1

Grow More

Virtual Internship Program

Author:- Prateek Kumars Singh

In [24]: iris.tail() Out[24]: 1.8 Iris-virginica iris.shape (150, 5)Out[25]: In [26]: iris.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 150 entries, 0 to 149 Data columns (total 5 columns): Non-Null Count Dtype # Column 0 sepal\_length 150 non-null float64 1 sepal\_width 150 non-null float64 float64 2 petal\_length 150 non-null 3 petal\_width 150 non-null float64 4 species 150 non-null object dtypes: float64(4), object(1) memory usage: 6.0+ KB In [35]: iris['species'].value\_counts() Iris-setosa 50 Out[35]: Iris-versicolor 50 Iris-virginica 50 Name: species, dtype: int64 3. Preprocessing The Dataset In [38]: # checking for missing values iris.isnull().sum() sepal\_length Out[38]: sepal\_width 0 petal\_length petal\_width species dtype: int64

4. Data Analysis In [46]: iris['sepal\_length'].hist() <AxesSubplot:> Out[46]: 25 20 15 10 5 5.0 5.5 6.0 6.5 iris['sepal\_width'].hist() <AxesSubplot:> Out[49]:

35 30 25 20 15 10

iris['petal\_length'].hist()

iris['petal\_width'].hist()

<AxesSubplot:>

<AxesSubplot:>

35

30

25

20

15

10

20 15 10

0.0

In [60]: **for** i **in** range(3):

plt.legend()

4.0

2.5

2.0

In [61]: **for** i **in** range(3):

plt.legend()

2.0

fj 1.5

In [62]: **for** i **in** range(3):

plt.legend()

petal length

In [63]: **for** i **in** range(3):

plt.legend()

2.5

2.0

0.5

0.0

2.0

plt.show()

plt.show()

Iris-setosa

Iris-versicolor

Iris-virginica

Iris-setosa

lris-versicolor

Iris-virginica

Iris-setosa

Iris-versicolor

Iris-virginica

Iris-setosa

Iris-versicolor

Iris-virginica

In [65]:

Out[67]:

Out[68]:

Out[70]:

In [73]:

In [83]:

Out[83]:

In [75]:

0.0

5. Level Encoding

le = LabelEncoder()

5.1

4.9

4.7

4.6

5.0

6.7

6.3

6.5

6.2

5.9

iris.head()

iris.tail()

2

145

146

147

148

149

2

3

iris.species

0

0

0

6. Model Training

Y = iris['species']

LogisticRegression()

In [77]: # accuracy on test data

wcss = []

**for** i **in** range(1, 11):

kmeans.fit(x)

plt.ylabel('WCSS')

plt.show()

800 700

100

In [85]:

Out[86]:

function over the prediction range

model = LogisticRegression()

model.fit(X\_train, Y\_train)

# accuracy on training data

Accuracy on Test data : 100.0

X\_train\_prediction = model.predict(X\_train)

Accuracy on Training data : 96.6666666666667

In [78]: print('Accuracy on Test data : ', test\_data\_accuracy\*100)

X\_test\_prediction = model.predict(X\_test)

x = iris.iloc[:, [0, 1, 2, 3, 4]].values

from sklearn.cluster import KMeans

wcss.append(kmeans.inertia\_)

plt.plot(range(1, 11), wcss) plt.title('The Elbow method') plt.xlabel('Number of clusters')

# we import KMeans algorithm using sklearn library

kmeans = KMeans(n\_clusters = i, init = 'k-means++',

# we use the very first method is Elbow Method

# WCSS means Within Cluster Sum of Squares

# plotting above result in line graph format

The Elbow method

Number of clusters

clusters and picking the elbow of the curve as the number of clusters to use.

8. Applying K-MEANS Method On Given Dataset

y\_kmeans = kmeans.fit\_predict(x)

print(y\_kmeans)

In [86]: kmeans.cluster\_centers\_

[5.006

plt.figure(figsize = (16,9)) from matplotlib import style

[5.9

In [89]: # Visualizing the clusters

plt.legend()

4.5

4.0

3.5

3.0 -

2.5

2.0

style.use('qqplot')

3 3]

kmeans = KMeans(n\_clusters=4,init='k-means++',max\_iter=300, n\_init=10, random\_state=5)

array([[7.42307692, 3.13076923, 6.26923077, 2.06923077, 2.

plt.scatter( $x[y\_kmeans == 0,0], x[y\_kmeans == 0,1], s = 100,$ 

plt.scatter( $x[y_kmeans == 1,0], x[y_kmeans == 1,1], s = 100,$ 

 $plt.scatter(x[y_kmeans == 2,0],x[y_kmeans == 2,1],s = 100,$ 

, 1.464

, 4.25

c = 'blue', label = 'Iris-setosa', marker='\*')

c = 'red', label = 'Iris-setosa', marker='\*')

c = 'yellow', label = 'Iris-versicolour', marker='\*')

