



NAME OF THE PROJECT

Micro-credit Defaulters



Submitted by:

NIDHI SINGH

ACKNOWLEDGMENT

This includes mentioning of all the references, research papers, data sources, professionals and other resources that helped me and guided me in completion of the project.

- I would like to thank FlipRobo Technologies for providing me this opportunity and guidance throughout the project and all the steps that are implemented.
- I have primarily referred to various articles scattered across various websites for the purpose of getting an idea on “*Micro-credit defaulters*” project.
- I would like to thank the technical support team also for helping me out and reaching out to me on clearing all my doubts as early as possible.
- I would like to thank my project SME Mr. Shwetank Mishra for providing the flexibility in time and for giving us guidance in creating the project.
- I have referred to various articles in Towards Data Science and Kaggle



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

- **Conceptual Background of the Domain Problem**

- We are working with client that is in Telecom Industry .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.
- Telecom Industry 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). The sample data is provided to us from our client database. It is hereby given to you 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.

- **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, which affect the regulations of profitability and impact assessment of MFIs have been summarized in the below report. We hope that the below report of 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

The various applications and methods which inspired us to build our project. We did a background survey regarding the basic ideas of our project and used those ideas for the collection of information like the technological stack, algorithms, and shortcomings of our project which led us to build a better project.

I have built 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.

- **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 as 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**

I am working with the micro credit defaulters dataset that contains various features and information about it. Using the data in form of 'read_csv' function provided by the Pandas package, which can import the data into our python environment. After importing the data, I have used the 'head' function to get a glimpse of our dataset.

In this label is used as my target column and it was having two classes Label '1' indicates that the loan has been paid i.e. Non-defaulter, whereas Label '0' indicates that the loan

has not been paid i.e. defaulter. It's clarify the binary classification problem, classification of algorithms for building model. There is no null values in the dataset and observed some unnecessary entries in some columns like in some columns it found more than 90% , zero values so dropped those columns. Those columns will create high skewness in the model.

To get better insight on the features uses plotting function like distribution plot, bar plot and count plot. With these plotting it is able to understand the relation between the features in better manner. Also outliers and skewness found in the dataset so it is removed outliers using percentile method and skewness using yeo-johnson method. classification algorithms while building model then tuned the best and saved the best model. Lastly predicted the label using saved model.

- **Data Sources and their formats**

Dataset has been provided by internship company – Flip Robo technologies in excel format. The sample data is provided from our client database. Data given is only for academic use, not for any commercial. 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. Also, dataset was having 209593 rows and 36 columns including target. In this particular datasets I have object, float and integer types of data. The dataset is in both numerical as well as categorical data. There may be some customers with no loan history. The dataset is imbalanced. Label '1' has approximately 87.5% records, while, label '0' has approximately 12.5% records.

Link for Dataset description: [Micro-Credit-Defaulter-Project/Data_Description.xlsx at main · DS0003/Micro-Credit-Defaulter-Project \(github.com\)](#)

- **Data Pre-processing Done**

- In order to get a better understanding of the data, we plotted a histogram of the data. We noticed that the dataset had many outliers, so removed outliers using percentile method .

We noticed skewness present in our dataset that we removed using yeo-johnson method. However, there were many data points that did not conform to this. This is because accident history and condition can have a significant effect of defaulter or non-defaulter, we pruned our dataset to standard deviations around the mean in order to remove outliers.

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

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

Exploratory Data Analysis (EDA)

- This section shows the exploration done on the dataset, which is what motivated the use of the algorithm. The following are the questions explored in this project and for the sake of writing I will only show some of the visuals here while I will provide the codes that shows the full visualization of all the questions explored.
- Is there a significant relationship between Non-defaulter & defaulter? It was used to check for this and we can see that there is a relationship between Label '1' as Non-defaulter, whereas Label '0' indicates that the loan has not been paid i.e. defaulter.
- Dataset is imbalanced. Label '1' has approximately 87.5% records, while, label '0' has approximately 12.5% records. Need to balance.
- There are two primary phases in the system:
 1. Training phase: The system is trained by using the data in the data set and fits a model (line/curve) based on the algorithm chosen accordingly.
 2. Testing phase: the system is provided with the inputs and is tested for its working. The accuracy is checked. And therefore, the data that is used to train the model or test it, has to be appropriate. The system is designed to detect and

predict and hence appropriate algorithms must be used to do the two different tasks. Before the algorithms are selected for further use, different algorithms were compared for its accuracy. The well-suited one for the task was chosen.

