



Project Report: **“Micro Credit Defaulter”**



Submitted by:
Rahul Nanduri
SME: Kushboo Gurg

ACKNOWLEDGMENT

Microfinance institutions play a major role in economic development in many developing countries. However, many of these microfinance institutions are faced with the problem of default because of the non-formal nature of the business and individuals they lend money. This study seeks to find the determinants of credit default in microfinance institutions.

First and foremost, I would like to warmly thank the “Flip Robo” team, who has given me this opportunity to deal with an interesting project on ML and it has helped me to improve my analysis skills.

INTRODUCTION

- **Business Problem Framing**

A Microfinance Institution (MFI) is an organization that offers financial services to low income populations. 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. Though, the MFI industries is primarily focusing on low income families and are very useful in such areas, the implementation of MFS has been uneven with both significant challenges and successes.

Today, microfinance is widely accepted as a poverty-reduction tool, representing \$70 billion in outstanding loans and a global outreach of 200 million clients.

We are working with one such client that is in Telecom Industry. They are a fixed wireless telecommunications network provider. They have launched various products and have developed its business and organization based on the budget operator model, offering better products at Lower Prices to all value conscious customers through a strategy of disruptive innovation that focuses on the subscriber.

- **Conceptual Background of the Domain Problem**

Microfinance is a proven tool for fighting poverty on a large scale. It provides very small loans, or micro-loans, to poor people, mostly women, to start or expand very small, self-sufficient businesses. Through their own ingenuity and drive, and with the support of the lending microfinance institution (MFI), poor women are able start their journey out of poverty

Unlike commercial loans, no collateral is required for a micro-loan and it is usually repaid within six months to a year. Those funds are then recycled as other loans, keeping money working and in the hands of borrowers.

The sample data is provided to us from our client database. It is hereby given for this exercise. In order to improve the selection of customers for the credit, the client wants some predictions that could help them in further investment and improvement in selection of customers.

We have to build a model which can be used to predict in terms of a probability for each loan transaction, whether the customer will be paying back the loaned amount within 5 days of insurance of loan. In this case, Label '1' indicates that the loan has been paid i.e. Non- defaulter, while, Label '0' indicates that the loan has not been paid i.e. defaulter.

- **Review of Literature**

An attempt has been made in this report to review the available literature in the area of microfinance. Approaches to microfinance, issues related to measuring social impact versus profitability of MFIs, issue of sustainability, variables impacting sustainability, the effect of regulations on profitability and impact assessment of MFIs have been summarized in the below report. We hope that the below report literature will provide a platform for further research and help the industry to combine theory and practice to take microfinance forward and contribute to alleviating the poor from poverty.

- **Motivation for the Problem Undertaken**

I have to model the micro-credit defaulters with the available independent variables. This model will then be used by the management to understand how the customer is considered a defaulter or non-defaulter based on the independent variables. They can accordingly manipulate the strategy of the firm and concentrate on areas that will yield high returns. Further, the model will be a good way for the management to understand whether the customer will be paying back the loaned amount within 5 days of insurance of loan. The **relationship between predicting defaulter and the economy** is an important motivating factor for predicting micro credit defaulter model.

Analytical Problem Framing

- **Mathematical/ Analytical Modeling of the Problem**

In this particular problem, I had label as my target column and it was having two classes Label '1' indicates that the loan has been paid i.e. Non- defaulter, while, Label '0' indicates that the loan has not been paid i.e. defaulter. So clearly it is a binary classification problem and I have to use all classification algorithms while building the model. There were no null values in the dataset. Also, I observed some unnecessary entries in some of the columns like in some columns I found more than 90% zero values so I decided to drop those columns. If I keep those columns as it is, it will create high skewness in the model. To get a better insight into the features I have used plotting like distribution plot, bar plot and count plot. With this plotting I was able to understand the relation between the features in a better manner. Also, I found outliers and skewness in the dataset so I removed outliers using the percentile method and I removed skewness using yeo-Johnson method. I have used all the classification algorithms while building model then tuned the best model and saved the best model. At last I have predicted the label using saved model.

- **Data Sources and their formats**

The data was collected for my internship company – Flip Robo technologies in excel format. The sample data is provided to us from our client database. It is hereby given to us for this exercise. In order to improve the selection of customers for the credit, the client wants some predictions that could help them in further investment and improvement in the selection of customers. Also, my dataset was having 209593 rows and 36 columns including the target. In this particular dataset I have object, float, and integer types of data. The information about features is as follows.

- **Features Information:**

1. label : Flag indicating whether the user paid back the credit amount within 5 days of issuing the loan{1:success, 0:failure}
2. msisdn : mobile number of user
3. aon : age on cellular network in days
4. daily_decr30 : Daily amount spent from main account, averaged over last 30 days (in Indonesian Rupiah)
5. daily_decr90 : Daily amount spent from main account, averaged over last 90 days (in Indonesian Rupiah)
6. rental30 : Average main account balance over last 30 days
7. rental90 : Average main account balance over last 90 days
8. last_rech_date_ma : Number of days till last recharge of main account
9. last_rech_date_da: Number of days till last recharge of data account
10. last_rech_amt_ma : Amount of last recharge of main account (in Indonesian Rupiah)
11. cnt_ma_rech30 : Number of times main account got recharged in last 30 days
12. fr_ma_rech30 : Frequency of main account recharged in last 30 days
13. sumamnt_ma_rech30 : Total amount of recharge in main account over last 30 days (in Indonesian Rupiah)
14. medianamnt_ma_rech30 : Median of amount of recharges done in main account over last 30 days at user level (in Indonesian Rupiah)
15. medianmarechprebal30 : Median of main account balance just before recharge in last 30 days at user level (in Indonesian Rupiah)
16. cnt_ma_rech90 : Number of times main account got recharged in last 90 days
17. fr_ma_rech90 : Frequency of main account recharged in last 90 days
18. sumamnt_ma_rech90 : Total amount of recharge in main account over last 90 days (in Indonesian Rupiah)
19. medianamnt_ma_rech90 : Median of amount of recharges done in main account over last 90 days at user level (in Indonesian Rupiah)

20. medianmarechprebal90 : Median of main account balance just before recharge in last 90 days at user level (in Indonesian Rupiah)
21. cnt_da_rech30 : Number of times data account got recharged in last 30 days
22. fr_da_rech30: Frequency of data account recharged in last 30 days
23. cnt_da_rech90 : Number of times data account got recharged in last 90 days
24. fr_da_rech90 : Frequency of data account recharged in last 90 days
25. cnt_loans30 : Number of loans taken by user in last 30 days
26. amnt_loans30: Total amount of loans taken by user in last 30 days
27. maxamnt_loans30 : maximum amount of loan taken by the user in last 30
28. days
29. medianamnt_loans30 : Median of amounts of loan taken by the user in last 30 days
30. cnt_loans90 : Number of loans taken by user in last 90 days
31. amnt_loans90 : Total amount of loans taken by user in last 90 days
32. maxamnt_loans90 : maximum amount of loan taken by the user in last 90 days
33. medianamnt_loans90 : Median of amounts of loan taken by the user in last 90 days
34. payback30 : Average payback time in days over last 30 days
35. payback90 : Average payback time in days over last 90 days
36. pcircle : telecom circle
37. pdate : date

• Data Preprocessing Done

- I. As a first step I have imported required libraries and I have imported the dataset which was in csv format.

- II. Then I did all the statistical analysis like checking shape, nunique, value counts, info etc.....
- III. Then while looking into the value counts, I found some columns with more than 90% zero values this creates skewness in the model and there are chances of getting model bias so I have dropped those columns with more than 90% zero values.
- IV. While checking for null values I found no null values in the dataset.
- V. I have also dropped Unnamed:0, msisdn and pcircle column as I found they are useless.
- VI. Next as a part of feature extraction I converted the pdate column to pyear, pmonth and pday. Thinking that this data will help us more than pdate.
- VII. In some columns I found negative values which were unrealistic so I have converted those negative values to positive using abs command.
- VIII. Also, I have converted all the float values in maxamnt_loans90 to zero as it is specified in the problem statement we can have only 0,6,12 as maximum amount of loan taken by the user in last 30 days. As well I have dropped all the data with amnt_loans90=0 as it gives the persons who have not taken any loans

• **Data Inputs- Logic- Output Relationships**

- I. Since I had all numerical columns, I have plotted dist plot to see the distribution of each column data.
- II. I have used box plot for each pair of categorical features that shows the relation between label and independent features. Also, we can observe wheather the person pays back the loan within the date based on features.
- III. In maximum features relation with target I observed Non-defaulter count is high compared to defaulters.

Importing important libraries:

```

import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.tree import DecisionTreeClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import SGDClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score as cvs
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_curve, auc
from sklearn.metrics import roc_auc_score
from scipy.stats import zscore
from sklearn.preprocessing import StandardScaler
from imblearn.over_sampling import SMOTE
import pickle
from sklearn.model_selection import cross_val_score
import warnings
warnings.filterwarnings('ignore')

```

Loading Data Set into a Variable:

Here I am loading the dataset into the variable `df_micro`.

```
df_micro = pd.read_csv("Data File.csv")
df_micro
```

	Unnamed: 0	label	msisdn	aon	daily_decr30	daily_decr90	rental30	rental90	last_rech_date_ma	last_rech_date_da	...	maxamnt_loans30	medianamnt_loans30	ci
0	1	0	21408170789	272.0	3055.050000	3065.150000	220.13	260.13	2.0	0.0	...	6.0	0.0	
1	2	1	76462170374	712.0	12122.000000	12124.750000	3691.26	3691.26	20.0	0.0	...	12.0	0.0	
2	3	1	17943170372	535.0	1398.000000	1398.000000	900.13	900.13	3.0	0.0	...	6.0	0.0	
3	4	1	55773170781	241.0	21.228000	21.228000	159.42	159.42	41.0	0.0	...	6.0	0.0	
4	5	1	03813182730	947.0	150.619333	150.619333	1098.90	1098.90	4.0	0.0	...	6.0	0.0	
...
209588	209589	1	22758185348	404.0	151.872333	151.872333	1089.19	1089.19	1.0	0.0	...	6.0	0.0	
209589	209590	1	95583184455	1075.0	36.936000	36.936000	1728.36	1728.36	4.0	0.0	...	6.0	0.0	
209590	209591	1	28556185350	1013.0	11843.111670	11904.350000	5861.83	8893.20	3.0	0.0	...	12.0	0.0	
209591	209592	1	59712182733	1732.0	12488.228330	12574.370000	411.83	984.58	2.0	38.0	...	12.0	0.0	
209592	209593	1	65061185339	1581.0	4489.362000	4534.820000	483.92	631.20	13.0	0.0	...	12.0	0.0	

209593 rows × 37 columns

Exploratory Data Analysis:

Exploratory Data Analysis refers to the critical process of performing initial investigations on data so as to discover patterns, to spot

anomalies, to test hypothesis and to check assumptions with the help of summary statistics and graphical representations. We performed some bi-variate analysis on the data to get a better overview of the data and to find outliers in our data-set. Outliers can occur due to some kind of errors while collecting the data and need to be removed so that it doesn't affect the performance of our model.

```
df_micro.shape
```

```
(209593, 37)
```

```
df_micro.info()
```

```
RangeIndex: 209593 entries, 0 to 209592
```

```
Data columns (total 37 columns):
```

#	Column	Non-Null Count	Dtype
0	Unnamed: 0	209593 non-null	int64
1	label	209593 non-null	int64
2	msisdn	209593 non-null	object
3	aon	209593 non-null	float64
4	daily_decr30	209593 non-null	float64
5	daily_decr90	209593 non-null	float64
6	rental30	209593 non-null	float64
7	rental90	209593 non-null	float64
8	last_rech_date_ma	209593 non-null	float64
9	last_rech_date_da	209593 non-null	float64
10	last_rech_amt_ma	209593 non-null	int64
11	cnt_ma_rech30	209593 non-null	int64
12	fr_ma_rech30	209593 non-null	float64
13	sumamnt_ma_rech30	209593 non-null	float64
14	medianamnt_ma_rech30	209593 non-null	float64
15	medianmarechprebal30	209593 non-null	float64
16	cnt_ma_rech90	209593 non-null	int64
17	fr_ma_rech90	209593 non-null	int64
18	sumamnt_ma_rech90	209593 non-null	int64
19	medianamnt_ma_rech90	209593 non-null	float64
20	medianmarechprebal90	209593 non-null	float64
21	cnt_da_rech30	209593 non-null	float64
22	fr_da_rech30	209593 non-null	float64
23	cnt_da_rech90	209593 non-null	int64
24	fr_da_rech90	209593 non-null	int64
25	cnt_loans30	209593 non-null	int64
26	amnt_loans30	209593 non-null	int64
27	maxamnt_loans30	209593 non-null	float64
28	medianamnt_loans30	209593 non-null	float64
29	cnt_loans90	209593 non-null	float64
30	amnt_loans90	209593 non-null	int64
31	maxamnt_loans90	209593 non-null	int64
32	medianamnt_loans90	209593 non-null	float64
33	payback30	209593 non-null	float64
34	payback90	209593 non-null	float64
35	pcircle	209593 non-null	object
36	pdate	209593 non-null	object

```
dtypes: float64(21), int64(13), object(3)
```

```
memory usage: 59.2+ MB
```

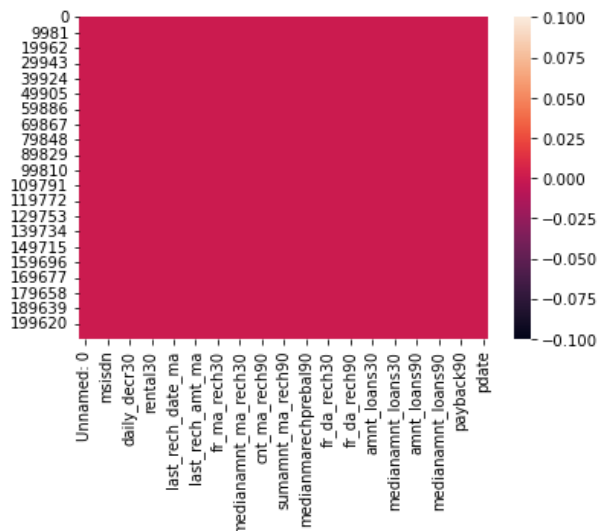
```
df_micro.columns
```

```
Index(['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'],
      dtype='object')
```

Checking the Null Values:

```
#checking if there are any null values present
df_micro.isnull().sum()
```

```
Unnamed: 0      0
label           0
msisdn          0
aon             0
daily_decr30    0
daily_decr90    0
rental30        0
rental90        0
last_rech_date_ma  0
last_rech_date_da  0
last_rech_amt_ma  0
cnt_ma_rech30    0
fr_ma_rech30     0
sumamnt_ma_rech30  0
medianamnt_ma_rech30  0
medianmarechprebal30  0
cnt_ma_rech90    0
fr_ma_rech90     0
sumamnt_ma_rech90  0
medianamnt_ma_rech90  0
medianmarechprebal90  0
cnt_da_rech30    0
fr_da_rech30     0
cnt_da_rech90    0
fr_da_rech90     0
cnt_loans30      0
amnt_loans30     0
maxamnt_loans30  0
medianamnt_loans30  0
cnt_loans90      0
amnt_loans90     0
maxamnt_loans90  0
medianamnt_loans90  0
payback30        0
payback90        0
pcircle          0
pdate            0
dtype: int64
```



We can see there are no null values in the dataset.

Checking the unique values of each feature:

After checking the unique values of each feature, I have the following conclusions-

Observations:

1. There are 209593 rows and 37 columns in the data sets which have different information in each attribute. There are no null values.
2. Basically, there are 2 type of observations made i.e, customer behavior for 30 days and 90 days. Two types of account held by customer main account, data account.
3. Target feature 'Label' has unbalanced data, we need to treat the target variable using sampling technique.
4. 'Unnamed: 0' attribute has all unique values as same as index columns which has no importance for analysis.
5. Approximately 90% of data in 'msisdn' has unique values, i.e, ID.
6. 'payback30','payback90' has nearly 50% of the values having 0.
7. More than 90% of 'last_rech_date_da', 'cnt_da_rech90', 'fr_da_rech90', 'medianamnt_loans30', 'medianamnt_loans90' has of values which is 0.
8. 'pcircle' has only 1 unique value through out column and 'pdate' is a categorical column we can drop this column.

```
#Checking unique values of each column  
df_micro.nunique()
```

```
Unnamed: 0      209593  
label           2  
msisdn          186243  
aon             4507  
daily_decr30    146328  
daily_decr90    155483  
rental30        131338  
rental90        139036  
last_rech_date_ma  1186  
last_rech_date_da  1174  
last_rech_amt_ma   70  
cnt_ma_rech30     71  
fr_ma_rech30     1083  
sumamnt_ma_rech30 15141  
medianamnt_ma_rech30 510  
medianmarechprebal30 23907  
cnt_ma_rech90     110  
fr_ma_rech90      89  
sumamnt_ma_rech90 31771  
medianamnt_ma_rech90 608  
medianmarechprebal90 22694  
cnt_da_rech30     1066  
fr_da_rech30     1072  
cnt_da_rech90     27  
fr_da_rech90      46  
cnt_loans30       40  
amnt_loans30      48  
maxamnt_loans30   1050  
medianamnt_loans30 6  
cnt_loans90       1110  
amnt_loans90      69  
maxamnt_loans90   3  
medianamnt_loans90 6  
payback30        1363  
payback90        2381  
pcircle           1  
pdate            82  
dtype: int64
```

In so many columns like

'last_rech_date_da', 'cnt_da_rech30', 'fr_da_rech30', 'cnt_da_rech90', 'fr_da_rech90', 'medianamnt_loans30', 'medianamnt_loans90' i found more than 90% zeros so they will create skewness in our dataset. So i have dropped these columns.

Checking description of data set:

```
#Dropping columns with more than 90% zeros
df_micro.drop(columns = ['last_rech_date_da', 'cnt_da_rech30', 'fr_da_rech30', 'cnt_da_rech90', 'fr_da_rech90', 'medianamnt_loans30', 'medianamnt_loans90'],
```

```
#Checking description of data set
df_micro.describe()
```

	label	aon	daily_decr30	daily_decr90	rental30	rental90	last_rech_date_ma	last_rech_amt_ma	cnt_ma_rech30	fr_ma_rech30	sumamnt_m
count	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000	209593.000000
mean	0.875177	8112.343445	5381.402289	6082.515068	2692.581910	3483.406534	3755.84780	2064.452797	3.978057	3737.355121	770.000000
std	0.330519	75696.082531	9220.623400	10918.812767	4308.586781	5770.461279	53905.89223	2370.786034	4.256090	53643.625172	1013.000000
min	0.000000	-48.000000	-93.012667	-93.012667	-23737.140000	-24720.580000	-29.000000	0.000000	0.000000	0.000000	0.000000
25%	1.000000	246.000000	42.440000	42.692000	280.420000	300.260000	1.000000	770.000000	1.000000	0.000000	154.000000
50%	1.000000	527.000000	1469.175667	1500.000000	1083.570000	1334.000000	3.000000	1539.000000	3.000000	2.000000	462.000000
75%	1.000000	982.000000	7244.000000	7802.790000	3356.940000	4201.790000	7.000000	2309.000000	5.000000	6.000000	1001.000000
max	1.000000	999860.755200	265926.000000	320630.000000	198926.110000	200148.110000	998650.37770	55000.000000	203.000000	999606.368100	81009.000000

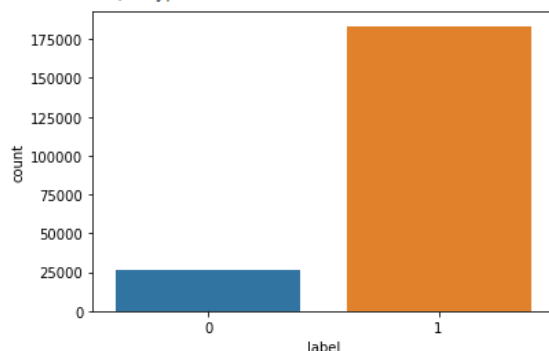
Univariate Analysis:

Uni means one, so in other words, the data has only one variable. Univariate data requires analyzing each variable separately. It doesn't deal with causes or relationships (unlike regression) and its major purpose is to describe; It takes data, summarizes that data and finds patterns in the data.

Analysing Target variable:

```
#Let us analyse target variable and check unique values
sns.countplot(df_micro['label'])
df_micro['label'].value_counts()
```

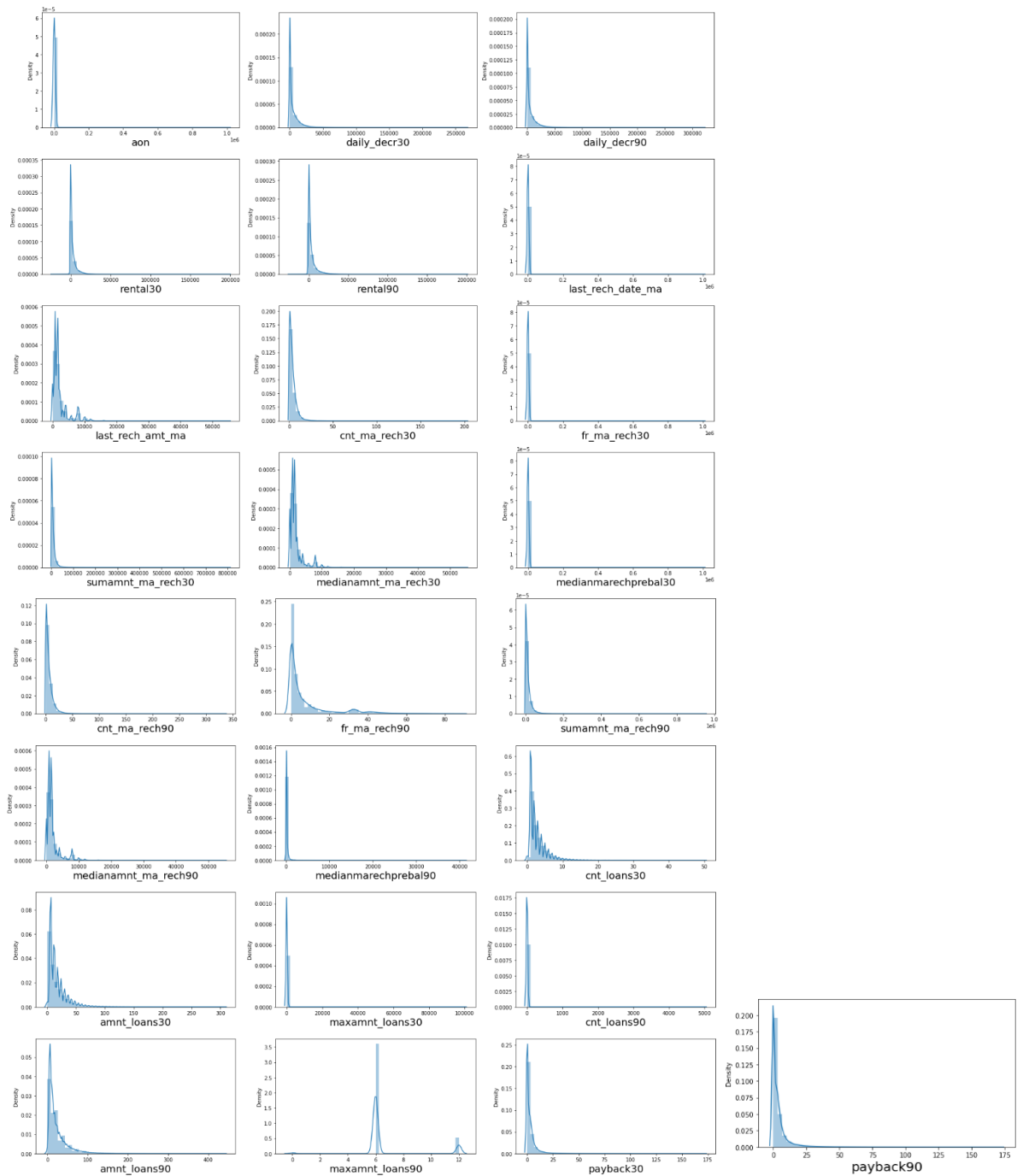
```
1    183431
0     26162
Name: label, dtype: int64
```



Here we can see that out of 209593, 183431 label are success which is around 87.5177% and 26162 labels are failure which is around 12.482%. Also Here we can see there is huge difference between two Categories of label. So the data is imbalanced. So we will apply SMOTE analysis before ML of final model.

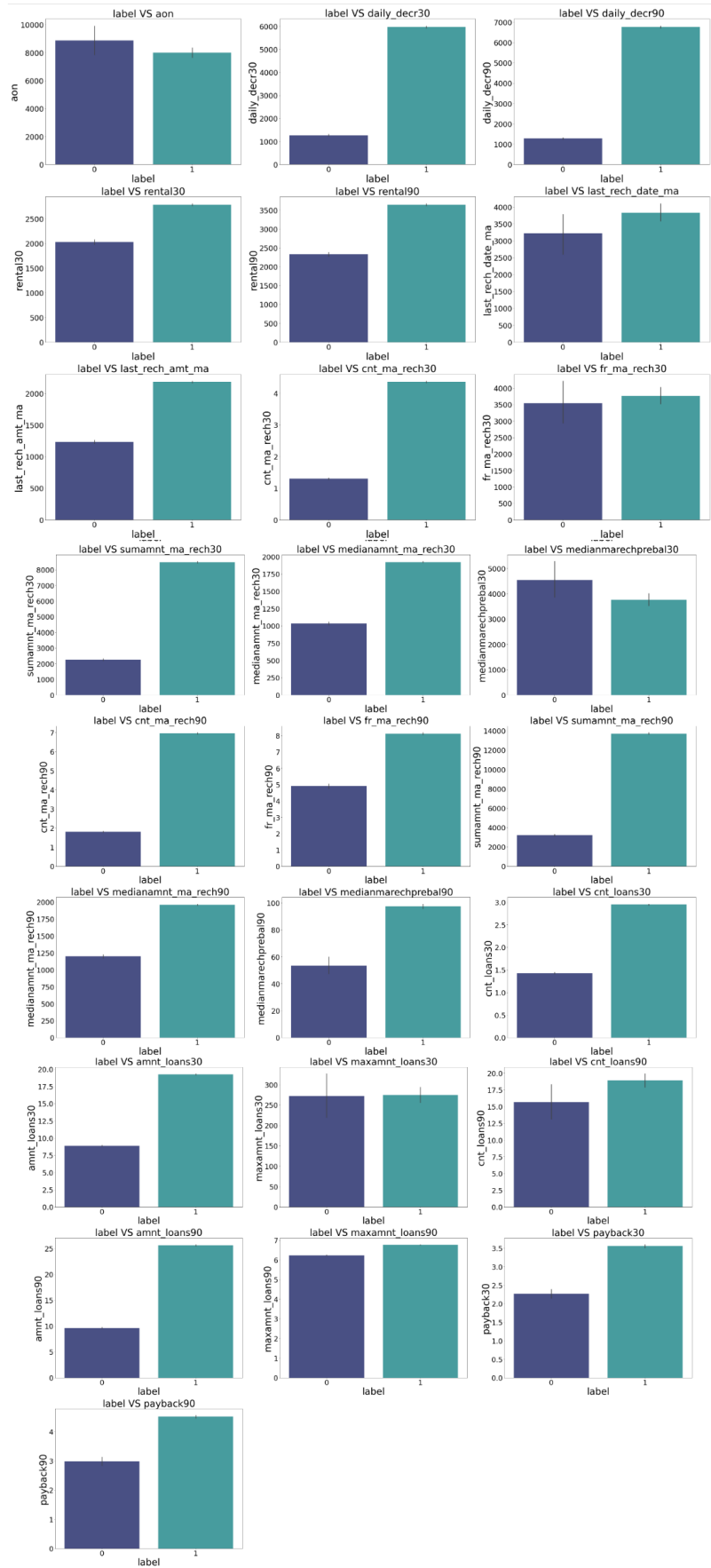
There are 87.5% of non-defaulters and 12.5% of defaulter customers, the data is unbalanced we will use SMOTE analysis technique to balance the target.

Analysing Numerical Columns:



Bivariate Analysis:

Bivariate analysis is finding some kind of empirical relationship between two variables. Specifically, the dependent vs independent Variables



OBSERVATIONS:

1. Customers with high value of Age on cellular network in days(aon) are maximum defaulters(who have not paid there loan amount-0).
2. Customers with high value of Daily amount spent from main account, averaged over last 30 days (in Indonesian Rupiah)(daily_decr30) are maximum Non-defaulters(who have paid there loan amount-1).
3. Customers with high value of Daily amount spent from main account, averaged over last 90 days (in Indonesian Rupiah)(daily_decr90) are maximum Non-defaulters(who have paid there loan amount-1).
4. Customers with high value of Average main account balance over last 30 days(rental30) are maximum Non-defaulters(who have paid there loan amount-1).
5. Customers with high value of Average main account balance over last 90 days(rental90) are maximum Non-defaulters(who have paid there loan amount-1).
6. Customers with high Number of days till last recharge of main account(last_rech_date_ma) are maximum Non-defaulters(who have paid there loan amount-1).
7. Customers with high value of Amount of last recharge of main account (in Indonesian Rupiah)(last_rech_amt_ma) are maximum Non-defaulters(who have paid there loan amount-1).
8. Customers with high value of Number of times main account got recharged in last 30 days(cnt_ma_rech30) are maximum Non-defaulters(who have paid there loan amount-1).
9. Customers with high value of Frequency of main account recharged in last 30 days(fr_ma_rech30) are maximum Non-defaulters(who have paid there loan amount-1) and also the count is high for defaulters comparitively Non-defaulters are more in number.
10. Customers with high value of Total amount of recharge in main account over last 30 days (in Indonesian Rupiah)(sumamnt_ma_rech30) are maximum Non-defaulters(who have paid there loan amount-1).
11. Customers with high value of Median of amount of recharges done in main account over last 30 days at user level (in Indonesian Rupiah)(medianamnt_ma_rech30) are maximum Non-defaulters(who have paid there loan amount-1).
12. Customers with high value of Median of main account balance just before recharge in last 30 days at user level (in Indonesian Rupiah)(medianmarechprebal30) are maximum defaulters(who have not paid there loan amount-0).
13. Customers with high value of Number of times main account got recharged in last 90 days(cnt_ma_rech90) are maximum Non-defaulters(who have paid there loan amount-1).

14. Customers with high value of Frequency of main account recharged in last 90 days(fr_ma_rech90) are maximum Non-defaulters(who have paid there loan amount-1).
15. Customers with high value of Total amount of recharge in main account over last 90 days (in Indonesian Rupiah)(sumamnt_ma_rech90) are maximum Non-defaulters(who have paid there loan amount-1).
16. Customers with high value of Median of amount of recharges done in main account over last 90 days at user level (in Indonesian Rupiah)(medianamnt_ma_rech90) are maximum Non-defaulters(who have paid there loan amount-1).
17. Customers with high value of Median of main account balance just before recharge in last 90 days at user level (in Indonesian Rupiah)(medianmarechprebal90) are maximum Non-defaulters(who have paid there loan amount-1).
18. Customers with high value of Number of loans taken by user in last 30 days(cnt_loans30) are maximum Non-defaulters(who have paid there loan amount-1).
19. Customers with high value of Total amount of loans taken by user in last 30 days(amnt_loans30) are maximum Non-defaulters(who have paid there loan amount-1).
20. Customers with high value of maximum amount of loan taken by the user in last 30 days(maxamnt_loans30) are maximum Non-defaulters(who have paid there loan amount-1).
21. Customers with high value of Number of loans taken by user in last 90 days(cnt_loans90) are maximum Non-defaulters(who have paid there loan amount-1).
22. Customers with high value of Total amount of loans taken by user in last 90 days(amnt_loans90) are maximum Non-defaulters(who have paid there loan amount-1).
23. Customers with high value of maximum amount of loan taken by the user in last 90 days(maxamnt_loans90) are maximum Non-defaulters(who have paid there loan amount-1).
24. Customers with high value of Average payback time in days over last 30 days(payback30) are maximum Non-defaulters(who have paid there loan amount-1).
25. Customers with high value of Average payback time in days over last 90 days(payback90) are maximum Non-defaulters(who have paid there loan amount-1).
26. In between 6th and 7th month maximum customers both defaulters and Non-defaulters have paid there loan amount.
27. Below 14th of each month all the customers have paid there loan amount.

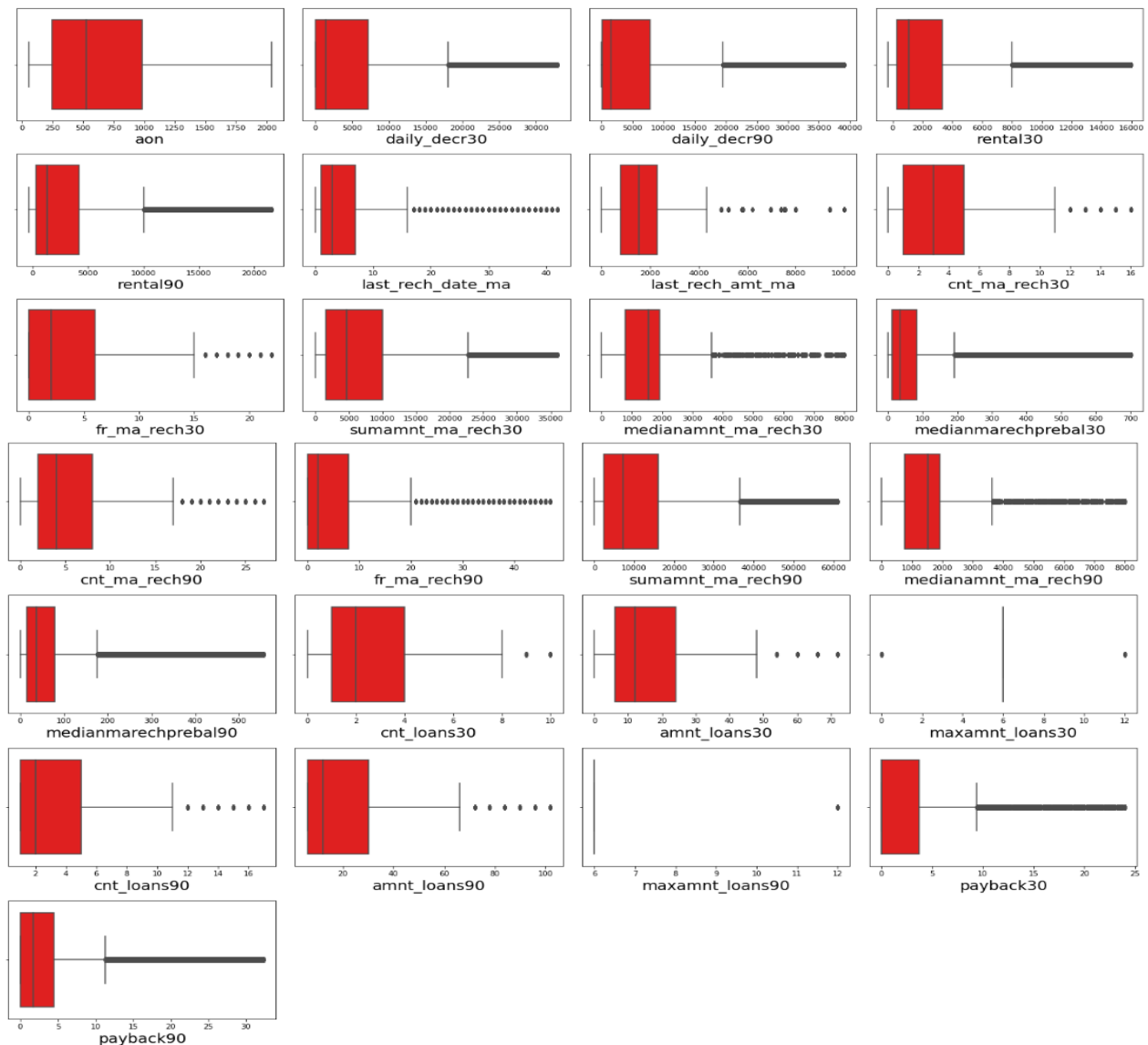
Removing outliers and skewness:

- To remove outliers, I have used percentile method.

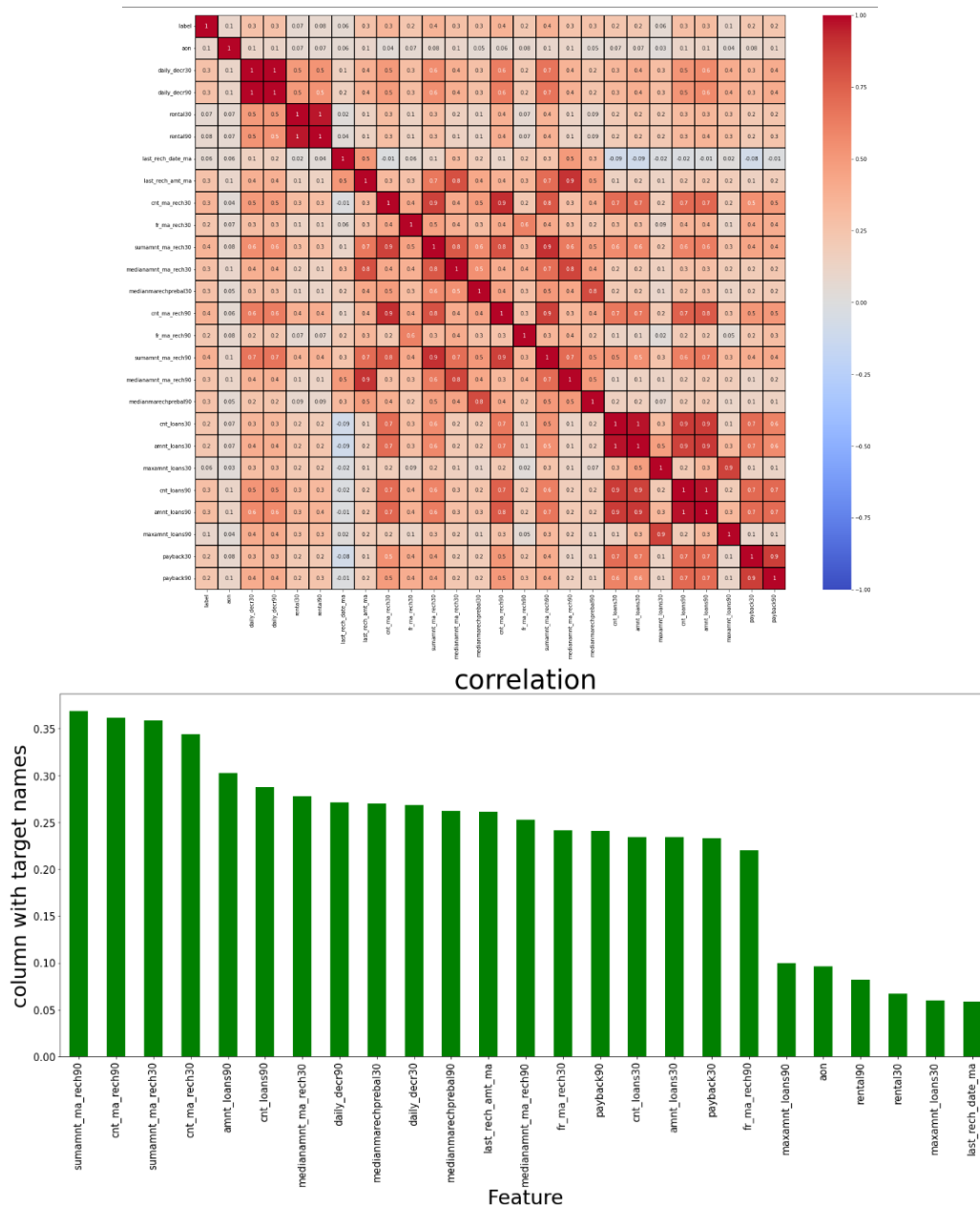
```
#Removing outliers using percentile method
for colu in features:
    if df_micro[colu].dtypes != 'object':
        percentile = df_micro[colu].quantile([0.01,0.99]).values
        df_micro[colu][df_micro[colu]<=percentile[0]]=percentile[0]
        df_micro[colu][df_micro[colu]>=percentile[1]]=percentile[1]
```

```
# Checking if the outliers is reduced or not

plt.figure(figsize=(20,25),facecolor='white')
plotnumber=1
for column in col:
    if plotnumber<=30:
        ax=plt.subplot(8,4,plotnumber)
        sns.boxplot(df_micro[column],color='red')
        plt.xlabel(column,fontsize=20)
        plotnumber+=1
plt.tight_layout()
```



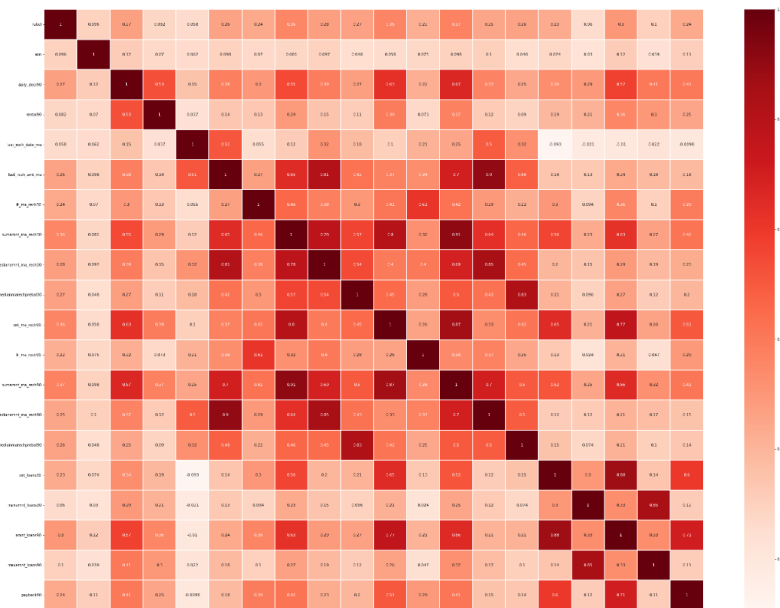
- To remove skewness I have used yeo-johnson method..
- Use of Pearson's correlation coefficient to check the correlation between dependent and independent features.



- There exist a lot of multicollinearity between 'daily_decr30' and 'daily_decr90'; 'rental30' and 'rental90' ; 'cnt_loans30' and 'amnt_loans30' ; 'cnt_loans90' and 'amnt_loans90' ; 'payback30' and 'payback90' ; . Hence we will drop either column from each group which is least correlated with the target column.

```
df_micro.drop(['daily_decr30', 'rental30', 'cnt_ma_rech30', 'amnt_loans30', 'cnt_loans90', 'payback30'], axis=1, inplace=True)
```

- Also, I have used Normalization to scale the data. After scaling we have to balance the target column using oversampling.



- I have used oversampling (SMOTE) to get rid of data imbalancing. The balanced output looks like this.

```
#Checking the value count of target column
y.value_counts()
```

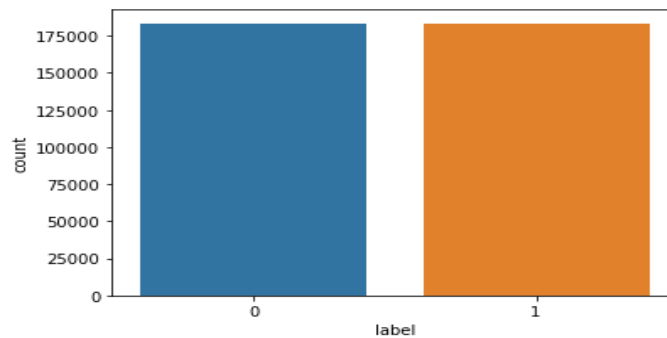
```
1    183431
0     26162
Name: label, dtype: int64
```

```
from imblearn.over_sampling import SMOTE
SM = SMOTE()
X, y = SM.fit_resample(X,y)
```

```
# Checking the value counts again
y.value_counts()
```

```
0    183431
1    183431
Name: label, dtype: int64
```

```
# Visualizing the target data after oversampling
sns.countplot(y)
```



Then followed by model building with all Classification algorithms.

1. Testing of Identified Approaches (Algorithms)

Since label was my target and it was a classification column with 0-defaulter and 1-Non-defaulter, so this particular problem was Classification problem. And I have used all Classification algorithms to build my model. By looking into the difference of accuracy score and cross validation score I found RandomForestClassifier as a best model with least difference. Also to get the best model we have to run through multiple models and to avoid the confusion of overfitting we have go through cross validation. Below are the list of classification algorithms I have used in my project.

- ☐ Logistic Regresson
- ☐ DecisionTreeClassifier
- ☐ K-Neighbour Classifier
- ☐ GaussianNB
- ☐ Random Forest Classifier
- ☐ AdaBoost Classifier

2. Key Metrics for success in solving problem under consideration

I have used the following metrics for evaluation

- **Precision** can be seen as a measure of quality, higher precision means that an algorithm returns more relevant results than irrelevant ones.□
- **Recall** is used as a measure of quantity and high recall means that an algorithm returns most of the relevant results.□
- **Accuracy score** is used when the True Positives and True negatives are more important. Accuracy can be used when the class distribution is similar.□
- **F1-score** is used when the False Negatives and False Positives are crucial. While F1-score is a better metric when there are imbalanced classes.□
- **Cross_val_score**: To run cross-validation on multiple metrics and also to return train scores, fit times and score times. Get predictions from each split of cross-validation for diagnostic purposes. Make a scorer from a performance metric or loss function.□
- **AUC_ROC_score**: ROC curve. It is a plot of the false positive rate (x-axis) versus the true positive rate (y-axis) for a number of different candidate threshold values between 0.0 and 1.0

I have used accuracy_score since I have balanced my data using oversampling.

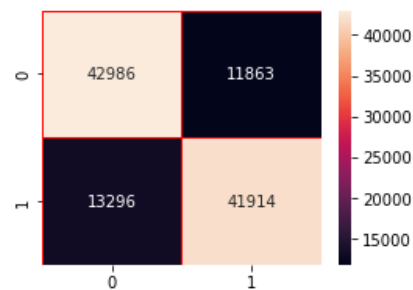
Run and Evaluate selected models:

1. Model Building:

1) Logistic Regression:

```
The accuracy score of train is : 77.09061031218482
The accuracy score test is : 77.14044285337864
The cross validation score is : 77.06997205080816
Confusion Matrix:
[[42986 11863]
 [13296 41914]]
Classification
```

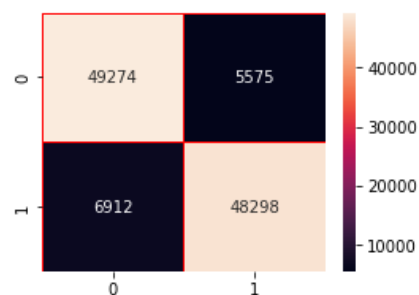
	precision	recall	f1-score	support
0	0.76	0.78	0.77	54849
1	0.78	0.76	0.77	55210
accuracy			0.77	110059
macro avg	0.77	0.77	0.77	110059
weighted avg	0.77	0.77	0.77	110059



2) Decision Tree Classifier:

```
The accuracy score of train is : 99.96729010175115
The accuracy score test is : 88.65426725665326
The cross validation score is : 87.68160826629368
Confusion Matrix:
[[49274 5575]
 [ 6912 48298]]
Classification
```

	precision	recall	f1-score	support
0	0.88	0.90	0.89	54849
1	0.90	0.87	0.89	55210
accuracy			0.89	110059
macro avg	0.89	0.89	0.89	110059
weighted avg	0.89	0.89	0.89	110059

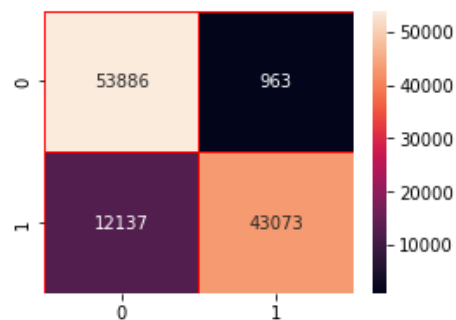


3) K-Neighbour Regressor:

The accuracy score of train is : 91.52268470383913
 The accuracy score test is : 88.09729326997338
 The cross validation score is : 87.31050636563702
 Confusion Matrix:
 [[53886 963]
 [12137 43073]]
 Classification precision recall f1-score support

 0 0.82 0.98 0.89 54849
 1 0.98 0.78 0.87 55210

 accuracy 0.88 110059
 macro avg 0.90 0.88 0.88 110059
 weighted avg 0.90 0.88 0.88 110059

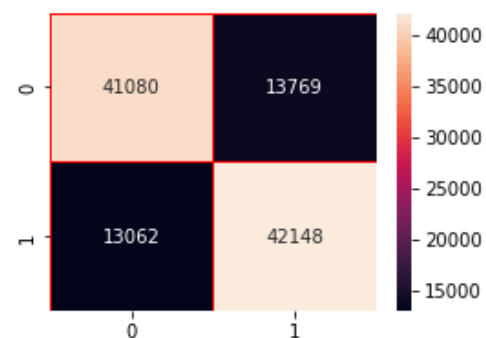


4) GaussianNB

The accuracy score of train is : 75.65215359633649
 The accuracy score test is : 75.62125768905769
 The cross validation score is : 75.6579947413947
 Confusion Matrix:
 [[41080 13769]
 [13062 42148]]
 Classification precision recall f1-score support

 0 0.76 0.75 0.75 54849
 1 0.75 0.76 0.76 55210

 accuracy 0.76 110059
 macro avg 0.76 0.76 0.76 110059
 weighted avg 0.76 0.76 0.76 110059



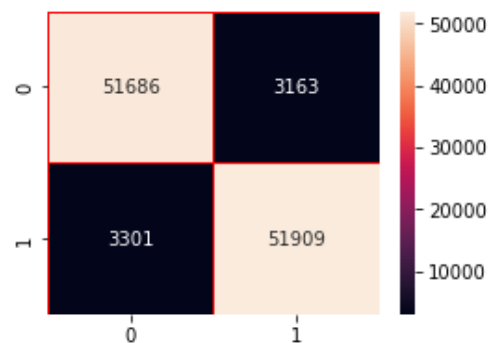
5) Random Forest Classifier


```

The accuracy score of train is : 99.9661218910994
The accuracy score test is : 94.12678654176396
The cross validation score is : 93.58808100968278
Confusion Matrix:
[[51686  3163]
 [ 3301 51909]]
Classification

```

	precision	recall	f1-score	support
0	0.94	0.94	0.94	54849
1	0.94	0.94	0.94	55210
accuracy			0.94	110059
macro avg	0.94	0.94	0.94	110059
weighted avg	0.94	0.94	0.94	110059



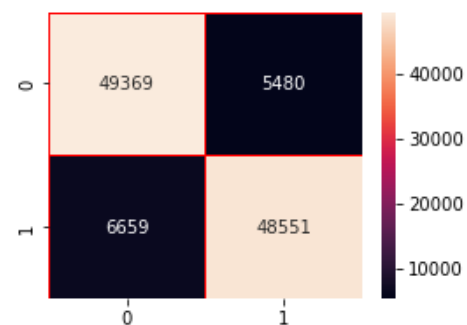
6) AdaBoost Classifier

```

The accuracy score of train is : 99.96729010175115
The accuracy score test is : 88.97046129803105
The cross validation score is : 88.03284997477692
Confusion Matrix:
[[49369  5480]
 [ 6659 48551]]
Classification

```

	precision	recall	f1-score	support
0	0.88	0.90	0.89	54849
1	0.90	0.88	0.89	55210
accuracy			0.89	110059
macro avg	0.89	0.89	0.89	110059
weighted avg	0.89	0.89	0.89	110059



Observation:

We have trained several models above for the dataset we had prepared and we got different results for different algorithm,

- Logistic regression model gave us train score of 77.1% and test score of 76.9% of accuracy and cross validation score of 77.1 % for the test model which is very near and also the precision, accuracy score are also high.
- DecisionTreeClassifier model gave us train score of 99.9% and test score of 89% of accuracy and cross validation score of 87.9 % for the test model. Here the model is overfitting as there is large difference between train score and test score.
- KNN classifier has given us 91.5% and 88% of accuracy and cv score of 87.5% for the test dataset.
- GaussianNB() model gave us train score of 75.7% and test score of 75.6% of accuracy and cross validation score of 75.7 % for the test model. Here the model accuracy is low as compared to other models..
- Random Forest classifier model gave us train score of 99.9% and test score of 94.1% of accuracy and cross validation score of 93.6% and also metric values are near to 1 which is very good score.
- AdaBoostClassifier model gave us train score of 99.9% and test score of 89.3% of accuracy and cross validation score of 93.6% .Here the model is overfitting as there is large difference between train score and test score.

