



Micro Credit Loan Defaulter Project Report



Submitted by:

Sweta Rai

ACKNOWLEDGMENT

I would like to express my deepest gratitude to my SME (Subject Matter Expert) Khushboo Garg as well as Flip Robo Technologies who gave me the opportunity to do this project on Surprise Housing Price Prediction, which also helped me in doing lots of research wherein I came to know about so many new things.

Also, I have utilized a few external resources that helped me to complete the project. I ensured that I learn from the samples and modify things according to my project requirement. All the external resources that were used in creating this project are listed below:

- 1) <https://www.google.com/>
- 2) <https://www.youtube.com/>
- 3) https://scikit-learn.org/stable/user_guide.html
- 4) <https://github.com/>
- 5) <https://www.kaggle.com/>
- 6) <https://medium.com/>
- 7) <https://towardsdatascience.com/>
- 8) <https://www.analyticsvidhya.com/>

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 industry 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.

- Conceptual Background of the Domain Problem

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. They understand the importance of communication and how it affects a person's life, thus, focusing on providing their services and products to low-income families and poor customers that can help them in the need of hour. They are

collaborating with an MFI to provide micro-credit on mobile balances to be paid back in 5 days. The Consumer is believed to be defaulter if he deviates from the path of paying back the loaned amount within the time duration of 5 days. For the loan amount of 5 (in Indonesian Rupiah), payback amount should be 6 (in Indonesian Rupiah), while, for the loan amount of 10 (in Indonesian Rupiah), the payback amount should be 12 (in Indonesian Rupiah).

- Review of Literature

1. What is Microfinance?

“Microfinance” is often seen as financial services for poor and low-income clients. In practice, the term is often used more narrowly to refer to loans and other services from providers that identify themselves as “microfinance institutions” (MFIs). Microfinance can also be described as a setup of a number of different operators focusing on the financially under-served people with the aim of satisfying their need for poverty alleviation, social promotion, emancipation, and inclusion. Microfinance institutions reach and serve their target market in very innovative ways. Microfinance operations differ in principle, from the standard disciplines of general and entrepreneurial finance. This difference can be attributed to the fact that the size of the loans granted with microcredit is typically too small to finance growth-oriented business projects. Some unique features of microfinance as follows:

- i. Delivery of very small loans to unsalaried workers.
- ii. Little or no collateral requirements.
- iii. Group lending and liability.
- iv. Pre-loan savings requirement.
- v. Gradually increasing loan sizes.

Implicit guarantee of ready access to future loans if present loans are repaid fully and promptly Microfinance is seen as a catalyst for

poverty alleviation, delivered in innovative and sustainable ways to assist the underserved poor, especially in developing countries.

2. Default in Microfinance

Default in microfinance is the failure of a client to repay a loan. The default could be in terms of the amount to be paid or the timing of the payment.

- Motivation for the Problem Undertaken

Our main objective of doing this project is to build a model to predict whether the users are paying the loan within the due date or not. We are going to predict by using Machine Learning algorithms.

The sample data is provided to us from our client database. 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.

Analytical Problem Framing

- Mathematical/ Analytical Modeling of the Problem

There are various analytics which I have done before moving forward with exploratory analysis, on the basis of accounts which got recharged in the last 30 days. I set the parameter that if the person is not recharging their main account within 3 months, I simply dropped their data because they are not valuable and they might be old customers, but there is no revenue rotating. Then I had checked the date columns and found that the data belongs to the year 2016. I extracted the month and day from the date, saved the data in separate columns, and tried to visualize the data on the basis of months and days.

I had checked the maximum amount of loan taken by the people and found that the data had more outliers. As per the description given by the client, the loan amount can be paid by the customer is either rupiah 6 or 12 so that I have dropped all the loan amount that shows the loan is taken more than 12 rupiah.

Then I separated the defaulter's data and checked the valuable customers in the network and we found that their monthly revenue is more than 10000 rupiah. Although the data is quite imbalanced and many columns doesn't have that expected maximum value, we dropped that columns. We checked the skewed data and try to treat the skewed data before model processing which caused NaN so avoided it.

When we try removing the unwanted data, i.e., the outliers, we found that almost 40000+ data has been chopped. Though the data given by the client had almost 37 columns and over 2 lakh columns I did not feel like losing on precious data so avoided the outlier removal part as well. After scaling my data, I have sent the data to various classification models and found that Extra Trees Classifier Algorithm is working well.

- Data Sources and their formats

The data is been provided by one of our clients from telecom industry. They are a fixed wireless telecommunications network provider and 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.

The data is been given by Indonesian telecom company and they gave it to us in a CSV file, with data description file in excel format. They also had provided the problem statement by explaining what they need from us and also the required criteria to be satisfied.

Let's check the data now. Below I have attached the snapshot below to give an overview.

```
import warnings
warnings.simplefilter("ignore")
warnings.filterwarnings("ignore")
import joblib

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

import missingno
import pandas_profiling
from imblearn.over_sampling import SMOTE
from sklearn.preprocessing import StandardScaler

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.neighbors import KNeighborsClassifier
import xgboost as xgb
import lightgbm as lgb

from sklearn import metrics
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from sklearn.metrics import classification_report
from sklearn.model_selection import GridSearchCV
from scikitplot.metrics import plot_roc_curve
from sklearn.metrics import roc_curve, auc, roc_auc_score
```

```
df = pd.read_csv("Data_File.csv")
```

I am importing the dataset comma separated values file and storing it into our dataframe for further usage.

```
df # checking the first 5 and last 5 rows
```

	Unnamed: 0	label	msisdn	aon	daily_decr30	daily_decr90	rental30	rental90	last_rech_date_ma	last_rech_date_da	...	maxamnt_loans30
0	1	0	21408170789	272.0	3055.050000	3065.150000	220.13	260.13	2.0	0.0	...	6.0
1	2	1	76462170374	712.0	12122.000000	12124.750000	3691.26	3691.26	20.0	0.0	...	12.0
2	3	1	17943170372	535.0	1398.000000	1398.000000	900.13	900.13	3.0	0.0	...	6.0
3	4	1	55773170781	241.0	21.228000	21.228000	159.42	159.42	41.0	0.0	...	6.0
4	5	1	03813182730	947.0	150.619333	150.619333	1098.90	1098.90	4.0	0.0	...	6.0
...
209588	209589	1	22758185348	404.0	151.872333	151.872333	1089.19	1089.19	1.0	0.0	...	6.0
209589	209590	1	95583184455	1075.0	36.936000	36.936000	1728.36	1728.36	4.0	0.0	...	6.0
209590	209591	1	28558185350	1013.0	11843.111667	11904.350000	5861.83	8893.20	3.0	0.0	...	12.0
209591	209592	1	59712182733	1732.0	12488.228333	12574.370000	411.83	984.58	2.0	38.0	...	12.0
209592	209593	1	65061185339	1581.0	4489.362000	4534.820000	483.92	631.20	13.0	0.0	...	12.0

209593 rows x 37 columns

Here we are taking a look at the first 5 and last 5 rows of our dataset. It shows that we have a total of 209593 rows and 37 columns present in our dataframe. We have the label column that stores the defaulter and non-defaulter values marked with 0 and 1 making this a Classification problem!

- Data Preprocessing Done

Checked for missing values to confirm the information of no null values present provided in the problem statement.