Data cleaning:

Remove outliers

```
# feature contain outliers  
featr=df[['aon', 'daily_decr30', 'daily_decr90', 'rental30', 'rental90', 'last_rech_date_ma', 'last_rech_amt_ma', 'cnt_ma_rech30'
```

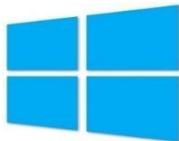
Steps :

- ✚ Importing the required packages into our python environment
- ✚ Importing the data to do some EDA on it
- ✚ Dataset having 209593 rows and 36 columns including target.
- ✚ Data Visualization
- ✚ Feature Selection & Data Split
- ✚ Modelling the data using the algorithms
- ✚ Evaluating the built model using the evaluation metrics

- State the set of assumptions (if any) related to the problem under consideration
- Finally, we conclude which model is best suitable for the given case by evaluating each of them using the evaluation metrics provided by the scikit-learn package. This 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
Technological stack, algorithms, and shortcomings of the project which led to build this project.

- Hardware and Software Requirements and Tools Used

- Listing down the hardware and software requirements along with the tools, libraries and packages used.
- Windows 10 64bit



- Anaconda 2021 / Python version – Python 3.9.5 (latest)
- Software: Jupyter notebook, Python, Panda library, numpy library, Matplotlib library, Seaborn library

- **Python:** Python is a general-purpose, and high-level programming language which is best known for its efficiency and powerful functions. Its ease to use, which makes it more accessible. Python provides data scientists with an extensive amount of tools and packages to build machine learning models. One of its special features is that we can build various machine learning with less-code.
- **Matplotlib** is a plotting library for the Python programming language and its numerical mathematics extension NumPy.
- **Seaborn** is a library for making statistical graphics in Python. It builds on top of matplotlib and integrates closely with pandas data structures. Seaborn helps you explore and understand your data.
- **NumPy** is a general-purpose array-processing package. it provides a high-performance multidimensional array object and tools for working with these arrays. It is the fundamental package for scientific computing with Python. Besides its obvious scientific uses, NumPy can also be used as an efficient multi-dimensional container of generic data. Arbitrary data-types can be defined using Numpy which allows NumPy to seamlessly and speedily integrate with a wide variety of databases.
- **Scikit-learn** provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python. It is licensed under a permissive simplified BSD license and is distributed under many Linux distributions, encouraging academic and commercial use. The library is built
- **Jupyter notebook:** The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations, and narrative text. It includes data cleaning and transformation, numerical simulation, statistical modelling, data visualization, machine learning.
The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. It includes data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning, and much more.

Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)
- The factors need to be found which can impact the micro credit. This can be done by analysing the various factors and the stores the respondent prefers. This will be done by checking each of the factors impacts the respondents in decision making.
- Machine Learning Algorithms:

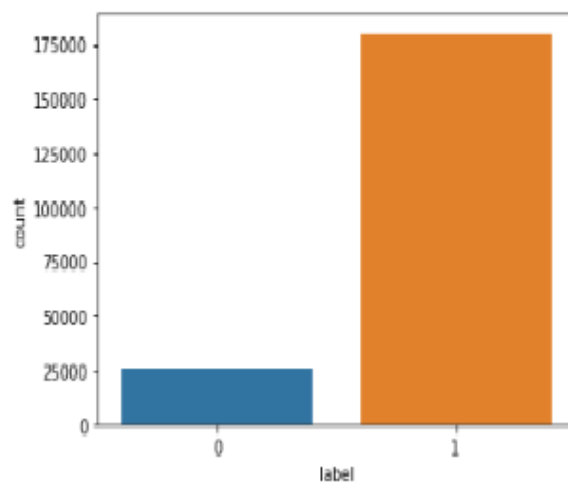
Machine learning-based systems are growing in popularity in research applications in most disciplines. Considerable decision-making knowledge from data has been acquired in the broad area of machine learning, in which decision-making tree-based ensemble techniques are recognized for supervised classification problems. Thus, classification is an essential form of data analysis in data mining that formulates models while describing significant data.