We are selecting random forest classifier model to increase the accuracy using Gridsearch CV method as it is giving the highest accuracy with least train and test score difference

3. Hyper Parameter Tuning:

```
1 parameters = {"n_estimators": [8, 4, 6],
2               "min_samples_leaf": [2, 3, 4, 5, ],
3               "criterion": ["gini", "entropy"],
4               "min_samples_split": [2, 4, 6, 8, 10]}
5
6 rfc = GridSearchCV(RandomForestClassifier(), parameters)
7 #fitting train and test data
8 rfc.fit(X_train, y_train)
9
10 #Best parameters
11 rfc.best_params_

{'criterion': 'entropy',
 'min_samples_leaf': 2,
 'min_samples_split': 4,
 'n_estimators': 8}
```

```

1 rfc = RandomForestClassifier(criterion='entropy',min_samples_leaf=2, min_samples_split=4,n_estimators=8,n_jobs=1,verbose=2)
2 rfc.fit(X_train, y_train)
3 y_pred_train = rfc.predict(X_train)
4 y_pred = rfc.predict(X_test)
5 print("The accuracy score of train is :", accuracy_score(y_train, y_pred_train)*100)
6 print("The accuracy score test is :", accuracy_score(y_test, y_pred)*100)
7 cv_score = cross_val_score(rfc,X_train, y_train,cv=5)
8 print("The cross validation score is :", cv_score.mean()*100)
9
10 print("Confusion Matrix: \n", confusion_matrix(y_test, y_pred))
11 print("Classification\n ", classification_report(y_test, y_pred))
12 print("*****")
13
14 plt.figure(figsize=(6,5))
15 sns.heatmap(confusion_matrix(y_test, y_pred),annot=True,fmt = "d",linecolor="r",linewidths=1)
16
17 plt.show()

```

The accuracy score of train is : 98.0027491890671
The accuracy score test is : 92.49493453511299

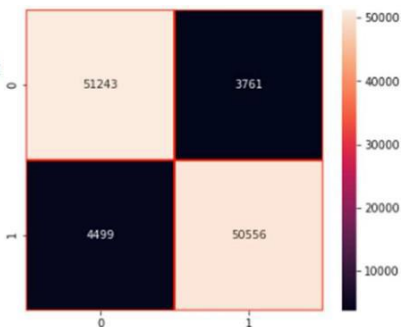
The cross validation score is : 92.00048255033806

Confusion Matrix:

```
[[51243  3761]
 [ 4499 50556]]
```

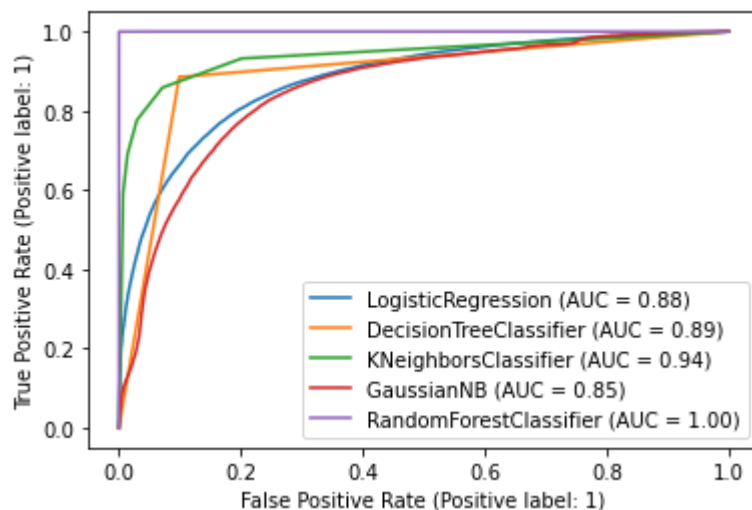
Classification

	precision	recall	f1-score	support
0	0.92	0.93	0.93	55004
1	0.93	0.92	0.92	55055
accuracy			0.92	110059
macro avg	0.93	0.92	0.92	110059
weighted avg	0.93	0.92	0.92	110059



We see that even after hyperparameter optimization, there is no much changes in the accuracy score, Hence saving the previous model for the reference.

AUC ROC Curve for Different model:



Here we can see that we are getting Random forest classifier accuracy score as best score also AUC ROC curve is having value 1.

Saving the model and Predictions:

We see that Random forest classifier model has given the highest AUC in graph, the accuracy score of 97% and CV score of 96% which is highest among all the models tested. also we see that evaluation metrics are high for this model. Hence we will be saving this model.

I have saved my best model using .pkl as follows.

```
import joblib
file = "MicroCredit_Defaulters_Project.joblib"
joblib.dump(rfc,file)
```

```
['MicroCredit_Defaulters_Project.joblib']
```

Interpretation of the Results

- a) The dataset was very challenging to handle it had 37 features with 30days and 90days information of customers.
- b) Firstly, the datasets were not having any null values.
- c) But there was huge number of zero entries in maximum columns so we have to be careful while going through the statistical analysis of the datasets.
- d) And proper plotting for proper type of features will help us to get better insight on the data. I found maximum numerical columns in the dataset so I have chosen bar plot to see the relation between target and features.
- e) I notice a huge amount of outliers and skewness in the data so we have choose proper methods to deal with the outliers and skewness. If we ignore this outliers and skewness we may end up with a bad model which has less accuracy.
- f) Then scaling dataset has a good impact like it will help the model not to get biased. Since we have not removed outliers and skewness completely from the dataset so we have to choose Normalization.
- g) We have to use multiple models while building model using dataset as to get the best model out of it.
- h) We have to use multiple metrics like F1_score, precision, recall and accuracy_score which will help us to decide the best model.
- i) I found Random Forest Classifier as the best model with 94.1% accuracy_score. Also I have improved the accuracy of the best model by running hyper parameter tuning.
- j) ☐At last I have predicted wheather the loan is paid back or not using saved model. It was good!! that I was able to get the predictions near to actual values.

CONCLUSION

▪ **Key Findings and Conclusions of the Study**

In this project report, we have used machine learning algorithms to predict the micro credit defaulters. We have mentioned the step by step procedure to analyze the dataset and finding the correlation between the features. Thus we can select the features which are correlated to each other and are independent in nature. These feature set were then given as an input to four algorithms and a hyper parameter tuning was done to the best model and the accuracy has been improved. Hence we calculated the performance of each model using different performance metrics and compared them based on these metrics. Then we have also saved the best model and predicted the label. It was good the the predicted and actual values were almost same.

▪ **Learning Outcomes of the Study in respect of Data Science**

I found that the dataset was quite interesting to handle as it contains all types of data in it. Improvement in computing technology has made it possible to examine social information that cannot previously be captured, processed and analyzed. New analytical techniques of machine learning can be used in property research. The power of visualization has helped us in understanding the data by graphical representation it has made me understand what data is trying to say. Data cleaning is one of the most important steps to remove unrealistic values and zero values. This study is an exploratory attempt to use four machine learning algorithms in estimating microcredit defaulter, and then compare their results.

To conclude, the application of machine learning in microcredit is still at an early stage. We hope this study has moved a small step ahead in providing some methodological and empirical contributions to crediting institutes, and presenting an alternative approach to the valuation of defaulters. Future direction of research may consider incorporating additional micro credit transaction data from a larger economical background with more features.

▪ **Limitations of this work and Scope for Future Work**

- First draw back is the length of the dataset it is very huge and hard to handle.
- Followed by more number of outliers and skewness these two will reduce our model accuracy.
- Also, we have tried best to deal with outliers, skewness and zero values. So it looks quite good that we have achieved a accuracy of 94.1% even after dealing all these drawbacks.
- Also, this study will not cover all Classification algorithms instead, it is focused on the chosen algorithm, starting from the basic ensembling techniques to the advanced ones.