## Stroke Healthcare Dataset

In this project, stroke disease dataset was loaded and some basic steps were performed for getting familiar with the dataset. We further inspected and vsiualised the dataset in detail to achieve a better understanding of it. This step also involved data-preprocessing to prepare the data for model building. Finally, we built predictive models and applied grid search validation to tune hyper-parameters in each model.

> Stroke dataset is provided by Kaggle: https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset (https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset)

```python
In [1]:
# Loading some of the libraries that will be use in this project

import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
from matplotlib import pyplot

import warnings
warnings.filterwarnings('ignore')
import copy

from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_curve, roc_auc_score
from sklearn import metrics
from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score,confusion_matrix
from sklearn.model_selection import cross_val_score
from sklearn.metrics import auc
from sklearn.metrics import precision_recall_curve
from imblearn.over_sampling import RandomOverSampler
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from xgboost.sklearn import XGBClassifier
import xgboost as xgb
```

## Introduction

Return Contents

According to the World Health Organization (WHO) stroke is the 2nd leading cause of death globally, responsible for approximately 11% of total deaths.

This dataset is used to predict whether a patient is likely to get stroke based on the input parameters like gender, age, various diseases, and smoking status. Each row in the data provides relevant information about the patient.

Attribute Information

1. id: unique identifier
2. gender: "Male", "Female" or "Other"
3. age: age of the patient
4. hypertension: 0 if the patient doesn't have hypertension, 1 if the patient has hypertension
5. heart_disease: 0 if the patient doesn't have any heart diseases, 1 if the patient has a heart disease
6. ever_married: "No" or "Yes"
7. work_type: "children", "Govt_jov", "Never_worked", "Private" or "Self-employed"
8. Residence_type: "Rural" or "Urban"
9. avg_glucose_level: average glucose level in blood
10. bmi: body mass index
11. smoking_status: "formerly smoked", "never smoked", "smokes" or "Unknown"*
12. stroke: 1 if the patient had a stroke or 0 if not

*Note: "Unknown" in smoking_status means that the information is unavailable for this patient

## Data loading

Return Contents

**Loading CSV file for dataset**

```
In [2]:   1  data = pd.read_csv('healthcare-dataset-stroke-data.csv')
          2  data
```

Out[2]:

|  | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smoked | 1 |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | NaN | never smoked | 1 |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked | 1 |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes | 1 |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24.0 | never smoked | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | Female | 80.0 | 1 | 0 | Yes | Private | Urban | 83.75 | NaN | never smoked | 0 |
| 5106 | 44873 | Female | 81.0 | 0 | 0 | Yes | Self-employed | Urban | 125.20 | 40.0 | never smoked | 0 |
| 5107 | 19723 | Female | 35.0 | 0 | 0 | Yes | Self-employed | Rural | 82.99 | 30.6 | never smoked | 0 |
| 5108 | 37544 | Male | 51.0 | 0 | 0 | Yes | Private | Rural | 166.29 | 25.6 | formerly smoked | 0 |
| 5109 | 44679 | Female | 44.0 | 0 | 0 | Yes | Govt_job | Urban | 85.28 | 26.2 | Unknown | 0 |

5110 rows × 12 columns

```
In [3]:   1  data.drop(['id'], axis =1, inplace = True)
```

**Note:**

- There are 5110 rows (instances) and 12 columns (features) in the given dataset
- Column "id" should be dropped as it could lead to overfitting because classifer might use that column to fit perfectly on the training set, ignoring all the other columns

## Exploratory data analysis

Return Contents

---

```
In [4]:   1  # Datatypes of columns
          2
          3  data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5110 entries, 0 to 5109
Data columns (total 11 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   gender             5110 non-null   object
 1   age                5110 non-null   float64
 2   hypertension       5110 non-null   int64
 3   heart_disease      5110 non-null   int64
 4   ever_married       5110 non-null   object
 5   work_type          5110 non-null   object
 6   Residence_type     5110 non-null   object
 7   avg_glucose_level  5110 non-null   float64
 8   bmi                4909 non-null   float64
 9   smoking_status     5110 non-null   object
 10  stroke             5110 non-null   int64
dtypes: float64(3), int64(3), object(5)
memory usage: 439.3+ KB
```

```
In [5]:   1  # Descriptive statsitics of continuous features
          2
          3  data.describe()
```

Out[5]:

|  | age | hypertension | heart_disease | avg_glucose_level | bmi | stroke |
|---|---|---|---|---|---|---|
| count | 5110.000000 | 5110.000000 | 5110.000000 | 5110.000000 | 4909.000000 | 5110.000000 |
| mean | 43.226614 | 0.097456 | 0.054012 | 106.147677 | 28.893237 | 0.048728 |
| std | 22.612647 | 0.296607 | 0.226063 | 45.283560 | 7.854067 | 0.215320 |
| min | 0.080000 | 0.000000 | 0.000000 | 55.120000 | 10.300000 | 0.000000 |
| 25% | 25.000000 | 0.000000 | 0.000000 | 77.245000 | 23.500000 | 0.000000 |
| 50% | 45.000000 | 0.000000 | 0.000000 | 91.885000 | 28.100000 | 0.000000 |
| 75% | 61.000000 | 0.000000 | 0.000000 | 114.090000 | 33.100000 | 0.000000 |
| max | 82.000000 | 1.000000 | 1.000000 | 271.740000 | 97.600000 | 1.000000 |

**Note:**

- Individuals are aged between less than a month old to 82 years
- BMI of indviduals ranges between 10.3 to 97.6 kg/m2

```
In [6]:   1  # Calculating number of diseased and non diseased individuals
          2
          3  total = len(data)
          4  diseased = len(data[data.stroke==1])
          5  non_diseased = len(data[data.stroke==0])
          6
          7  print ("Total number of individuals : ", total)
          8  print ("Number of individuals with stroke :", diseased)
          9  print ("Number of individuals without stroke : ", non_diseased)
```

```
Total number of individuals :  5110
Number of individuals with stroke : 249
Number of individuals without stroke :  4861
```

```python
In [7]:    1  # Calculating number of individuals in each gender type
           2
           3  print(data.gender.unique())
           4
           5  print("Other individuals : ", (data['gender']=='Other').sum())
           6  print("Number of males : ", (data['gender']=='Male').sum())
           7  print("Number of females : ", (data['gender']=='Female').sum())
```

```
['Male' 'Female' 'Other']
Other individuals :  1
Number of males :  2115
Number of females :  2994
```

```python
In [8]:    1  # Dropping individual with gender type "Other"
           2
           3  data.drop(data.index[data['gender'] == 'Other'], inplace = True)
```

```python
In [9]:    1  # Calculating number of individuals by smoking status
           2
           3  print(data.smoking_status.unique())
           4
           5  print("Unknown status : ", data.smoking_status.value_counts()['Unknown'])
           6  print("Formerly smoked : ",data.smoking_status.value_counts()['formerly smoked'])
           7  print("Never smoked : ", data.smoking_status.value_counts()['never smoked'])
           8  print("Smokes : ", data.smoking_status.value_counts()['smokes'])
```

```
['formerly smoked' 'never smoked' 'smokes' 'Unknown']
Unknown status :  1544
Formerly smoked :  884
Never smoked :  1892
Smokes :  789
```

