

IE406 Machine Learning

Lab Assignment - 3

Group 39

201901201: Sharvil Sheth
201901216: Prabhav Shah
201901301: Mohil Desai
201901462: Dhaval Vaidya
201901464: Vishvesh Patel

Question 1

In this question we will use the data given in file “Social_Network_Ads.csv” which is a categorical dataset to determine whether a user purchased a product or not by using three features to determine user’s decision. Visualize the data by 3D plotting features using different colors for label 0 and 1. Use data in files “Social_Network_Ads.csv” to perform logistic regression by implementing logistic function and with available library function and compare your results. Use 90% data points from each set for training and remaining 10% for testing the accuracy of classification. Using confusion matrix find accuracy, precision, F1 score and recall.

Answer

code

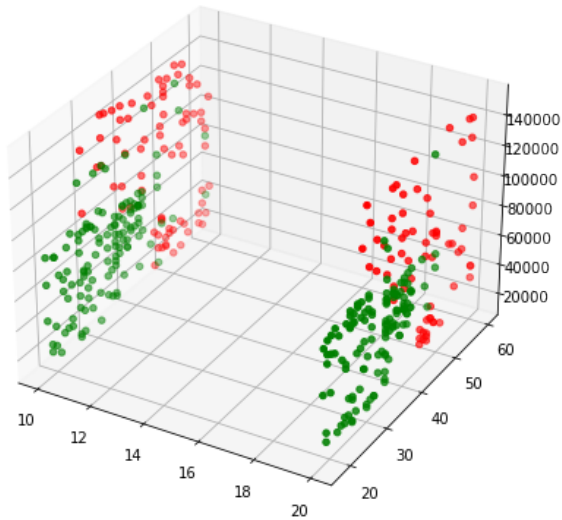
```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn import preprocessing
5 import scipy
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.metrics import confusion_matrix
9
10 df = pd.read_csv('Social_Network_Ads.csv')
11 df['Gender'] = df['Gender'].replace("Male", 20)
12 df['Gender'] = df['Gender'].replace("Female", 10)
13 df.info()
14 <class 'pandas.core.frame.DataFrame'>
15 RangeIndex: 400 entries, 0 to 399
16 Data columns (total 5 columns):
17 #   Column          Non-Null Count  Dtype
18 ---  ---
19 0    User ID         400 non-null    int64
20 1    Gender          400 non-null    int64
21 2    Age             400 non-null    int64
22 3    EstimatedSalary 400 non-null    int64
23 4    Purchased       400 non-null    int64
24 dtypes: int64(5)
25 memory usage: 15.8 KB
26
27 grouped = df.groupby(df.Purchased)
28 bias_0_plot = grouped.get_group(0)
29 bias_1_plot = grouped.get_group(1)
30
31 X_0_plot = bias_0_plot[['Gender', 'Age', 'EstimatedSalary']].values
32
33 X_1_plot = bias_1_plot[['Gender', 'Age', 'EstimatedSalary']].values
34
35 X_0_norm0 = preprocessing.normalize([X_0_plot[:,0]])
36 X_0_norm1 = preprocessing.normalize([X_0_plot[:,1]])
37 X_0_norm2 = preprocessing.normalize([X_0_plot[:,2]])
```

```

39
40 X_1_norm0 = preprocessing.normalize([X_1_plot[:,0]])
41 X_1_norm1 = preprocessing.normalize([X_1_plot[:,1]])
42 X_1_norm2 = preprocessing.normalize([X_1_plot[:,2]])
43
44 fig = plt.figure(figsize = (10, 7))
45 ax = plt.axes(projection = "3d")
46
47 ax.scatter3D(X_0_plot[:,0], X_0_plot[:,1], X_0_plot[:,2], color = "green")
48 ax.scatter3D(X_1_plot[:,0], X_1_plot[:,1], X_1_plot[:,2], color = "red")
49 plt.show()

```

Listing 1: Question 1

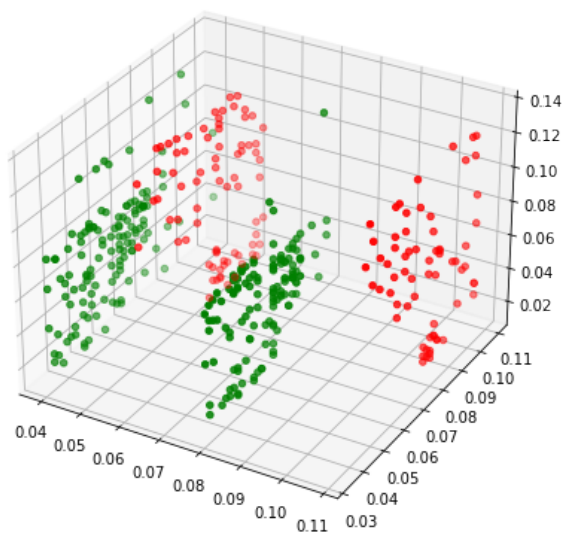


```

1 fig = plt.figure(figsize = (10, 7))
2 ax = plt.axes(projection = "3d")
3
4 ax.scatter3D(X_0_norm0, X_0_norm1, X_0_norm2, color = "green")
5 ax.scatter3D(X_1_norm0, X_1_norm1, X_1_norm2, color = "red")
6 plt.show()