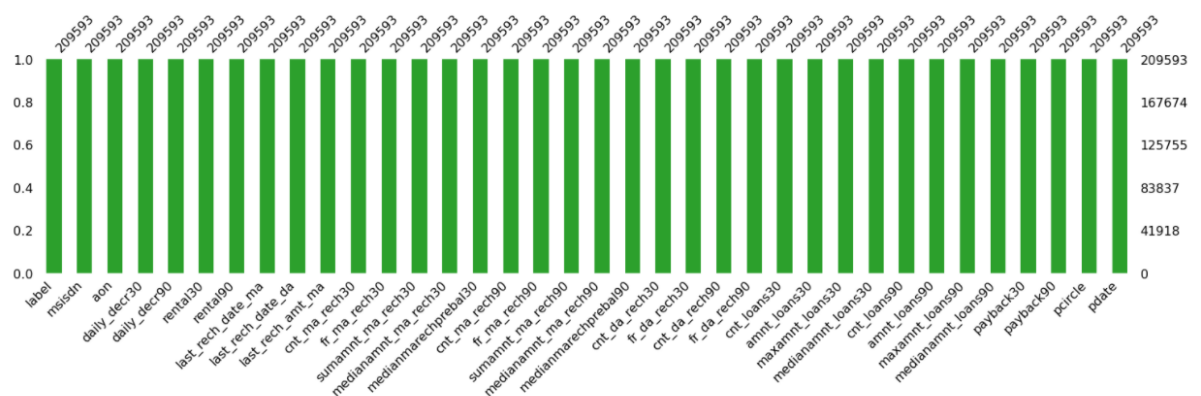

```
df.isna().sum() # checking for missing values
```

```
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
```

Took a visual on the missing data information as well.

```
missingno.bar(df, figsize = (25,5), color="tab:green")
```

<AxesSubplot:>



Using the info method, we are able to confirm the non-null count details as well as the datatype information. We have 21 float/decimal datatype, 12 integer datatype and 3 object/categorical datatype columns. We will need to convert the object datatype

columns to numerical data before we input the information in our machine learning models.

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 209592 entries, 0 to 209592
Data columns (total 36 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   label                                209592 non-null  int64
1   msisdn                             209592 non-null  object
2   aon                                 209592 non-null  float64
3   daily_decr30                       209592 non-null  float64
4   daily_decr90                       209592 non-null  float64
5   rental30                           209592 non-null  float64
6   rental90                           209592 non-null  float64
7   last_rech_date_ma                  209592 non-null  float64
8   last_rech_date_da                  209592 non-null  float64
9   last_rech_amt_ma                   209592 non-null  int64
10  cnt_ma_rech30                      209592 non-null  int64
11  fr_ma_rech30                       209592 non-null  float64
12  sumamnt_ma_rech30                  209592 non-null  float64
13  medianamnt_ma_rech30               209592 non-null  float64
14  medianmarechprebal30               209592 non-null  float64
15  cnt_ma_rech90                      209592 non-null  int64
16  fr_ma_rech90                       209592 non-null  int64
17  sumamnt_ma_rech90                  209592 non-null  int64
18  medianamnt_ma_rech90               209592 non-null  float64
19  medianmarechprebal90               209592 non-null  float64
20  cnt_da_rech30                      209592 non-null  float64
21  fr_da_rech30                       209592 non-null  float64
22  cnt_da_rech90                      209592 non-null  int64
23  fr_da_rech90                       209592 non-null  int64
24  cnt_loans30                        209592 non-null  int64
25  amnt_loans30                       209592 non-null  int64
26  maxamnt_loans30                    209592 non-null  float64
27  medianamnt_loans30                 209592 non-null  float64
28  cnt_loans90                        209592 non-null  float64
29  amnt_loans90                       209592 non-null  int64
30  maxamnt_loans90                    209592 non-null  int64
31  medianamnt_loans90                 209592 non-null  float64
32  payback30                          209592 non-null  float64
33  payback90                          209592 non-null  float64
34  pcircle                            209592 non-null  object
35  pdate                             209592 non-null  object
dtypes: float64(21), int64(12), object(3)
```

- Data Inputs- Logic- Output Relationships

Data description on each column present in our dataset.

label : Flag indicating whether the user paid back the credit amount within 5 days of issuing the loan {1: success, 0: failure}

msisdn : Mobile number of users

aon : Age on cellular network in days

daily_decr30 : Daily amount spent from main account, averaged over last 30 days (in Indonesian Rupiah)

daily_decr90 : Daily amount spent from main account, averaged over last 90 days (in Indonesian Rupiah)

rental30 : Average main account balance over last 30 days

rental90 : Average main account balance over last 90 days

last_rech_date_ma : Number of days till last recharge of main account

last_rech_date_da : Number of days till last recharge of data account

last_rech_amt_ma : Amount of last recharge of main account (in Indonesian Rupiah)

cnt_ma_rech30 : Number of times main account got recharged in last 30 days

fr_ma_rech30 : Frequency of main account recharged in last 30 days

sumamnt_ma_rech30 : Total amount of recharge in main account over last 30 days (in Indonesian Rupiah)

medianamnt_ma_rech30 : Median of amount of recharges done in main account over last 30 days at user level (in Indonesian Rupiah)

medianmarechprebal30 : Median of main account balance just before recharge in last 30 days at user level (in Indonesian Rupiah)

cnt_ma_rech90 : Number of times main account got recharged in last 90 days

fr_ma_rech90 : Frequency of main account recharged in last 90 days

sumamnt_ma_rech90 : Total amount of recharge in main account over last 90 days (in Indonesian Rupiah)

medianamnt_ma_rech90 : Median of amount of recharges done in main account over last 90 days at user level (in Indonesian Rupiah)

medianmarechprebal90 : Median of main account balance just before recharge in last 90 days at user level (in Indonesian Rupiah)

cnt_da_rech30 : Number of times data account got recharged in last 30 days

fr_da_rech30 : Frequency of data account recharged in last 30 days

cnt_da_rech90 : Number of times data account got recharged in last 90 days

fr_da_rech90 : Frequency of data account recharged in last 90 days

cnt_loans30 : Number of loans taken by user in last 30 days

amnt_loans30 : Total amount of loans taken by user in last 30 days

maxamnt_loans30 : Maximum amount of loan taken by the user in last 30 days

medianamnt_loans30: Median of amounts of loan taken by the user in last 30 days

cnt_loans90 : Number of loans taken by user in last 90 days

amnt_loans90 : Total amount of loans taken by user in last 90 days

maxamnt_loans90 : Maximum amount of loan taken by the user in last 90 days

medianamnt_loans90: Median of amounts of loan taken by the user in last 90 days

payback30 : Average payback time in days over last 30 days

payback90 : Average payback time in days over last 90 days

pcircle : Telecom circle

pdate : Date

Data description in a tabular format:

Column Names	Column Definition
label	Flag indicating whether the user paid back the credit amount within 5 days of issuing the loan {1:success, 0:failure}
msisdh	Mobile number of user
aon	Age on cellular network in days
daily_decr30	Daily amount spent from main account, averaged over last 30 days (in Indonesian Rupiah)
daily_decr90	Daily amount spent from main account, averaged over last 90 days (in Indonesian Rupiah)
rental30	Average main account balance over last 30 days
rental90	Average main account balance over last 90 days
last_rech_date_ma	Number of days till last recharge of main account
last_rech_date_da	Number of days till last recharge of data account
last_rech_amt_ma	Amount of last recharge of main account (in Indonesian Rupiah)
cnt_ma_rech30	Number of times main account got recharged in last 30 days
fr_ma_rech30	Frequency of main account recharged in last 30 days
sumamnt_ma_rech30	Total amount of recharge in main account over last 30 days (in Indonesian Rupiah)
medianamnt_ma_rech30	Median of amount of recharges done in main account over last 30 days at user level (in Indonesian Rupiah)
medianmarechprebal30	Median of main account balance just before recharge in last 30 days at user level (in Indonesian Rupiah)
cnt_ma_rech90	Number of times main account got recharged in last 90 days
fr_ma_rech90	Frequency of main account recharged in last 90 days
sumamnt_ma_rech90	Total amount of recharge in main account over last 90 days (in Indonesian Rupiah)
medianamnt_ma_rech90	Median of amount of recharges done in main account over last 90 days at user level (in Indonesian Rupiah)
medianmarechprebal90	Median of main account balance just before recharge in last 90 days at user level (in Indonesian Rupiah)
cnt_da_rech30	Number of times data account got recharged in last 30 days
fr_da_rech30	Frequency of data account recharged in last 30 days
cnt_da_rech90	Number of times data account got recharged in last 90 days
fr_da_rech90	Frequency of data account recharged in last 90 days
cnt_loans30	Number of loans taken by user in last 30 days
amnt_loans30	Total amount of loans taken by user in last 30 days
maxamnt_loans30	Maximum amount of loan taken by the user in last 30 days
medianamnt_loans30	Median of amounts of loan taken by the user in last 30 days
cnt_loans90	Number of loans taken by user in last 90 days
amnt_loans90	Total amount of loans taken by user in last 90 days
maxamnt_loans90	Maximum amount of loan taken by the user in last 90 days
medianamnt_loans90	Median of amounts of loan taken by the user in last 90 days
payback30	Average payback time in days over last 30 days
payback90	Average payback time in days over last 90 days
pcircle	Telecom circle
pdate	Date

- State the set of assumptions (if any) related to the problem under consideration

I had made an assumption that any telecom company keeps the data of customer within 3 months so I have chopped off my data on basis of that.

I have dropped the 2016 year from pdate columns because the data is from the year 2016, only the date and months are different. We separated months and days to different columns.

Then I separately checked the defaulter's data and found that many valuable users are defaulters as they might have forgotten to pay or they are having a busy life. I separated them so that company can deal politely, because we cannot lose these customers.

- Hardware and Software Requirements and Tools Used

Hardware technology being used.

RAM : 8 GB

CPU : AMD Ryzen 5 3550H with Radeon Vega Mobile Gfx 2.10 GHz

GPU : AMD Radeon™ Vega 8 Graphics and NVIDIA GeForce GTX 1650 Ti

Software technology being used.

Programming language : Python

Distribution : Anaconda Navigator

Browser based language shell : Jupyter Notebook

Libraries/Packages specifically being used.

Pandas , NumPy, matplotlib, seaborn, scikit-learn, pandas-profiling, missingno

Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)

We have used the describe method to check the numerical data details. There are 33 columns which have numerical values in them and it looks like the count, mean, standard deviation, minimum value, 25% quartile, 50% quartile, 75% quartile and maximum value are all mostly properly distributed in terms of data points but I do see some abnormality that we will confirm with a visual on it.

label	0.8800	0.3300	0.0000	1.0000	1.0000	1.0000	1.0000
aon	8112.3800	75696.2600	-48.0000	246.0000	527.0000	982.0000	999860.7600
daily_decr30	5381.4100	9220.6400	-93.0100	42.4400	1469.0900	7244.1000	265926.0000
daily_decr90	6082.5300	10918.8400	-93.0100	42.6900	1500.0000	7802.8000	320630.0000
rental30	2692.5800	4308.6000	-23737.1400	280.4200	1083.5400	3356.9400	198926.1100
rental90	3483.4100	5770.4800	24720.5800	300.2600	1334.0000	4201.7900	200148.1100
last_rech_date_ma	3755.8700	53906.0200	-29.0000	1.0000	3.0000	7.0000	998650.3800
last_rech_date_da	3712.2200	53374.9600	-29.0000	0.0000	0.0000	0.0000	999171.8100
last_rech_amt_ma	2064.4600	2370.7900	0.0000	770.0000	1539.0000	2309.0000	55000.0000
cnt_ma_rech30	3.9800	4.2600	0.0000	1.0000	3.0000	5.0000	203.0000
fr_ma_rech30	3737.3700	53643.7500	0.0000	0.0000	2.0000	6.0000	999606.3700
sumamnt_ma_rech30	7704.5000	10139.6500	0.0000	1540.0000	4628.0000	10010.0000	810096.0000
medianamnt_ma_rech30	1812.8200	2070.8700	0.0000	770.0000	1539.0000	1924.0000	55000.0000
medianmarechprebal30	3851.9500	54006.5000	-200.0000	11.0000	33.9000	83.0000	999479.4200
cnt_ma_rech90	6.3200	7.1900	0.0000	2.0000	4.0000	8.0000	336.0000
fr_ma_rech90	7.7200	12.5900	0.0000	0.0000	2.0000	8.0000	88.0000
sumamnt_ma_rech90	12396.2400	16857.8300	0.0000	2317.0000	7226.0000	16000.0000	953036.0000
medianamnt_ma_rech90	1864.6000	2081.6900	0.0000	773.0000	1539.0000	1924.0000	55000.0000
medianmarechprebal90	92.0300	369.2200	-200.0000	14.6000	36.0000	79.3100	41456.5000
cnt_da_rech30	262.5800	4183.9100	0.0000	0.0000	0.0000	0.0000	99914.4400
fr_da_rech30	3749.5100	53885.5400	0.0000	0.0000	0.0000	0.0000	999809.2400
cnt_da_rech90	0.0400	0.4000	0.0000	0.0000	0.0000	0.0000	38.0000
fr_da_rech90	0.0500	0.9500	0.0000	0.0000	0.0000	0.0000	64.0000
cnt_loans30	2.7600	2.5500	0.0000	1.0000	2.0000	4.0000	50.0000
amnt_loans30	17.9500	17.3800	0.0000	6.0000	12.0000	24.0000	306.0000
maxamnt_loans30	274.6600	4245.2700	0.0000	6.0000	6.0000	6.0000	99864.5600
medianamnt_loans30	0.0500	0.2200	0.0000	0.0000	0.0000	0.0000	3.0000
cnt_loans90	18.5200	224.8000	0.0000	1.0000	2.0000	5.0000	4997.5200
amnt_loans90	23.6500	26.4700	0.0000	6.0000	12.0000	30.0000	438.0000
maxamnt_loans90	6.7000	2.1000	0.0000	6.0000	6.0000	6.0000	12.0000
medianamnt_loans90	0.0500	0.2000	0.0000	0.0000	0.0000	0.0000	3.0000
payback30	3.4000	8.8100	0.0000	0.0000	0.0000	3.7500	171.5000
payback90	4.3200	10.3100	0.0000	0.0000	1.6700	4.5000	171.5000
	mean	std	min	25%	50%	75%	max

In the above report we can see that the maximum value for columns aon, daily_decr30, daily_decr90, rental30, rental90, last_rech_date_ma, last_rech_date_da, fr_ma_rech30, sumamnt_ma_rech30, medianmarechprebal30, sumamnt_ma_rech90 and fr_da_rech30 have quite a high number than the other column values.

- Testing of Identified Approaches (Algorithms)

Listing down all the 8 classification machine learning algorithms used for the training and testing.

```
LR = LogisticRegression()
ETC = ExtraTreesClassifier()
SVCM = SVC(C=1.0, kernel='rbf', gamma='auto', random_state=42)
DTC = DecisionTreeClassifier(max_depth=15, random_state=21)
RFC = RandomForestClassifier(max_depth=15, random_state=111)
KNN = KNeighborsClassifier(n_neighbors=15)
XGB = xgb.XGBClassifier(verbosity=0)
LGBM = lgb.LGBMClassifier()

models = {'Logistic Regression' : LR,
          'Extra Trees Classifier' : ETC,
          'Support Vector Classifier' : SVCM,
          'Decision Tree Classifier' : DTC,
          'Random Forest Classifier' : RFC,
          'K Nearest Neighbors Classifier' : KNN,
          'XGB Classifier' : XGB,
          'LGBM Classifier' : LGBM}
```

- Run and Evaluate selected models

I created a Classification Model function incorporating the evaluation metrics so that we can get the required data for all the models.

Machine Learning Model for Classification with Evaluation Metrics

```
# Classification Model Function
def classify(model_func):
    for model_name, model in model_func.items():
        # Training the model
        model.fit(X_train, Y_train)

        # Predicting Y_test
        pred = model.predict(X_test)

        print('\n#####', model_name, '#####')

        # Classification Report
        class_report = classification_report(Y_test, pred)
        print("\nClassification Report for {}: \n".format(model_name), class_report)

        # Accuracy Score
        acc_score = (accuracy_score(Y_test, pred))*100
        print("Accuracy Score for {}:".format(model_name), acc_score)

        # Cross Validation Score
        cv_score = (cross_val_score(model, X, Y, cv=5).mean())*100
        print("Cross Validation Score for {}:".format(model_name), cv_score)

        # Result of accuracy minus cv scores
        result = acc_score - cv_score
        print("\nAccuracy Score - Cross Validation Score is", result)
```


- Key Metrics for success in solving problem under consideration

The key metrics used here were accuracy_score, cross_val_score, classification report, auc_score and confusion matrix. We tried to find out the best parameters and also to increase our scores by using Hyperparameter Tuning and we will be using GridSearchCV method.

1. Cross Validation:

Cross-validation helps to find out the over fitting and under fitting of the model. In the cross validation the model is made to run on different subsets of the dataset which will get multiple measures of the model. If we take 5 folds, the data will be divided into 5 pieces where each part being 20% of full dataset. While running the Cross-validation the 1st part (20%) of the 5 parts will be kept out as a holdout set for validation and everything else is used for training data. This way we will get the first estimate of the model quality of the

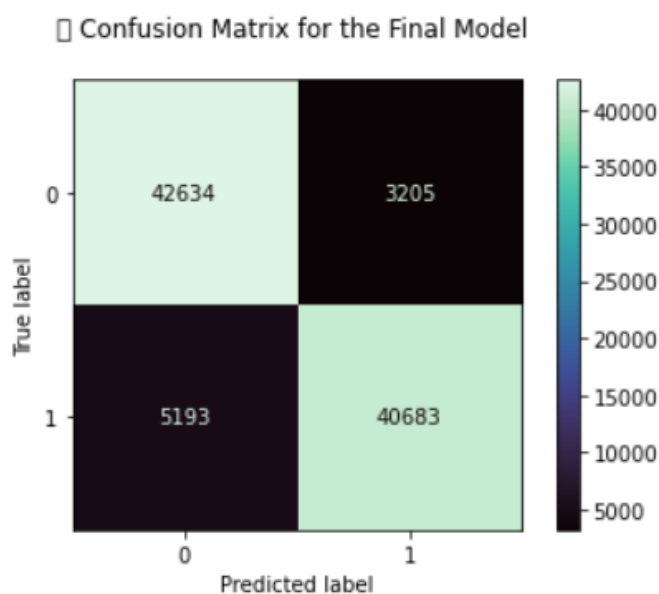
dataset. In the similar way further iterations are made for the second 20% of the dataset is held as a holdout set and remaining 4 parts are used for training data during process. This way we will get the second estimate of the model quality of the dataset. These steps are repeated during the cross-validation process to get the remaining estimate of the model quality.

2. Confusion Matrix: A confusion matrix, also known as an error matrix, is a specific table layout that allows visualization of the performance of an algorithm, typically a supervised learning one (in unsupervised learning it is usually called a matching matrix). Each row of the matrix represents the instances in a predicted class, while each column represents the instances in an actual class (or vice versa). The name stems from the fact that it makes it easy to see whether the system is confusing two classes (i.e., commonly mislabelling one as another).

It is a special kind of contingency table, with two dimensions ("actual" and "predicted"), and identical sets of "classes" in both dimensions (each combination of dimension and class is a variable in the contingency table).

Confusion Matrix

```
metrics.plot_confusion_matrix(Classifier, X_test, Y_test, cmap='mako')  
plt.title('\t Confusion Matrix for the Final Model \n')  
plt.show()
```



3. Classification Report: The classification report visualizer displays the precision, recall, F1, and support scores for the model. There are four ways to check if the predictions are right or wrong: 1. TN / True Negative: the case was negative and predicted negative 2. TP / True Positive: the case was positive and predicted positive 3. FN / False Negative: the case was positive but predicted negative 4. FP / False Positive: the case was negative but predicted positive

Precision: Precision is the ability of a classifier not to label an instance positive that is actually negative. For each class, it is defined as the ratio of true positives to the sum of a true positive and false positive. It is the accuracy of positive predictions. The formula of precision is given below: $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$

Recall: Recall is the ability of a classifier to find all positive instances. For each class it is defined as the ratio of true positives to the sum of true positives and false negatives. It is also the fraction of positives that were correctly identified. The formula of recall is given below:

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

F1 score: The F1 score is a weighted harmonic mean of precision and recall such that the best score is 1.0 and the worst is 0.0. F1 scores are lower than accuracy measures as they embed precision and recall into their computation. As a rule of thumb, the weighted average of F1 should be used to compare classifier models, not global accuracy. The formula is: $\text{F1 Score} = 2 * (\text{Recall} * \text{Precision}) / (\text{Recall} + \text{Precision})$

Support: Support is the number of actual occurrences of the class in the specified dataset. Imbalanced support in the training data may indicate structural weaknesses in the reported scores of the classifier and could indicate the need for stratified sampling or rebalancing. Support doesn't change between models but instead diagnoses the evaluation process.

Extra Trees Classifier

Classification Report for Extra Trees Classifier:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	45930
1	0.95	0.95	0.95	45785
accuracy			0.95	91715
macro avg	0.95	0.95	0.95	91715
weighted avg	0.95	0.95	0.95	91715

Accuracy Score for Extra Trees Classifier: 95.1207545112577

Cross Validation Score for Extra Trees Classifier: 94.83481437060458

Accuracy Score - Cross Validation Score is 0.2859401406531248

4. AUC-ROC Curve and score:

AUC (Area Under the Curve) - ROC (Receiver Operating Characteristics) curve is a performance measurement for the classification problems at various threshold settings. ROC is a probability curve and AUC represent the degree or measure of

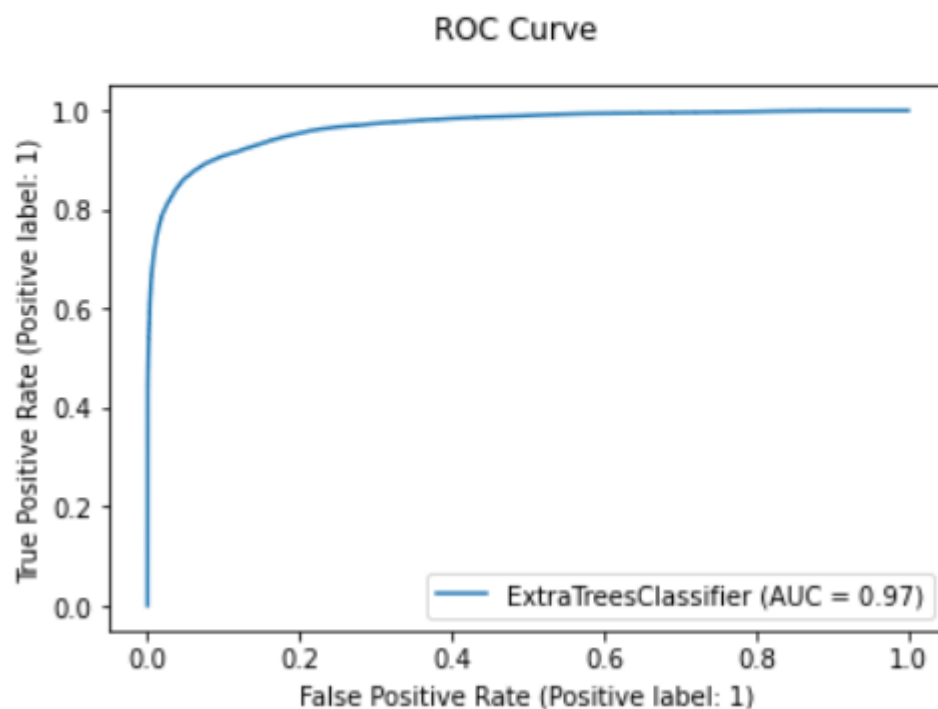
separability. It tells how much the model is capable of distinguishing between classes. Higher the AUC, the better the model is at predicting 0s as 0s and 1s as 1s. By analogy, the Higher the AUC, the better the model is at distinguishing between patients with the disease and no disease.

The ROC curve is plotted with TPR against the FPR where TPR is on the y-axis and FPR is on the x-axis.

Score is the area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores.

AUC ROC Curve

```
disp = metrics.plot_roc_curve(Final_Model, X_test, Y_test)
disp.figure_.suptitle("ROC Curve")
plt.show()
```



5. Hyperparameter Tuning: There is a list of different machine learning models. They all are different in some way or the other, but what makes them different is nothing but input parameters for the model. These input parameters are named as Hyperparameters.

These hyperparameters will define the architecture of the model, and the best part about these is that you get a choice to select these for your model. You must select from a specific list of hyperparameters for a given model as it varies from model to model. We are not aware of optimal values for hyperparameters which would generate the best model output. So, what we tell the model is to explore and select the optimal model architecture automatically. This selection procedure for hyperparameter is known as Hyperparameter Tuning. We can do tuning by using GridSearchCV. GridSearchCV is a function that comes in Scikit-learn (or SK-learn) model selection package. An important point here to note is that we need to have Scikit-learn library installed on the computer. This function helps to loop through predefined hyperparameters and fit your estimator (model) on your training set. So, in the end, we can select the best parameters from the listed hyperparameters.

Hyper parameter tuning on the best Classification ML Model

```
# Choosing Extra Trees Classifier
```

```
fmod_param = {'criterion' : ["gini", "entropy"],  
              'max_depth' : [30, 40],  
              'n_estimators' : [300, 350],  
              'min_samples_split' : [3, 4],  
              'random_state' : [42, 72]  
              }
```

```
GSCV = GridSearchCV(ExtraTreesClassifier(), fmod_param, cv=5)  
GSCV.fit(X_train, Y_train)  
GSCV.best_params_
```

```
Final_Model = ExtraTreesClassifier(criterion="entropy", max_depth=30, min_samples_split=3,  
                                  n_estimators=350, random_state=72)  
Classifier = Final_Model.fit(X_train, Y_train)  
fmod_pred = Final_Model.predict(X_test)  
fmod_acc = (accuracy_score(Y_test, fmod_pred))*100  
print("Accuracy score for the Best Model is:", fmod_acc)
```

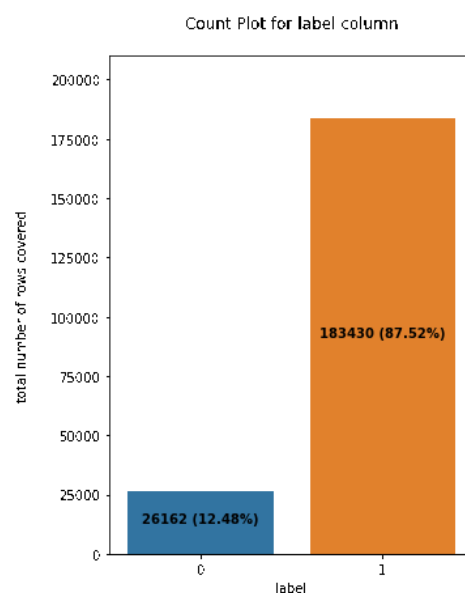
```
Accuracy score for the Best Model is: 90.8433734939759
```

- Visualizations

Now, we will see the different plots done with this dataset in order to know the insight of the data present. Below are the codes given for the plots and the output obtained:

Univariate Analysis

```
try:
    x = 'label'
    k=0
    plt.figure(figsize=[5,7])
    axes = sns.countplot(df[x])
    for i in axes.patches:
        ht = i.get_height()
        mr = len(df[x])
        st = f"{ht} ({round(ht*100/mr,2)}%)"
        plt.text(k, ht/2, st, ha='center', fontweight='bold')
        k += 1
    plt.ylim(0,210000)
    plt.title(f'Count Plot for {x} column\n')
    plt.ylabel(f'total number of rows covered\n')
    plt.show()
except Exception as e:
    print("Error:", e)
    pass
```



Bivariate Analysis

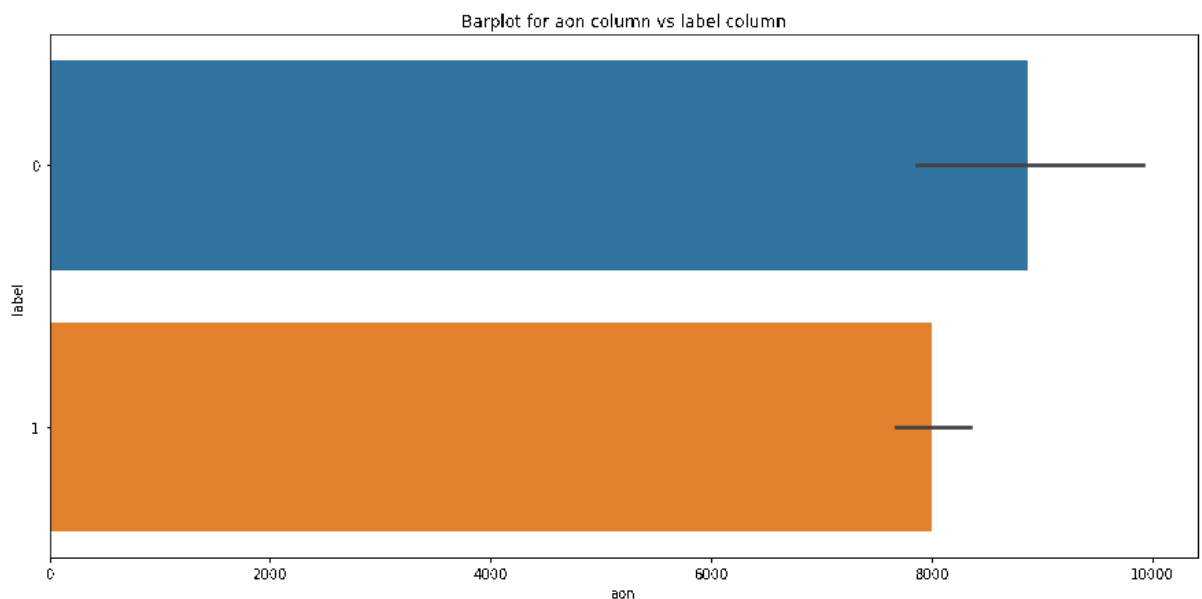
```
y = 'label'

x = 'aon'
plt.figure(figsize=[15,7])
sns.barplot(x,y,data=df,orient='h')
plt.title(f"Barplot for {x} column vs {y} column")
plt.show()

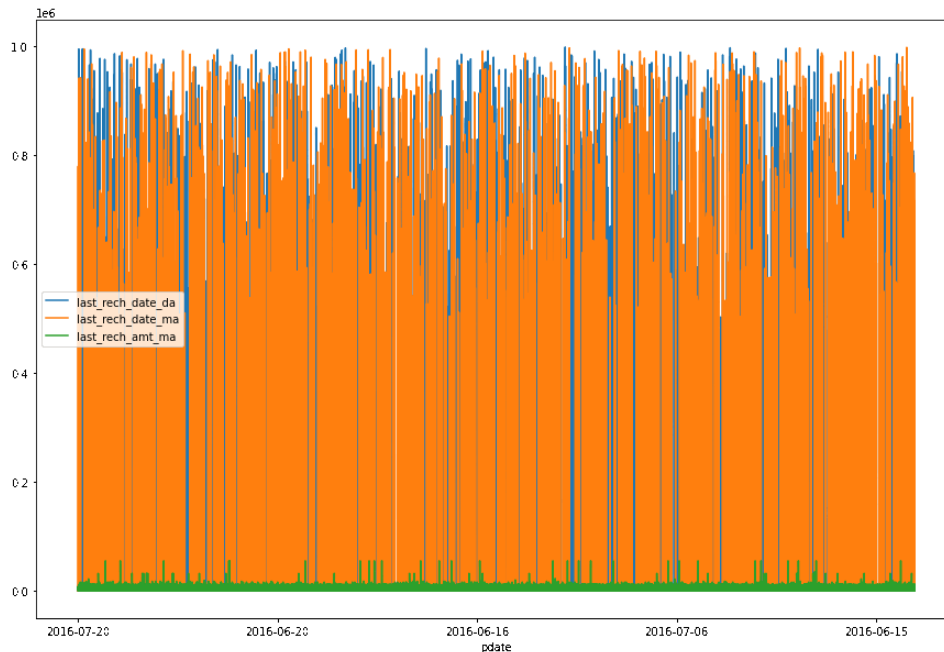
x = 'last_rech_date_da'
plt.figure(figsize=[15,7])
sns.barplot(x,y,data=df,orient='h')
plt.title(f"Barplot for {x} column vs {y} column")
plt.show()

x = 'last_rech_date_ma'
plt.figure(figsize=[15,7])
sns.barplot(x,y,data=df,orient='h')
plt.title(f"Barplot for {x} column vs {y} column")
plt.show()

x = 'last_rech_amt_ma'
plt.figure(figsize=[15,7])
sns.barplot(x,y,data=df,orient='h')
plt.title(f"Barplot for {x} column vs {y} column")
plt.show()
```



```
df.plot(kind="line", x="pdate", y=["last_rech_date_da", "last_rech_date_ma", "last_rech_amt_ma"], figsize=[15,10])
df.plot(kind="line", x="msisdn", y=["last_rech_date_da", "last_rech_date_ma", "last_rech_amt_ma"], figsize=[15,10])
```



```
plt.figure(figsize=(15,5))
sns.scatterplot(x='medianamnt_loans30', y='medianamnt_loans90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='maxamnt_loans30', y='maxamnt_loans90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='cnt_da_rech30', y='cnt_da_rech90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='cnt_loans30', y='cnt_loans90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='amnt_loans30', y='amnt_loans90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='cnt_ma_rech30', y='cnt_ma_rech90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='fr_da_rech30', y='fr_da_rech90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='fr_ma_rech30', y='fr_ma_rech90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='medianamnt_ma_rech30', y='medianamnt_ma_rech90', data=df, hue='label')

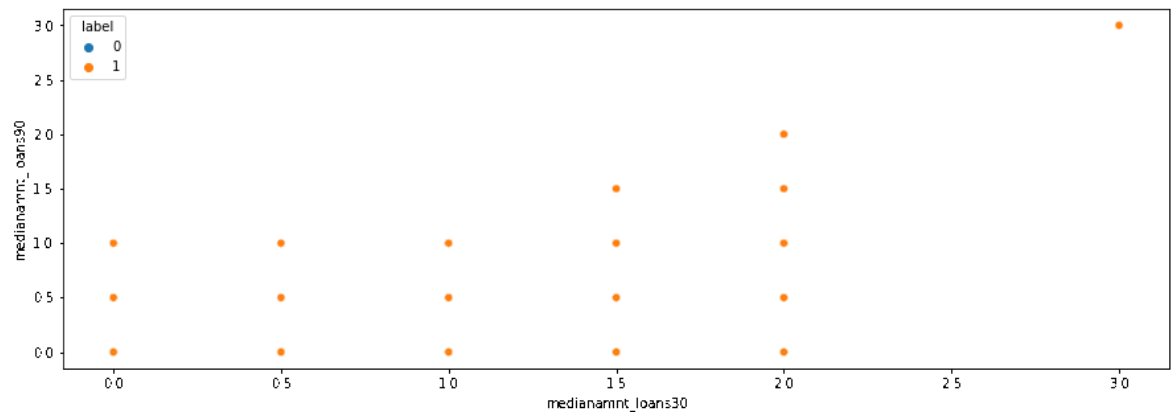
plt.figure(figsize=(15,5))
sns.scatterplot(x='daily_decr30', y='daily_decr90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='rental30', y='rental90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='payback30', y='payback90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='medianmarechprebal30', y='medianmarechprebal90', data=df, hue='label')

plt.figure(figsize=(15,5))
sns.scatterplot(x='sumamnt_ma_rech30', y='sumamnt_ma_rech90', data=df, hue='label')
```

● Interpretation of the Results

for feature aon:

Data ranges from -48 to 999860 with Mean value of 8112.34.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature daily_descr30:

Data ranges from -93 to 265926 with Mean value of 5381.4.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature daily_descr90:

Data ranges from -93 to 320630 with Mean value of 6082.52.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature rental30:

Data ranges from -23737.14 to 198926 with Mean value of 2692.58.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature rental90:

Data ranges from -24720 to 200148 with Mean value of 3483.41.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature last_rech_date_ma:

Data ranges from -29 to 998650 with Mean value of 3755.85.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature last_rech_date_da:

Data ranges from -29 to 999178 with Mean value of 3712.2.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature last_rech_amt_ma:

Data ranges from 0 to 55000 with Mean value of 2064.45.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature cnt_ma_rech30:

Data ranges from 0 to 203 with Mean value of 3.98.

Data is not distributed normally or in well curve.

Data is spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature fr_ma_rech30:

Data ranges from 0 to 999606 with Mean value of 3737.36.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature sumamnt_ma_rech30:

Data ranges from 0 to 810096 with Mean value of 7704.5.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature medianamnt_ma_rech30:

Data ranges from 0 to 55000 with Mean value of 1812.82.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature medianmarechprebal30:

Data ranges from -200 to 999479 with Mean value of 3851.93.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature cnt_ma_rech90:

Data ranges from 0 to 336 with Mean value of 6.32.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature fr_ma_rech90:

Data ranges from 0 to 88 with Mean value of 7.72.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature sumamnt_ma_rech90:

Data ranges from 0 to 953036 with Mean value of 12396.22.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature medianamnt_ma_rech90:

Data ranges from 0 to 55000 with Mean value of 1864.6.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature medianmarechprebal90:

Data ranges from -200 to 41456 with Mean value of 92.03.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature cnt_da_rech30:

Data ranges from 0 to 99914 with Mean value of 262.58.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature fr_da_rech30:

Data ranges from 0 to 999809 with Mean value of 3749.49.

Data is not distributed normally or in well curve.

Data is highly spreaded and needs to be treated accordingly.

Data is positively skewed and needs to be treated accordingly.

for feature cnt_da_rech90:

Data ranges from 0 to 38 with Mean value of 0.04.

Data is distributed normally but not in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature fr_da_rech90:

Data ranges from 0 to 64 with Mean value of 0.05.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature cnt_loans30:

Data ranges from 0 to 50 with Mean value of 2.76.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature amnt_loans30:

Data ranges from 0 to 306 with Mean value of 17.95.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature maxamnt_loans30:

Data ranges from 0 to 99864 with Mean value of 274.66.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature medianamnt_loans30:

Data ranges from 0 to 3 with Mean value of 0.05.

Data is not distributed normally or in well curve and it is understandable as feature has only limited set of values.

Data is positively skewed and needs to be treated accordingly.

for feature cnt_loans90:

Data ranges from 0 to 4997.52 with Mean value of 18.52.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature amnt_loans90:

Data ranges from 0 to 438 with Mean value of 23.65.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature maxamnt_loans90:

Data ranges from 0 to 12 with Mean value of 6.7.

Data is not distributed normally or in well curve and it understandable as user has two option for loans i.e., 5 and 10 for with 6 and 12 has to be paid.

Data is positively skewed and needs to be treated accordingly.

for feature medianamnt_loans90:

Data ranges from 0 to 3 with Mean value of 0.05.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature payback30:

Data ranges from 0 to 171.5 with Mean value of 3.4.

Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

for feature payback90:

Data ranges from 0 to 171.5 with Mean value of 4.32.

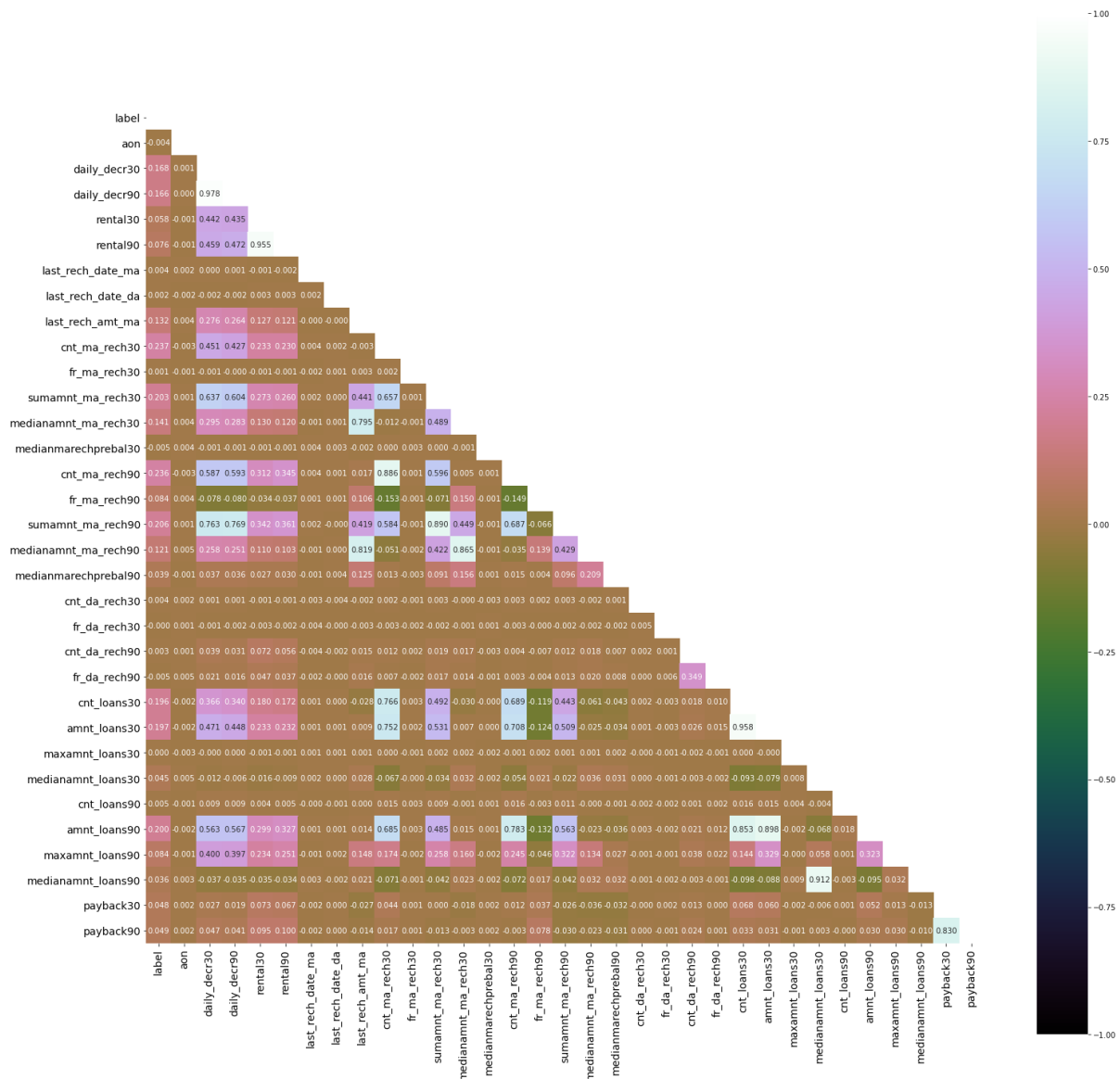
Data is not distributed normally or in well curve.

Data is positively skewed and needs to be treated accordingly.

Correlation using a Heatmap

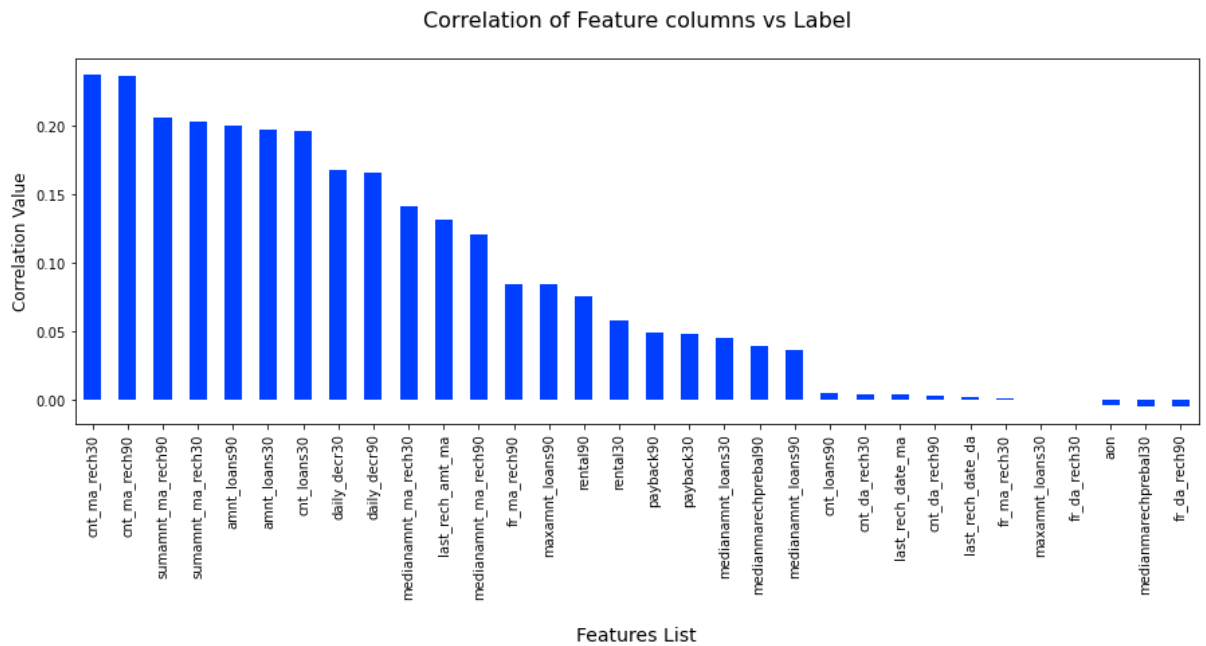
- Positive correlation - A correlation of +1 indicates a perfect positive correlation, meaning that both variables move in the same direction together.
- Negative correlation - A correlation of -1 indicates a perfect negative correlation, meaning that as one variable goes up, the other goes down.