**Note:**

- There are a total of 5110 indiviuals out of which only 249 have stroke. This suggests that dataset is highly imbalanced
- There are more females than males in the dataset
- There is only 1 individual with gender category as "Other". Hence we can remove it and conduct analyses only on males and females
- There are more individuals in never smoked category as compared to those in formerly smoked and currently smokes catgeory

```python
In [10]:   1  # Encode Categorical Columns
           2
           3  categorical_columns = ['gender','ever_married','work_type', 'Residence_type', 'smoking_status']
           4  le = LabelEncoder()
           5  data[categorical_columns] = data[categorical_columns].apply(le.fit_transform)
           6  data.head()
```

Out[10]:

| | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 67.0 | 0 | 1 | 1 | 2 | 1 | 228.69 | 36.6 | 1 | 1 |
| 1 | 0 | 61.0 | 0 | 0 | 1 | 3 | 0 | 202.21 | NaN | 2 | 1 |
| 2 | 1 | 80.0 | 0 | 1 | 1 | 2 | 0 | 105.92 | 32.5 | 2 | 1 |
| 3 | 0 | 49.0 | 0 | 0 | 1 | 2 | 1 | 171.23 | 34.4 | 3 | 1 |
| 4 | 0 | 79.0 | 1 | 0 | 1 | 3 | 0 | 174.12 | 24.0 | 2 | 1 |

```python
In [11]:   1  # Checking data for missing values
           2
           3  print('Missing data: ')
           4  print(data.isnull().sum())
```

```
Missing data:
gender                 0
age                    0
hypertension           0
heart_disease          0
ever_married           0
work_type              0
Residence_type         0
avg_glucose_level      0
bmi                  201
smoking_status         0
stroke                 0
dtype: int64
```

```
1  # Checking distribution and skewness of BMI column
2
3  data.hist(column = 'bmi')
4  print(data.skew(axis=0))
```

```
gender              0.349410
age                -0.137430
hypertension        2.715026
heart_disease       3.946786
ever_married       -0.658345
work_type          -0.308679
Residence_type     -0.032506
avg_glucose_level   1.572815
bmi                 1.055063
smoking_status     -0.039430
stroke              4.192807
dtype: float64
```

```
1  # Replacing missing valeus with median
2
3  data['bmi'].fillna(data['bmi'].median(), inplace=True)
```

```
1  data.isnull().sum()
```

```
gender              0
age                 0
hypertension        0
heart_disease       0
ever_married        0
work_type           0
Residence_type      0
avg_glucose_level   0
bmi                 0
smoking_status      0
stroke              0
dtype: int64
```

**Note:**

- BMI column has 201 missing values. Instead of dropping the individuals with missing data, we can impute missing BMI values
- The distribution of BMI is positively skewed with value for skewness equals 1.05. Hence, we can impute missing values with median of BMI column.

## Data visualisation

Return Contents

```
1  labels_gender = ['Female', 'Male']
2  labels_disease = ['Not diseased', 'Diseased']
3  labels_smoking = ['Formerly smoked', 'Never smoked', 'Smokes', 'Unknown']
4  labels_hypertension = ['No', 'Yes']
5  labels_heart = ['No', 'Yes']
```

```
In [16]:    1   # Displaying gender and stroke distribution in the dataset
            2
            3   fig = plt.subplots(figsize=(12, 8))
            4
            5   colors = ['#ff9999','#66b3ff']
            6
            7   ax1 = plt.subplot2grid((1,2),(0,0))
            8   ax1.pie(data['gender'].value_counts(), labels=labels_gender, autopct='%.1f%%',  colors = colors, shadow = True, startangle=90)
            9   ax1.set_title('Gender Distribution')
           10   plt.tight_layout()
           11
           12   ax2 = plt.subplot2grid((1,2),(0,1))
           13   ax2.pie(data['stroke'].value_counts(), labels=labels_disease, autopct='%.1f%%', colors = colors,  shadow = True, startangle=90)
           14   ax2.set_title('Stroke Distribution')
           15   plt.tight_layout()
```



**Important findings from pie chart**

- Proportion of females is 58.6% and males is 41.4%
- Class label is highly imbalanced with 5% diseased individuals and 95% non diseased individuals

```
In [17]:    1   # Displaying distribution of diseased and non-diseased individuals by gender
            2
            3   plt.figure(figsize=(10,6))
            4   ax3 = sns.countplot(data['gender'], hue=data['stroke'], palette=['#00CED1', '#FF7F50'], saturation=0.8)
            5   ax3.set_xticklabels(labels_gender, fontsize=12)
            6   plt.xlabel('Gender')
            7   plt.ylabel('Count')
            8   plt.title('Stroke by gender', fontsize=16)
            9   plt.legend(loc='upper left', fontsize=12, labels=['Non diseased', 'Diseased'])
           10   plt.show()
```

```
1  # Displaying distribution of diseased and non-diseased individuals by smoking status
2
3  plt.figure(figsize=(10,6))
4  ax4 = sns.countplot(data['smoking_status'], hue=data['stroke'], palette=['#00CED1', '#FF7F50'], saturation=0.8)
5  ax4.set_xticklabels(labels_smoking, fontsize=12)
6  plt.xlabel('Smoking status')
7  plt.ylabel('Count')
8  plt.title('Stroke by smoking status', fontsize=16)
9  plt.legend(loc='upper left', fontsize=12, labels=['Non diseased', 'Diseased'])
10 plt.show()
```

```
1  # Displaying distribution of diseased and non-diseased individuals by heart disease
2
3  plt.figure(figsize=(10,6))
4  ax5 = sns.countplot(data['heart_disease'], hue=data['stroke'], palette=['#00CED1', '#FF7F50'], saturation=0.8)
5  ax5.set_xticklabels(labels_heart, fontsize=12)
6  plt.xlabel('Heart disease')
7  plt.ylabel('Count')
8  plt.title('Stroke by heart disease', fontsize=16)
9  plt.legend(loc='upper left', fontsize=12, labels=['Non diseased', 'Diseased'])
10 plt.show()
```

```
1  # Displaying distribution of diseased and non-diseased individuals by hypertension
2
3  plt.figure(figsize=(10,6))
4  ax6 = sns.countplot(data['hypertension'], hue=data['stroke'], palette=['#00CED1', '#FF7F50'], saturation=0.8)
5  ax6.set_xticklabels(labels_hypertension, fontsize=12)
6  plt.xlabel('Hypertension')
7  plt.ylabel('Count')
8  plt.title('Stroke by hypertension', fontsize=16)
9  plt.legend(loc='upper left', fontsize=12, labels=['Non diseased', 'Diseased'])
10 plt.show()
```