```

Listing 2: Question 1



```

1
2 def sigmoid(x):
3     # Activation function used to map any real value between 0 and 1
4     return 1 / (1 + np.exp(-x))
5

```

```

6 def net_input(theta, x):
7     # Computes the weighted sum of inputs
8     return np.dot(x, theta)
9
10 def probability(theta, x):
11     # Returns the probability after passing through sigmoid
12     return sigmoid(net_input(theta, x))
13
14     def cost_function(theta, x, y):
15         # Computes the cost function for all the training samples
16         m = x.shape[0]
17         total_cost = -(1 / m) * np.sum(y * np.log(probability(theta, x)) + (1 - y) * np.log(1 -
18             probability(theta, x)))
19         return total_cost
20
21 def gradient(theta, x, y):
22     # Computes the gradient of the cost function at the point theta
23     m = x.shape[0]
24     return (1 / m) * np.dot(x.T, sigmoid(net_input(theta, x)) - y)
25
26     def fit(x, y, theta):
27         opt_weights = scipy.optimize.fmin_tnc(func=cost_function, x0=theta,
28             fprime=gradient, args=(x, y.flatten()))
29         return opt_weights[0]
30
31     X = df.iloc[:, 1:4]
32     X = np.c_[np.ones((X.shape[0], 1)), X]
33     Y = df[['Purchased']].values
34     theta = np.zeros((X.shape[1], 1))
35
36     X_train_own, X_test_own, Y_train_own, Y_test_own = train_test_split(X, Y, test_size=0.2)
37
38     parameters = fit(X_train_own, Y_train_own, theta)
39     parameters
40     output = array([-1.54919397e+01,  2.61066273e-02,  2.97935984e-01,  3.83055362e-05])
41
42 def predict(x):
43     theta = parameters[:, np.newaxis]
44     return probability(theta, x)
45
46 def accuracy(x, actual_classes, probab_threshold=0.5):
47     predicted_classes = (predict(x) >=
48         probab_threshold).astype(int)
49     predicted_classes = predicted_classes.flatten()
50     accuracy = np.mean(predicted_classes == actual_classes)
51     return accuracy
52
53 accuracy(X_test_own, Y_test_own.flatten())
54 output- 0.7875

```

Listing 3: Logistic Regression our own implementation

```

1 X = df[['Gender', 'Age', "EstimatedSalary"]].values
2 Y = df[['Purchased']].values
3
4 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2)
5
6 clf = LogisticRegression(random_state=0, max_iter=1000).fit(X_train, Y_train)
7 Y_pred = clf.predict(X_test)
8 clf.score(X_test, Y_test)
9 output- 0.65
10
11 tn, fp, fn, tp = confusion_matrix(Y_test, Y_pred).ravel()
12
13 acc = (tn+tp)/(tn+tp+fp+fn)
14 prec = tp/(tp+fp)
15 recall = tp/(tp+fn)
16 f1_score = 2*(prec*recall)/(prec+recall)
17 print('Accuracy - ', acc)
18 print('Precision - ', prec)
19 print('Recall - ', recall)
20 print('F1 Score - ', f1_score)
21
22 Accuracy - 0.65
23 Precision - 0.5151515151515151

```

```
24 Recall - 0.5862068965517241
25 F1 Score - 0.5483870967741935
```

Listing 4: Implementing sklearn logistic regression

Observation/ Justification

In this question we observed that the logistic regression implementation that we did on our own gave better performance than the sklearn library

Question 2

You will work with a widely used Iris dataset. The Iris Dataset contains four features (sepal length, sepal width, petal length, and petal width) of 50 samples of three species of Iris (Iris setosa, Iris virginica, and Iris versicolor). Plot features' histogram. Compute pdf and compare it with histogram. perform the exploratory data analysis by plotting the basic statistics like mean, median, min, and max value of each feature (sepal and petal lengths and widths) for each of the three classes (setosa, virginica, and versicolor).

