# CAPSTONE PROJECT : 3
## CREDIT CARD DEFAULT PREDECTION

## BY
## ABHAY DEEP SINGH

AI

## PROBLEM DESCRIPTION:

This project is aimed at predicting the case of customers default payments in Taiwan. From the perspective of risk management, the result of predictive accuracy of the estimated probability of default will be more valuable than the binary result of classification - credible or not credible clients. We can use the K-S chart to evaluate which customers will default on their credit card payments.

# DATA SUMMARY

[ ]

| ID | LIMIT_BAL | SEX | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 | PAY_4 | ... | BILL_AMT4 | BILL_AMT5 | BILL_A |
|----|-----------|-----|-----------|----------|-----|-------|-------|-------|-------|-----|-----------|-----------|--------|
| 1 | 20000 | 2 | 2 | 1 | 24 | 2 | 2 | -1 | -1 | ... | 0 | 0 | |
| 2 | 120000 | 2 | 2 | 2 | 26 | -1 | 2 | 0 | 0 | ... | 3272 | 3455 | 3 |
| 3 | 90000 | 2 | 2 | 2 | 34 | 0 | 0 | 0 | 0 | ... | 14331 | 14948 | 15 |
| 4 | 50000 | 2 | 2 | 1 | 37 | 0 | 0 | 0 | 0 | ... | 28314 | 28959 | 29 |
| 5 | 50000 | 1 | 2 | 1 | 57 | -1 | 0 | -1 | 0 | ... | 20940 | 19146 | 19 |

vs × 25 columns

# **DATA SUMMARY**

- This data set contains 30000 rows and 23 columns of six month.
- There are nine categorical features in our dataset.
- There are no missing value.
- There are no duplicate value.
- There are no null value.

# Feature Summary

X1: Amount of the given credit, includes both individual and family credit.

- X2: Gender(1=Male and 2=Female)
- X3: Education(1=graduate, 2= university, 3= high school and 4= others)
- X4: Marital status (1= Married, 2 = single, 3= others)
- X5: Age in year.
- X6-X11: History of past payment from April to September
- X12-17: Amount of bill statement fro April to September
- X18-X23: Amount of previous payment from April to September
- Y: Default payment

# Approach to analyze the dataset

**AI**

### DATA CLEANING

- Find information on document columns value.
- Clear data for analysis.

### DATA EXPLORATION

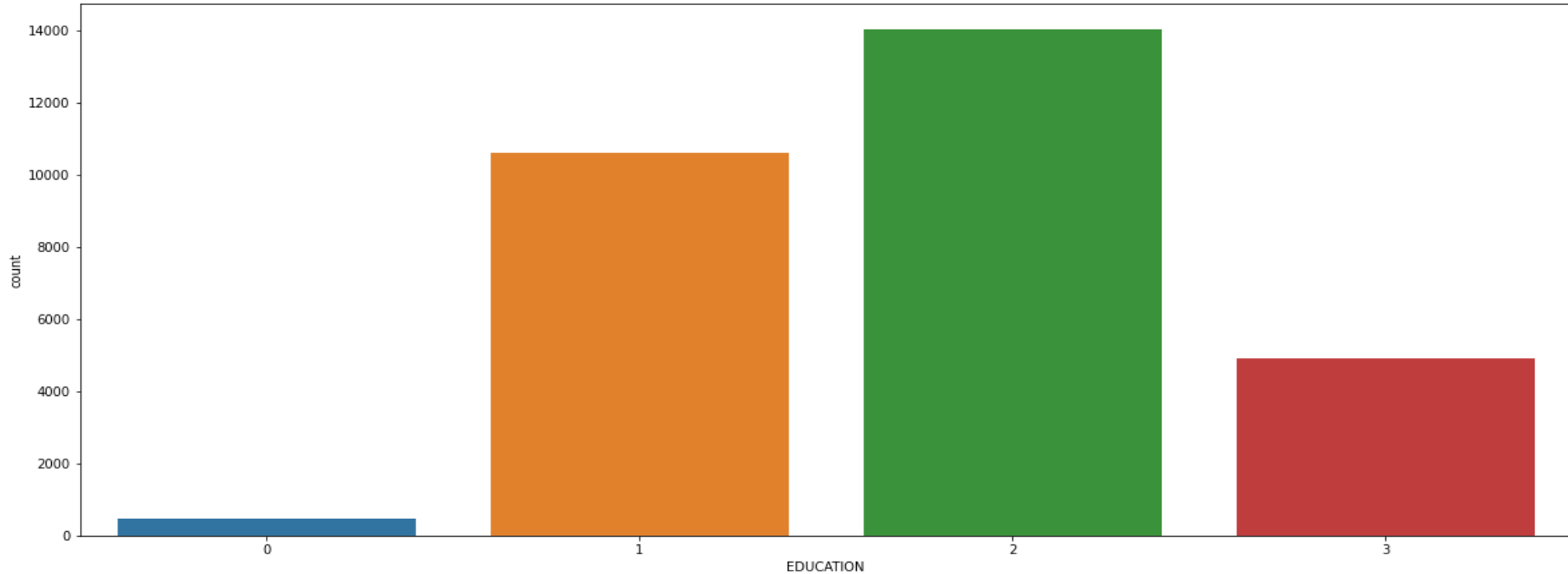- Analyze the data with exploratory data analysis.