**Important findings from count plots**

- Count of males and females with stroke is approximately similar
- The number of individuals with stroke are slightly more in currently smokes category as compared to other smoking categories
- There are more individuals without hypertension and heart disease who suffer from stroke as compared to those with hypertension and stroke

```
1  # Displaying histogram for Age, BMI and Average glucose level
2
3  fig, axes = plt.subplots(1, 3, figsize=(15, 5))
4
5  # Histogram Plot for Age
6
7  ax8 = sns.histplot(x='age', data=data, ax = axes[0])
8  ax8.set_title("Age Histogram")
9
10 # Histogram Plot for BMI
11
12 ax9 = sns.histplot(x='bmi', data=data, ax =axes[1])
13 ax9.set_title("BMI Histogram")
14
15 # Histogram Plot for Average glucose level
16
17 ax10 = sns.histplot(x='avg_glucose_level', data=data, ax =axes[2])
18 ax10.set_title("Average glucose level Histogram")
19
```

Out[21]: Text(0.5, 1.0, 'Average glucose level Histogram')



**Important findings from histogram**

- Distribution of BMI and average glucose level is highly positively skewed
- Distribution of age is approximately normal

In [22]:

```
1  # Displaying histogram for Age, BMI and Average glucose level
2
3  # Box Plot for Age
4  fig, axes = plt.subplots(1, 3, figsize=(15, 5))
5
6  ax14 = sns.boxplot(y= 'age', data=data, ax = axes[0])
7  ax14.set_title("Age Box Plot")
8
9  # Box Plot for BMI
10
11 ax15 = sns.boxplot(y= 'bmi', data=data, ax =axes[1])
12 ax15.set_title("BMI Box Plot")
13
14 # Box Plot for Average glucose level
15
16 ax16 = sns.boxplot(y='avg_glucose_level', data=data, ax =axes[2])
17 ax16.set_title("Average glucose level Box Plot")
```

Out[22]: Text(0.5, 1.0, 'Average glucose level Box Plot')

```
In [23]:    1  # Displaying boxplot for Age, BMI and Average glucose level by disease status
            2
            3
            4  # Box Plot for Age
            5  fig, axes = plt.subplots(1, 3, figsize=(15, 5))
            6
            7  ax11 = sns.boxplot(x='stroke', y= 'age', data=data, ax = axes[0])
            8  ax11.set_title("Age Box Plot")
            9  ax11.set_xticklabels(['Non diseased', 'Diseased'])
           10
           11  # Box Plot for BMI
           12
           13  ax12 = sns.boxplot(x='stroke', y= 'bmi', data=data, ax =axes[1])
           14  ax12.set_title("BMI Box Plot")
           15  ax12.set_xticklabels(['Non diseased', 'Diseased'])
           16
           17  # Box Plot for Average glucose level
           18
           19  ax13 = sns.boxplot(x='stroke', y= 'avg_glucose_level', data=data, ax =axes[2])
           20  ax13.set_title("Average glucose level Box Plot")
           21  ax13.set_xticklabels(['Non diseased', 'Diseased'])
```

Out[23]: [Text(0, 0, 'Non diseased'), Text(1, 0, 'Diseased')]



**Important findings from box plots**

- Distribution of all 3 features age, bmi and average glucose level is highly dispersed with high range score
- There is a greater variability in non diseased individuals for bmi and average glucose level as well as larger outliers.
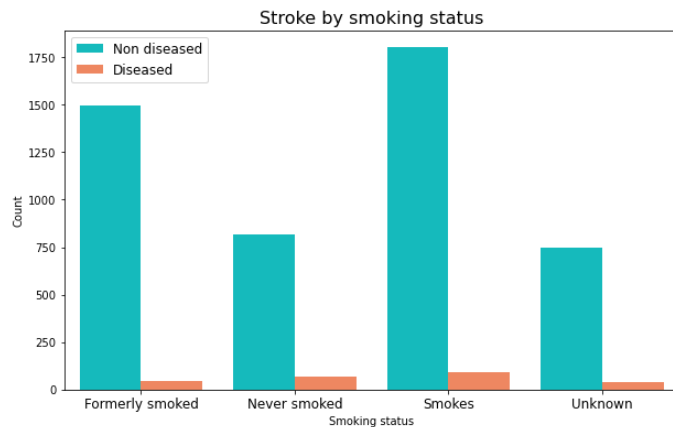
```
In [24]:    1  # Displaying distribution of diseased and non-diseased individuals by gender and age
            2
            3  plt.figure(figsize=(16, 6))
            4  palette=['#1E90FF', '#ee638f']
            5  labels_gender = ['Female', 'Male']
            6  ax7=sns.boxenplot(x=data.gender, y=data.age, hue=data.stroke, palette=palette, linewidth=3)
            7  ax7.set_xticklabels(labels_gender, fontsize=12)
            8  handles = ax7.get_legend_handles_labels()[0]
            9  ax7.legend(handles, ['Non diseased', 'Diseased'], loc = 'upper left')
           10  ax7.set_title("Boxen plot for Gender, Age and Disease status",fontsize=16)
           11  plt.show()
```



**Important findings from boxen plot**

- Older individuals suffer from stroke more as compared to younger individuals. However, there are some outliers in stroke group with childhood ages. Although, this is very rare but paediatric stroke is still possible.
- Men suffer from stroke at a slightly older age as compared to women

```
1  # Correlation matrix
2
3  plt.figure(figsize=(20, 8))
4  cmap = sns.diverging_palette(2, 165, as_cmap=True)
5  heatmap =sns.heatmap(data.corr(),  linewidths=1, cmap=cmap, annot = True, center=0)
6  heatmap.set_title('Correlation Matrix for Stroke Dataset', fontdict={'fontsize':18}, pad=16);
```

### Correlation Matrix for Stroke Dataset

| | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|
| gender | 1 | -0.028 | 0.021 | 0.086 | -0.03 | 0.057 | -0.0061 | 0.055 | -0.026 | -0.062 | 0.0091 |
| age | -0.028 | 1 | 0.28 | 0.26 | 0.68 | -0.36 | 0.014 | 0.24 | 0.32 | 0.27 | 0.25 |
| hypertension | 0.021 | 0.28 | 1 | 0.11 | 0.16 | -0.052 | -0.008 | 0.17 | 0.16 | 0.11 | 0.13 |
| heart_disease | 0.086 | 0.26 | 0.11 | 1 | 0.11 | -0.028 | 0.003 | 0.16 | 0.037 | 0.048 | 0.13 |
| ever_married | -0.03 | 0.68 | 0.16 | 0.11 | 1 | -0.35 | 0.006 | 0.16 | 0.33 | 0.26 | 0.11 |
| work_type | 0.057 | -0.36 | -0.052 | -0.028 | -0.35 | 1 | -0.0073 | -0.05 | -0.3 | -0.31 | -0.032 |
| Residence_type | -0.0061 | 0.014 | -0.008 | 0.003 | 0.006 | -0.0073 | 1 | -0.0048 | -0.00044 | 0.0082 | 0.015 |
| avg_glucose_level | 0.055 | 0.24 | 0.17 | 0.16 | 0.16 | -0.05 | -0.0048 | 1 | 0.17 | 0.063 | 0.13 |
| bmi | -0.026 | 0.32 | 0.16 | 0.037 | 0.33 | -0.3 | -0.00044 | 0.17 | 1 | 0.22 | 0.036 |
| smoking_status | -0.062 | 0.27 | 0.11 | 0.048 | 0.26 | -0.31 | 0.0082 | 0.063 | 0.22 | 1 | 0.028 |
| stroke | 0.0091 | 0.25 | 0.13 | 0.13 | 0.11 | -0.032 | 0.015 | 0.13 | 0.036 | 0.028 | 1 |