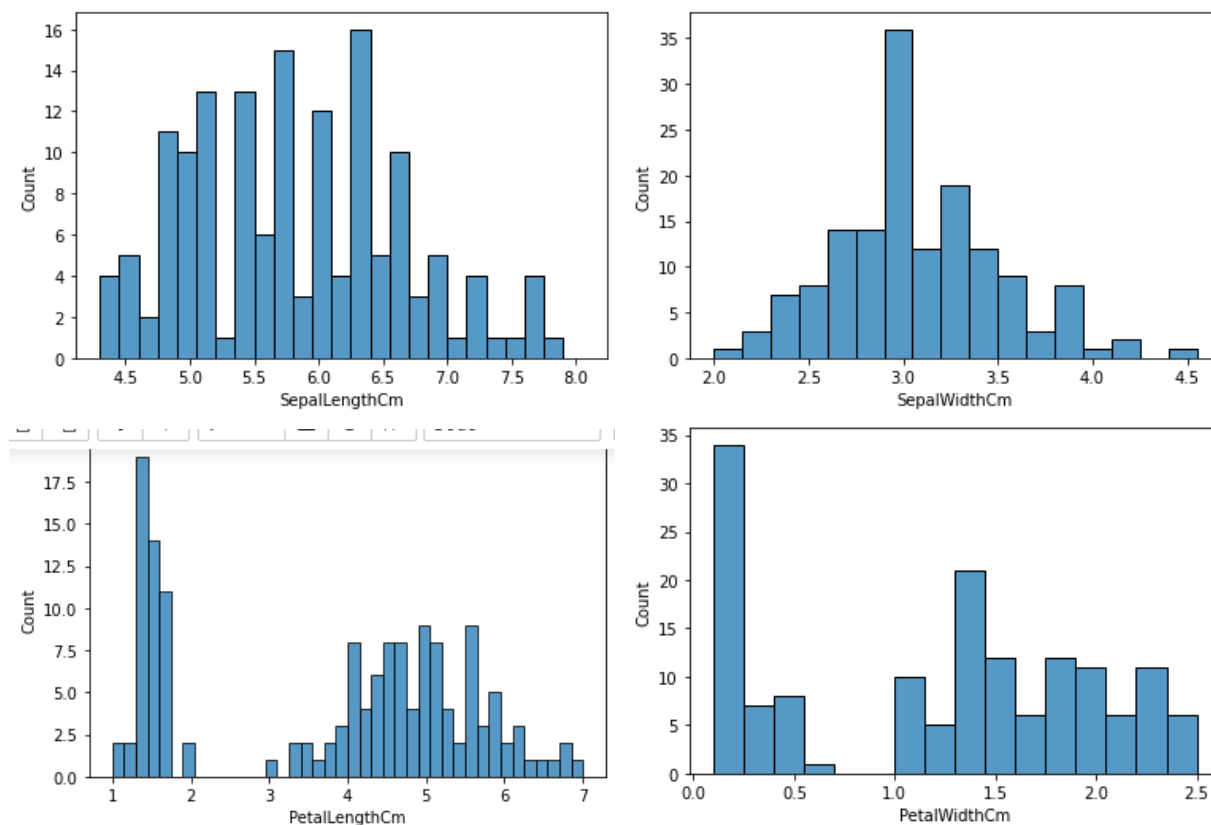
Answer

code

```
1 #Code for plotting Histograms
2 import seaborn as sns
3 plt . figure ()
4 sns . histplot ( data =df[df. columns [1]] , binwidth =0.15)
5 plt . figure ()
6 sns . histplot ( data =df[df. columns [2]] , binwidth =0.15)
7 plt . figure ()
8 sns . histplot ( data =df[df. columns [3]] , binwidth =0.15)
9 plt . figure ()
10 sns . histplot ( data =df[df. columns [4]] , binwidth =0.15) ;
```

Listing 5: Question 2

Result



code

```

1 #Code for basic statistics
2 iris=df.iloc[:,1:]
3 iris.groupby('Species').agg(['mean','median','min','max'])
4
5 #Code for boxplot to plot and analyse the basic stats
6 import seaborn as sns
7 iris=df.iloc[:,1:]
8 sns.set(style="ticks")
9 plt.figure(figsize=(12,12))
10 plt.subplot(2,2,1)
11 sns.boxplot(x='Species',y='SepalLengthCm',data=iris)
12 plt.subplot(2,2,2)
13 sns.boxplot(x='Species',y='SepalWidthCm',data=iris)
14 plt.subplot(2,2,3)
15 sns.boxplot(x='Species',y='PetalLengthCm',data=iris)
16 plt.subplot(2,2,4)
17 sns.boxplot(x='Species',y='PetalWidthCm',data=iris)
18 plt.show()

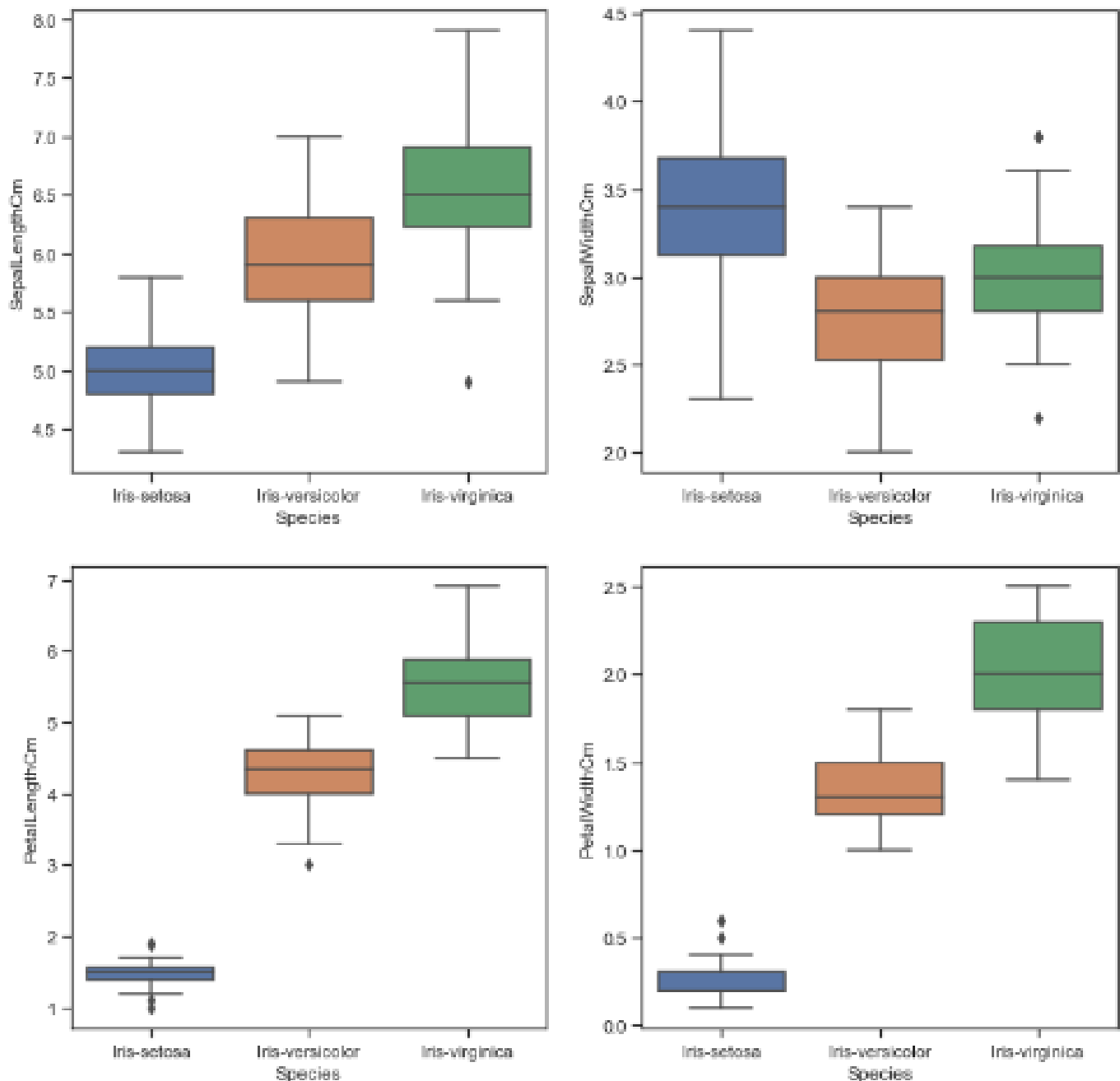
```

