

IE406 Machine Learning

Lab Assignment - 03

Group 16

Student ID: 201801046 Name: Ravi Patel

Student ID: 201801057 Name: Ayan Khokhar

Student ID: 201801432 Name: Prakhar Maheshwari

Student ID: 201801452 Name: Jaydeep Machhi

Student ID: 201801459 Name: Aakash Panchal

Question 1

In this question we will use the data given in file “Social_Network_Ads.csv” which is a categorical data-set 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 # Scaling the dataset values from 0 to 1
2 scaler = MinMaxScaler()
3 dropped_df = df.drop(['User ID'], axis = 1)
4 column_names = dropped_df.columns
5
6 scaled_df = scaler.fit_transform(dropped_df)
7 scaled_df = pd.DataFrame(scaled_df, columns= column_names)
8 x = scaled_df[['Gender', 'Age', 'EstimatedSalary']].values
9 y = scaled_df['Purchased'].values
10 print(scaled_df.head()) # Scaled Dataset

1 # Logistic Regression implementation from scratch
2 class MyLogisticRegression:
3     def __init__(self, alpha=0.01, iterations=100000, fit_intercept=True):
4         self.alpha = alpha #Learning Rate
```

```

5         self.iterations = iterations #Finite Iterations for algorithm
convergence
6         self.fit_intercept = fit_intercept #Set if to true if intercept
needed
7
8         def _conc_intercept(self, X): # Append Intercept
9             intercept = np.ones((X.shape[0], 1))
10            return np.concatenate((intercept, X), axis=1)
11
12        def _logistic(self, z): # Logistic(Sigmoid) Function for Logistic
Regression
13            return 1 / (1 + np.exp(-z))
14
15        def _hypothesis(self, h, y): # Binary Cross Entropy Function
16            return (-y * np.log(h) - (1 - y) * np.log(1 - h)).mean()
17
18        def fit(self, X, y): # Using Gradient Descent to calculate minima of
the function
19            if self.fit_intercept:
20                X = self._conc_intercept(X)
21
22            self._features = np.zeros(X.shape[1])
23            for i in range(self.iterations):
24                z = np.dot(X, self._features)
25                h = self._logistic(z)
26                gradient = np.dot(X.T, (h - y)) / y.size
27                self._features -= self.alpha * gradient
28
29
30        def predict(self, X, threshold = 0.5): # Prediction Function if prob
>= thresh then one class else another class
31            if self.fit_intercept:
32                X = self._conc_intercept(X)
33            return self._logistic(np.dot(X, self._features)) >= threshold

1 plt.figure(figsize = (8,8))
2 ax = plt.axes(projection = '3d')
3
4 df0 = df[df.Purchased == 0]
5 df1 = df[df.Purchased == 1]
6
7 # Classifying Purchased as blue and not purchased as red data points in
form of a scatter plot
8 ax.scatter(df0['Gender'].values, df0['Age'].values, df0['EstimatedSalary'].
values, c = 'red')
9 ax.scatter(df1['Gender'].values, df1['Age'].values, df1['EstimatedSalary'].
values, c = 'blue')
10 plt.legend(['Not Purchased', 'Purchased'])
11 ax.set_title('Scatter plot of Features')
12 ax.set_xlabel('Gender')
13 ax.set_ylabel('Age')
14 ax.set_zlabel('EstimatedSalary')

1 regressor = MyLogisticRegression()

```

```

2 X = x
3 # print(X,y)
4
5 X_train, X_test, y_train, y_test = train_test_split(
6     X, y, test_size=0.1, random_state=4)
7 regressor.fit(X_train, y_train)
8
9 y_pred = regressor.predict(X_test)
10 score = accuracy_score(y_pred, y_test)
11 print(score)
12
13 confusion_mat = confusion_matrix(y_test, y_pred)
14
15 # print(y_predicted)
16 graph_cm_matrix = pd.DataFrame(confusion_mat, columns=['Actual True:1', '
    Actual False:0'],
17                                index=['Predicted Positive:1', 'Predicted
    Negative:0'])
18
19 sns.heatmap(graph_cm_matrix, annot=True, fmt='d')

1 # Function to calculate Classification metrics from the confusion matrix
2 def calculate_classification_metrics(confusion_mat):
3     TP = confusion_mat[0,0]
4     FP = confusion_mat[0,1]
5     FN = confusion_mat[1,0]
6     TN = confusion_mat[1,1]
7     accuracy = (TP + TN)/(TP + TN + FP + FN)
8     precision = (TP)/(TP + FP)
9     recall = (TP)/(TP + FN)
10    F1 = (2 * precision * recall)/(precision + recall)
11    print(f'Accuracy: {accuracy}')
12    print(f'Precision: {precision}')
13    print(f'Recall: {recall}')
14    print(f'F1 Score: {F1}')

1 #Using scikit-learn's implementation of Logistic Regression
2 X = x
3
4 X_train, X_test, y_train, y_test = train_test_split(
5     X, y, test_size=0.1, random_state=4)
6
7 clf = LogisticRegression(random_state=10, max_iter = 100000).fit(X_train,
    y_train)
8 y_predicted = clf.predict(X_test)
9 confusion_mat = confusion_matrix(y_test, y_predicted)
10
11 # print(y_predicted)
12 graph_cm_matrix = pd.DataFrame(confusion_mat, columns=['Actual True:1', '
    Actual False:0'],
13                                index=['Predicted Positive:1', 'Predicted
    Negative:0'])
14
15 sns.heatmap(graph_cm_matrix, annot=True, fmt='d')

```

Result

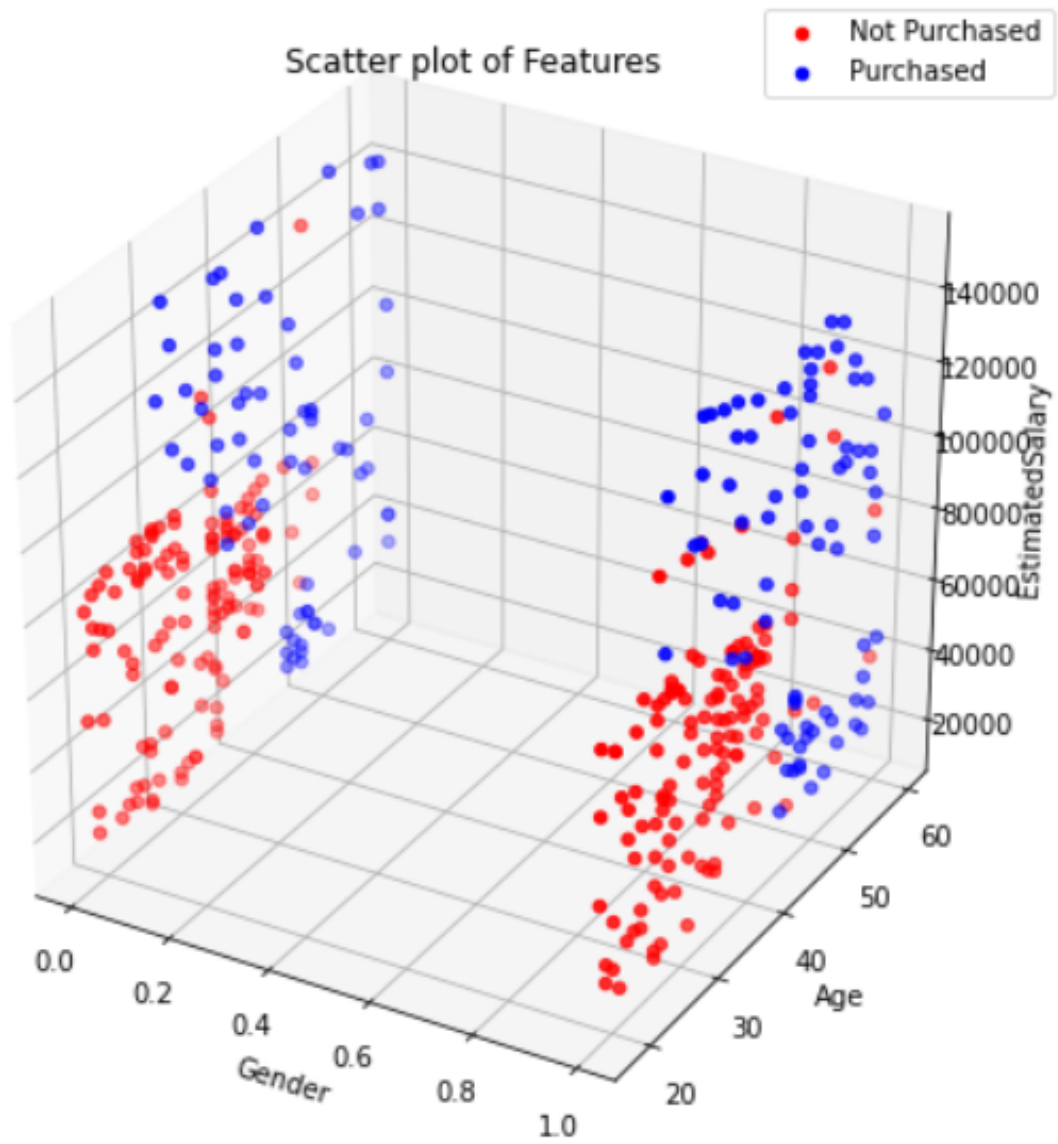


Figure 1: 3D plot for the data-set

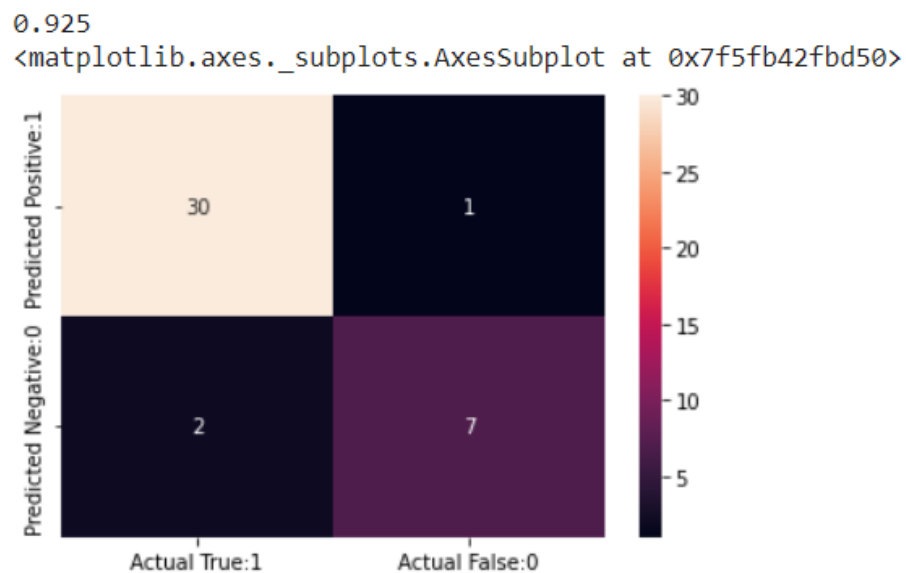


Figure 2: Confusion matrix for our own logistic function

Accuracy: 0.925
Precision: 0.967741935483871
Recall: 0.9375
F1 Score: 0.9523809523809523

Figure 3: Accuracy, precision, recall and F1 Score for our own logistic function

<matplotlib.axes._subplots.AxesSubplot at 0x7f5fb389b090>



Figure 4: Confusion matrix for ScikitLearn's regression function

Accuracy: 0.95
Precision: 1.0
Recall: 0.9393939393939394
F1 Score: 0.96875

Figure 5: Accuracy, precision, recall and F1 Score for ScikitLearn's function

Observation/ Justification

- We have scaled the values of the dataset from 0 to 1 using the MinMaxScaler function of sci-kit-learn. To improve the classification performance.
- By changing the learning rate(alpha) to a smaller value and iterations to a larger number can increase the accuracy of our algorithm.
- Moreover, gradient descent is not the most optimal method for calculating the minima. By making use of more sophisticated algorithms we can increase the performance of our implementation.
- From the above analysis of the confusion matrix of our own implementation of Logistic Regression with that of the Scikit Learn's Implementation. It can be clearly seen that both perform pretty well in terms of binary classification of Purchased/Not Purchased.

- Moreover from the analysis of classification metrics, our implementation performs pretty well in terms of accuracy, precision, recall and F1 score.

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 iris_df = pd.read_csv('/content/Iris.csv')
2 print(iris_df.head())
3
4 print(iris_df.shape)
```