**Important findings from correlation matrix**

- There are no strong correlations observed between features as well as between features and class label

## Model building

Return Contents

```
1  # Dividing data into input features and output class
2  X, y = data.iloc[:, :-1], data.iloc[:, -1:]
```

```
1  # Input data
2
3  X.head()
```

| | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 67.0 | 0 | 1 | 1 | 2 | 1 | 228.69 | 36.6 | 1 |
| 1 | 0 | 61.0 | 0 | 0 | 1 | 3 | 0 | 202.21 | 28.1 | 2 |
| 2 | 1 | 80.0 | 0 | 1 | 1 | 2 | 0 | 105.92 | 32.5 | 2 |
| 3 | 0 | 49.0 | 0 | 0 | 1 | 2 | 1 | 171.23 | 34.4 | 3 |
| 4 | 0 | 79.0 | 1 | 0 | 1 | 3 | 0 | 174.12 | 24.0 | 2 |

```
1  # Class label
2
3  y.head()
```

| | stroke |
|---|---|
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | 1 |
| 4 | 1 |

```
1  print(X.shape)
2  print(y.shape)
```

```
(5109, 10)
(5109, 1)
```

```
In [30]:    1  #  Splitting data into train and test
            2
            3  X_train_tune, X_test_df, y_train_tune, y_test_df = train_test_split(X, y, test_size=0.2, random_state=1)
            4
            5  X_train_df, X_val_df, y_train_df, y_val_df = train_test_split(X_train_tune, y_train_tune, test_size=0.15, random_state=1) # 0.25 x 0.8 = 0
            6
            7  # X_train_df, X_test_df, y_train_df, y_test_df = train_test_split(X, y, test_size=0.3, random_state= 0)
```

```
In [31]:    1  print("Train Shape: {}".format(X_train_df.shape))
            2  print("Validation Shape: {}".format(X_val_df.shape))
            3  print("Test Shape: {}".format(X_test_df.shape))
```

```
Train Shape: (3473, 10)
Validation Shape: (614, 10)
Test Shape: (1022, 10)
```

```
In [32]:    1  # Performing oversampling
            2
            3  OS = RandomOverSampler(random_state=0)
            4  osx, osy = OS.fit_resample(X_train_df, y_train_df)
            5  print("X_train shape after oversampling: ", osx.shape)
            6  print("y_train shape after oversampling: ", osy.shape)
```

```
X_train shape after oversampling:  (6606, 10)
y_train shape after oversampling:  (6606, 1)
```

```
In [33]:    1  # Performing standardization
            2
            3  sc = StandardScaler()
            4  X_train = sc.fit_transform(osx)
            5  X_val = sc.fit_transform(X_val_df)
            6  X_test = sc.fit_transform(X_test_df)
```

```
In [34]:    1  X_train
```

```
Out[34]: array([[ 1.20617142,  0.31939861, -0.46826179, ..., -0.61591723,
                   1.01154191, -0.42030218],
                 [-0.82906955, -0.71710541, -0.46826179, ..., -0.71883621,
                  -0.98632002,  0.5467881 ],
                 [ 1.20617142, -2.06906718, -0.46826179, ...,  0.04870741,
                  -1.57392647, -1.38739245],
                 ...,
                 [ 1.20617142, -0.04112452, -0.46826179, ..., -0.5664509 ,
                   0.24765353,  1.51387837],
                 [-0.82906955,  0.86018332, -0.46826179, ..., -0.80001174,
                  -0.41340373,  0.5467881 ],
                 [-0.82906955,  1.04044489,  2.13555755, ...,  1.53577786,
                   2.39241707,  0.5467881 ]])
```

```
In [35]:    1  X_val
```

```
Out[35]: array([[-0.85130435,  0.38147098, -0.3105295 , ...,  0.03177528,
                   0.25078602,  0.62127172],
                 [-0.85130435, -1.42357819, -0.3105295 , ..., -0.79802466,
                  -1.88840489, -1.2578457 ],
                 [-0.85130435,  0.33634475, -0.3105295 , ..., -0.3935435 ,
                   0.42035603,  0.62127172],
                 ...,
                 [-0.85130435, -0.65643229, -0.3105295 , ..., -0.71652645,
                  -0.3100994 ,  1.56083043],
                 [ 1.17466802, -0.61130606,  3.22030594, ..., -0.61581134,
                  -0.10139785, -0.31828699],
                 [-0.85130435, -0.61130606, -0.3105295 , ..., -0.74685675,
                  -0.16661709,  0.62127172]])
```

```
In [36]:    1  X_test
```

```
Out[36]: array([[ 1.21926946, -1.55443213, -0.3401772 , ..., -1.05401379,
                  -1.26468711, -1.35160688],
                 [ 1.21926946, -1.81639089, -0.3401772 , ..., -0.34378401,
                  -0.73833399, -1.35160688],
                 [-0.82016325, -1.77273109, -0.3401772 , ..., -0.36412003,
                  -1.72837916, -1.35160688],
                 ...,
                 [ 1.21926946, -1.51077234, -0.3401772 , ..., -0.24997595,
                  -0.73833399, -1.35160688],
                 [-0.82016325,  0.10463996, -0.3401772 , ..., -0.320168  ,
                  -0.53781851,  0.5419397 ],
                 [ 1.21926946, -0.20097858, -0.3401772 , ...,  0.16964564,
                   1.39214295, -0.40483359]])
```

```
In [37]:    1  # Converting dataframe to numpy array
            2
            3  y_train = osy.to_numpy().ravel()
            4  y_train
```

```
Out[37]: array([0, 0, 0, ..., 1, 1, 1])
```

```
In [38]:    1  y_val = y_val_df.to_numpy().ravel()
            2  y_val
```

```
Out[38]: array([0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
```

```
In [39]:    1  y_test = y_test_df.to_numpy().ravel()
            2  y_test
```

```
Out[39]: array([0, 0, 0, ..., 0, 0, 0])
```