Listing 6: Question 2

Result

Species	SepalLengthCm				SepalWidthCm				PetalLengthCm				PetalWidthCm			
	mean	median	min	max	mean	median	min	max	mean	median	min	max	mean	median	min	max
Iris-setosa	5.006	5.0	4.3	5.8	3.418	3.4	2.3	4.4	1.464	1.50	1.0	1.9	0.244	0.2	0.1	0.6
Iris-versicolor	5.936	5.9	4.9	7.0	2.770	2.8	2.0	3.4	4.260	4.35	3.0	5.1	1.326	1.3	1.0	1.8
Iris-virginica	6.588	6.5	4.9	7.9	2.974	3.0	2.2	3.8	5.552	5.55	4.5	6.9	2.026	2.0	1.4	2.5

Species	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
	mean	median	min	max
Iris-setosa	0.352490	0.381024	0.173511	0.107210
Iris-versicolor	0.516171	0.313798	0.469911	0.197753
Iris-virginica	0.635880	0.322497	0.551895	0.274650



Question 3

Visualize the data in the Iris Dataset by considering maximum combinations of two features in a 2D plot. Use red, green, and blue colors for labeling the three classes: Iris setosa, Iris virginica, and Iris versicolor, respectively. Comment on whether any two classes among the three can be separated by a line? Report your observations for each case.

Answer

code

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import scipy.linalg as la
5 import math
6 from scipy.fftpack import fft, fftfreq
7 from scipy.linalg import toeplitz
8 from matplotlib import animation
```

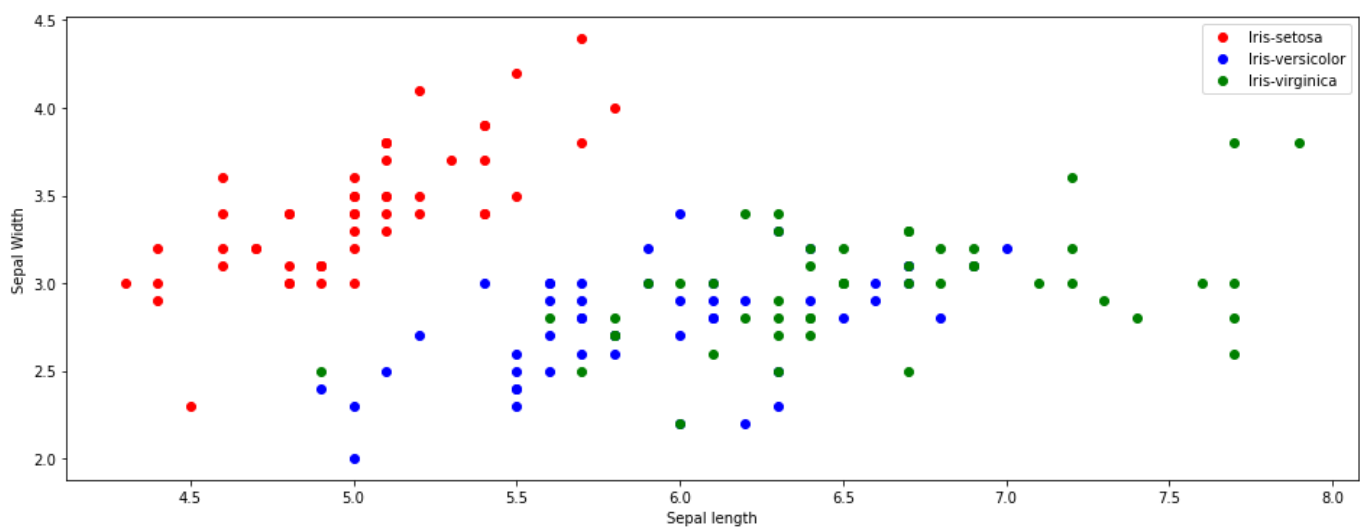
```

