

---

# Micro-Credit Defaulter Model

---



**FLIP ROBO**

---

ABHISHEK MISHRA



# **Micro-Credit Defaulter Model**

**SUBMITTED BY**  
**ABHISHEK MISHRA**

---

## **ACKNOWLEDGEMENT**

First of all, I extol the Almighty for pouring his blessings on me and giving me potentiality and opportunity to carry the work to its end with a success.

“It’s impossible to prepare a project report without the help and fillip of some people and certainly this project report is of no exception.”

At the commencement of this project report I would like to evince my deepest sense of gratitude to Ms. Astha Mishra, my honored mentor. Without her guidance, insightful decision, valuable comments and correction it would not have possible to reach up to this mark.

I would like to draw my gratitude to Flip Robo and Data Trained for providing me a suitable environment and guidance to complete my work. Last but not least thanks to the brilliant authors from where I have got the idea to carry out the project.

References were taken from various articles from Medium, KDnuggets, Towards Data Science, Machine Learning Mastery, Analytics Vidya, American Statistical Association, Research Gate and documentations of Python and Sklearn.

Abhishek Mishra

---

# CONTENTS

<b>ACKNOWLEDGEMENT</b>	<b>03</b>
<b>CHAPTER 1: INTRODUCTION</b>	<b>06</b>
<b>CHAPTER 2: ANALYTICAL PROBLEM FRAMING</b>	<b>09</b>
<b>CHAPTER 3: MODEL DEVELOPMENT</b>	<b>14</b>
<b>CHAPTER 4: CONCLUSION</b>	<b>26</b>

# **LIST OF FIGURES**

<b><u>FIGURES</u></b>	<b><u>PAGE NO</u></b>
<b>Fig.1: NULL VALUES IN THE DATASET</b>	<b>09</b>
<b>Fig.2: NEGATIVE VALUES IN DATASET</b>	<b>10</b>
<b>Fig.3: INFO OF THE DATASET</b>	<b>11</b>
<b>Fig.5: DATASET</b>	<b>13</b>
<b>Fig.6: PROCESSED DATASET</b>	<b>13</b>
<b>Fig.7: DATA PREPARATION</b>	<b>14</b>
<b>Fig.8: MODELLING LOGISTIC REGRESSION</b>	<b>16</b>
<b>Fig.9: CROSS VAL SCORE LOGISTIC REGRESSION</b>	<b>16</b>
<b>Fig.10: RANDOMIZED SEARCH CV RFC</b>	<b>17</b>
<b>Fig.11: AUC OF RANDOM FOREST CLASSIFIER</b>	<b>17</b>
<b>Fig.12: MODELLING XGBOOST</b>	<b>18</b>
<b>Fig.13: HEATMAP OF CONFUSION MATRIX OF XGBOOST</b>	<b>18</b>
<b>Fig.14: RESULTS</b>	<b>19</b>
<b>Fig.15: NULL</b>	<b>20</b>
<b>Fig.16: DEFAULTERS AND NON-DEFAULTERS</b>	<b>20</b>
<b>Fig.17: AGE ON NETWORK</b>	<b>21</b>
<b>Fig.18: MAIN ACCOUNT BALANCE OF LAST 30 DAYS</b>	<b>21</b>
<b>Fig.19: LOAN TAKEN IN 30 DAYS</b>	<b>22</b>
<b>Fig.20: RECHARGE FREQUENCY IN 90 DAYS</b>	<b>22</b>
<b>Fig.21: FREQUENCY DISTRIBUTION</b>	<b>23</b>
<b>Fig.22: PROBABILITY DENSITY</b>	<b>24</b>
<b>Fig.23: MODEL SAVING</b>	<b>25</b>

---

# CHAPTER-1

## INTRODUCTION

### 1.1 BUSINESS PROBLEM

Our client is a Telecom Industry having fixed wireless telecommunication network; they are keen to provide better services at low price range.

They are collaborating with an MFI to provide micro-credit on mobile balances to be paid back in 5 days. The consumer will be considered as a defaulter if he doesn't pay back the loaned amount within 5 days. The consumer has to pay 6 Indonesian Rupiah for a loan of 5 Indonesian Rupiah and for a loan of 10 Indonesian Rupiah the payback amount should be 12 Indonesian Rupiah.

### 1.2 BACKGROUND

As a result of incapacity of development and traditional banks to effectively finance the low-income population of the world, microfinance seems like a continuum between pure capitalism and socialism economies (World Bank, 2008). Access to formal financial services has been limited for many, if not most, of the world's poorest: more than 2.5 billion people do not use formal financial services.

In developing countries, the problems triggered by informational lopsidedness that are distinctive to credit markets are intensified since low-income people lack collateral that can be provided against loans and because of the weak legal system enforcement cannot be possible,

MFS becomes very useful when targeting especially the unbanked poor families living in remote areas with not much sources of income. The Microfinance services (MFS) provided by MFI are Group Loans, Agricultural Loans, Individual Business Loans and so on.

---

Many microfinance institutions (MFI), experts and donors are supporting the idea of using mobile financial services (MFS) which they feel are more convenient and efficient, and cost saving, than the traditional high-touch model used since long for the purpose of delivering microfinance services. In current world scenario microfinance is widely accepted as a poverty-reduction tool, and it represents \$70 billion in outstanding loans and a global outreach of 200 million clients.

In current world scenario the importance of communication in a person life is known to all. So here our client having fixed wireless telecommunication network trying to introduce services low price in collaboration with a MFI. The low-income group are the main target of this new launch.

As we know that the telecom sector is one of the most competitive fields so this data is very helpful in understanding the problem for the lower-class people specially by providing them the facility of network and the credit amount provided by the help of MFI and MFS. From this data we get to know that what the criteria to become defaulters and successor are. And the useful information from the data to know how much amount people spend on data recharge or on the main balance recharge.

## **1.3 MOTIVATION FOR THE PROBLEM UNDERTAKEN**

The project was the first provided to me by FlipRobo as a part of the internship programme. The exposure to real world data and the opportunity to deploy my skillset in solving a real time problem has been the primary motivation. Further diving into the dataset, the motive is to help the poor or low-income band to have continuous access to their mobile accounts and to make emergency calls even when they do not have account balance making use of the loan facility. Alternatively, it is also important to ensure that the provider does not incur a loss for providing the facility. Hence, the entire focused is on building a model that can effectively predicts a defaulter by using the historical data which would help in approval process of the loans to end users with a clean sheet.