### MODELING

- Logistic regression
- XG Boost
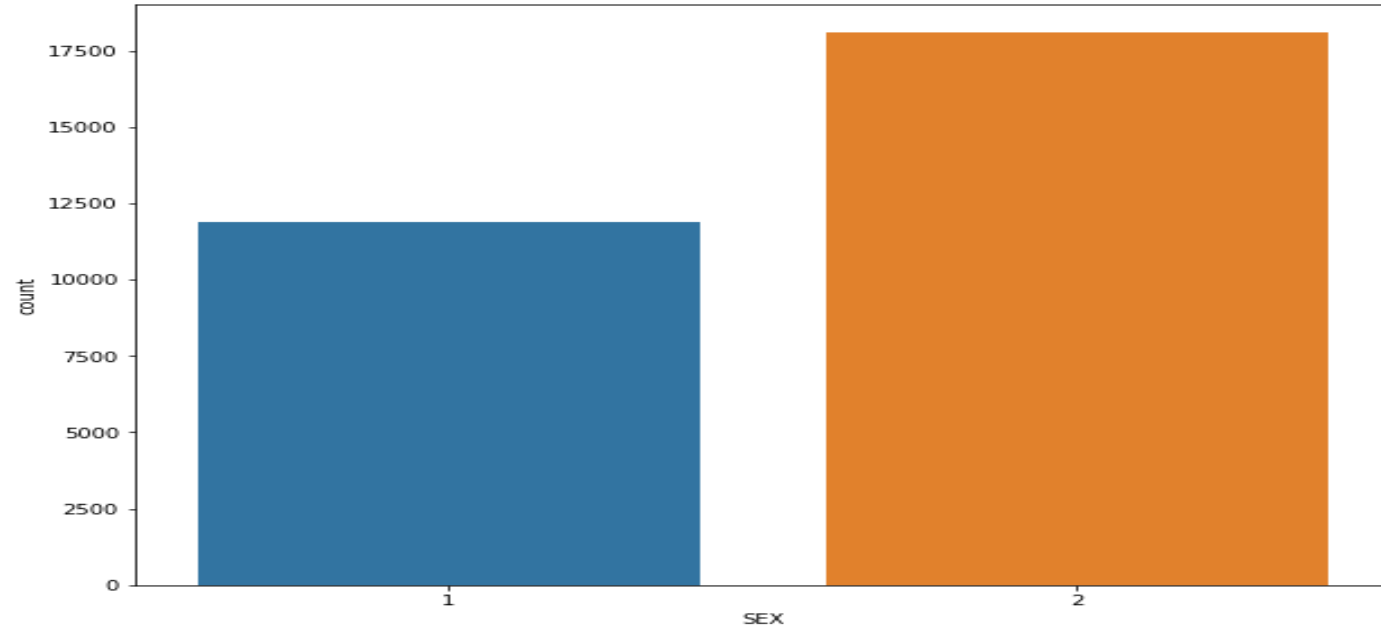- Random forest
- SVC