9
10 df = pd.read_csv(r'C:\Users\Vatsal\Downloads\Iris.csv')
11 X=df.iloc[:,1:3]
12 Y=df.iloc[:,5]
13 X=np.array(X)
14 Y=np.array(Y)
15 a=X[np.where(Y=='Iris-setosa')]
16 b=X[np.where(Y=='Iris-versicolor')]
17 c=X[np.where(Y=='Iris-virginica')]
18 plt.figure(figsize=(16,20))
19 plt.subplot(3, 1, 1)
20 plt.scatter(a[:,0],a[:,1],c='red',alpha=1,label='Iris-setosa')
21 plt.scatter(b[:,0],b[:,1],c='blue',alpha=1,label='Iris-versicolor')
22 plt.scatter(c[:,0],c[:,1],c='green',alpha=1,label='Iris-virginica')
23 plt.legend()
24 plt.xlabel('Sepal length')
25 plt.ylabel('Sepal Width')
26 plt.show()

```

Listing 7: Question 3

Result



Observation

Class setosa and Class versicolor can be separated by drawing a line and we can get 100% training accuracy as both class dots are not overlapping with each other in majority of area. Similarly, class setosa and class virginica can be separated by drawing a line and can achieve a 100% train accuracy. Here, class versicolor and class virginica are overlapping and cannot be separated by a line.

code

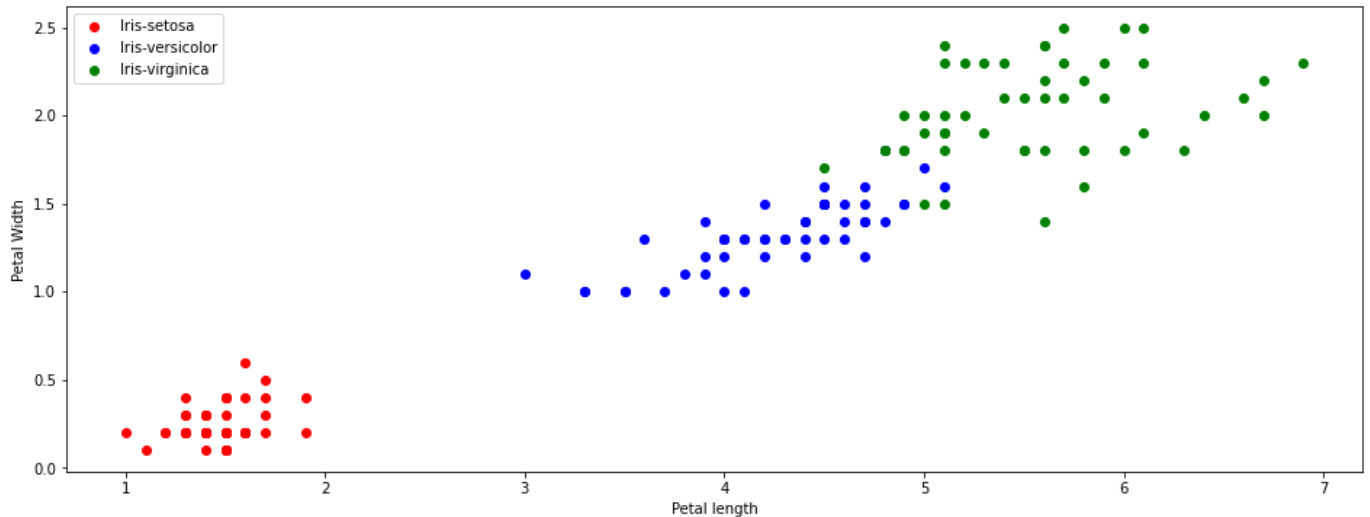
```