• Visualizations

- As the value counts observation I find imbalance dataset in which defaulter values is less and Non defaulter values is high. About to 15% and 85% respectively
- Dataset is imbalanced. Label '1' has approximately 87.5% records, while, label '0' has approximately 12.5% records. Need to balance.

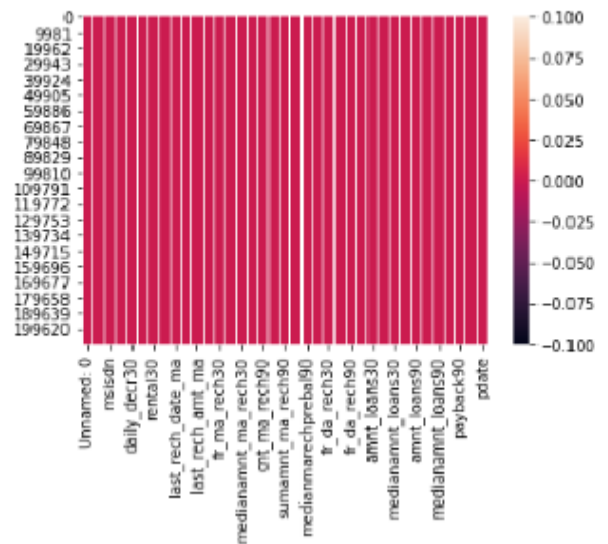
```
#count for target column
sns.countplot(df['label'])
```

```
<AxesSubplot: xlabel='label', ylabel='count'>
```