Listing 1: Loading data

```
1 #Plotting probability density function and histogram for each feature
2 features = iris_df.columns
3
4 for i in range(1,5):
5     grid = sns.FacetGrid(iris_df, hue='Species', height=5)
6     grid1 = sns.FacetGrid(iris_df, hue='Species', height=5)
7     grid.map(sns.kdeplot, features[i]).add_legend();
8     grid1.map(sns.histplot, features[i]).add_legend();
```

Listing 2: Plotting probability density function and histogram for each feature

```
1 # Mapping Species name to integer for prediction
2 species = {'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-virginica' : 2}
3 iris_df['Species'] = [species[item] for item in iris_df['Species']]
4
5 # Performing Exploratory Data Analysis on the Dataset
6 df_dict = {g: d for g, d in iris_df.groupby('Species')}
7
8 SepalL = []
9 SepalW = []
10 PetalL = []
11 PetalW = []
12
13 x = list(species.keys())
14 for i in range(3):
15     SepalL.append(df_dict[i].mean()['SepalLengthCm'])
16     SepalW.append(df_dict[i].mean()['SepalWidthCm'])
17     PetalL.append(df_dict[i].mean()['PetalLengthCm'])
```

```

18     PetalW.append(df_dict[i].mean()['PetalWidthCm'])
19
20
21 print(SepalL)
22 fig = plt.figure()
23 fig.set_figwidth(10)
24 fig.set_figheight(8)
25 plt.subplot(2,2,1)
26 plt.bar(x, SepalL)
27 plt.title('Mean of SepalLength')
28 plt.subplot(2,2,2)
29
30 plt.bar(x, SepalW)
31
32 plt.title('Mean of SepalWidth')
33 plt.subplot(2,2,3)
34
35 plt.bar(x, PetalL)
36
37 plt.title('Mean of PetalLength')
38 plt.subplot(2,2,4)
39 plt.bar(x, SepalW)
40
41 plt.title('Mean of PetalWidth')

```

Listing 3: Mean for each features

```

1
2 SepalL = []
3 SepalW = []
4 PetalL = []
5 PetalW = []
6 for i in range(3):
7     SepalL.append(df_dict[i].median()['SepalLengthCm'])
8     SepalW.append(df_dict[i].median()['SepalWidthCm'])
9     PetalL.append(df_dict[i].median()['PetalLengthCm'])
10    PetalW.append(df_dict[i].median()['PetalWidthCm'])
11
12 print(SepalL)
13 fig = plt.figure()
14 fig.set_figwidth(10)
15 fig.set_figheight(8)
16 plt.subplot(2,2,1)
17 plt.bar(x, SepalL)
18 plt.title('Median of SepalLength')
19
20 plt.subplot(2,2,2)
21 plt.bar(x, SepalW)
22 plt.title('Median of SepalWidth')
23
24 plt.subplot(2,2,3)
25 plt.bar(x, PetalL)
26 plt.title('Median of PetalLength')
27

```



```

28 plt.subplot(2,2,4)
29 plt.bar(x, SepalW)
30 plt.title('Median of PetalWidth')

```

Listing 4: Median for each features

```

1
2
3 SepalL = []
4 SepalW = []
5 PetalL = []
6 PetalW = []
7 for i in range(3):
8     SepalL.append(df_dict[i].min()['SepalLengthCm'])
9     SepalW.append(df_dict[i].min()['SepalWidthCm'])
10    PetalL.append(df_dict[i].min()['PetalLengthCm'])
11    PetalW.append(df_dict[i].min()['PetalWidthCm'])
12
13 print(SepalL)
14 fig = plt.figure()
15 fig.set_figwidth(10)
16 fig.set_figheight(8)
17 plt.subplot(2,2,1)
18 plt.bar(x, SepalL)
19 plt.title('Minimum of SepalLength')
20
21 plt.subplot(2,2,2)
22 plt.bar(x, SepalW)
23 plt.title('Minimum of SepalWidth')
24
25 plt.subplot(2,2,3)
26 plt.bar(x, PetalL)
27 plt.title('Minimum of PetalLength')
28
29 plt.subplot(2,2,4)
30 plt.bar(x, SepalW)
31 plt.title('Minimum of PetalWidth')

```

Listing 5: Min value for each features

```

1
2
3 SepalL = []
4 SepalW = []
5 PetalL = []
6 PetalW = []
7 for i in range(3):
8     SepalL.append(df_dict[i].max()['SepalLengthCm'])
9     SepalW.append(df_dict[i].max()['SepalWidthCm'])
10    PetalL.append(df_dict[i].max()['PetalLengthCm'])
11    PetalW.append(df_dict[i].max()['PetalWidthCm'])
12
13 print(SepalL)
14 fig = plt.figure()

```

```

15 fig.set_figwidth(10)
16 fig.set_figheight(8)
17 plt.subplot(2,2,1)
18 plt.bar(x, SepalL)
19 plt.title('Maximum of SepalLength')
20
21 plt.subplot(2,2,2)
22 plt.bar(x, SepalW)
23 plt.title('Maximum of SepalWidth')
24
25 plt.subplot(2,2,3)
26 plt.bar(x, PetalL)
27 plt.title('Maximum of PetalLength')
28
29 plt.subplot(2,2,4)
30 plt.bar(x, SepalW)
31 plt.title('Maximum of PetalWidth')