**Logistic Regression**

```
In [40]:    1  # Logistic regression
            2
            3  lor = LogisticRegression(max_iter = 20000)
            4  lor.fit(X_train, y_train)
            5  y_val_pred_lor = lor.predict(X_val)
```

```
In [41]:    1  # Printing scores for accuracy, precision, recall and F1 with threshold 0.5
            2
            3  print("Accuracy : ", metrics.accuracy_score(y_val, y_val_pred_lor))
            4  print("Precision : ", metrics.precision_score(y_val, y_val_pred_lor))
            5  print("Recall : ", metrics.recall_score(y_val, y_val_pred_lor))
            6  print("F1 score : ", metrics.f1_score(y_val, y_val_pred_lor))
```

```
Accuracy :  0.5553745928338762
Precision :  0.059027777777777776
Recall :  0.8947368421052632
F1 score :  0.11074918566775245
```

```
In [42]:    1  val_lor_proba = lor.predict_proba(X_val)
            2  y_val_prob_lor = val_lor_proba[:, 1]
            3  y_val_prob_lor
```

```
Out[42]: array([0.5901919 , 0.07348657, 0.55158436, 0.8692085 , 0.81749494,
                 0.03563532, 0.93889825, 0.04000919, 0.13656255, 0.21132713,
                 0.13315282, 0.14863129, 0.8012939 , 0.20443922, 0.50997239,
                 0.09294311, 0.07758153, 0.43499084, 0.5294772 , 0.34350958,
                 0.17778339, 0.84533251, 0.12316908, 0.0367498 , 0.13534342,
                 0.28062785, 0.10598128, 0.97459354, 0.16509053, 0.96784532,
                 0.59833853, 0.76919572, 0.80635983, 0.46907555, 0.94735302,
                 0.6283407 , 0.70578857, 0.92609312, 0.57164407, 0.63196893,
                 0.07222791, 0.40361193, 0.74390412, 0.42652412, 0.25841966,
                 0.60342324, 0.61226367, 0.79849076, 0.92239013, 0.05378858,
                 0.69976722, 0.75610241, 0.58323382, 0.04803207, 0.06110172,
                 0.91880072, 0.38966455, 0.90732755, 0.1102547 , 0.26207897,
                 0.76322798, 0.2811851 , 0.11135666, 0.38189743, 0.76340337,
                 0.11210117, 0.81468001, 0.62315942, 0.84396673, 0.60669488,
                 0.3587874 , 0.18129746, 0.29918093, 0.54519342, 0.61435371,
                 0.95067833, 0.75847122, 0.87901685, 0.15169192, 0.91614148,
                 0.38326964, 0.29508097, 0.18540214, 0.13517798, 0.02784977,
                 0.52395234, 0.96421653, 0.24765595, 0.52412349, 0.39798421,
                 0.24177847, 0.0629371 , 0.16683823, 0.95039168, 0.54809554,
```

```
In [43]:    1  def lor_bestThreshold(y_true, y_pred):
            2
            3  # Calculating best threshold for maximum F1 score
            4      P, R, T = precision_recall_curve(y_true, y_pred.round(3))
            5      F1index, = np.where( (2*(P*R)/(P+R)) == max((2*(P*R)/(P+R))))
            6      global lor_f1index
            7      lor_f1index = T[F1index][0]
            8      print("Best Threshold for maximum F1 score: ", lor_f1index)
```

```python
In [44]:   1  def lor_performance(y_true, y_pred):
           2
           3      # Predicting class label based on best threshold
           4
           5      y_pred_new = np.where(y_pred >= lor_f1index, 1, 0)
           6
           7  # Calculating performace metrices
           8
           9      pf1 = metrics.precision_score(y_true, y_pred_new)
          10      rf1 = metrics.recall_score(y_true,y_pred_new)
          11      f1_1 = (2 * (pf1*rf1) / ( pf1 + rf1))
          12
          13  # Printing accuracy, precision, recall and F1 score:
          14
          15      print('Accuracy: %.3f' % metrics.accuracy_score(y_true, y_pred_new.round(2)))
          16      print("Precision at Best F1 :", pf1)
          17      print("Recall at Best F1 :", rf1)
          18      print("Best F1 Score :", f1_1)
          19
          20  # Printing confusion matrix
          21
          22      lor_cm = confusion_matrix(y_true, y_pred_new)
          23      lor_df = pd.DataFrame(lor_cm,
          24                      index = [0, 1],
          25                      columns = [0, 1])
          26      plt.figure(figsize=(10,5))
          27      sns.heatmap(lor_df, fmt="d", cmap = 'BuGn', annot=True)
          28      plt.title('Confusion Matrix')
          29      plt.ylabel('Actual Values')
          30      plt.xlabel('Predicted Values')
          31      plt.show()
```

```python
In [45]:   1  # Best threshold for F1 score and performance on validaton set
           2
           3  lor_bestThreshold(y_val, y_val_prob_lor)
           4  lor_performance(y_val, y_val_prob_lor)
```

```
Best Threshold for maximum F1 score:  0.894
Accuracy: 0.888
Precision at Best F1 : 0.1527777777777778
Recall at Best F1 : 0.5789473684210527
Best F1 Score : 0.2417582417582418
```



```python
In [46]:   1  test_lor_proba = lor.predict_proba(X_test)
           2  y_test_prob_lor = test_lor_proba[:, 1]
           3  y_test_prob_lor
```

```
Out[46]: array([0.04513708, 0.03643784, 0.04723962, ..., 0.06081127, 0.47787425,
            0.31931657])
```

```python
In [47]:   1  # Test set peformance: Best threshold  F1
           2
           3  lor_performance(y_test, y_test_prob_lor)
```

```
Accuracy: 0.875
Precision at Best F1 : 0.22131147540983606
Recall at Best F1 : 0.45
Best F1 Score : 0.2967032967032967
```