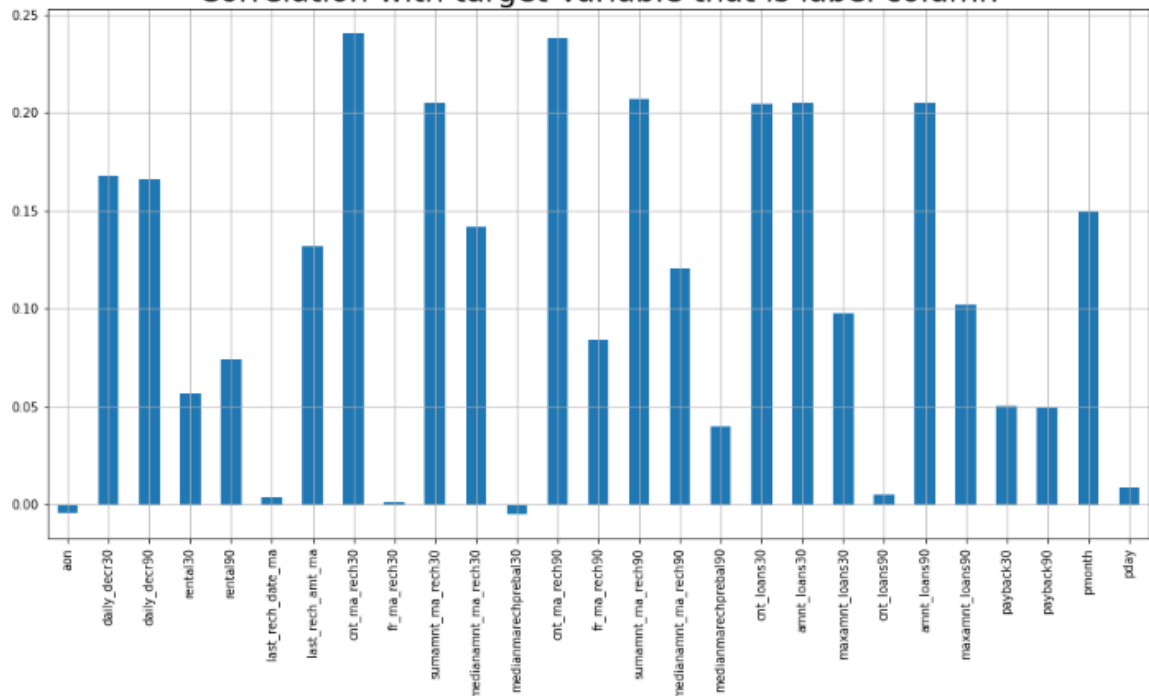


```
sns.heatmap(df.isnull())
```

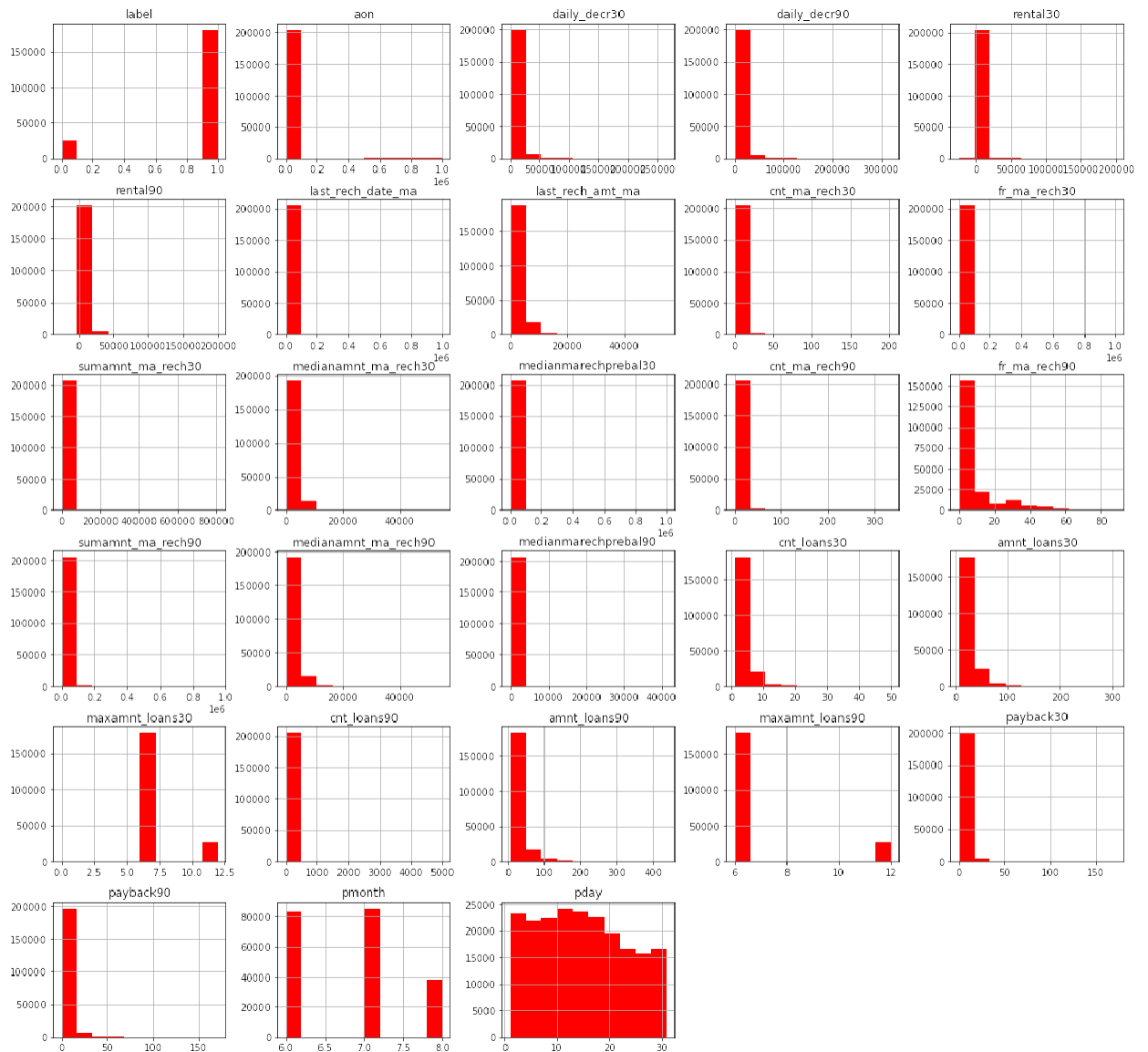
```
<AxesSubplot: >
```



Correlation with target Variable that is label column

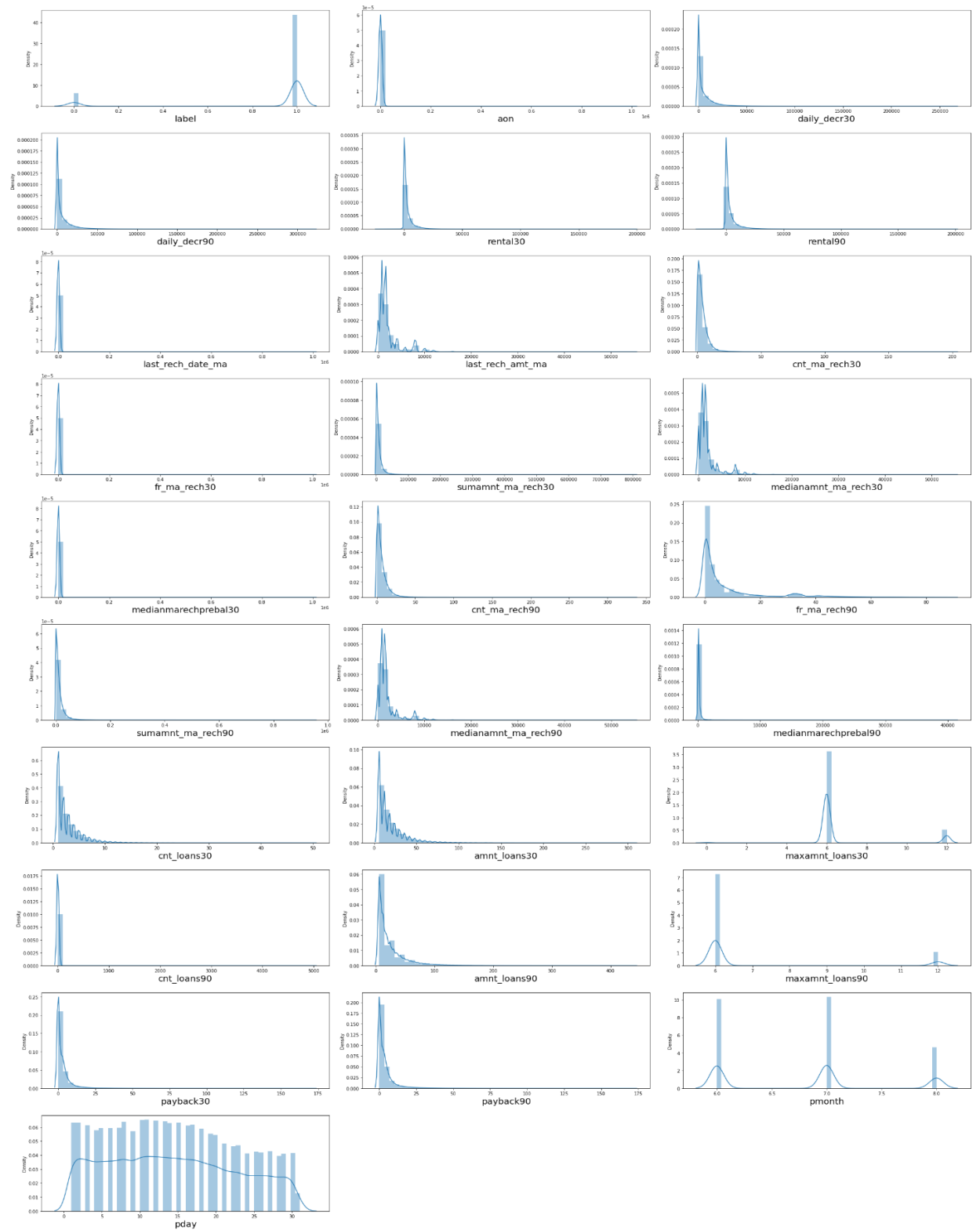


Plotting the Histogram



- To remove outliers I have used percentile method. And to remove skewness I have used yeo-johnson method. We have dropped all the unnecessary columns in the dataset according to our understanding. Use of Pearson's correlation coefficient to check the correlation between dependent and independent features. Also I have used Normalization to scale the data.

Bar plots to see the relation of numerical feature with target and 2 types of plots for numerical columns like distribution plot for univariate and bar plot for bivariate analysis.