# CHAPTER-2

## ANALYTICAL PROBLEM FRAMING

### 2.1 Analytical Modeling of the Problem

The dataset provided has a shape of (209593, 37). Here the target or the dependent variable named “Label” have two different distinct values 0 and 1. Where 0 represents the defaulter & 1 represents the non-defaulter category of people. As the target is giving binary output so, it is a classification-based problem. Here the dataset has no null values and no duplicate values.

Different values of the dataset like the count, mean, standard deviation, min, 25%,50%,75%, max can be obtained using df.describe() function. The values obtained from this function shows that maximum values have high standard deviation.

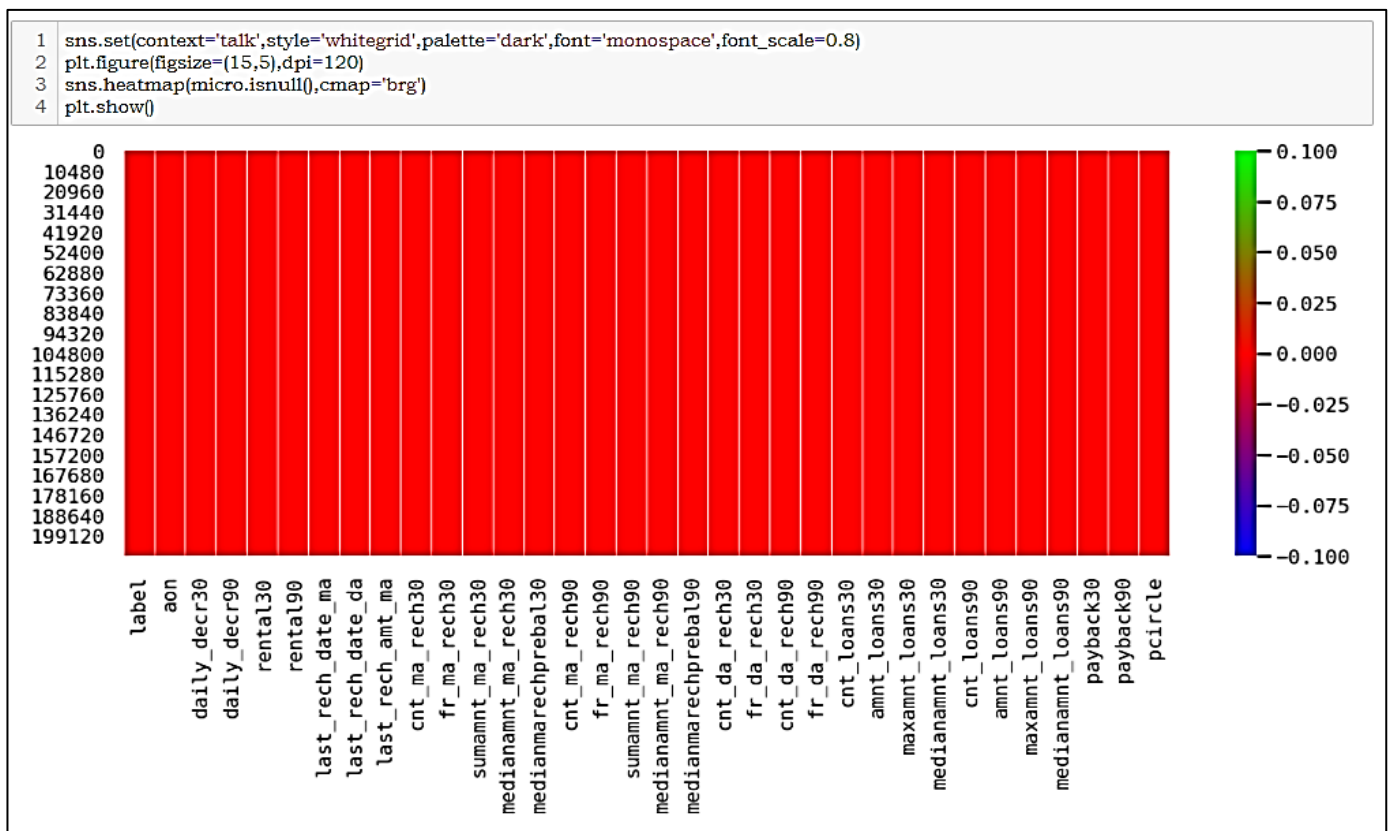


Fig.1: NULL VALUES IN THE DATASET

The describe function showed that few columns have values less than zero which is absurd and the of boxplot have shown the outliers in the dataset but we can't remove all of them using Zscore method as this will lead to loss of more than 10% of data so all the outliers were removed using proper condition.

```
1 mfi.drop(mfi.index[mfi['aon']<0], inplace = True)
2 mfi.drop(mfi.index[mfi['daily_decr30']<0], inplace = True)
3 mfi.drop(mfi.index[mfi['daily_decr90']<0], inplace = True)
4 mfi.drop(mfi.index[mfi['rental30']<0], inplace = True)
5 mfi.drop(mfi.index[mfi['rental90']<0], inplace = True)
6 mfi.drop(mfi.index[mfi['last_rech_date_ma']<0], inplace = True)
7 mfi.drop(mfi.index[mfi['last_rech_date_da']<0], inplace = True)
8 mfi.drop(mfi.index[mfi['medianmarechprebal30']<0], inplace = True)
9 mfi.drop(mfi.index[mfi['cnt_ma_rech90']<0], inplace = True)
10 mfi.drop(mfi.index[mfi['medianmarechprebal90']<0], inplace = True)
```

999860.75 days or 2739.34 year on cellular network is impossible so removing any data where aon is greater than 7300 days or 20 years.

```
1 mfi.drop(mfi.index[mfi['aon']>7301], inplace = True)
```

Fig.2: NEGATIVE VALUES IN DATASET

## 2.2 DATA SOURCES AND THEIR FORMATS

The data source obtained was of CSV form and have 209593 rows and 37 columns. The df.info() function shows the data types of the columns.