```

Listing 6: Max value for each features

Result

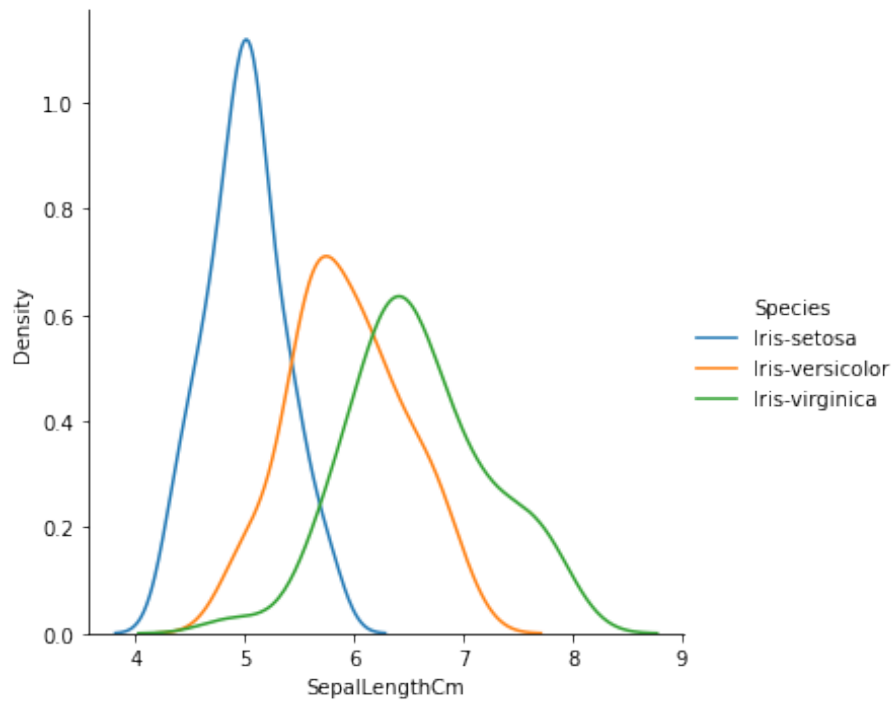


Figure 6: PDF of Sepal Length

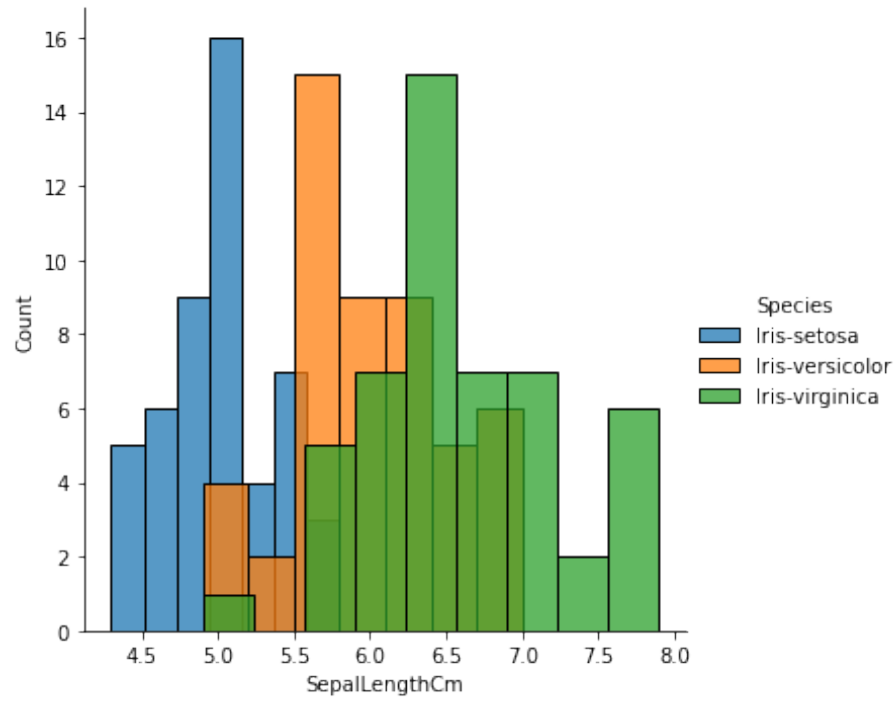


Figure 7: Histogram for Sepal Length

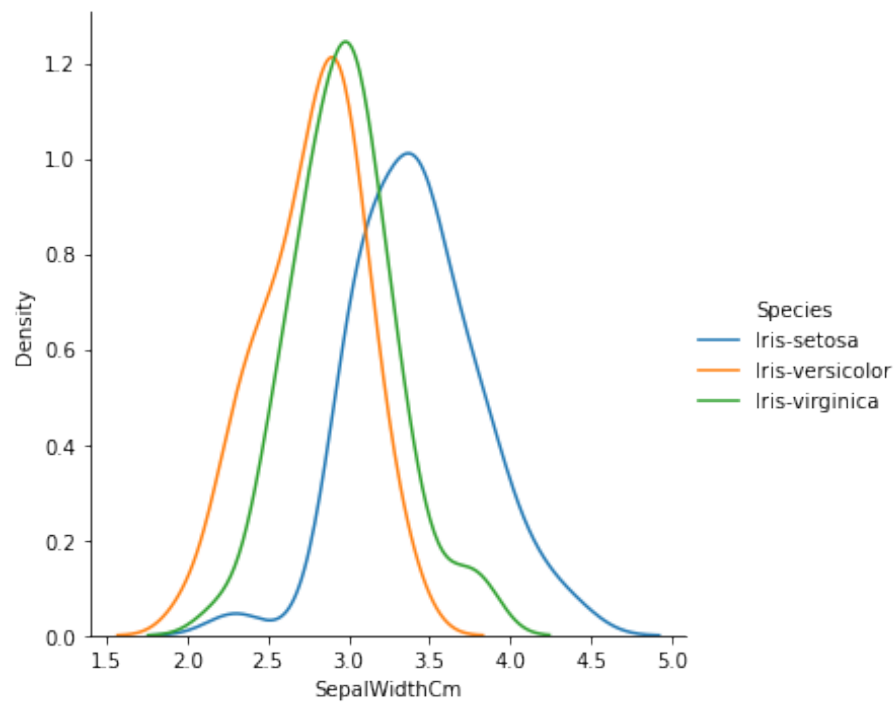


Figure 8: PDF for Sepal Width

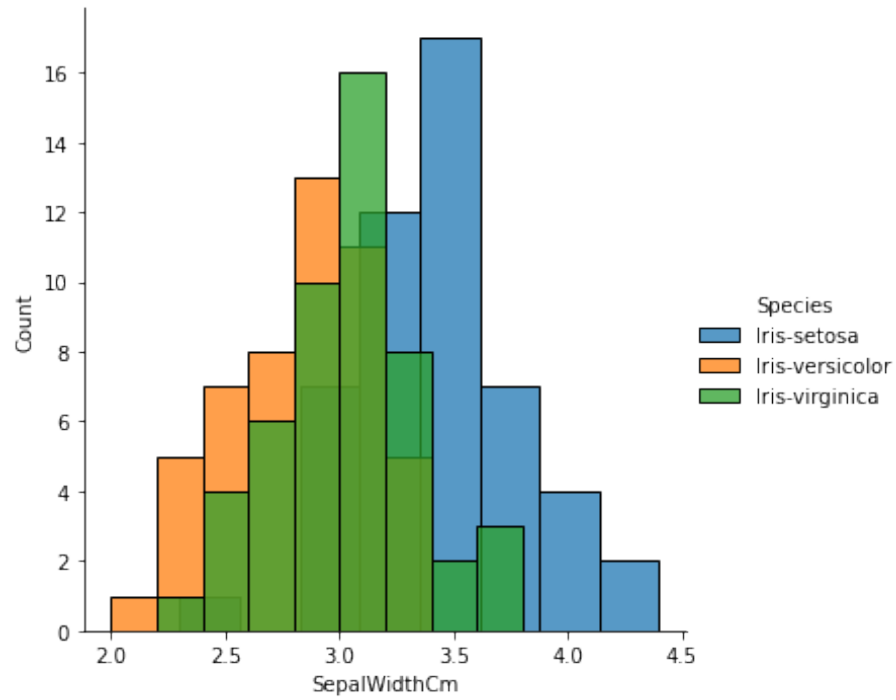


Figure 9: Histogram for Sepal Width

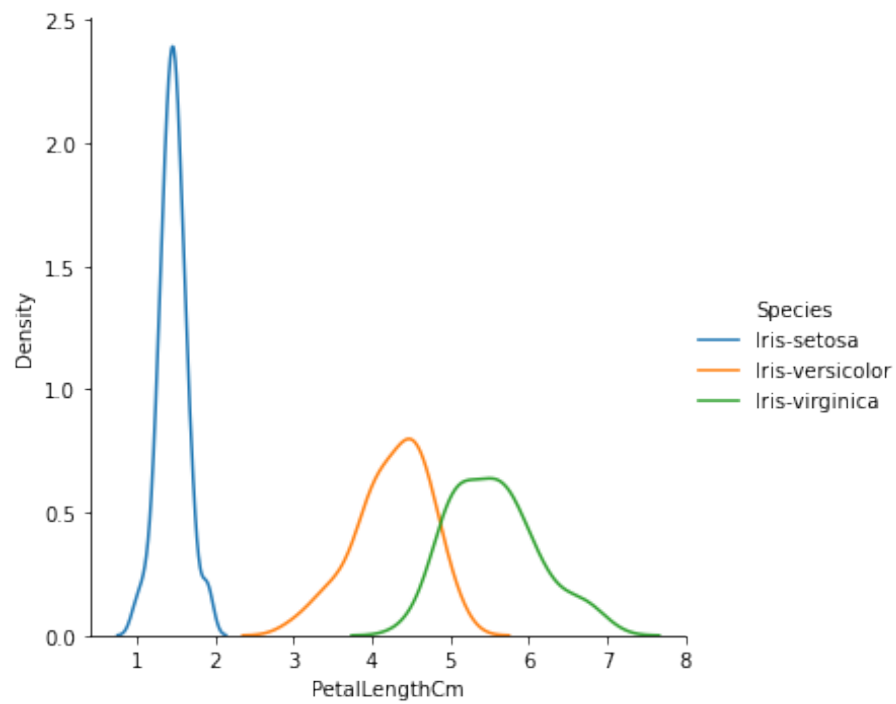


Figure 10: PDF for Petal Length

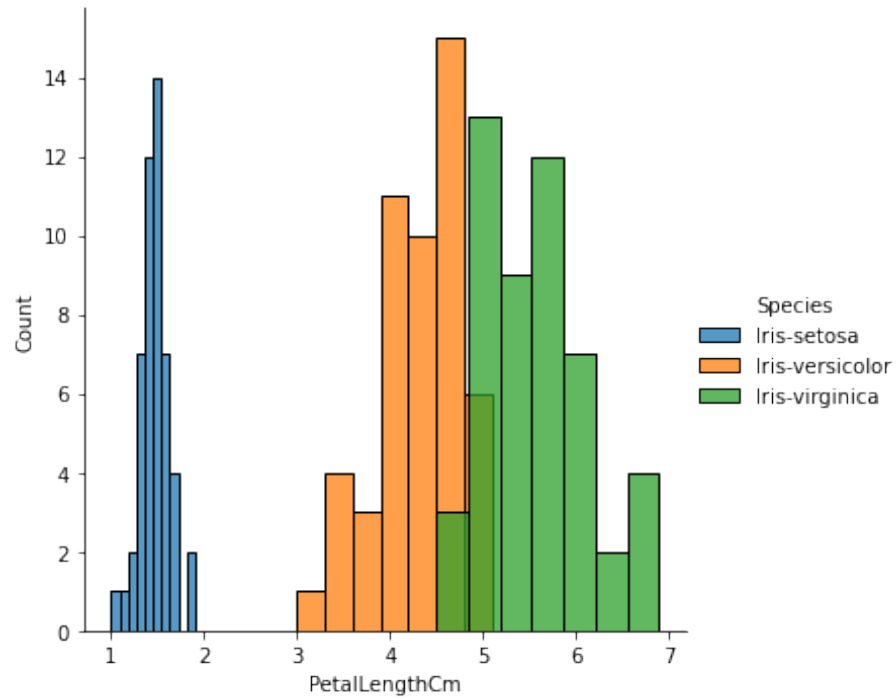


Figure 11: Histogram for Petal Length

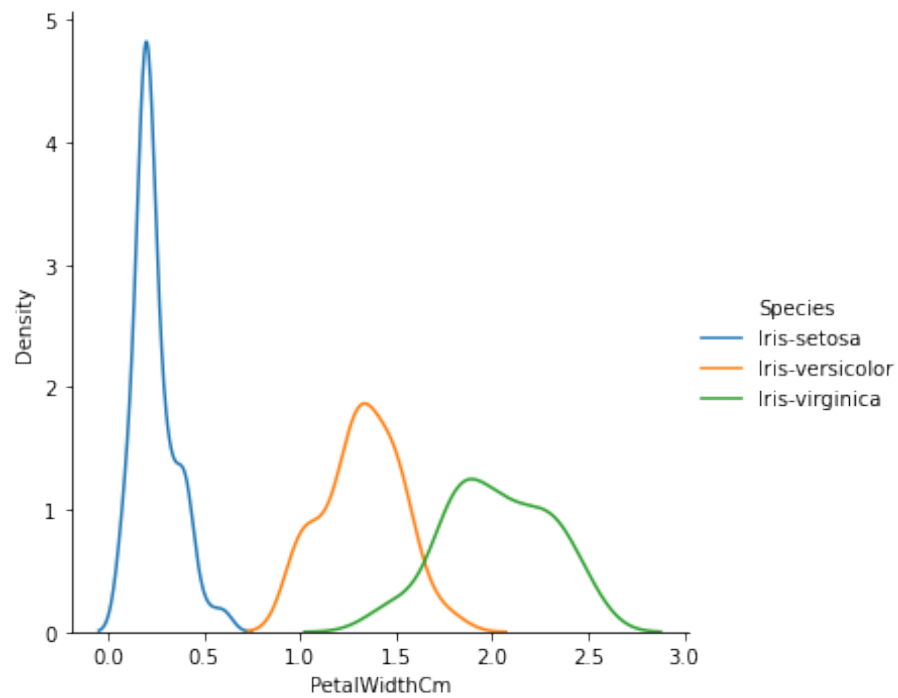


Figure 12: Probability for Petal Width

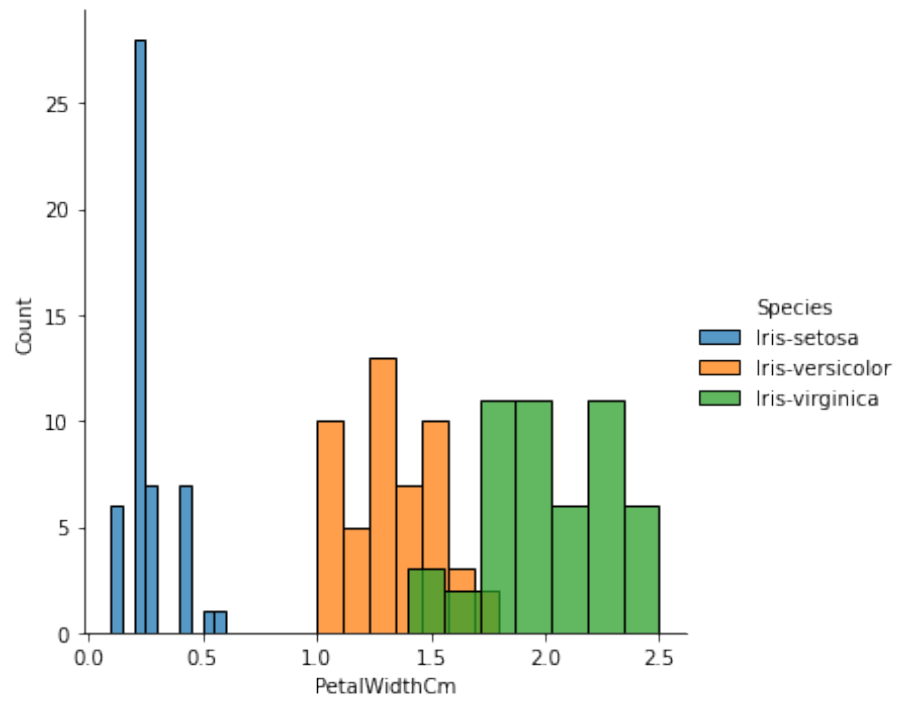


Figure 13: Histogram for Petal Width

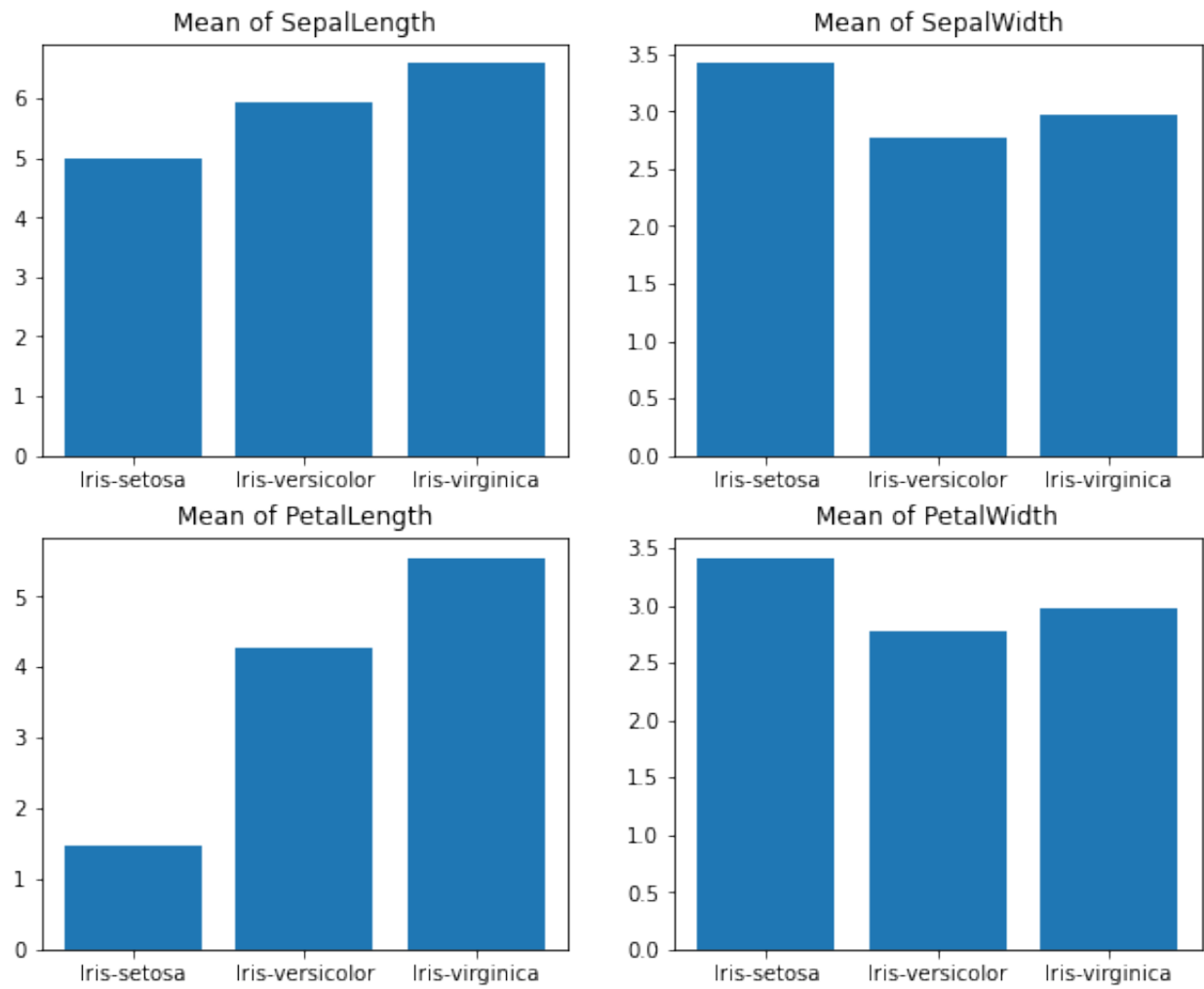


Figure 14: Mean of features

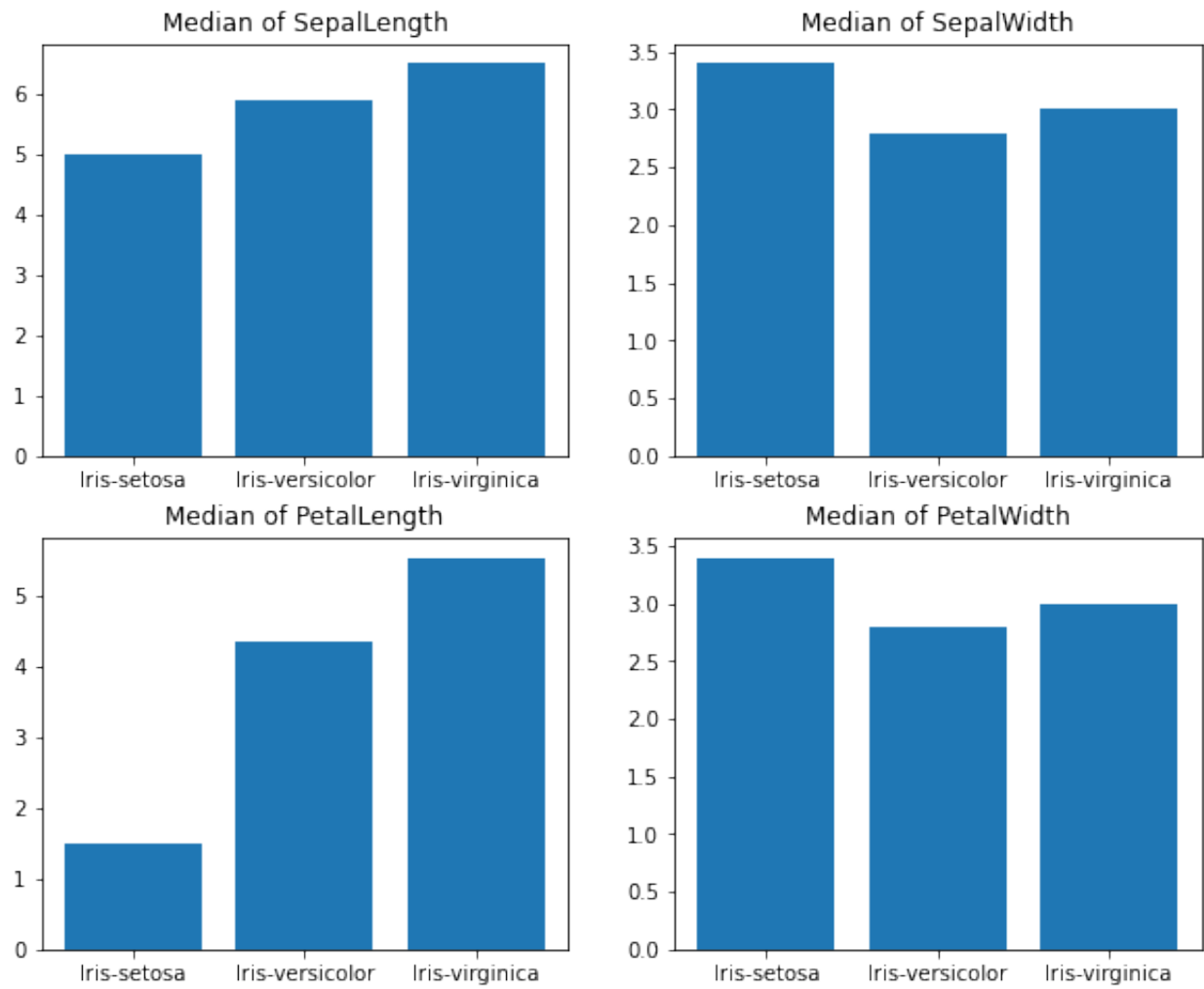


Figure 15: Median of features

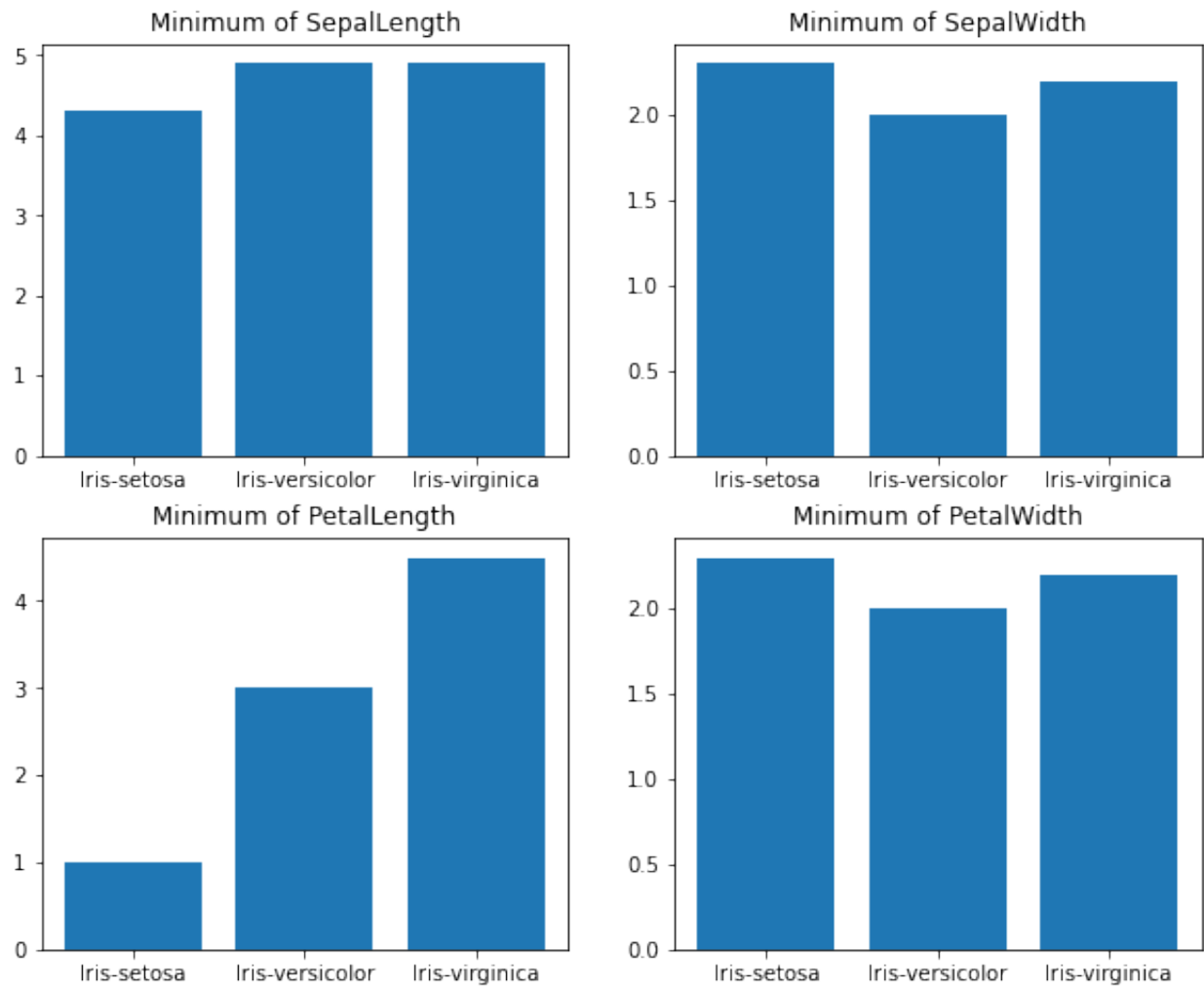


Figure 16: Min for features

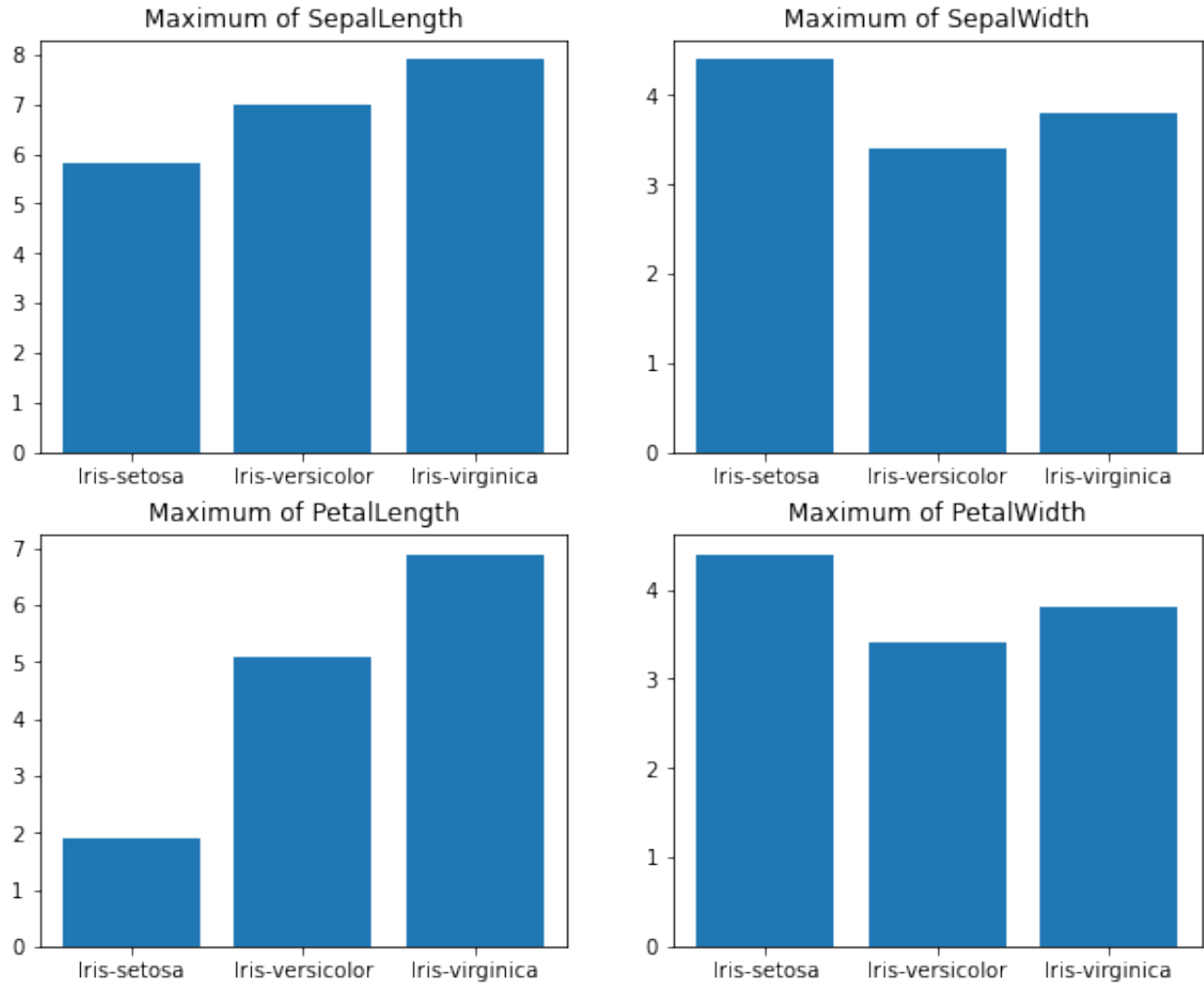


Figure 17: Max for features

Observation/ Justification

PDF for each of the features follows a Gaussian distribution. Looking at the PDF we can get a high level overview of range of each feature depending on the species.