Outliers in most of the columns so we have to treat them using suitable methods.



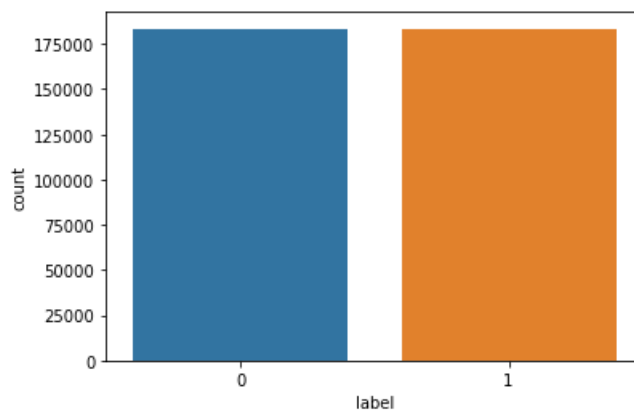
- After scaling we have to balance the target column using oversampling. Then followed by model building with all Classification algorithms. I have used oversampling (SMOTE) to get rid of data unbalancing.

SMOTE (STANDARDISING TARGET COLUMN)

```
] : #Oversampling the data
    from imblearn.over_sampling import SMOTE
    SM = SMOTE()
    x, y = SM.fit_resample(x,y)

] : #Visualize the data after balancing
    sns.countplot(y)

] : <AxesSubplot:xlabel='label', ylabel='count'>
```



MACHINE LEARNING

Selecting best Random State with Best Accuracy

```
69]: #Lets find best random state and accuracy score
max_acc=0
max_rs=0
for i in range(1,200):
    x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.30, random_
    mod=RandomForestClassifier()
    mod.fit(x_train,y_train)
    pred=mod.predict(x_test)
    acc=accuracy_score(y_test,pred)
    if acc>max_acc:
        max_acc=acc
        max_rs=i

print("Best Accuracy is ", max_acc, "on random_state ",max_rs)
```

Best Accuracy is 0.9550604675673957 on random state 0

- Run and Evaluate selected models
- Random Forest Classifier, Decision Tree Classifier, Extra Tree Classifier, K Neighbor Classifier, Logistic Regressor were our baseline methods. For most of the model implementations, the open-source Scikit-Learn package was used.