1	micro.info()
---	--------------

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 209593 entries, 0 to 209592
Data columns (total 34 columns):
#   Column                               Non-Null Count  Dtype
---  -
0   label                                209593 non-null  int64
1   aon                                  209593 non-null  float64
2   daily_decr30                        209593 non-null  float64
3   daily_decr90                        209593 non-null  float64
4   rental30                            209593 non-null  float64
5   rental90                            209593 non-null  float64
6   last_rech_date_ma                   209593 non-null  float64
7   last_rech_date_da                   209593 non-null  float64
8   last_rech_amt_ma                    209593 non-null  int64
9   cnt_ma_rech30                       209593 non-null  int64
10  fr_ma_rech30                        209593 non-null  float64
11  sumamnt_ma_rech30                   209593 non-null  float64
12  medianamnt_ma_rech30                209593 non-null  float64
13  medianmarechprebal30                209593 non-null  float64
14  cnt_ma_rech90                       209593 non-null  int64
15  fr_ma_rech90                        209593 non-null  int64
16  sumamnt_ma_rech90                   209593 non-null  int64
17  medianamnt_ma_rech90                209593 non-null  float64
18  medianmarechprebal90                209593 non-null  float64
19  cnt_da_rech30                       209593 non-null  float64
20  fr_da_rech30                        209593 non-null  float64
21  cnt_da_rech90                       209593 non-null  int64
22  fr_da_rech90                        209593 non-null  int64
23  cnt_loans30                         209593 non-null  int64
24  amnt_loans30                        209593 non-null  int64
25  maxamnt_loans30                     209593 non-null  float64
26  medianamnt_loans30                  209593 non-null  float64
27  cnt_loans90                         209593 non-null  float64
28  amnt_loans90                        209593 non-null  int64
29  maxamnt_loans90                     209593 non-null  int64
30  medianamnt_loans90                  209593 non-null  float64
31  payback30                           209593 non-null  float64
32  payback90                           209593 non-null  float64
33  pcircle                             209593 non-null  object
dtypes: float64(21), int64(12), object(1)
memory usage: 54.4+ MB

```

**Fig.3: INFO OF THE DATASET**

<b>Label</b>	<b>Flag indicating whether the user paid back the credit amount within 5 days of issuing the loan {1: success, 0: failure}</b>
<b>Msisdn</b>	<b>mobile number of users</b>
<b>Aon</b>	<b>age on cellular network in days</b>
<b>daily_decr30</b>	<b>Daily amount spent from main account, averaged over last 30 days (in Indonesian Rupiah)</b>
<b>daily_decr90</b>	<b>Daily amount spent from main account, averaged over last 90 days (in Indonesian Rupiah)</b>
<b>rental30</b>	<b>Average main account balance over last 30 days</b>
<b>rental90</b>	<b>Average main account balance over last 90 days</b>
<b>last_rech_date_ma</b>	<b>Number of days till last recharge of main account</b>
<b>last_rech_date_da</b>	<b>Number of days till last recharge of data account</b>
<b>last_rech_amt_ma</b>	<b>Amount of last recharge of main account (in Indonesian Rupiah)</b>
<b>cnt_ma_rech30</b>	<b>Number of times main account got recharged in last 30 days</b>
<b>fr_ma_rech30</b>	<b>Frequency of main account recharged in last 30 days</b>
<b>sumamnt_ma_rech30</b>	<b>Total amount of recharge in main account over last 30 days (in Indonesian Rupiah)</b>
<b>medianamnt_ma_rech30</b>	<b>Median of amount of recharges done in main account over last 30 days at user level (in Indonesian Rupiah)</b>

<b>medianmarechprebal30</b>	<b>Median of main account balance just before recharge in last 30 days at user level (in Indonesian Rupiah)</b>
<b>cnt_ma_rech90</b>	<b>Number of times main account got recharged in last 90 days</b>
<b>fr_ma_rech90</b>	<b>Frequency of main account recharged in last 90 days</b>
<b>sumamnt_ma_rech90</b>	<b>Total amount of recharge in main account over last 90 days (in Indonesian Rupee)</b>
<b>medianamnt_ma_rech90</b>	<b>Median of amount of recharges done in main account over last 90 days at user level (in Indonesian Rupee)</b>
<b>medianmarechprebal90</b>	<b>Median of main account balance just before recharge in last 90 days at user level (in Indonesian Rupee)</b>
<b>cnt_da_rech30</b>	<b>Number of times data account got recharged in last 30 days</b>
<b>fr_da_rech30</b>	<b>Frequency of data account recharged in last 30 days</b>
<b>cnt_da_rech90</b>	<b>Number of times data account got recharged in last 90 days</b>
<b>fr_da_rech90</b>	<b>Frequency of data account recharged in last 90 days</b>
<b>cnt_loans30</b>	<b>Number of loans taken by user in last 30 days</b>
<b>amnt_loans30</b>	<b>Total amount of loans taken by user in last 30 days</b>
<b>maxamnt_loans30</b>	<b>maximum amount of loan taken by the user in last 30 days</b>
<b>medianamnt_loans30</b>	<b>Median of amounts of loan taken by the user in last 30 days</b>
<b>cnt_loans90</b>	<b>Number of loans taken by user in last 90 days</b>
<b>amnt_loans90</b>	<b>Total amount of loans taken by user in last 90 days</b>
<b>maxamnt_loans90</b>	<b>maximum amount of loan taken by the user in last 90 days</b>
<b>medianamnt_loans90</b>	<b>Median of amounts of loan taken by the user in last 90 days</b>
<b>payback30</b>	<b>Average payback time in days over last 30 days</b>
<b>payback90</b>	<b>Average payback time in days over last 90 days</b>
<b>Pcircle</b>	<b>telecom circle</b>
<b>Pdate</b>	<b>Date</b>

**TABLE 1: METADATA**

## 2.3 DATA PREPROCESSING

After loading the dataset, the null and duplicate values were checked, the absurd values and skewed values were removed with proper conditions. The columns named 'msisdn', 'pcircle' and 'date' were dropped as they served no purpose in modelling. After preprocessing the dataset shape has been reduced to 189339,33.

```
1 print('Shape of the dataset - ',micro.shape)
2 print('\nColumns in the dataset-\n\n',micro.columns.values)

Shape of the dataset - (209593, 37)

Columns in the dataset-

['Unnamed: 0' 'label' 'msisdn' 'aon' 'daily_decr30' 'daily_decr90'
'rental30' 'rental90' 'last_rech_date_ma' 'last_rech_date_da'
'last_rech_amt_ma' 'cnt_ma_rech30' 'fr_ma_rech30' 'sumamnt_ma_rech30'
'medianamnt_ma_rech30' 'medianmarechprebal30' 'cnt_ma_rech90'
'fr_ma_rech90' 'sumamnt_ma_rech90' 'medianamnt_ma_rech90'
'medianmarechprebal90' 'cnt_da_rech30' 'fr_da_rech30' 'cnt_da_rech90'
'fr_da_rech90' 'cnt_loans30' 'amnt_loans30' 'maxamnt_loans30'
'medianamnt_loans30' 'cnt_loans90' 'amnt_loans90' 'maxamnt_loans90'
'medianamnt_loans90' 'payback30' 'payback90' 'pcircle' 'pdate']
```

Fig.5: DATASET

```
1 print(' Earlier the shape of dataset with outliers was:',micro.shape,'\n The shape of the dataset after outlier removal is:',mfi.shape,
2 '\n Percentage of data_loss:', data_loss)