Histogram also serves the same purpose as the PDF.

Below is an ordering based on decreasing order of mean for each features.

Sepal Length: Iris virginica, Iris versicolor, Iris setosa
 Sepal Width: Iris setosa, Iris versicolor, Iris virginica
 Petal Length: Iris virginica, Iris versicolor, Iris setosa
 Petal Width: Iris setosa, Iris virginica, Iris versicolor

Below is an ordering based on decreasing order of median for each features.

Sepal Length: Iris virginica, Iris versicolor, Iris setosa
 Sepal Width: Iris setosa, Iris virginica, Iris versicolor

Petal Length: Iris virginica, Iris versicolor, Iris setosa
Petal Width: Iris setosa, Iris virginica, Iris versicolor

Below is an ordering based on decreasing order of min for each features.

Sepal Length: Iris virginica, Iris versicolor, Iris setosa
Sepal Width: Iris setosa, Iris virginica , Iris versicolor
Petal Length: Iris virginica, Iris versicolor, Iris setosa
Petal Width: Iris setosa, Iris virginica, Iris versicolor

Below is an ordering based on decreasing order of max for each features.

Sepal Length: Iris virginica, Iris versicolor, Iris setosa
Sepal Width: Iris setosa, Iris virginica, Iris versicolor
Petal Length: Iris virginica, Iris versicolor, Iris setosa
Petal Width: Iris setosa, Iris virginica, Iris versicolor

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 df = pd.read_csv('/content/Iris.csv')
2
3 species = {'Iris-setosa' : 0, 'Iris-versicolor' : 1, 'Iris-virginica' : 2}
4 df['Species'] = [species[item] for item in df['Species']]
5 print(df.head())
6
7 color_map = {0 : 'r', 1 : 'b', 2 : 'g'}
8 sps = df['Species'].unique()
9 rev_species = { 0 : 'Iris-setosa', 1 : 'Iris-versicolor', 2 : 'Iris-
    virginica'}
10 marker_map = {0:'o', 1:'x', 2:'v'}
11
12 feature = df.columns[1:-1]
13 plt.figure(figsize = (14,14))
14 print(feature)
15 print(sps)
16
17 p = 1
18 for i in range(4):
19     for j in range(i+1, 4):
20
21         plt.subplot(3, 2, p)
```

```

22     for k in sps:
23
24         plt.scatter(df[df['Species'] == k][feature[i]], df[df['Species
25         ''] == k][feature[j]], c = color_map[k], marker = marker_map[k], label =
            rev_species[k]) # Use this if you want to use markers
26         # plt.scatter(df[df['Species'] == k][feature[i]], df[df['
27         Species'] == k][feature[j]], c = color_map[k], label = k) # Use this if
28         don't want to use markers
29         plt.xlabel(feature[i])
30         plt.ylabel(feature[j])
31         plt.legend()
32
33     p += 1
34
35 plt.show()

```

Listing 7: Combination of two features and their separation

Result

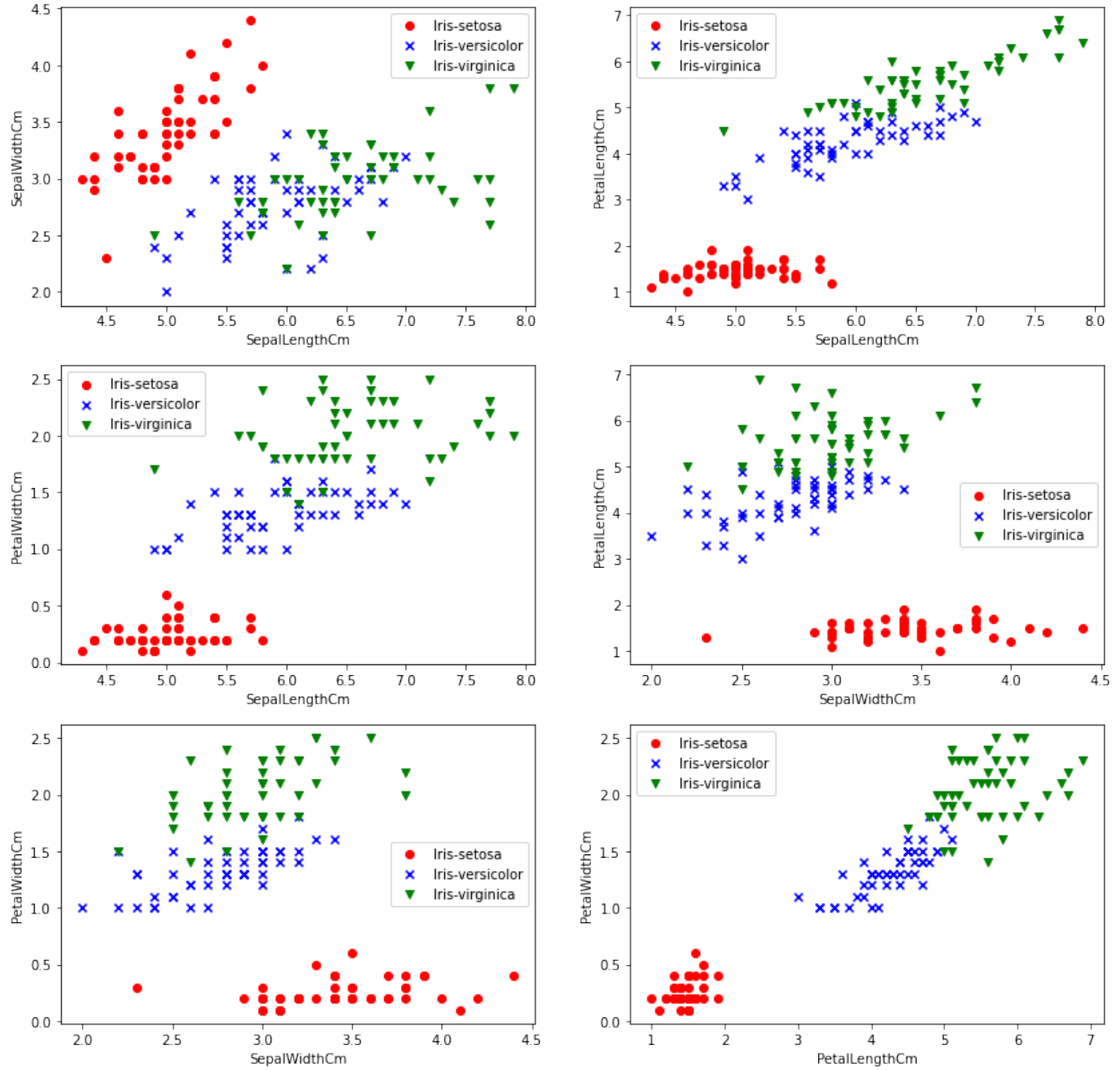


Figure 18: Combination of two features and their separation

Observation/ Justification

We can clearly see from the figure that using one vs all method for separation, we can easily separate Iris Setosa from the other two. Infact, Iris virginica can also be separated from other two but with some error in separation. But, it is almost impossible to separate Iris - versicolor from other two using line as a boundry. For this we need some complex boundaries(i.e. circle).

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 iris_df = pd.read_csv('/content/Iris.csv')
2 print(iris_df.head())
3 species = {'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-virginica' : 2}
4 iris_df['Species'] = [species[item] for item in iris_df['Species']]
5 dropped_df = iris_df.drop(['Id'], axis=1)
6 column_names = dropped_df.columns
7 scaled_df = scaler.fit_transform(dropped_df)
8 scaled_df = pd.DataFrame(scaled_df, columns =column_names)
```

Listing 8: Q4: Loading data

```
1 species = {0: 0, 0.5: 1, 1 : 2}
2 scaled_df['Species'] = [species[item] for item in scaled_df['Species']]
3 X = scaled_df[['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', '
    PetalWidthCm']]
4 y = scaled_df['Species']
```

Listing 9: Q4: Extracting useful data

```
1 X_train, X_test, y_train, y_test = train_test_split(
2     X, y, test_size=0.2, random_state=2)
3
4 clf = LogisticRegression(random_state=10, n_jobs=1).fit(X_train, y_train)
5 y_predicted = clf.predict(X_test)
6 confusion_mat = confusion_matrix(y_test, y_predicted)
7
8 # print(y_predicted)
9 graph_cm_matrix = pd.DataFrame(confusion_mat, columns=['Iris-setosa: 0', '
    Iris-versicolor: 1', 'Iris-virginica : 2'],
10                                index=['Iris-setosa: 0', 'Iris-versicolor:
    1', 'Iris-virginica : 2'])
11
12 sns.heatmap(graph_cm_matrix, annot=True, fmt='d')
```

Listing 10: Q4: Confusion Matrix

```
1 print('Macro')
2 # Classification Report
3 print(classification_report(y_test, y_predicted))
4
5 macro = precision_recall_fscore_support(y_test, y_predicted, average='
    macro')
6 print(f'Precision = {macro[0]}\nRecall = {macro[1]}\nF1Score = {macro[2]}'
7     )
```

```

8 print('\nMicro')
9 micro = precision_recall_fscore_support(y_test, y_predicted, average='
    micro')
10 print(f'Precision = {micro[0]}\nRecall = {micro[1]}\nF1Score = {micro[2]}',
    )

```

Listing 11: Q4: accuracy

Result

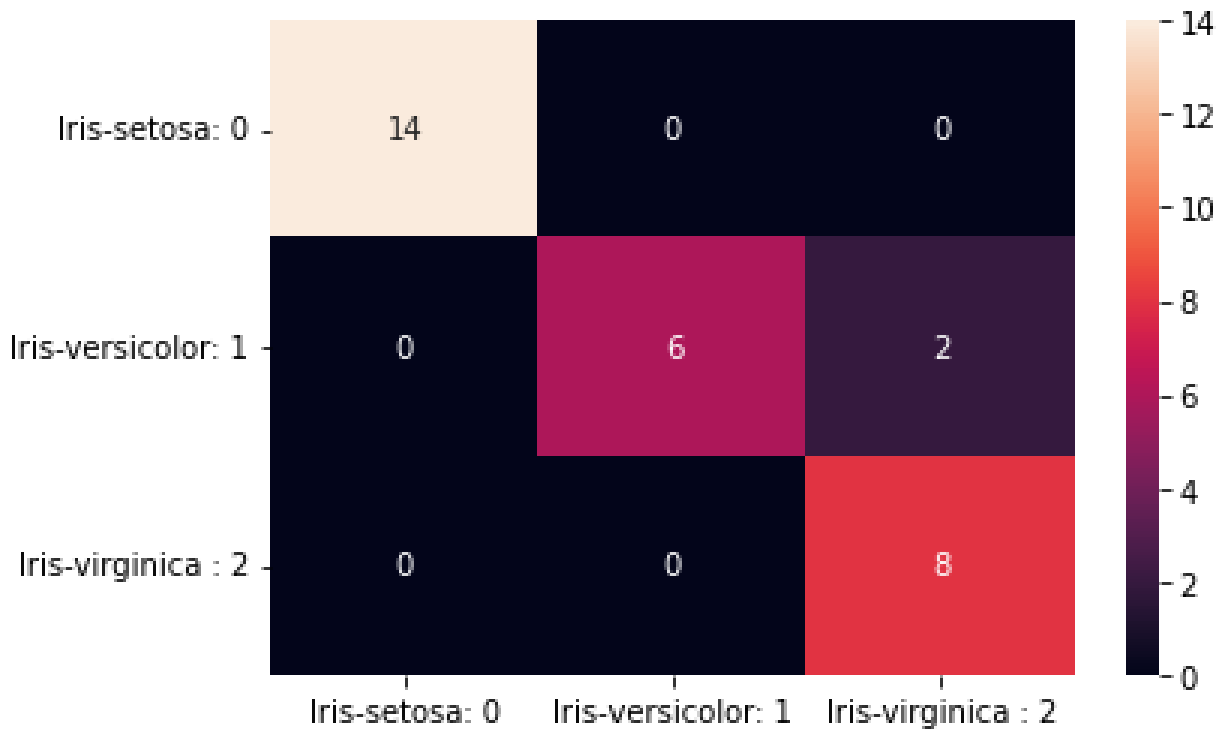


Figure 19: Confusion Matrix

Result

	precision	recall	f1-score	support
0	1.00	1.00	1.00	14
1	1.00	0.75	0.86	8
2	0.80	1.00	0.89	8
accuracy			0.93	30
macro avg	0.93	0.92	0.92	30
weighted avg	0.95	0.93	0.93	30

Figure 20: Classification report

Result

```
Macro
Precision = 0.9333333333333332
Recall = 0.9166666666666666
F1Score = 0.9153439153439153

Micro
Precision = 0.9333333333333333
Recall = 0.9333333333333333
F1Score = 0.9333333333333333

Weighted
Precision = 0.9466666666666667
Recall = 0.9333333333333333
F1Score = 0.9322751322751323
```

Figure 21: Macro, Micro and Weighted parameters

Observation/ Justification

The standard classification metrics help us to find out how accurate our model is and how much better predictions it can make for the different ranges of the values provided.

Confusion matrix has almost zero entries in cells which are not a main diagonal. This shown that our model is accurate. non zero entry of 2 in a cell (1,2) indicates that species should be Iris-verginica but it is predicted as Iris-versicolor.

Almost every parameters(accuracy, precision, recall and F1 score) is close to 1.00 which shows the correctness of our model.

Question 5

Imbalanced dataset typically refers to a dataset where the classes are not represented equally. Classification problems having multiple classes with imbalanced dataset present a different challenge than a binary classification problem. The skewed distribution makes the machine learning algorithms less effective, especially in predicting minority class examples.

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 you observations at the end.