1 X=df.iloc[:,3:5]
2 Y=df.iloc[:,5]
3 X=np.array(X)
4 Y=np.array(Y)
5 a=X[np.where(Y=='Iris-setosa')]
6 b=X[np.where(Y=='Iris-versicolor')]
7 c=X[np.where(Y=='Iris-virginica')]
8 plt.figure(figsize=(16,20))
9 plt.subplot(3, 1, 1)
10 plt.scatter(a[:,0],a[:,1],c='red',alpha=1,label='Iris-setosa')
11 plt.scatter(b[:,0],b[:,1],c='blue',alpha=1,label='Iris-versicolor')
12 plt.scatter(c[:,0],c[:,1],c='green',alpha=1,label='Iris-virginica')
13 plt.legend()
14 plt.xlabel('Petal length')
15 plt.ylabel('Petal Width')
16 plt.show()

```

Listing 8: Question 3

Result



Observation/ Justification

Class setosa and Class versicolor can be separated by drawing a line and we can achieve 100% train accuracy as both class dots are not overlapping with majority area of each other. Similarly, Class setosa and Class virginica can be separated by drawing a line and can achieve a 100% train accuracy. Here, Class versicolor and Class Virginica are separable with a line, but can't guarantee 100% train accuracy since there are some dots overlapping. Hence this dots near decision boundary will not guarantee completely accurate predictions.

Question 4

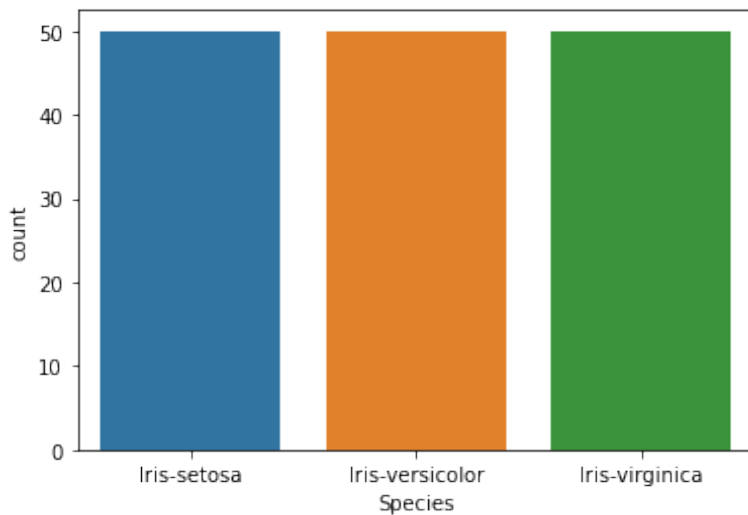
Perform logistic regression on IRIS Dataset and plot confusion matrix. Using confusion matrix find accuracy, precision, F1 score and recall.

Answer

code

```
1 import matplotlib.pyplot as p
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.metrics import classification_report
8 from sklearn.metrics import f1_score, precision_score, recall_score, accuracy_score
9 from sklearn.metrics import confusion_matrix
10
11 data = pd.read_csv('Iris.csv')
12 p.figure(1)
13 sns.countplot(x = 'Species', data = data)
14
15
16 data.drop('Id', axis=1, inplace=True)
```

Listing 9: Question 4



```

1 # Changing
2 # Iris-setosa to 0
3 # Iris-versicolor to 1
4 # Iris-virginica to 2
5
6
7 for i in range (len(data['Species'])):
8     if (data['Species'][i] == 'Iris-setosa'):
9         data['Species'][i] = 0
10    elif (data['Species'][i] == 'Iris-versicolor'):
11        data['Species'][i] = 1
12    elif (data['Species'][i] == 'Iris-virginica'):
13        data['Species'][i] = 2

```

Listing 10: Question 4

```

1
2 X = data.drop(['Species'], axis = 1)
3 y = data['Species']
4 y = y.astype('int')
5
6 ## CHANGE TEST SIZE HERE
7 testSize = 0.5
8
9
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=testSize, random_state= 40
    )

```

Listing 11: Data Splitting to Test and Training

```

1
2 logmodel = LogisticRegression(max_iter=1000)
3 logmodel.fit(X_train, y_train)
4
5 predictions = logmodel.predict(X_test)

```

Listing 12: Training Data and Predictions

```

1
2 report = classification_report(y_test, predictions)
3 print(report)
4
5 acc = accuracy_score(y_test, predictions)
6 precision = precision_score(y_test, predictions, average='micro')
7 recall = recall_score(y_test, predictions, average='micro')
8 f1 = f1_score(y_test, predictions, average='micro')
9 print("Accuracy: ", acc)
10 print("Precision: ", precision)
11 print("Recall: ", recall)
12 print("f1 score: ", f1)
13
14
15 conf_matrix = confusion_matrix(y_test, predictions)

```

```

16 print('\n\nConfusion Matrix :-\n', confusion_matrix(y_test, predictions))
17
18 precision    recall  f1-score   support
19
20      0       1.00      1.00      1.00        27
21      1       0.88      1.00      0.94        22
22      2       1.00      0.88      0.94        26
23
24   accuracy                0.96        75
25  macro avg              0.96      0.96      0.96        75
26 weighted avg              0.96      0.96      0.96        75
27
28 Accuracy:  0.96
29 Precision:  0.96
30 Recall:    0.96
31 f1 score:   0.96
32
33
34 Confusion Matrix :-
35 [[27  0  0]
36 [ 0 22  0]
37 [ 0  3 23]]
38

```

Listing 13: Confusion Matrix