Earlier the shape of dataset with outliers was: (209593, 33)
The shape of the dataset after outlier removal is: (189339, 33)
Percentage of data_loss: 9.663
```

Fig.6: PROCESSED DATASET

## 2.4 HARDWARE & TOOL USED

In this project the below mentioned machine, IDE and packages were used;

HARDWARE	LAPTOP: ASUS TUF A17 OS: WIN 10 HOME BASIC PROCESSOR: AMD RYZEN 7 4800H RAM: 16GB VRAM: 6GB NVIDIA GTX 1660Ti
LANGUAGE	Python 3.8
IDE	JUPYTER NOTEBOOK 6.0.3
PACKAGES	PANDAS, NUMPY, SCIPY, SKLEARN, MATPLOTLIB, SEABORN

TABLE 2: DEVICE AND TOOLS

# CHAPTER-3

## DEVELOPMENT AND EVALUATION

### 3.1 IDENTIFICATION OF POSSIBLE PROBLEM-SOLVING APPROACHES

After loading the dataset, the absurd values, unwanted columns and the skewness were removed. Which leads to removal of approximately 9% of data. After this the dataset is split into two-part, one part named 'x' contains all the independent columns and one part 'y' which contains the target column.

The x part then normalized using standard scaler and y part is converted into numpy array and reshaped so that the shape of x and y remain same. Once this part is completed, they can be sent for modelling.

```
1 x=mfi.drop(['label'],axis=1)
2 y=mfi.label
```

```
1 ss=StandardScaler()
2 x=ss.fit_transform(x)
3 print(x)
```

```
[[-0.63497471  0.53116441  0.50874106 ... -0.23671909  2.6639062
  2.35395517]
 [ 0.46850625  1.00573954  0.97600258 ... -0.23671909 -0.85094496
 -0.92930934]
 [ 0.14051447  0.26206303  0.24207332 ... -0.23671909 -0.85094496
 -0.92930934]
 ...
 [ 0.87331506  0.99772413  0.96976859 ... -0.23671909  0.81227363
  0.40891675]
 [ 1.48936289  1.01598956  0.9883761 ... -0.23671909 -0.85094496
  1.38544149]
 [ 1.38457446  0.66369256  0.64182028 ... -0.23671909 -0.85094496
 -0.92930934]]
```

```
1 y=np.array(y)
2 y=y.reshape(-1,1)
```

```
1 print('shape of x:',x.shape,'\nshape of y:',y.shape)
```

```
shape of x: (189339, 32)
shape of y: (189339, 1)
```

**Fig.7: DATA PREPARATION**

## 3.2 TESTING OF IDENTIFIED APPROACHES

Once the normalized data were obtained, they can be sent for modelling. Here the output column 'label' is generating binary output so it's a classification-based problem and we can use the following algorithms for modelling and the highest performing algorithms to get our final model;

- Logistic Regression
- Decision Tree Classifier
- Gaussian NB
- Random Forest Classifier
- XgBoost Classifier

With the help of RandomizedSearchCV hyper parameter tuning will be done and the best parameters for each model will be found.

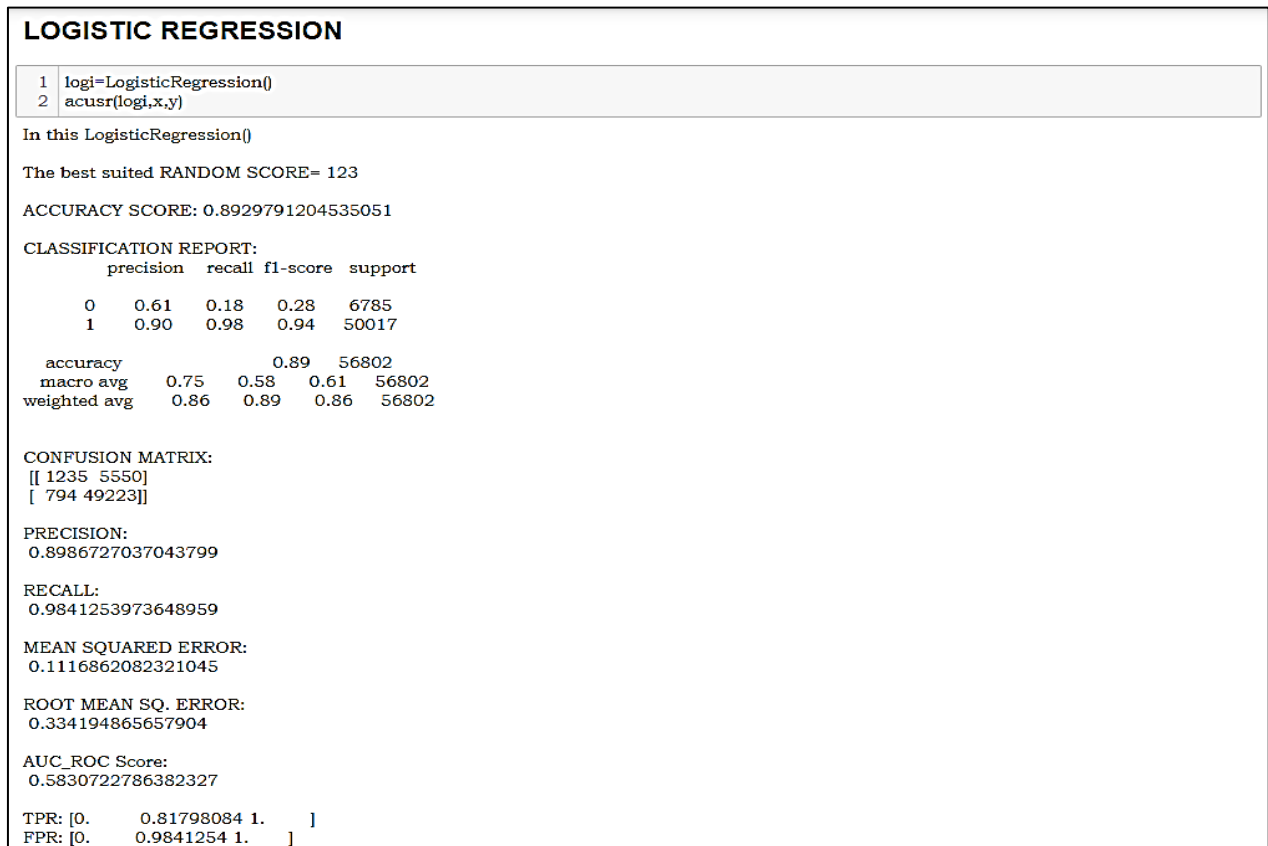
During modelling various metrics like f1 score, confusion matrix, accuracy score, classification report, roc curve, auc, roc auc score, mean squared error, precision score, recall score will be used to determine the performance of the model.

To check whether the model suffering from over fitting or underfitting cross val score will be used. To view the best performing model AUC Curve will be used. At the end the best model will be saved using Joblib library.

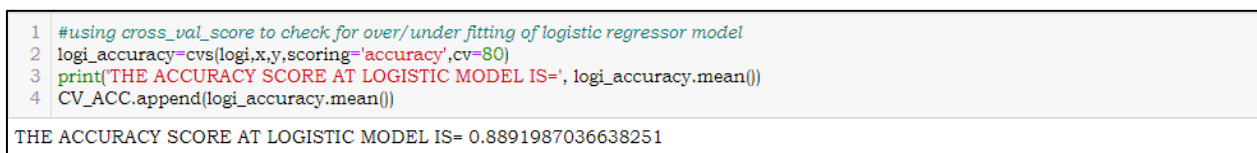
## 3.3 RUNNING AND EVALUATION

Here I have generated a function which will find the best random score for the model in a range of 25 to 180. It'll also show the confusion matrix, accuracy score, classification report, roc curve, auc, roc auc score, mean squared error, precision score, recall score, tpr, fpr at that random score.

The values so obtained will be then append to their respective lists. This function will also draw the AUC Curve and heat map of the confusion matrix. Below are the few images of modelling, heatmap and auc curve.



**Fig.8: MODELLING LOGISTIC REGRESSION**



**Fig.9: CROSS VAL SCORE LOGISTIC REGRESSION**



## RANDOM FOREST CLASSIFIER