# Count of credit card on basis of age



**People from age 24 to 36 uses more credit and as the age increase the count decrease.**

# Count of credit card basis of education



- **Most number of credit card are used by university students and university student shown by 2.**
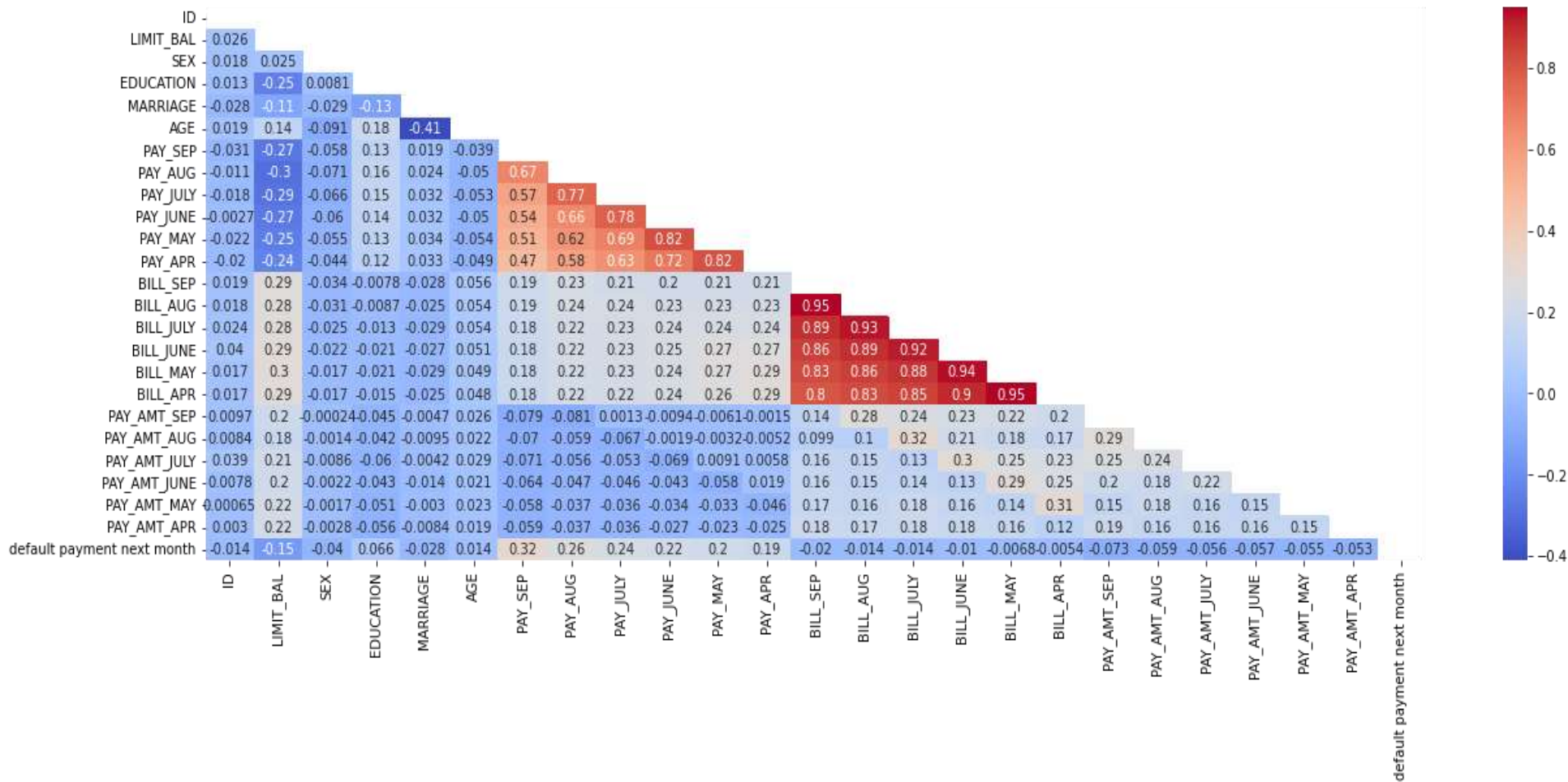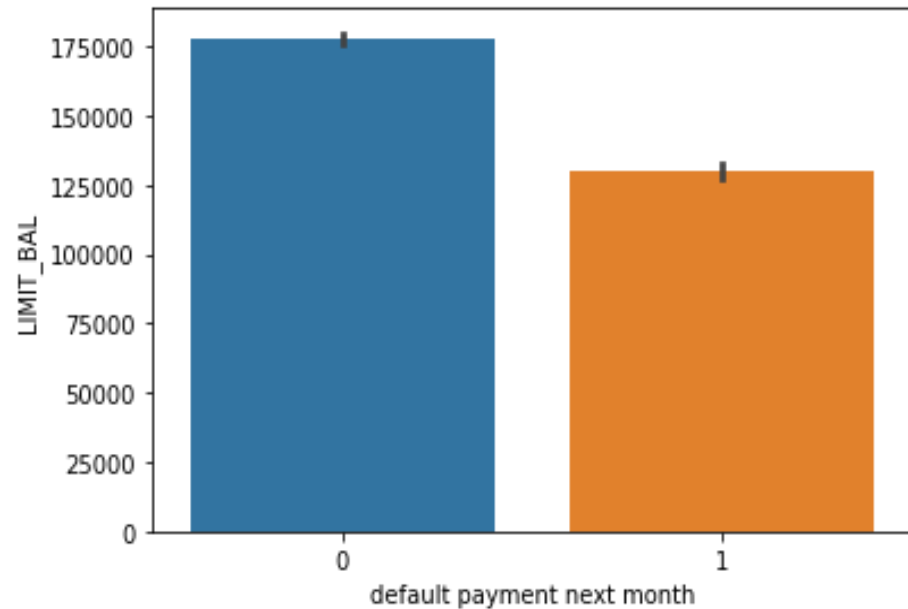
# Count of credit card on basis of Gender