```

1
2
3 p.figure(3)
4
5 fig, ax = p.subplots(figsize=(7.5, 7.5))
6 ax.matshow(conf_matrix, cmap=p.cm.Blues, alpha=0.3)
7 for i in range(conf_matrix.shape[0]):
8     for j in range(conf_matrix.shape[1]):
9         ax.text(x=j, y=i, s=conf_matrix[i, j], va='center', ha='center', size='xx-large')
10
11
12 p.xlabel('Predictions', fontsize=18)
13 p.ylabel('Actuals', fontsize=18)
14 p.title('Confusion Matrix for Test Size: ' + str(testSize) , fontsize=18)
15
16
17 p.show()

```

Listing 14: Confusion Matrix

Result



Observation/ Justification

Question 5

In this question you will perform logistic regression for multiclass classification on the

20 News groups dataset. Since this dataset is a balanced one, you will perform the pre- processing to create an imbalanced version of the dataset (by removing some news

articles from some groups). One example is given below. Perform multiclass classification using logistic regression on both the balanced and the imbalanced version of the dataset. Compare the performance in each case by obtaining the confusion matrix and accuracy. Report your observations at the end.

Answer

code

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, classification_report
7 from sklearn.metrics import f1_score, precision_score, recall_score, accuracy_score
8
9 df = pd.read_csv('Iris.csv')
10 # df['Species'] = df['Species'].replace('Iris-setosa', 0)
11 # df['Species'] = df['Species'].replace('Iris-versicolor', 1)
12 # df['Species'] = df['Species'].replace('Iris-virginica', 2)
13 df.info()
14 <class 'pandas.core.frame.DataFrame'>
15 RangeIndex: 150 entries, 0 to 149
16 Data columns (total 6 columns):
17 #   Column          Non-Null Count  Dtype
18 ---  ---
19 0    Id              150 non-null    int64
20 1    SepalLengthCm   150 non-null    float64
21 2    SepalWidthCm    150 non-null    float64
22 3    PetalLengthCm   150 non-null    float64
23 4    PetalWidthCm    150 non-null    float64
24 5    Species         150 non-null    object
25 dtypes: float64(4), int64(1), object(1)
26 memory usage: 7.2+ KB
```

Listing 15: Question 5

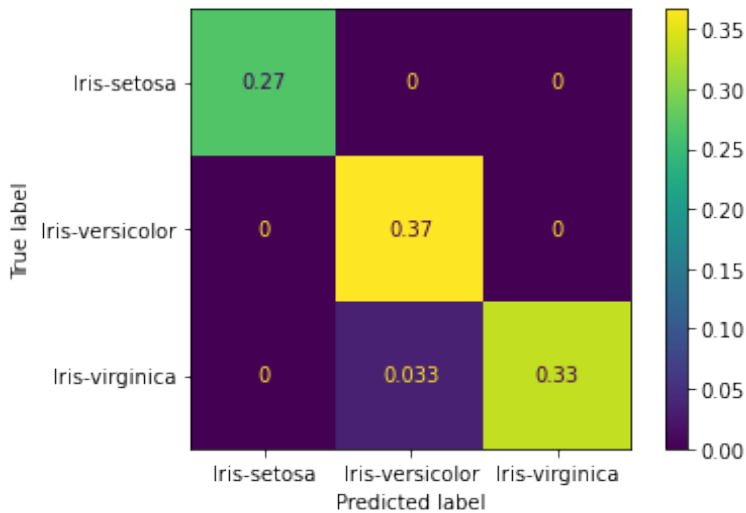
```
1 #loading data on X and Y
2 X = df.iloc[:, 1:5]
3 X = np.c_[np.ones((X.shape[0], 1)), X]
4 Y = df[['Species']].values
5 theta = np.zeros((X.shape[1], 1))
6
7 np.array(np.unique(Y, return_counts=True)).T
8
9 array([[ 'Iris-setosa', 50],
10        [ 'Iris-versicolor', 50],
11        [ 'Iris-virginica', 50]], dtype=object)
12
13 #splitting dataset for training and testing
14 X_train_bal, X_test_bal, Y_train_bal, Y_test_bal = train_test_split(X, Y, test_size=0.2)
15
16 #training model
17 clf_bal = LogisticRegression(random_state=0, max_iter=1000).fit(X_train_bal, Y_train_bal)
18 Y_pred_bal = clf_bal.predict(X_test_bal)
19
20 acc_bal = accuracy_score(Y_test_bal, Y_pred_bal)
21 precision_bal = precision_score(Y_test_bal, Y_pred_bal, average='micro')
22 recall_bal = recall_score(Y_test_bal, Y_pred_bal, average='micro')
23 f1_bal = f1_score(Y_test_bal, Y_pred_bal, average='micro')
24 print("Accuracy: ", acc_bal)
25 print("Precision: ", precision_bal)
26 print("Recall: ", recall_bal)
27 print("f1 score: ", f1_bal)
```

```

28
29 report_bal = classification_report(Y_test_bal, Y_pred_bal)
30 print(report_bal)
31
32 Accuracy:  0.9666666666666667
33 Precision:  0.9666666666666667
34 Recall:    0.9666666666666667
35 f1 score:  0.9666666666666667
36
37      precision    recall  f1-score   support
38
39  Iris-setosa      1.00      1.00      1.00        8
40  Iris-versicolor  0.92      1.00      0.96       11
41  Iris-virginica   1.00      0.91      0.95       11
42
43   accuracy
44  macro avg   0.97      0.97      0.97       30
45  weighted avg 0.97      0.97      0.97       30
46
47 #generating confusion matrix
48 cm_bal = confusion_matrix(Y_test_bal, Y_pred_bal, labels=clf_bal.classes_, normalize='all')
49 disp_bal = ConfusionMatrixDisplay(confusion_matrix=cm_bal, display_labels=clf_bal.classes_)
50 disp_bal.plot()