- Our primary packages for this project are going to be pandas for data processing, NumPy to work with arrays, matplotlib & seaborn for data visualizations, and finally scikit-learn for building an evaluating our ML model.

```
[82]: #Logistic Regression
model = LogisticRegression()
classifier(model, x,y)

Accuracy Score: 76.68614106978984

Classification Report:
      precision    recall  f1-score   support

     0       0.76       0.77       0.77       54730
     1       0.77       0.76       0.77       55329

 accuracy         0.77         0.77         0.77       110059
 macro avg       0.77         0.77         0.77       110059
weighted avg       0.77         0.77         0.77       110059

Cross Validation Score: 76.68496642183295

Accuracy Score - Cross Validation Score is 0.0011746479568870427
```

```
#DECISION TREE CLASSIFIER
model = DecisionTreeClassifier()
classifier(model, x, y)

Accuracy Score: 91.11203990586867

Classification Report:
      precision    recall  f1-score   support

     0       0.90       0.92       0.91       54730
     1       0.92       0.90       0.91       55329

 accuracy         0.91         0.91         0.91       110059
 macro avg       0.91         0.91         0.91       110059
weighted avg       0.91         0.91         0.91       110059

Cross Validation Score: 91.0416045947137

Accuracy Score - Cross Validation Score is 0.07043531115496648
```

```
[87]: #K NEIGHBORS CLASSIFIER
model = KNeighborsClassifier()
classifier(model, x, y)

Accuracy Score: 89.57831708447287

Classification Report:
      precision    recall  f1-score   support

     0       0.83       0.99       0.90       54730
     1       0.99       0.80       0.89       55329

 accuracy         0.91         0.90         0.90       110059
 macro avg       0.91         0.90         0.89       110059
weighted avg       0.91         0.90         0.89       110059

Cross Validation Score: 90.08346460057861

Accuracy Score - Cross Validation Score is -0.5051475161057368
```

```
[88]: #EXTRATREES CLASSIFIER
model = ExtraTreesClassifier()
classifier(model, x, y)

Accuracy Score: 96.03303682570258

Classification Report:
      precision    recall  f1-score   support

     0       0.95       0.98       0.96       54730
     1       0.98       0.94       0.96       55329

 accuracy         0.96         0.96         0.96       110059
 macro avg       0.96         0.96         0.96       110059
weighted avg       0.96         0.96         0.96       110059

Cross Validation Score: 96.42290802217832

Accuracy Score - Cross Validation Score is -0.3898711964757382
```

```
: #RANDOM FOREST CLASSIFIER
model = RandomForestClassifier(random_state=0)
classifier(model, x, y)

Accuracy Score: 95.30524536839332

Classification Report:
      precision    recall  f1-score   support

     0       0.95       0.96       0.95       54730
     1       0.96       0.95       0.95       55329

 accuracy         0.95         0.95         0.95       110059
 macro avg       0.95         0.95         0.95       110059
weighted avg       0.95         0.95         0.95       110059

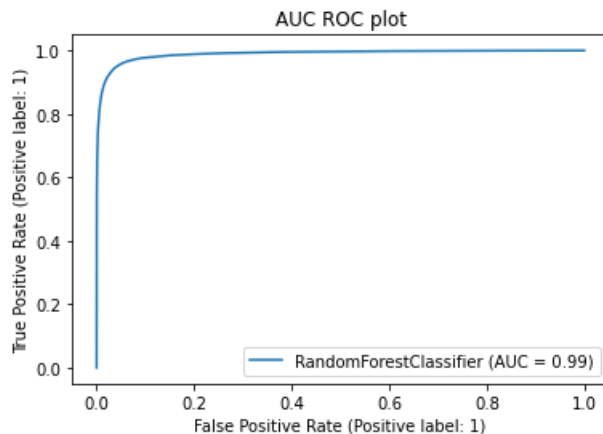
Cross Validation Score: 95.20338679081331

Accuracy Score - Cross Validation Score is 0.10185857758000338
```

- **Key Metrics for success in solving problem under consideration**

- Accuracy score is used when the True Positives and True negatives are more important. Accuracy can be used when the class distribution is similar.
- 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

Using Hyper-parameter : model parameters are estimated from data automatically and model hyper-parameters are set manually and are used in processes to help estimate model and Grid search is a basic method for hyper-parameter tuning. It performs an exhaustive search on the hyper-parameter set specified by users.