, 3.418

, 2.76

<matplotlib.legend.Legend at 0x21417fdfdc0>

[6.34324324, 2.93243243, 5.31351351, 2.01081081, 1.97297297]])

, 0.244

5.5

5.5

5.0

plt.scatter( $x[y_kmeans == 0,0], x[y_kmeans == 0,1], s = 100,$ 

plt.scatter( $x[y_kmeans == 1,0], x[y_kmeans == 1,1], s = 100,$ 

plt.scatter( $x[y\_kmeans == 2,0], x[y\_kmeans == 2,1], s = 100,$ 

s=100, c='red', label='Cluster\_Centroids')

c = 'green',label = 'Iris-setosa',marker='\*')

c = 'blue', label = 'Iris-versicolour', marker='\*')

c = 'yellow', label = 'Iris-setosa', marker='\*')

plt.scatter(kmeans.cluster\_centers\_[:,0],kmeans.cluster\_centers\_[:,1],

In [90]: # plotting Centroids of the clusters in above graph

<matplotlib.legend.Legend at 0x21417fcba30>

plt.figure(figsize = (16,9))

style.use('ggplot')

plt.legend()

4.5 -

4.0

3.5

3.0

2.5

Thank You

6.0

, 1.326

, 0.

, 1.02

training\_data\_accuracy = accuracy\_score(X\_train\_prediction, Y\_train)

In [76]: print('Accuracy on Training data : ', training\_data\_accuracy\*100)

test\_data\_accuracy = accuracy\_score(X\_test\_prediction, Y\_test)

7. Optimal Number Of Clustring For K-Means Classification Algoritham

 $max_iter = 300$ ,  $n_init = 10$ ,  $random_state = 0$ )

In [81]: # Finding the optimal number of clusters for k-means classification Algoritham

In [71]: X = iris.drop(columns='species', axis=1)

0.5

from sklearn.preprocessing import LabelEncoder

In [67]: iris['species'] = le.fit\_transform(iris['species'])

3.5

3.0

3.2

3.1

3.6

3.0

2.5

3.0

3.4

3.0

Name: species, Length: 150, dtype: int32

sepal\_length sepal\_width petal\_length

1.0

1.5

petal\_width

1.4

1.4

1.3

1.5

1.4

5.0

5.2

5.1

sepal\_length sepal\_width petal\_length petal\_width species

2.0

2.5

3.0

sepal\_width

petal\_length

3.5

4.0

4.5

2.5

implement label encoding in Python using the scikit-learn library and also understand the challenges with label encoding

0

0

0

0

0

2

2

2

petal\_width species

0.2

0.2

0.2

2.3

1.9

2.0

2.3

1.8

In [72]: X\_train, X\_test, Y\_train, Y\_test = train\_test\_split(X, Y, test\_size=0.20, stratify=Y, random\_state=2)

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

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

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

Iris-setosa Iris-versicolour ★ Iris-setosa

Iris-setosa Iris-versicolour Iris-setosa Cluster Centroids

],

],

],

Out[63]:

4.5

5.0

plt.xlabel('petal length') plt.ylabel('petal width')

> Iris-virginica Iris-versicolor

plt.xlabel('sepal length') plt.ylabel('petal length')

> Iris-setosa Iris-virginica Iris-versicolor

plt.xlabel('sepal width') plt.ylabel('pepal width')

2.5

3.0

sepal width

3.5

In [59]: # scatterplot

Out[60]:

Out[61]:

0.5

plt.xlabel('sepal length') plt.ylabel('sepal width')

1.0

x = iris[iris['species'] == species[i]]

<matplotlib.legend.Legend at 0x214161d2bb0>

sepal length

petal length

x = iris[iris['species'] == species[i]]

<matplotlib.legend.Legend at 0x214160bf970>

sepal length

x = iris[iris['species'] == species[i]]

<matplotlib.legend.Legend at 0x21415fe0700>

x = iris[iris['species'] == species[i]]

<matplotlib.legend.Legend at 0x214159e8250>

colours = ['red', 'yellow', 'violet']

1.5

species = ['Iris-setosa', 'Iris-virginica', 'Iris-versicolor']

2.0

plt.scatter(x['sepal\_length'], x['sepal\_width'], c = colours[i], label=species[i])

Iris-setosa Iris-virginica

Iris-versicolor

7.5

plt.scatter(x['petal\_length'],x['petal\_width'], c = colours[i], label=species[i])

plt.scatter(x['sepal\_length'], x['petal\_length'], c = colours[i], label=species[i])

7.5

plt.scatter(x['sepal\_width'], x['petal\_width'], c = colours[i], label=species[i])

Iris-setosa Iris-virginica Iris-versicolor

4.0

sns.violinplot(y='species', x='sepal\_length', data=iris, inner='quartile')

sns.violinplot(y='species', x='sepal\_width', data=iris, inner='quartile')

sns.violinplot(y='species', x='petal\_length', data=iris, inner='quartile')

sns.violinplot(y='species', x='petal\_width', data=iris, inner='quartile')

sepal\_length

4.5

Out[51]:

Out[53]: