**ANTIM**     **2021ume0204@iitjammu.ac.in**     **IIT JAMMU**

**LINK FOR THE COOLAB FILE : https://colab.research.google.com/drive/1g5v8JE_pGZlTp0-zruGiPZPOoqTn93eX?usp=sharing**

# <u>Machine Learning Intern Assessment Assignment</u>

## Customer Churn Prediction

**TASKS:**

- Analysed the given customer dataset.
- Made the proper visualization of the dataset using pie chart, bar graphs, histograms, curves etc
- DATA Preprocessing
- Developed a machine learning model to predict customer churn based on historical customer data.
- Model Evaluation And Prediction

**Data**:  In this assignment we were provided with a dataset containing historical customer information, including customer attributes

- Customer ID ,
-  Name , Age ,
-  Gender ,
- Location,
-  Subscription_Length_Months ,
- Monthly_Bill ,
- Total_Usage_GB ,
- Churn.

**Description of Work:**

1: **DATA LOADING:** First we converted the xlsx format data into csv format. Then converted the dataset into pandas dataframe. We got the understanding of the data using various python libraries pandas , numpy , missingno , seaborn , plotly , sklearn , xgboost , catboost.

All the dependencies are mentioned and imported under the Dependencies Section.

2: **DATA VISUALIZATION**: The basic information about the dataset was retrieved by printing shape of the data, attributes , values of the attributes and many statistical measures were analysed

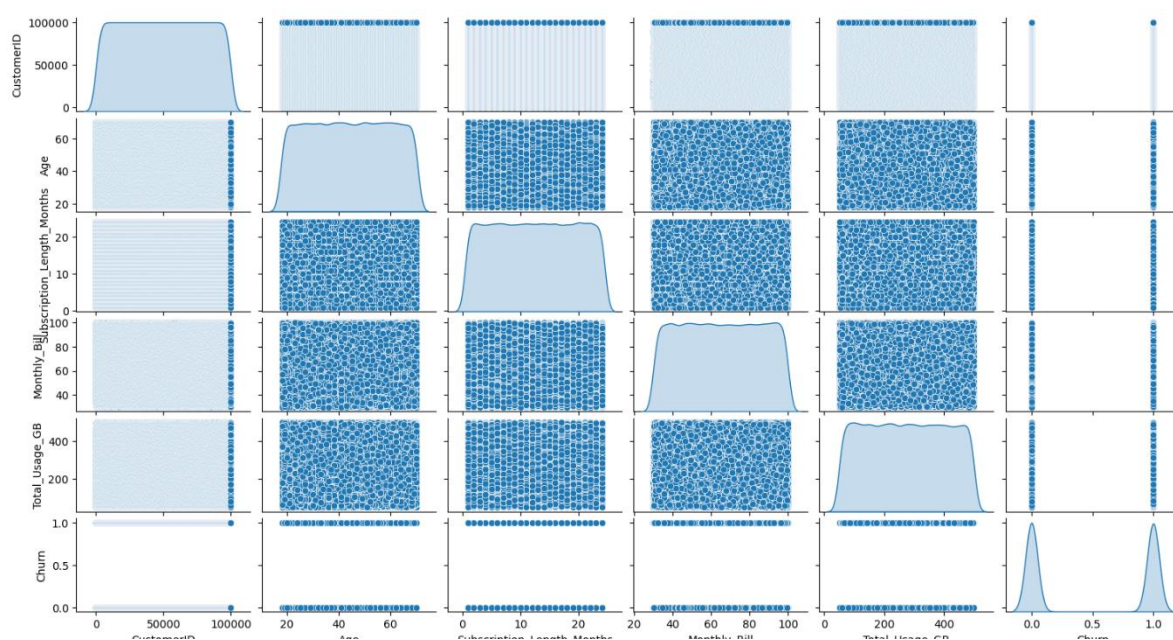All the basic visualization of the data is performed in the first cell of the VISUALIZING THE DATA section.

**ANTIM**          [2021ume0204@iitjammu.ac.in](mailto:2021ume0204@iitjammu.ac.in)          **IIT JAMMU**
**LINK FOR THE COOLAB FILE :** [https://colab.research.google.com/drive/1g5v8JE_pGZlTp0-zruGiPZPOoqTn93eX?usp=sharing](https://colab.research.google.com/drive/1g5v8JE_pGZlTp0-zruGiPZPOoqTn93eX?usp=sharing)

For more Illustration and Analysis many plots are also provided in the subsequent cells of this section.

It includes Pi chart, Bar graph, Histograms, Curves plotted between many pair of attributes in the Customer Data set.

A Correlation plot between various Attributes was also analysed and I found out that no Two different attributes were following any standard relation i.e. No correlation was found between the Attributes.

Here is the Figure which I analysed.



As we can observe there is not any nice curve between the pair of attributes except the diagonal elements of the plot.

Which is also affecting the accuracy of the model.

3: **DATA PREPROCESSING**:

Under this Section I performed the following Tasks:

- Wrote a function to convert the column with dtype='object' to int using the LabelEncoder() function.
- Wrote a function to Split the Input data into Training and Test Data with Test size of 30% and applying the Standard Scaler Transform .
- Wrote a function to displaying  the Distribution plot of the features.

4: **MODEL TRAINING AND EVALUATION:**

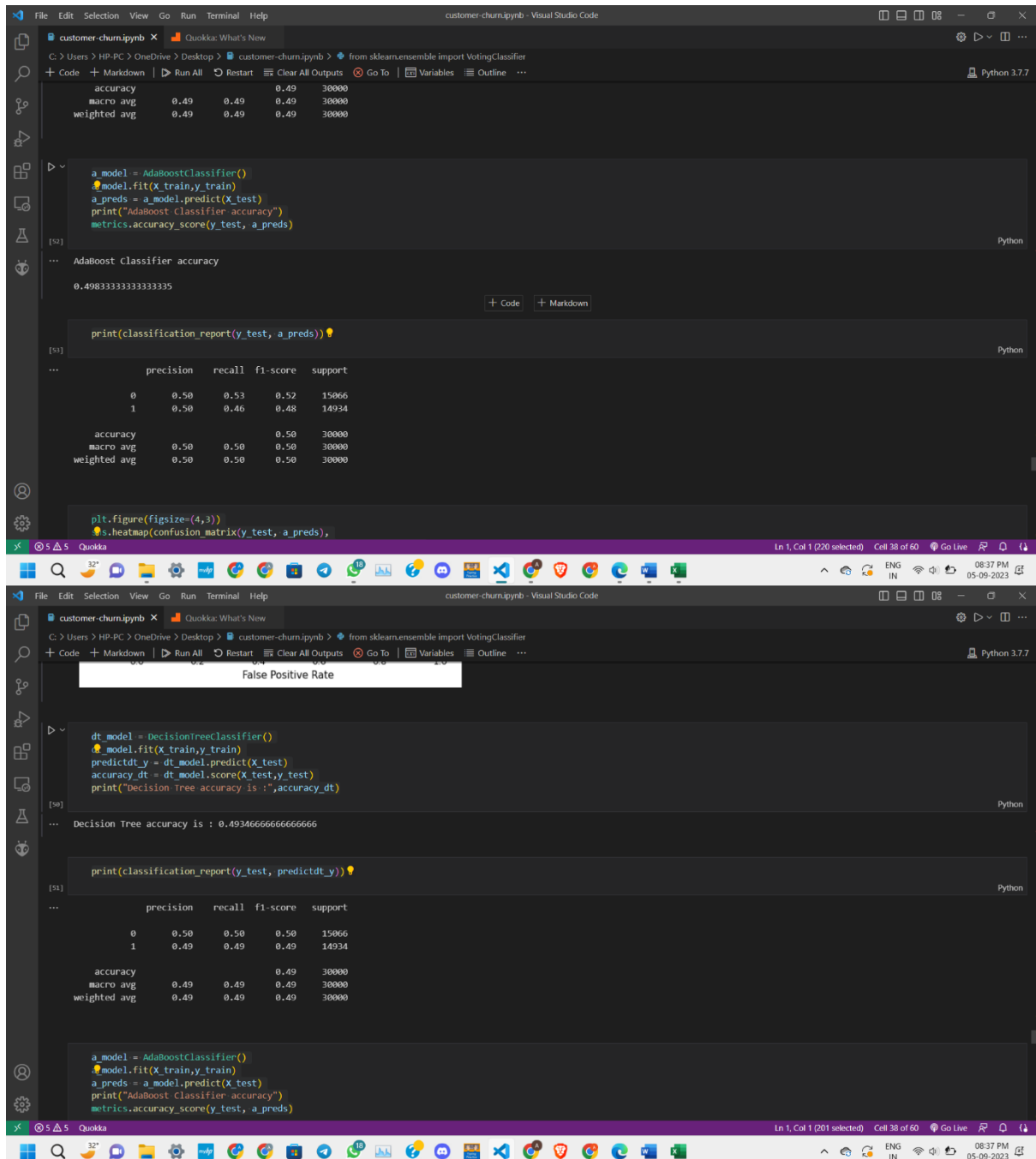**I tried to fit various classifiers on this Dataset and trained the model,namely**

**KNN classifier,AdaBoost,Random Forest, Logistic Regression,SVM,SVC;**

**ANTIM**          **2021ume0204@iitjammu.ac.in**          **IIT JAMMU**

**LINK FOR THE COOLAB FILE : https://colab.research.google.com/drive/1g5v8JE_pGZlTp0-zruGiPZPOoqTn93eX?usp=sharing**

**Out of all these classifiers None of them was giving the accuracy of more than 50 % ;**

**But the Support Vector Classifier was giving the best accuracy of 51% , so I used that classifier in my Model.**

**Here are attachments of the obtained accuracies for the different Classifiers:**

```python
gb = GradientBoostingClassifier()
gb.fit(X_train, y_train)
gb_pred = gb.predict(X_test)
print("Gradient Boosting Classifier", accuracy_score(y_test, gb_pred))
```

Gradient Boosting Classifier 0.49516666666666664

```python
print(classification_report(y_test, gb_pred))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.50      | 0.57   | 0.53     | 15066   |
| 1            | 0.49      | 0.42   | 0.46     | 14934   |
| accuracy     |           |        | 0.50     | 30000   |
| macro avg    | 0.49      | 0.49   | 0.49     | 30000   |
| weighted avg | 0.49      | 0.50   | 0.49     | 30000   |

```python
plt.figure(figsize=(4,3))
sns.heatmap(confusion_matrix(y_test, gb_pred),
            annot=True,fmt = "d",linecolor="k",linewidths=3)
plt.title("Gradient Boosting Classifier Confusion Matrix",fontsize=14)
plt.show()
```

Gradient Boosting Classifier Confusion Matrix



```python
scaler= StandardScaler()

X_train[num_cols] = scaler.fit_transform(X_train[num_cols])
X_test[num_cols] = scaler.transform(X_test[num_cols])
```

```python
knn_model = KNeighborsClassifier(n_neighbors = 50)
knn_model.fit(X_train,y_train)
predicted_y = knn_model.predict(X_test)
accuracy_knn = knn_model.score(X_test,y_test)
print("KNN accuracy:",accuracy_knn)
```
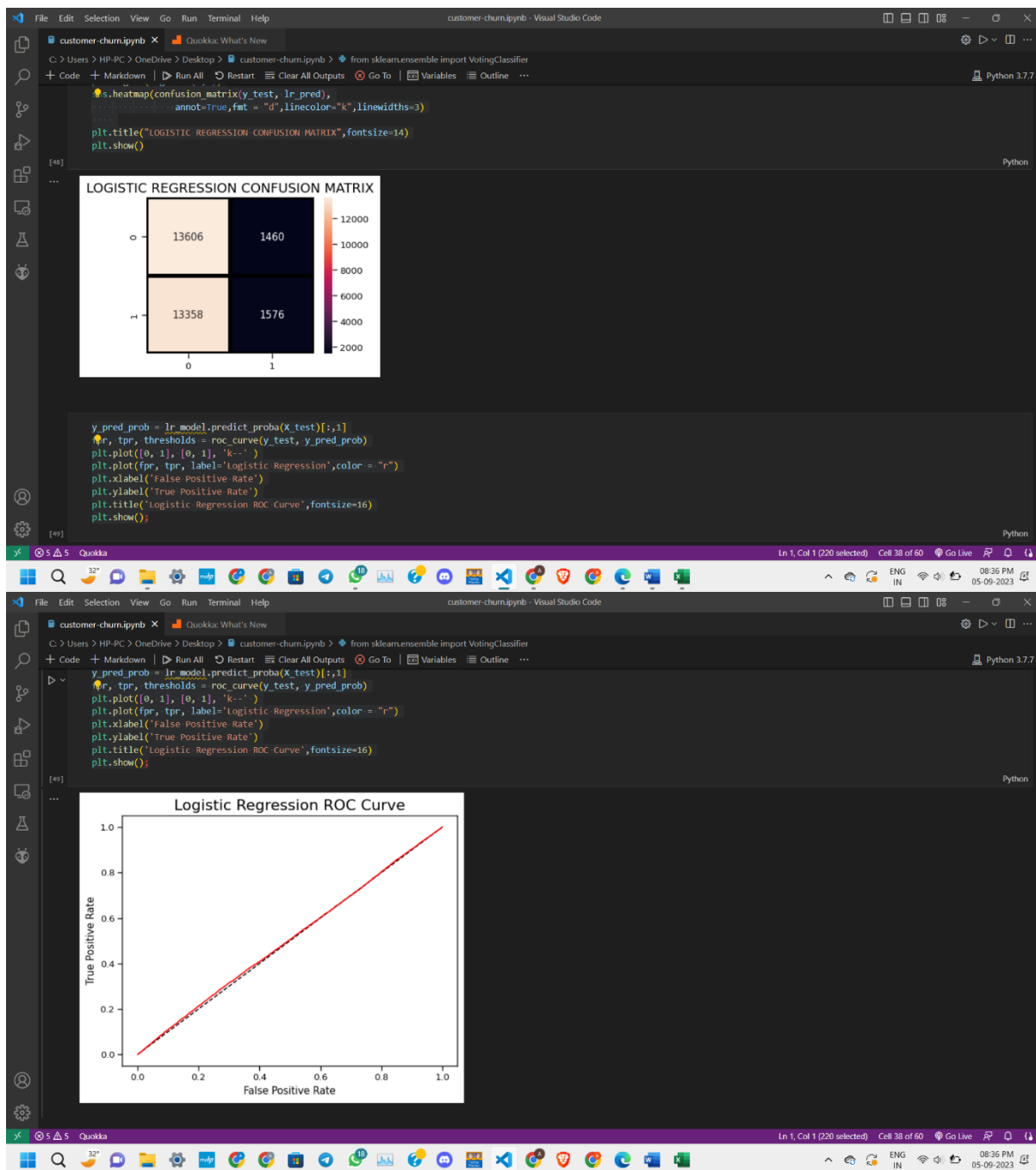
KNN accuracy: 0.49473333333333336

```python
print(classification_report(y_test, predicted_y))
```
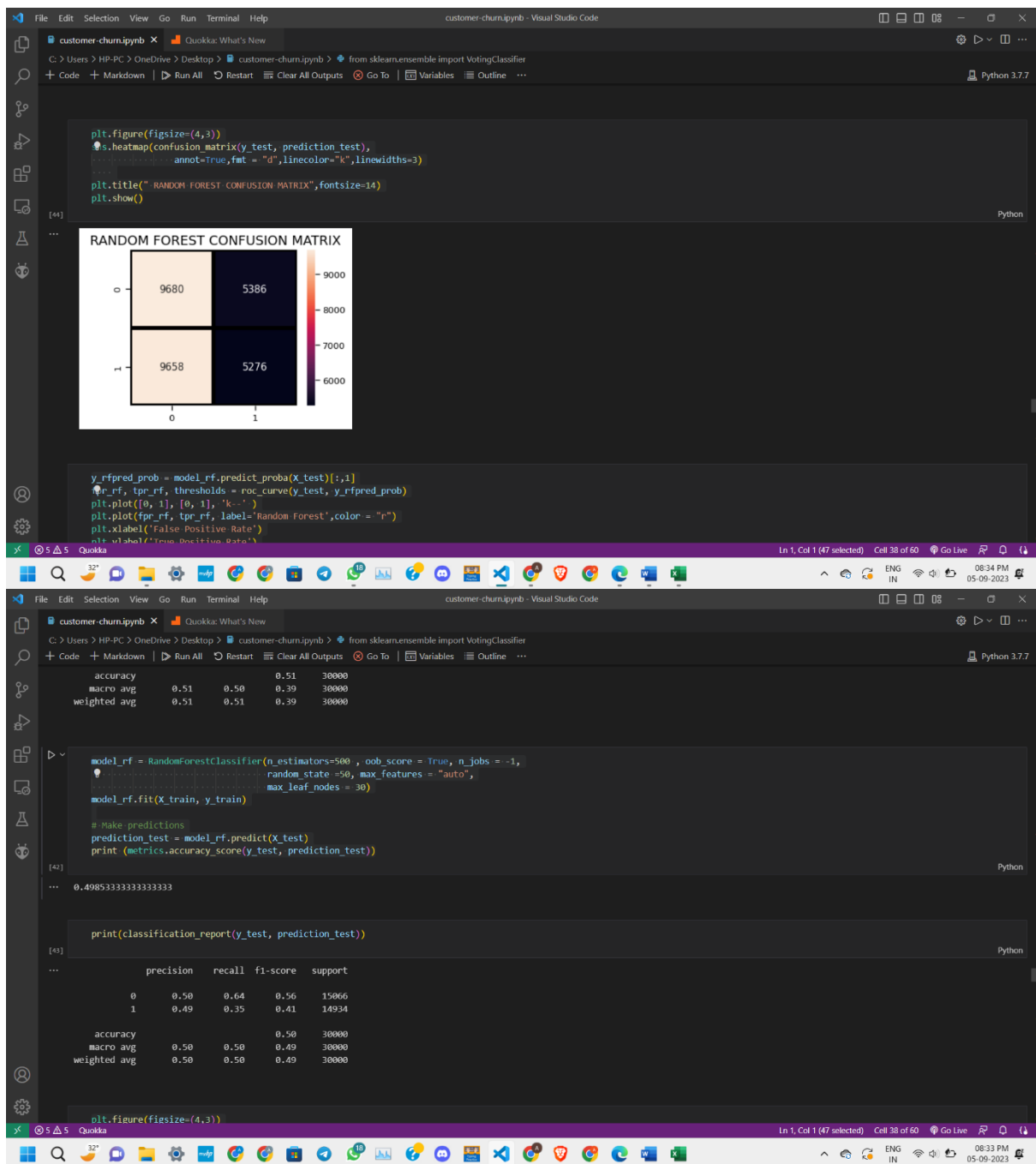
|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.50      | 0.56   | 0.53     | 15066   |
| 1            | 0.49      | 0.43   | 0.46     | 14934   |
| accuracy     |           |        | 0.49     | 30000   |
| macro avg    | 0.49      | 0.49   | 0.49     | 30000   |
| weighted avg | 0.49      | 0.49   | 0.49     | 30000   |