```
upper_triangle = np.triu(df.corr())
plt.figure(figsize=(25,25))
sns.heatmap(df.corr(), vmin=-1, vmax=1, annot=True, square=True, fmt='.3f',
            annot_kws={'size':10}, cmap="cubehelix", mask=upper_triangle)
plt.xticks(fontsize=14)
plt.yticks(fontsize=14)
plt.show()
```



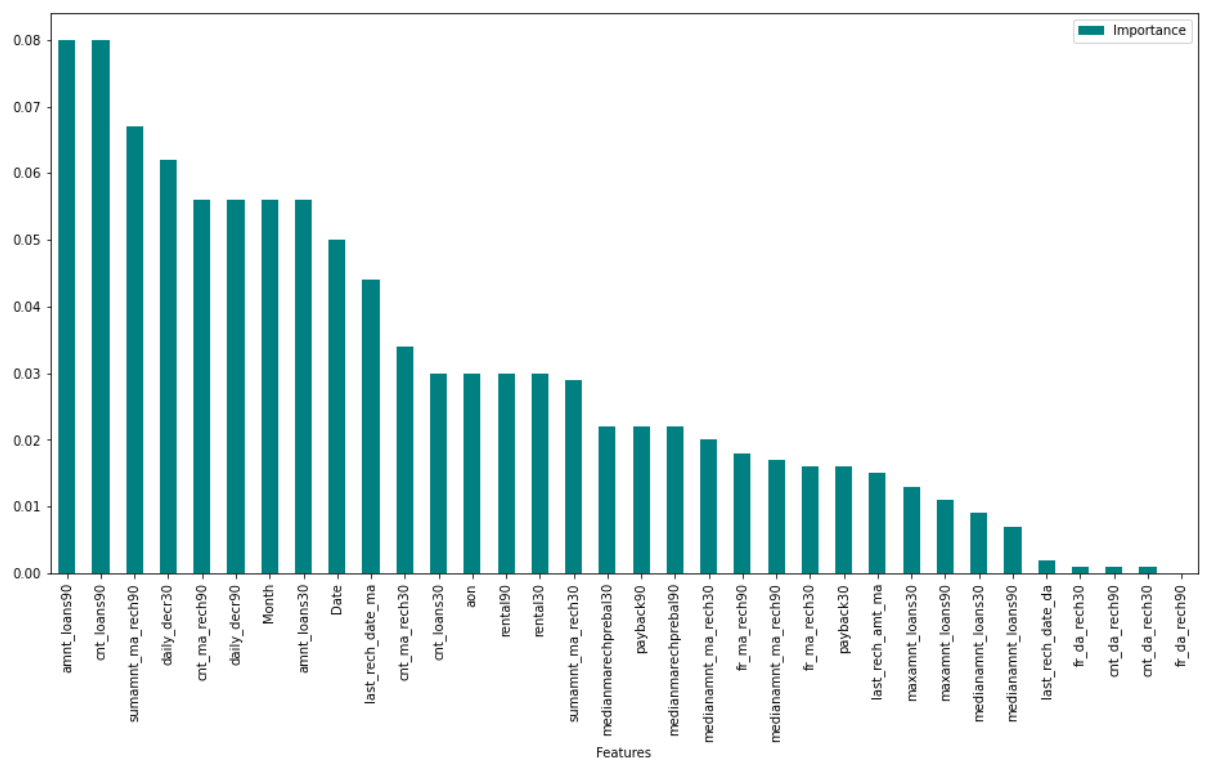
Correlation Bar Plot comparing Gender column with the remaining columns

```
df_corr = df.corr()
plt.figure(figsize=(15,5))
df_corr['label'].sort_values(ascending=False).drop('label').plot.bar()
plt.title("Correlation of Feature columns vs Label\n", fontsize=16)
plt.xlabel("\nFeatures List", fontsize=14)
plt.ylabel("Correlation Value", fontsize=12)
plt.show()
```



Feature importance bar graph

```
rf=RandomForestClassifier()
rf.fit(X_train, Y_train)
importances = pd.DataFrame({'Features':x.columns, 'Importance':np.round(rf.feature_importances_,3)})
importances = importances.sort_values('Importance', ascending=False).set_index('Features')
plt.rcParams["figure.figsize"] = (16,8)
importances.plot.bar(color='teal')
importances
```



CONCLUSION

- Key Findings and Conclusions of the Study

From the final model MFI can find if a person will return money or not and should an MFI provide a loan to that person or not judging from the various features taken into consideration.

- Learning Outcomes of the Study in respect of Data Science

I built multiple classification models and did not rely on one single model for getting better accuracy and using cross validation comparison I ensured that the model does not fall into overfitting and underfitting issues. I picked the best one and performed hyper parameter tuning on it to enhance the scores.

- Limitations of this work and Scope for Future Work

Limitation is it will only work for this particular use case and will need to be modified if tried to be utilized on a different scenario but on a similar scale. Scope is that we can use it in companies to find whether we should provide loan to a person or not and we can also make prediction about a person buying an expensive service on the basis of their personal details that we have in this dataset like number of times data account got recharged in last 30 days and daily amount spent from main account, averaged over last 30 days (in Indonesian Rupiah) so even a marketing company can also use this.