HYPER PARAMETER TUNING

```
: # Lets import all required libraries for hyper parameter tuning.
from sklearn.model_selection import GridSearchCV

: parameters={'criterion':['gini', 'entropy'],
              'n_estimators':[100,200,300],
              'max_features':['auto', 'sqrt', 'log2'],
              }

: gcv=GridSearchCV(RandomForestClassifier(),parameters,cv=5)

: gcv.fit(x_train,y_train)

: 
  > GridSearchCV
  > estimator: RandomForestClassifier
    > RandomForestClassifier

: gcv.best_params_

: {'criterion': 'entropy', 'max_features': 'log2', 'n_estimators': 300}
```


Roc & Auc

- Present the receiver operating characteristic (ROC) curves and their respective areas under the curve (AUCs). ROC curves and AUCs are used to measure the quality of a classifier's output; thus, they measure how correctly a classifier has been tuned. Movement along the ROC curve is typically a trade-off between the classifier's sensitivity (true positive rate (TPR)) and specificity (TNR), and the steeper the curve, the better. For the ROC curve, sensitivity increases as we move up, and specificity decreases as we move right. The ROC curve along a 45_ angle
- Interpretation of the Results
 - In this research, two experiments were performed, the first experiment was validating and filtering data using all the variables available in the dataset after pre-processing, while the second experiment was conducted using most important variables and the goal of this is to be able to improve the model's performance using fewer variables.
 - Requirement of train and test and building of many models to get accuracy of the model.
 - There are multiple of matrix which decide the best fit model like as : R-squared ,RMSE value, VIF, CDF & PDF Z-score , Roc & Auc and etc.
 - Database helped in making perfect model and will help in understanding Indonesian micro finance services (MFS) And use multiple metrics like F1_score, precision, recall and accuracy_score which will help to decide the best model.
 - Random forest Classifier as the best model with 95.44% accuracy_score..
 - Lastly predicted wheather the loan is paid back or not using saved model. It was good!! that was able to get the predictions near to actual values.

FINAL MODEL

```
9]: f_model=RandomForestClassifier(criterion='entropy',n_estimators=
f_model.fit(x_train,y_train)
pred=f_model.predict(x_test)
acc=accuracy_score(y_test,pred)
print('Accuracy Score: ',(accuracy_score(y_test,pred)*100))
print('Confusion Matrix: ',confusion_matrix(y_test,pred))
print(classification_report(y_test,pred))
```

```
Accuracy Score: 95.41336919288746
Confusion Matrix: [[52529 2201]
 [2847 52482]]
      precision    recall  f1-score   support

0         0.95      0.96      0.95      54730
1         0.96      0.95      0.95      55329

 accuracy          0.95      0.95      0.95      110059
 macro avg         0.95      0.95      0.95      110059
weighted avg         0.95      0.95      0.95      110059
```

SAVING MODEL IN PICKLE FORMAT

```
11]: # pickeling or serialization of a file
import pickle
filename = 'MicroCredit_Final_Mode.pkl'
pickle.dump(f_model, open(filename, 'wb'))
```

CONCLUSION



PREDICTION Conclusion

The Predicted values were almost similar to the original values.

t[102]:

	original	predicted
0	1	1
1	1	1
2	0	0
3	1	1
4	0	0
...
110054	1	1
110055	0	0
110056	0	0
110057	0	0
110058	0	0

110059 rows × 2 columns

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

This research evaluated individuals' credit risk performance in a micro-finance environment using machine learning and deep learning techniques. While traditional methods utilizing models such as linear regression are commonly adopted to estimate reasonable accuracy nowadays, these models have been succeeded by extensive employment of machine and deep learning models that have been broadly applied and produce prediction outcomes with greater precision. Using real data, we compared the various machine learning algorithms' accuracy by performing detailed experimental analysis while classifying individuals' requesting a loan into three classes, namely, good, average, and poor.

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.

Calculated the performance of each model using different performance metrics and compared them based on these metrics. Then we have also saved the best fit model and predicted the label. This is interesting that predicted and actual values were almost same.

Learning Outcomes of the Study in respect of Data Science

- Dataset is imbalanced. Label '1' has approximately 87.5% records, while, label '0' has approximately 12.5% records, and defaulter are higher.
 - This 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
-
- **Limitations of this work and Scope for Future Work**
 - The length of the dataset it is very huge and hard to handle.
 - 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 95.4% even after dealing all these drawbacks.
 - This study will not cover all Classification algorithms instead, it is focused on the chosen algorithm, starting from the basic assembling techniques to the advanced ones.

THANK YOU
-NIDHI SINGH