```

Listing 16: Multiclass classification on Balanced dataset



```

1 #making imbalanced dataset from balanced one
2 df_imbalance = df.copy()
3 df_imbalance = df_imbalance.append(df[0:50])
4 df_imbalance = df_imbalance.append(df[0:50])
5 df_imbalance = df_imbalance.append(df[50:100])
6
7 #loading data on X and Y
8 X = df_imbalance.iloc[:, 1:5]
9 X = np.c_[np.ones((X.shape[0], 1)), X]
10 Y = df_imbalance[['Species']].values
11 theta = np.zeros((X.shape[1], 1))
12
13 np.array(np.unique(Y, return_counts=True)).T
14
15 array([[ 'Iris-setosa', 150],
16        [ 'Iris-versicolor', 100],
17        [ 'Iris-virginica', 50]], dtype=object)
18
19 #splitting dataset for training and testing
20 X_train_imbal, X_test_imbal, Y_train_imbal, Y_test_imbal = train_test_split(X, Y, test_size
    =0.2)
21
22 #training model and calculating accuracy
23 clf_imbal = LogisticRegression(random_state=0, max_iter=1000).fit(X_train_imbal, Y_train_imbal
    )
24 Y_pred_imbal = clf_imbal.predict(X_test_imbal)
25
26 acc_imbal = accuracy_score(Y_test_imbal, Y_pred_imbal)
27 precision_imbal = precision_score(Y_test_imbal, Y_pred_imbal, average='micro')

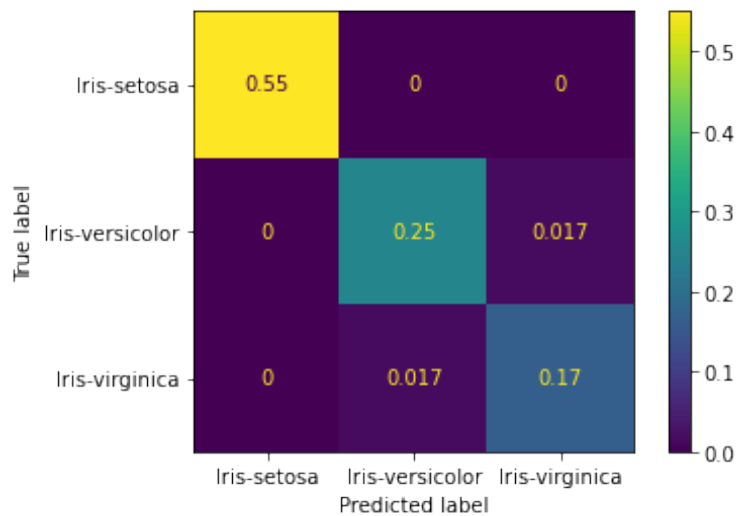
```

```

28 recall_imbal = recall_score(Y_test_imbal, Y_pred_imbal, average='micro')
29 f1_imbal = f1_score(Y_test_imbal, Y_pred_imbal, average='micro')
30 print("Accuracy: ", acc_imbal)
31 print("Precision: ", precision_imbal)
32 print("Recall: ", recall_imbal)
33 print("f1 score: ", f1_imbal)
34
35 report_imbal = classification_report(Y_test_imbal, Y_pred_imbal)
36 print(report_imbal)
37
38 Accuracy:  0.9666666666666667
39 Precision:  0.9666666666666667
40 Recall:    0.9666666666666667
41 f1 score:  0.9666666666666667
42
43      precision    recall  f1-score   support
44
45  Iris-setosa      1.00      1.00      1.00        33
46  Iris-versicolor  0.94      0.94      0.94        16
47  Iris-virginica   0.91      0.91      0.91        11
48
49   accuracy
50   macro avg   0.95      0.95      0.95        60
51   weighted avg   0.97      0.97      0.97        60
52
53 #generating confusion matrix
54 cm_imbal = confusion_matrix(Y_test_imbal, Y_pred_imbal, labels=clf_imbal.classes_, normalize='
55 all')
56 disp_imbal = ConfusionMatrixDisplay(confusion_matrix=cm_imbal, display_labels=clf_imbal.
57 classes_)
58 disp_imbal.plot()

```

Listing 17: Multiclass classification on imbalanced data set



```

1 f, axes = plt.subplots(1, 2, figsize=(15,5), sharey=True)
2
3 disp_bal.plot(ax=axes[0], xticks_rotation=45)
4 disp_bal.ax_.set_title("Balanced dataset")
5 disp_bal.im_.colorbar.remove()
6 disp_bal.ax_.set_ylabel('True label', fontsize=15)
7 disp_bal.ax_.set_xlabel('')
8
9 disp_imbal.plot(ax=axes[1], xticks_rotation=45)
10 disp_imbal.ax_.set_title("Imbalanced dataset")
11 disp_imbal.ax_.set_ylabel('')
12 disp_imbal.ax_.set_xlabel('')
13
14 f.text(0.47, -0.1, 'Predicted label', fontsize=15)
15 plt.show()

```

Listing 18: Observation