```python
svc_model = SVC(random_state = 1)
svc_model.fit(X_train,y_train)
```
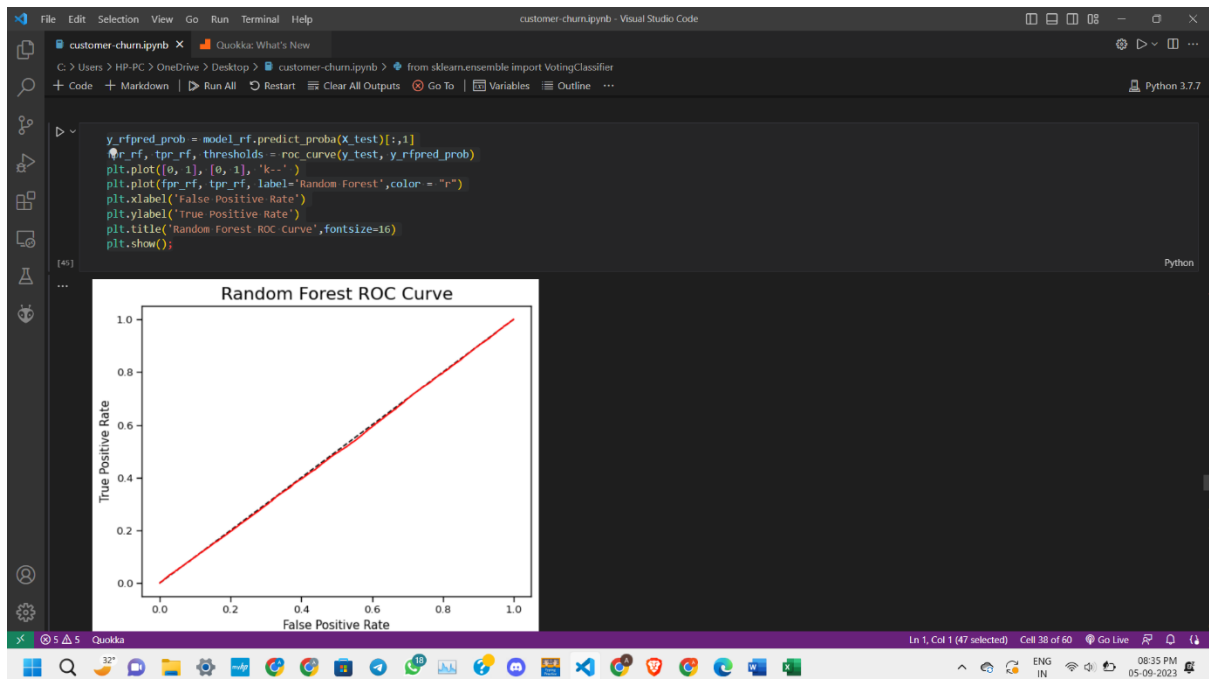
**ANTIM**          **2021ume0204@iitjammu.ac.in**          **IIT JAMMU**

LINK FOR THE COOLAB FILE : https://colab.research.google.com/drive/1g5v8JE_pGZlTp0-zruGiPZPOoqTn93eX?usp=sharing

As it is clear that all of these classifiers are giving the accuracy less than 50%.

But SVC classifier gave the accuracy of 51% so SVC is used to train the model.

ALL the Code for the Complete Model is given in Detail in the provided Google collab file along with Detailed Data visualization.

**ANTIM**          2021ume0204@iitjammu.ac.in          **IIT JAMMU**

LINK FOR THE COOLAB FILE : https://colab.research.google.com/drive/1g5v8JE_pGZlTp0-zruGiPZPOoqTn93eX?usp=sharing