```
1 rfc=RandomForestClassifier()
2 rfc_para={'n_estimators':[300,350,400],'max_depth':[3,6,9],'criterion':['gini','entropy']}
3 rfc_rsv=rsv(rfc,rfc_para,cv=30,n_jobs=-1)
4 rfc_rsv.fit(x,y)
5 print(rfc_rsv)
6 print('\nbest score=',rfc_rsv.best_score_)
7 print('\nbest parameters for RFC=\n',rfc_rsv.best_params_)
```

```
RandomizedSearchCV(cv=30, estimator=RandomForestClassifier(), n_jobs=-1,
    param_distributions={'criterion': ['gini', 'entropy'],
        'max_depth': [3, 6, 9],
        'n_estimators': [300, 350, 400]})
```

best score= 0.909126004969696

best parameters for RFC=  
{'n\_estimators': 300, 'max\_depth': 9, 'criterion': 'gini'}

Fig.10: RANDOMIZED SEARCH CV RFC

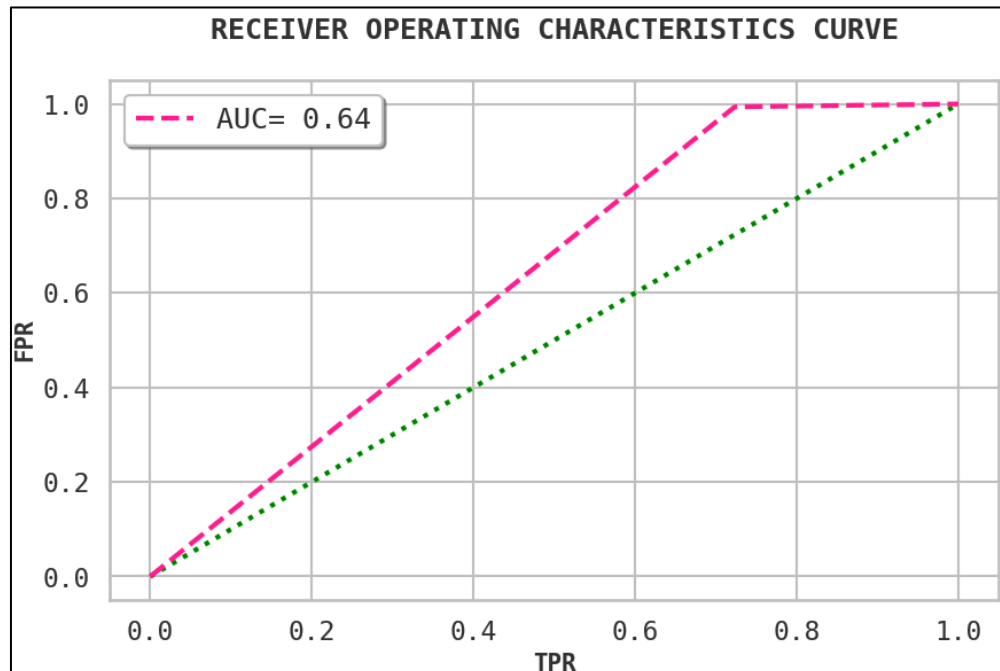


Fig.11: AUC OF RANDOM FOREST CLASSIFIER

## XGBOOST CLASSIFIER

```

1 XGB = XGBClassifier()
2 acusr(XGB,x,y)

```

In this XGBClassifier(base\_score=0.5, booster='gbtree', colsample\_bylevel=1, colsample\_bynode=1, colsample\_bytree=1, gamma=0, gpu\_id=-1, importance\_type='gain', interaction\_constraints='', learning\_rate=0.300000012, max\_delta\_step=0, max\_depth=6, min\_child\_weight=1, missing=nan, monotone\_constraints='()', n\_estimators=100, n\_jobs=0, num\_parallel\_tree=1, random\_state=0, reg\_alpha=0, reg\_lambda=1, scale\_pos\_weight=1, subsample=1, tree\_method='exact', validate\_parameters=1, verbosity=None)

The best suited RANDOM SCORE= 123

ACCURACY SCORE: 0.9191753811485511

CLASSIFICATION REPORT:

	precision	recall	f1-score	support
0	0.76	0.43	0.55	6785
1	0.93	0.98	0.95	50017
accuracy			0.92	56802
macro avg	0.84	0.70	0.75	56802
weighted avg	0.91	0.92	0.90	56802

CONFUSION MATRIX:

```

[[ 2886 3899]
 [ 905 49112]]

```

PRECISION:

```

0.9264492275188169

```

RECALL:

```

0.9819061519083512

```

MEAN SQUARED ERROR:

```

0.08457448681384458

```

ROOT MEAN SQ. ERROR:

```

0.2908169300674302

```

AUC\_ROC Score:

```

0.7036280943771676

```

TPR: [0. 0.57464996 1. ]

FPR: [0. 0.98190615 1. ]

Fig.12: MODELLING XGBOOST

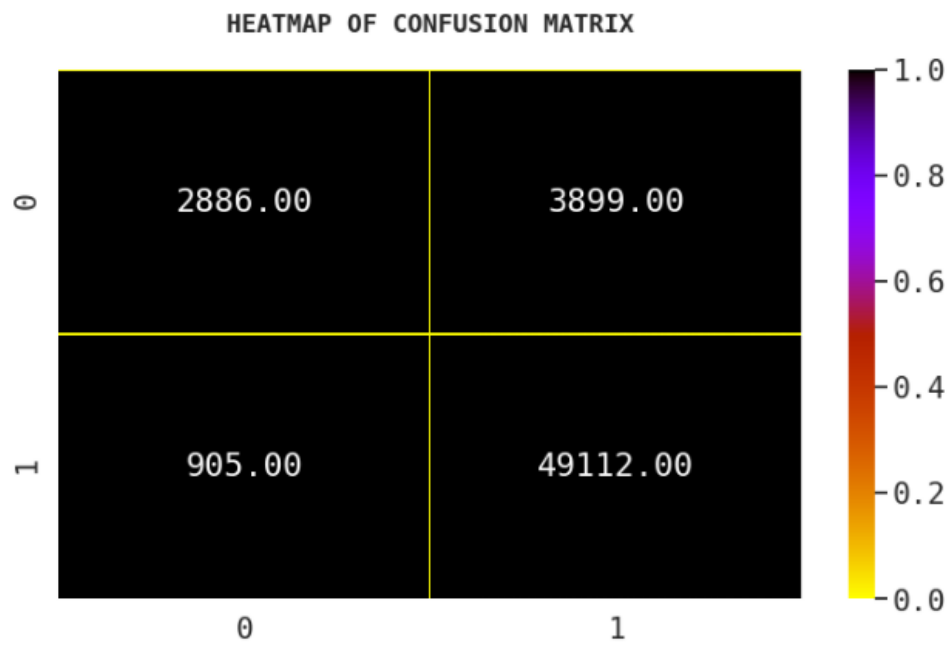


Fig.13: HEATMAP OF CONFUSION MATRIX OF XGBOOST

### 3.3 METRICE OF EVALUATION

In the modeling I have chosen metrics like Precision, Recall, Mean Squared Error, Root Mean Square Error, Classification Report, Accuracy score, Confusion Matrix, AUC, tpr, fpr and Cross val Score as my evaluation criteria. All the values were stored in a list and later they were saved in form of a DataFrame for proper evaluation and visualization of the values. Basing on the values the best model has been selected.

```
1 model=["LOGISTIC","DTC","GAUSSIAN-NB","RANDOM FOREST","XGBOOST"]
2 results = pd.DataFrame({'MODEL':model,'Acuracy':ACCURACY,'Precision': PRECESION , 'Recall': RECALL,
3                          'RMSE':RMSE,'MSE':MSE,'AUC':AUC,'TPR':TPR,'FPR':FPR,'CV_ACCURACY':CV_ACC})
4
5
6
7 results.style.set_properties(**{'background-color':'midnightblue','color': 'gold','border-color': 'darkorange'})
```

	MODEL	Acuracy	Precision	Recall	RMSE	MSE	AUC	TPR	FPR	CV_ACCURACY
0	LOGISTIC	0.892979	0.898673	0.984125	0.334195	0.111686	0.583072	[0. 0.81798084 1. ]	[0. 0.9841254 1. ]	0.889199
1	DTC	0.911975	0.922459	0.977548	0.303524	0.092127	0.685900	[0. 0.60574797 1. ]	[0. 0.97754763 1. ]	0.910911
2	GAUSSIAN-NB	0.797278	0.947967	0.807166	0.456961	0.208813	0.740281	[0. 0.3266028 1. ]	[0. 0.80716556 1. ]	0.792890
3	RANDOM FOREST	0.912855	0.910144	0.993922	0.302915	0.091757	0.635281	[0. 0.72336035 1. ]	[0. 0.99392207 1. ]	0.908909
4	XGBOOST	0.919175	0.926449	0.981906	0.290817	0.084574	0.703628	[0. 0.57464996 1. ]	[0. 0.98190615 1. ]	0.916293

Fig.14: RESULTS

### 3.4 VISUALIZATION

Visualization is a part of EDA which helps to understand the data better. Here to understand the data various plots like Heatmap, Countplot, Bar graph were plotted.

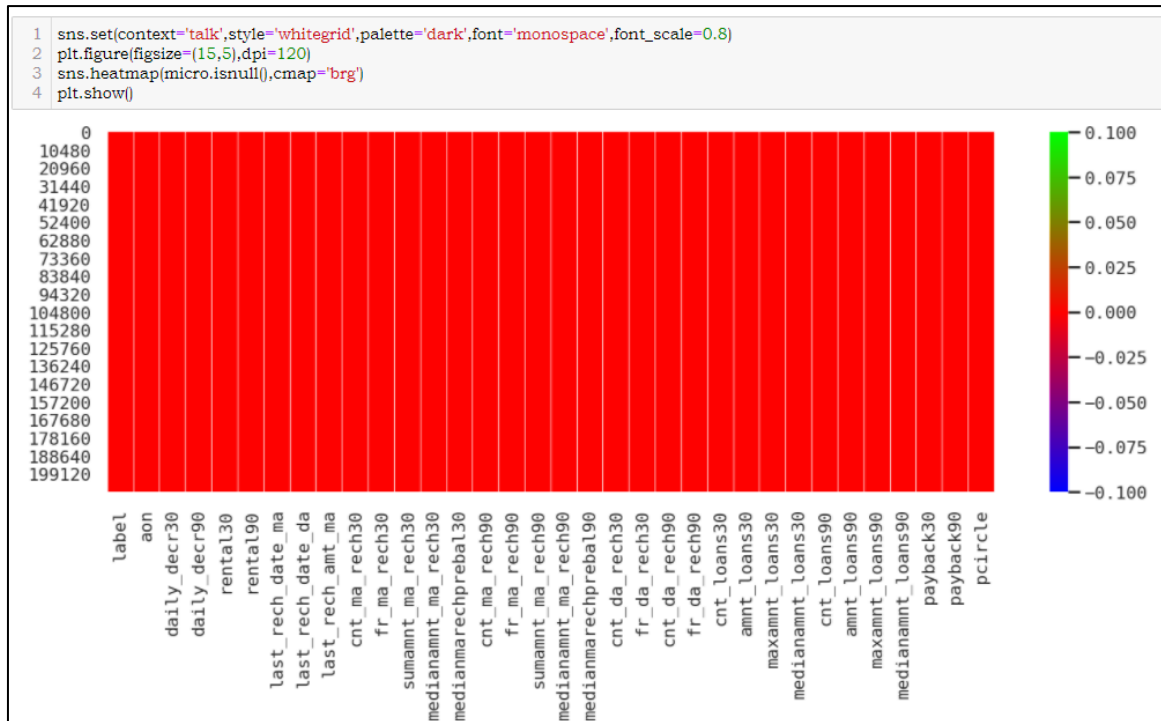


Fig.15: NULL

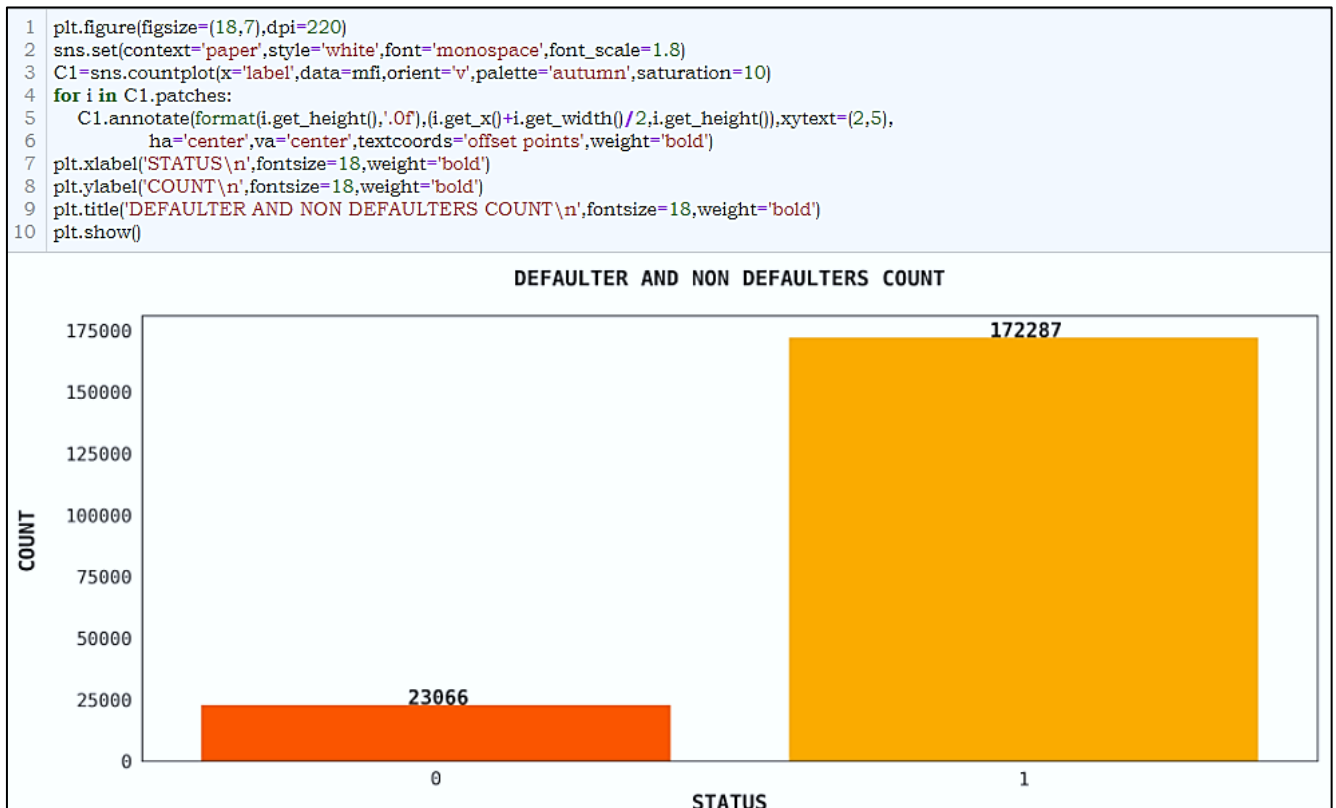


Fig.16: DEFAULTERS AND NON-DEFAULTERS

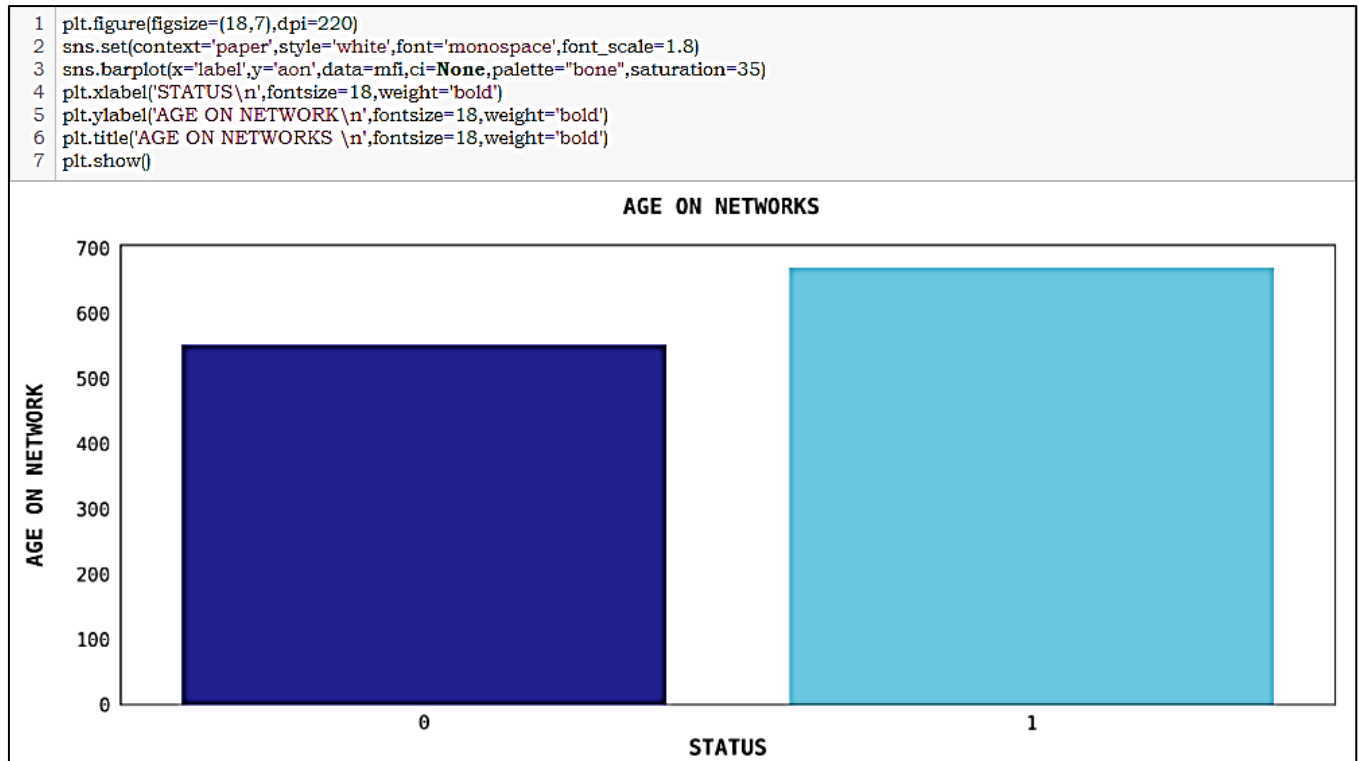


Fig.17: AGE ON NETWORK

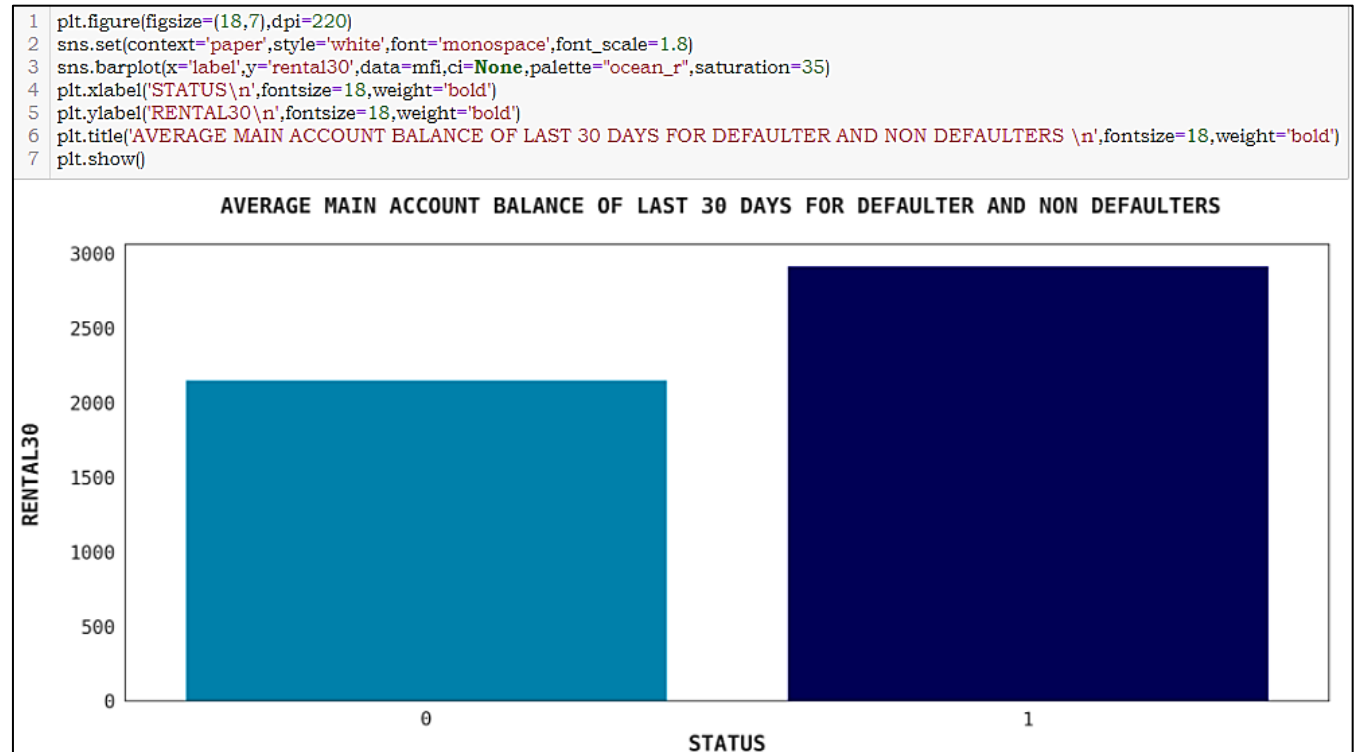


Fig.18: MAIN ACCOUNT BALANCE OF LAST 30 DAYS



Fig.19: LOAN TAKEN IN 30 DAYS



Fig.20: RECHARGE FREQUENCY IN 90 DAYS