Answer

We have taken Iris dataset. First of all, we will split the dataset based on (same) species, and then with the help of those chunks we can pre-process the data to make unbalanced or balanced versions of the dataset.

```
1 iris_df = pd.read_csv('/content/Iris.csv')
2
3 # Picking Unique Species among the Species column
4 spec = iris_df['Species'].unique()
5
6 # Splitting the dataframe into small dataframes having data from the same
   species
7 iris_df0 = iris_df[iris_df['Species'] == spec[0]]
8 iris_df1 = iris_df[iris_df['Species'] == spec[1]]
9 iris_df2 = iris_df[iris_df['Species'] == spec[2]]
```

Balanced Version

code

```
1 # Creating Balanced Dataset from the Given Iris dataset by selecting a
   fraction of the rows from the splitted dataframes
2 new_subset_iris_df0 = iris_df0.sample(frac = 0.8)
3 new_subset_iris_df1 = iris_df1.sample(frac = 0.8)
4 new_subset_iris_df2 = iris_df2.sample(frac = 0.8)
5
6 # Merging the Dataframes back to generate a Balanced Dataframe
7 new_subset = new_subset_iris_df0.append(new_subset_iris_df1).append(
   new_subset_iris_df2)
8 cols = new_subset.columns[1:-1]
9 species = {'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-virginica' : 2}
10 new_subset['Species'] = [species[item] for item in new_subset['Species']]
11 X = new_subset[cols]
12 y = new_subset[new_subset.columns[-1]]
13
14 # Performing Logistic Regression and predicting values
15 clf = LogisticRegression(max_iter = 10000)
16 X_train, X_test, y_train, y_test = train_test_split(
17     X, y, test_size = 0.2, random_state = 4)
18 clf.fit(X_train, y_train)
19 pred_y = clf.predict(X_test)
20
21 # Generating confusion matrix and its heat-map
22 conf_matrix = confusion_matrix(y_test, pred_y)
23 graph_cm_matrix = pd.DataFrame(conf_matrix, columns=['Iris-setosa: 0', '
   Iris-versicolor: 1', 'Iris-virginica : 2'],
24                                index=['Iris-setosa: 0', 'Iris-versicolor:
   1', 'Iris-virginica : 2'])
25 sns.heatmap(graph_cm_matrix, annot = True, fmt = 'd')
26
27
28 # Classification Report
```

```

29 print(classification_report(y_test, pred_y))
30
31 print('Macro')
32 macro = precision_recall_fscore_support(y_test, pred_y, average='macro')
33 print(f'Precision = {macro[0]}\nRecall = {macro[1]}\nF1Score = {macro[2]}',
34       )
35
36 print('\nMicro')
37 micro = precision_recall_fscore_support(y_test, pred_y, average='micro')
38 print(f'Precision = {micro[0]}\nRecall = {micro[1]}\nF1Score = {micro[2]}',
39       )
40
41 print('\nWeighted')
42 weighted = precision_recall_fscore_support(y_test, pred_y, average='
    weighted')
43 print(f'Precision = {weighted[0]}\nRecall = {weighted[1]}\nF1Score = {
    weighted[2]}')

```

Result

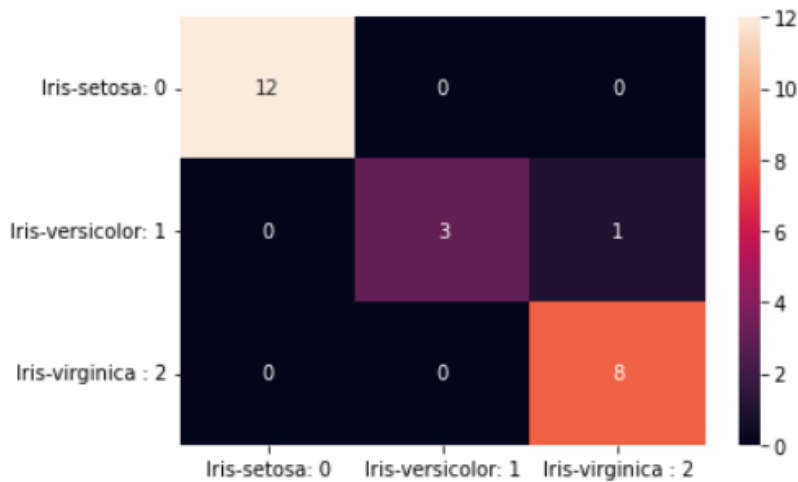


Figure 22: Heat-map of confusion matrix

	precision	recall	f1-score	support
0	1.00	1.00	1.00	12
1	1.00	0.75	0.86	4
2	0.89	1.00	0.94	8
accuracy			0.96	24
macro avg	0.96	0.92	0.93	24
weighted avg	0.96	0.96	0.96	24

Macro
Precision = 0.9629629629629629
Recall = 0.9166666666666666
F1Score = 0.9327731092436974

Micro
Precision = 0.9583333333333334
Recall = 0.9583333333333334
F1Score = 0.9583333333333334

Weighted
Precision = 0.9629629629629629
Recall = 0.9583333333333334
F1Score = 0.9565826330532213

Figure 23: Results

Unbalanced Version

code

```

1 # Creating Imbalanced Dataset from the Given Iris dataset by selecting a
  fraction of the rows from the splitted dataframes
2 subset_iris_df0 = iris_df0.sample(frac = 0.5)
3 subset_iris_df1 = iris_df1.sample(frac = 0.9)
4 subset_iris_df2 = iris_df2.sample(frac = 0.3)
5
6 # Merging the Dataframes back to generate the Imbalanced Dataframe
7 new_subset = subset_iris_df0.append(subset_iris_df1).append(
  subset_iris_df2)
8 cols = new_subset.columns[1:-1]
9 species = {'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-virginica' : 2}
10 new_subset['Species'] = [species[item] for item in new_subset['Species']]
11 X = new_subset[cols]
12 y = new_subset[new_subset.columns[-1]]
13
14 # Performing Logistic Regression and predicting values
15 clf = LogisticRegression(max_iter = 10000)
16 X_train, X_test, y_train, y_test = train_test_split(
17     X, y, test_size = 0.2, random_state = 4)
18 clf.fit(X_train, y_train)
19 pred_y = clf.predict(X_test)

```

```

20
21 # Generating confusion matrix and its heat-map
22 conf_matrix = confusion_matrix(y_test, pred_y)
23 graph_cm_matrix = pd.DataFrame(conf_matrix, columns=['Iris-setosa: 0',
24     Iris-versicolor: 1', 'Iris-virginica : 2'],
25     index=['Iris-setosa: 0', 'Iris-versicolor:
26     1', 'Iris-virginica : 2'])
27 sns.heatmap(graph_cm_matrix, annot = True, fmt = 'd')
28
29 # Classification Report
30 print(classification_report(y_test, pred_y))
31
32 print('Macro')
33 macro = precision_recall_fscore_support(y_test, pred_y, average='macro')
34 print(f'Precision = {macro[0]}\nRecall = {macro[1]}\nF1Score = {macro[2]}')
35
36 print('\nMicro')
37 micro = precision_recall_fscore_support(y_test, pred_y, average='micro')
38 print(f'Precision = {micro[0]}\nRecall = {micro[1]}\nF1Score = {micro[2]}')
39
40 print('\nWeighted')
41 weighted = precision_recall_fscore_support(y_test, pred_y, average='
42     weighted')
43 print(f'Precision = {weighted[0]}\nRecall = {weighted[1]}\nF1Score = {
44     weighted[2]}')

```

Result

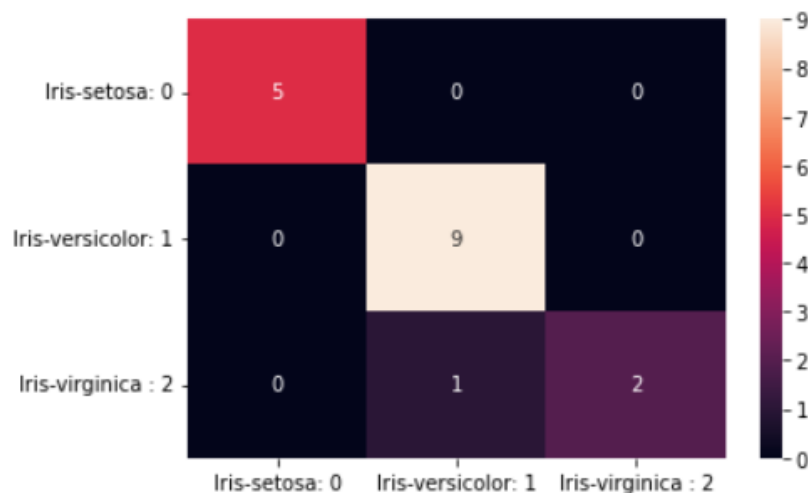


Figure 24: Heat-map of confusion matrix

	precision	recall	f1-score	support
0	1.00	1.00	1.00	5
1	0.90	1.00	0.95	9
2	1.00	0.67	0.80	3
accuracy			0.94	17
macro avg	0.97	0.89	0.92	17
weighted avg	0.95	0.94	0.94	17
Macro				
Precision =	0.9666666666666667			
Recall =	0.8888888888888888			
F1Score =	0.9157894736842106			
Micro				
Precision =	0.9411764705882353			
Recall =	0.9411764705882353			
F1Score =	0.9411764705882353			
Weighted				
Precision =	0.9470588235294118			
Recall =	0.9411764705882353			
F1Score =	0.9368421052631579			

Figure 25: Results

Observation/Justification

1. The standard classification metrics help us to find out how accurate our model is and how much better predictions it can make for the different ranges of the values provided.
2. But in case of imbalance or skew in the class distribution, the standard classification metrics does not perform pretty well as the dataset becomes biased towards the features having larger samples. And hence, the model ends up giving bad accuracy towards minority classes.
3. The basic problem with these metrics is, they assume that the skewness or the bias that is present in the training dataset will be present throughout the dataset. Thus they become inaccurate in classifying the values.