**Decision Tree**

```python
In [48]:   1  #  Decision Tree
           2
           3  dtc = DecisionTreeClassifier(criterion='gini', max_depth=3, min_samples_split=3, min_samples_leaf=1)
           4  dtc.fit(X_train, y_train)
           5  y_val_pred_dtc = dtc.predict(X_val)
```

```
In [49]:    1  # Printing scores for accuracy, precision, recall and F1 with threshold 0.5
            2
            3  print("Accuracy : ", metrics.accuracy_score(y_val, y_val_pred_dtc))
            4  print("Precision : ", metrics.precision_score(y_val, y_val_pred_dtc))
            5  print("Recall : ", metrics.recall_score(y_val, y_val_pred_dtc))
            6  print("F1 score : ", metrics.f1_score(y_val, y_val_pred_dtc))
```

```
Accuracy :  0.49185667752442996
Precision :  0.04923076923076923
Recall :  0.8421052631578947
F1 score :  0.0930232558139535
```

```
In [50]:    1  val_dtc_proba = dtc.predict_proba(X_val)
            2  y_val_prob_dtc = val_dtc_proba[:, 1]
            3  y_val_prob_dtc
```

```
Out[50]: array([0.62885906, 0.01351351, 0.62885906, 0.7522604 , 0.62885906,
         0.01351351, 0.85365854, 0.01351351, 0.01351351, 0.21116139,
         0.01351351, 0.01351351, 0.62885906, 0.21116139, 0.32404181,
         0.01351351, 0.01351351, 0.62885906, 0.62885906, 0.21116139,
         0.01351351, 0.7522604 , 0.01351351, 0.01351351, 0.01351351,
         0.21116139, 0.01351351, 0.85365854, 0.01351351, 0.85365854,
         0.62885906, 0.85365854, 0.62885906, 0.62885906, 0.85365854,
         0.7522604 , 0.85365854, 0.7522604 , 0.62885906, 0.62885906,
         0.01351351, 0.32404181, 0.62885906, 0.62885906, 0.21116139,
         0.62885906, 0.7522604 , 0.85365854, 0.85365854, 0.01351351,
         0.62885906, 0.85365854, 0.62885906, 0.01351351, 0.01351351,
         0.85365854, 0.32404181, 0.7522604 , 0.01351351, 0.21116139,
         0.85365854, 0.21116139, 0.01351351, 0.32404181, 0.7522604 ,
         0.01351351, 0.7522604 , 0.62885906, 0.7522604 , 0.62885906,
         0.62885906, 0.21116139, 0.21116139, 0.62885906, 0.32404181,
         0.7522604 , 0.7522604 , 0.62885906, 0.01351351, 0.7522604 ,
         0.21116139, 0.21116139, 0.01351351, 0.01351351, 0.01351351,
         0.62885906, 0.85365854, 0.21116139, 0.32404181, 0.62885906,
         0.21116139, 0.01351351, 0.01351351, 0.85365854, 0.62885906,
         0.85365854, 0.01351351, 0.7522604 , 0.7522604 , 0.62885906,
```

```
In [51]:    1  def dtc_bestThreshold(y_true, y_pred):
            2
            3  # Calculating best threshold for maximum F1 score
            4      P, R, T = precision_recall_curve(y_true, y_pred.round(3))
            5      F1index, = np.where( (2*(P*R)/(P+R)) == max((2*(P*R)/(P+R))))
            6      global dtc_f1index
            7      dtc_f1index = T[F1index][0]
            8      print("Best Threshold for maximum F1 score: ", dtc_f1index)
```

```
In [52]:    1  def dtc_performance(y_true, y_pred):
            2
            3  # Predicting class label based on best threshold
            4
            5      y_pred_new = np.where(y_pred >= dtc_f1index, 1, 0)
            6
            7  # Calculating performace metrices
            8
            9      pf1 = metrics.precision_score(y_true, y_pred_new)
           10      rf1 = metrics.recall_score(y_true,y_pred_new)
           11      f1_1 = (2 * (pf1*rf1) / ( pf1 + rf1))
           12
           13  # Printing accuracy, precision, recall and F1 score:
           14
           15      print('Accuracy: %.3f' % metrics.accuracy_score(y_true, y_pred_new.round(2)))
           16      print("Precision at Best F1 :", pf1)
           17      print("Recall at Best F1 :", rf1)
           18      print("Best F1 Score :", f1_1)
           19
           20  # Printing confusion matrix
           21
           22      dtc_cm = confusion_matrix(y_true, y_pred_new)
           23      dtc_df = pd.DataFrame(dtc_cm,
           24                      index = [0, 1],
           25                      columns = [0, 1])
           26      plt.figure(figsize=(10,5))
           27      sns.heatmap(dtc_df, fmt="d", cmap = 'BuGn', annot=True)
           28      plt.title('Confusion Matrix')
           29      plt.ylabel('Actual Values')
           30      plt.xlabel('Predicted Values')
           31      plt.show()
```

```
In [53]:    1  # Best threshold for F1 score and performance on validaton set
            2
            3  dtc_bestThreshold(y_val, y_val_prob_dtc)
            4  dtc_performance(y_val, y_val_prob_dtc)
```

```
Best Threshold for maximum F1 score:  0.752
Accuracy: 0.705
Precision at Best F1 : 0.07368421052631578
Recall at Best F1 : 0.7368421052631579
Best F1 Score : 0.13397129186602869
```



Confusion Matrix
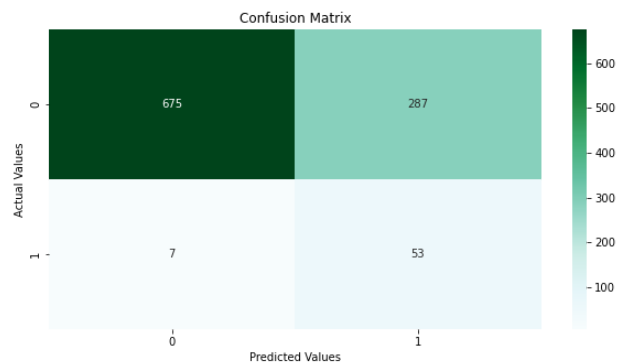
```python
In [54]:   1  test_dtc_proba = dtc.predict_proba(X_test)
           2  y_test_prob_dtc = test_dtc_proba[:, 1]
           3  y_test_prob_dtc
```

```
Out[54]: array([0.01351351, 0.01351351, 0.01351351, ..., 0.01351351, 0.62885906,
                0.62885906])
```

```python
In [55]:   1  # Test set peformance: Best threshold  F1
           2  dtc_performance(y_test, y_test_prob_dtc)
```

```
Accuracy: 0.712
Precision at Best F1 : 0.15588235294117647
Recall at Best F1 : 0.8833333333333333
Best F1 Score : 0.265
```



Confusion Matrix

**Random Forest**

```python
In [56]:   1  # Random Forest
           2
           3  rfc =  RandomForestClassifier(n_estimators=50, max_depth=3, min_samples_split=4, min_samples_leaf = 1)
           4  rfc.fit(X_train, y_train)
           5  y_val_pred_rfc = rfc.predict(X_val)
```

```python
In [57]:   1  # Printing scores for accuracy, precision, recall and F1 with threshold 0.5
           2  print("Accuracy : ", metrics.accuracy_score(y_val, y_val_pred_rfc))
           3  print("Precision : ", metrics.precision_score(y_val, y_val_pred_rfc))
           4  print("Recall : ", metrics.recall_score(y_val, y_val_pred_rfc))
           5  print("F1 score : ", metrics.f1_score(y_val, y_val_pred_rfc))
```

```
Accuracy :  0.511400651465798
Precision :  0.056782334384858045
Recall :  0.9473684210526315
F1 score :  0.10714285714285715
```

```python
In [58]:   1  val_rfc_proba = rfc.predict_proba(X_val)
           2  y_val_prob_rfc = val_rfc_proba[:, 1]
           3  y_val_prob_rfc
```

```
Out[58]: array([0.59297157, 0.04094033, 0.57607762, 0.66224072, 0.74110192,
                0.03505334, 0.73664409, 0.0334643 , 0.1919148 , 0.25479207,
                0.20719544, 0.23448558, 0.62123817, 0.26580574, 0.5529246 ,
                0.04051241, 0.06928641, 0.512068  , 0.55490021, 0.27702061,
                0.21483538, 0.66255897, 0.03763732, 0.03482201, 0.20229804,
                0.27561875, 0.21055901, 0.75089014, 0.19342663, 0.7350246 ,
                0.59200824, 0.61872996, 0.63648861, 0.46726754, 0.67850677,
                0.5970046 , 0.58022316, 0.60471499, 0.59033394, 0.69389336,
                0.19971047, 0.48370525, 0.63375816, 0.45178967, 0.24991452,
                0.5956311 , 0.60950397, 0.659911  , 0.63458148, 0.07824234,
                0.6166636 , 0.67612425, 0.61981484, 0.03839601, 0.03839601,
                0.74959841, 0.43308421, 0.62148192, 0.20026379, 0.27190205,
                0.63194384, 0.29664841, 0.192634  , 0.43031146, 0.60082409,
                0.21271897, 0.62046825, 0.591809  , 0.66735782, 0.59548848,
                0.47048845, 0.26869351, 0.27827135, 0.60220737, 0.56976963,
                0.68153672, 0.6031155 , 0.64177233, 0.21733885, 0.61990574,
                0.26574834, 0.30214481, 0.19012038, 0.21658015, 0.0885531 ,
                0.5777183 , 0.70575606, 0.27993626, 0.55400065, 0.4896141 ,
                0.27629352, 0.06928641, 0.21259379, 0.6440106 , 0.59189552,
```

```python
In [59]:   1  def rfc_bestThreshold(y_true, y_pred):
           2
           3  # Calculating best threshold for maximum F1 score
           4      P, R, T = precision_recall_curve(y_true, y_pred.round(3))
           5      F1index, = np.where( (2*(P*R)/(P+R)) == max((2*(P*R)/(P+R))))
           6      global rfc_f1index
           7      rfc_f1index = T[F1index][0]
           8      print("Best Threshold for maximum F1 score: ", rfc_f1index)
```

```
In [60]:    1  def rfc_performance(y_true, y_pred):
            2
            3      # Predicting class label based on best threshold
            4
            5      y_pred_new = np.where(y_pred >= rfc_f1index, 1, 0)
            6
            7      # Calculating performace metrics
            8
            9      pf1 = metrics.precision_score(y_true, y_pred_new)
           10      rf1 = metrics.recall_score(y_true,y_pred_new)
           11      f1_1 = (2 * (pf1*rf1) / ( pf1 + rf1))
           12
           13      # Printing accuracy, precision, recall and F1 score:
           14
           15      print('Accuracy: %.3f' % metrics.accuracy_score(y_true, y_pred_new.round(2)))
           16      print("Precision at Best F1 :", pf1)
           17      print("Recall at Best F1 :", rf1)
           18      print("Best F1 Score :", f1_1)
           19
           20      # Printing confusion matrix
           21
           22      rfc_cm = confusion_matrix(y_true, y_pred_new)
           23      rfc_df = pd.DataFrame(rfc_cm,
           24                    index = [0, 1],
           25                    columns = [0, 1])
           26      plt.figure(figsize=(10,5))
           27      sns.heatmap(rfc_df, fmt="d", cmap = 'BuGn', annot=True)
           28      plt.title('Confusion Matrix')
           29      plt.ylabel('Actual Values')
           30      plt.xlabel('Predicted Values')
           31      plt.show()
```

```
In [61]:    1  # Best threshold for F1 score and performance on validaton set
            2
            3  rfc_bestThreshold(y_val, y_val_prob_rfc)
            4  rfc_performance(y_val, y_val_prob_rfc)
```

```
Best Threshold for maximum F1 score:  0.726
Accuracy: 0.923
Precision at Best F1 : 0.13157894736842105
Recall at Best F1 : 0.2631578947368421
Best F1 Score : 0.17543859649122803
```



```
In [62]:    1  test_rfc_proba = rfc.predict_proba(X_test)
            2  y_test_prob_rfc = test_rfc_proba[:, 1]
            3  y_test_prob_rfc
```

```
Out[62]: array([0.04400927, 0.07936583, 0.03848498, ..., 0.07936583, 0.53176916,
                0.45801838])
```

```
In [63]:    1  # Test set peformance: Best threshold  F1
            2
            3  rfc_performance(y_test, y_test_prob_rfc)
```

```
Accuracy: 0.900
Precision at Best F1 : 0.23076923076923078
Recall at Best F1 : 0.3
Best F1 Score : 0.2608695652173913
```