```

1 sns.set_context('talk',font_scale=1.9)
2 sns.set(context='notebook',style='whitegrid',palette='dark',font='monospace',font_scale=1.3)
3 plt.figure(figsize=(12,3),dpi=120)
4 mfi.hist(figsize=(80,80),grid=True)
5 plt.show()

```



Fig.21: FREQUENCY DISTRIBUTION

```

1 plt.figure(figsize=(70,60))
2 sns.set(style='darkgrid')
3 sns.set_context('talk',font_scale=1.1)
4 for i in range(len(mf1cn)):
5     plt.subplot(7,5,i+1)
6     sns.distplot(mf1[mf1cn[i]], color='dodgerblue', bins=80,kde=True,rug=True)
7     plt.tight_layout()

```



Fig.22: PROBABILITY DENSITY



### 3.5 INTERPRETATION

Basing on the result obtained XgBoost Classifier have performed well and has given better result as compared to other models. Here the accuracy score is 91.91%, Precision is 92.64%, AUC value is 70% and after cross validation the accuracy becomes 91.62%. As this model have all the desired characters and performance so XgBoost has been selected as final model and it will be saved using joblib library.

```
1 joblib.dump(XGB,'MFI.obj')  
['MFI.obj']
```

**Fig.23: MODEL SAVING**

# CHAPTER-4

## CONCLUSION

### 4.1 KEY FINDINGS

From the above analysis the below mentioned results were achieved which depicts the chances and conditions of a user being a defaulter

- Users who are new to the network have a chance of being a defaulter.
- Users maintaining a low average of main account balance for last 30 days have a chance of being a defaulter.
- Users with low frequency of recharge in last 90 days have a chance of being a defaulter.
- Users who have taken low amount of loan in last 30 days have a chance of being a defaulter.

### 4.2 LEARNING OUTCOMES OF THE STUDY

The dataset obtained was imbalanced and have values which are absurd. Few values were negative while few are too high to believe. So those value were treated using proper conditions.

The count plots, bar plot, heatmap gave a vivid idea to understand the behavioral patterns which differentiate the defaulters and non-defaulters.

In this classification problem the model created using XgBoost classifier worked efficiently as compared to other models.

### 4.3 LIMITATIONS OF THIS WORK AND SCOPE FOR FUTURE WORK

There were certain limitations found in this dataset.

- The metadata provided could have been more precise.

- 
- The 'pcircle' column have shown a single network circle which put a limitation to the analysis.
  - Dataset have a large number of absurd values which must be taken into consideration.