**XGBoost**

```
In [64]:    1  # XGBoost
            2
            3  xgb = XGBClassifier(booster = 'gbtree', learning_rate = 0.2,  max_depth = 3, min_child_weight = 4, n_estimators  = 100)
            4  xgb.fit(X_train, y_train)
            5  y_val_pred_xgb = xgb.predict(X_val)
```

```
In [65]:   1  # Printing scores for accuracy, precision, recall and F1 with threshold 0.5
           2  print("Accuracy : ", metrics.accuracy_score(y_val, y_val_pred_xgb))
           3  print("Precision : ", metrics.precision_score(y_val, y_val_pred_xgb))
           4  print("Recall : ", metrics.recall_score(y_val, y_val_pred_xgb))
           5  print("F1 score : ", metrics.f1_score(y_val, y_val_pred_xgb))
```

```
Accuracy :  0.737785016286645
Precision :  0.05625
Recall :  0.47368421052631576
F1 score :  0.1005586592178771
```

```
In [66]:   1  val_xgb_proba = xgb.predict_proba(X_val)
           2  y_val_prob_xgb = val_xgb_proba[:, 1]
           3  y_val_prob_xgb
```

```
Out[66]: array([6.78817153e-01, 2.98074214e-03, 3.70944083e-01, 7.87746072e-01,
                 4.77646738e-01, 1.98925799e-03, 5.62546790e-01, 4.13316506e-04,
                 9.00347997e-03, 3.64418253e-02, 1.21001888e-03, 2.25282721e-02,
                 4.02316898e-02, 4.70763594e-02, 3.32985967e-01, 5.70106134e-03,
                 2.30093766e-03, 6.05080366e-01, 6.83031827e-02, 4.31910791e-02,
                 1.33705884e-02, 1.68648228e-01, 3.02909000e-04, 1.49806775e-03,
                 2.99105770e-04, 8.07479247e-02, 4.54469398e-02, 5.08826733e-01,
                 2.56600732e-04, 7.44657040e-02, 9.39358249e-02, 5.69034040e-01,
                 3.43420744e-01, 4.61541861e-02, 4.77598786e-01, 2.09523112e-01,
                 5.61590254e-01, 3.15018177e-01, 5.36622286e-01, 7.00215101e-01,
                 9.56199598e-04, 5.67753434e-01, 1.17981724e-01, 8.37576240e-02,
                 3.31274047e-02, 5.41981578e-01, 4.91242260e-01, 5.63326955e-01,
                 6.90959632e-01, 3.37825762e-03, 5.37757754e-01, 6.28371656e-01,
                 2.46840447e-01, 6.62178500e-04, 3.83408013e-04, 8.05537879e-01,
                 4.15352322e-02, 5.15885472e-01, 2.82293232e-03, 9.95177496e-03,
                 4.18184280e-01, 2.12080926e-01, 1.79896061e-03, 4.76122871e-02,
                 7.49996185e-01, 2.89719389e-03, 8.67162943e-01, 5.49740732e-01,
                 2.25281745e-01, 5.26018560e-01, 4.78884697e-01, 2.71404423e-02,
                 1.09763034e-01, 1.38265774e-01, 1.88607216e-01, 6.98699415e-01,
```

```
In [67]:   1  def xgb_bestThreshold(y_true, y_pred):
           2
           3      # Calculating best threshold for maximum F1 score
           4      P, R, T = precision_recall_curve(y_true, y_pred.round(3))
           5      F1index, = np.where( (2*(P*R)/(P+R)) == max((2*(P*R)/(P+R))))
           6      global xgb_f1index
           7      xgb_f1index = T[F1index][0]
           8      print("Best Threshold for maximum F1 score: ", xgb_f1index)
```

```
In [68]:   1  def xgb_performance(y_true, y_pred):
           2
           3      # Predicting class label based on best threshold
           4
           5      y_pred_new = np.where(y_pred >= xgb_f1index, 1, 0)
           6
           7      # Calculating performace metrices
           8
           9      pf1 = metrics.precision_score(y_true, y_pred_new)
          10      rf1 = metrics.recall_score(y_true,y_pred_new)
          11      f1_1 = (2 * (pf1*rf1) / ( pf1 + rf1))
          12
          13      # Printing accuracy, precision, recall and F1 score:
          14
          15      print('Accuracy: %.3f' % metrics.accuracy_score(y_true, y_pred_new.round(2)))
          16      print("Precision at Best F1 :", pf1)
          17      print("Recall at Best F1 :", rf1)
          18      print("Best F1 Score :", f1_1)
          19
          20      # Printing confusion matrix
          21
          22      xgb_cm = confusion_matrix(y_true, y_pred_new)
          23      xgb_df = pd.DataFrame(xgb_cm,
          24                     index = [0, 1],
          25                     columns = [0, 1])
          26      plt.figure(figsize=(10,5))
          27      sns.heatmap(xgb_df, fmt="d", cmap = 'BuGn', annot=True)
          28      plt.title('Confusion Matrix')
          29      plt.ylabel('Actual Values')
          30      plt.xlabel('Predicted Values')
          31      plt.show()
```

```
In [69]:   1  # Best threshold for F1 score and performance on validaton set
           2
           3  xgb_bestThreshold(y_val, y_val_prob_xgb)
           4  xgb_performance(y_val, y_val_prob_xgb)
```

```
Best Threshold for maximum F1 score:  0.439
Accuracy: 0.699
Precision at Best F1 : 0.06315789473684211
Recall at Best F1 : 0.631578947368421
Best F1 Score : 0.11483253588516747
```
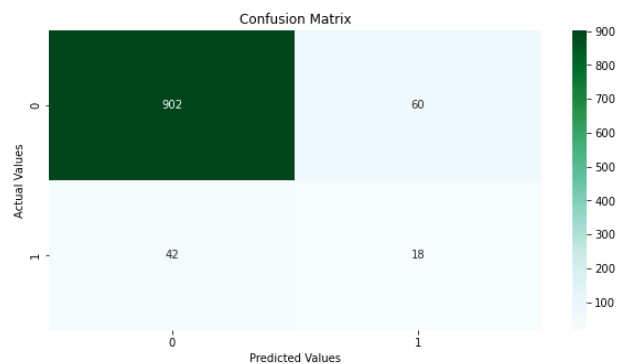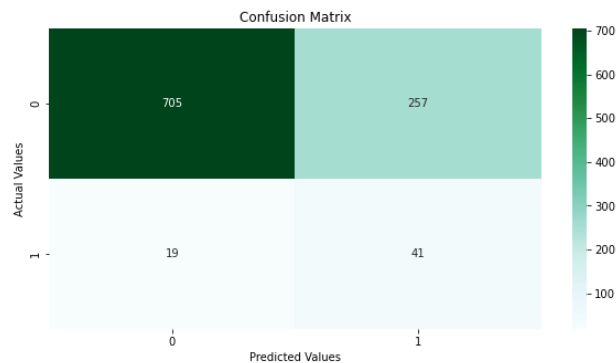

Confusion Matrix

```
In [70]:    1  test_xgb_proba = xgb.predict_proba(X_test)
            2  y_test_prob_xgb = test_xgb_proba[:, 1]
            3  y_test_prob_xgb
```

Out[70]:  array([0.00295475, 0.00448489, 0.00062596, ..., 0.0042313 , 0.32490686,
          0.09692448], dtype=float32)

```
In [71]:    1  # Test set peformance: Best threshold  F1
            2
            3  xgb_performance(y_test, y_test_prob_xgb)
```

Accuracy: 0.730
Precision at Best F1 : 0.13758389261744966
Recall at Best F1 : 0.6833333333333333
Best F1 Score : 0.22905027932960892



```
In [75]:    1  fpr_lor , tpr_lor, threshold_lor = roc_curve(y_test, y_test_prob_lor)
            2
            3  fpr_dtc , tpr_dtc, threshold_dtc = roc_curve(y_test, y_test_prob_dtc)
            4
            5  fpr_rfc , tpr_rfc, threshold_rfc = roc_curve(y_test, y_test_prob_rfc)
            6
            7  fpr_xgb , tpr_xgb, threshold_xgb = roc_curve(y_test, y_test_prob_xgb)
            8
            9  lor_AUC = roc_auc_score(y_test, y_test_prob_lor)
           10  dtc_AUC = roc_auc_score(y_test, y_test_prob_dtc)
           11  rfc_AUC = roc_auc_score(y_test, y_test_prob_rfc)
           12  xgb_AUC = roc_auc_score(y_test, y_test_prob_xgb)
```

```
In [76]:    1  print('AUC for Logistic Regression: ' , lor_AUC)
            2  print('AUC for Decision Tree: ' , dtc_AUC)
            3  print('AUC for Random Forest: ' , rfc_AUC)
            4  print('AUC for XGBoost: ' , xgb_AUC)
```

AUC for Logistic Regression:  0.8507796257796258
AUC for Decision Tree:  0.8311850311850312
AUC for Random Forest:  0.8430959805959806
AUC for XGBoost:  0.7785343035343035

```
In [78]:    1  plt.plot([0,1], ls = '--', linewidth=3, color = 'black')
            2  plt.plot([0,0],[1,0], c='.5')
            3  plt.plot([1,1],c='.5')
            4  plt.plot(fpr_lor, tpr_lor, linewidth=3, color = '#33CC33', label= "Logistic Regression Classifier ")
            5  plt.plot(fpr_dtc, tpr_dtc, linewidth=3, color = '#005aff', label= "Decision Tree Classifier")
            6  plt.plot(fpr_rfc, tpr_rfc, linewidth=3, color = '#ffa500', label= "Random Forest Classifier")
            7  plt.plot(fpr_xgb, tpr_xgb, linewidth=3, color = '#cc0000', label= "XGBoost Classifier")
            8  plt.legend()
            9  plt.xlabel("False Positive Rate")
           10  plt.ylabel("True Positive Rate")
           11  plt.title('ROC curve')
           12  plt.show()
```



## Conclusion

Return Contents

- For classification problems that have a severe class imbalance, the default threshold can result in poor performance. As such, a simple spproach to improving the performance of a classifier that predicts probabilities on an imbalanced classification problem is to tune the threshold used to map probabilities to class labels