From above plot we can see that number of female credit card holder are more than male.

- **1: Male**
- **2: Female**

# Correlation Matrix

# SMOTE(Synthetic Minority Oversampling Technique)



- Defaulters are less than non defaulters.
- We have solve the imbalance by SMOTE.
- We can see we have imbalanced dataset

# Confusion Matrix

- **A confusion matrix is a table that is used to define the performance of a classification algorithm**.

- **It is a table with 4 different combinations of predicted and actual values.**

# Modeling Overview

## Models Used

- Logistic regression

- XG Boost

- Random forest

- SVC

# Modeling Approach

**DATA PREPROCESSING**
- Feature selection
- Feature engineering
- Train test split
- SMOTE over sampeling

**DATA FITTING AND TUNING**
- Start with default parameter
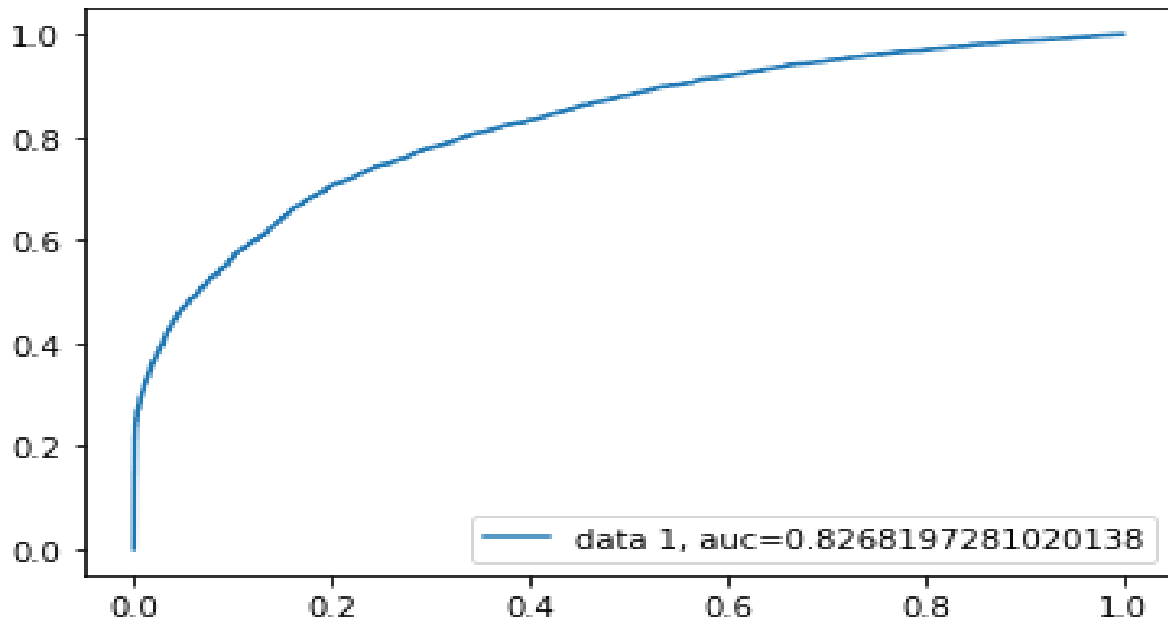- Hyperparameter tuning
- Measure RUC and AUC on training data

**MODEL EVALUATION**
- Model testing
- Precision score
- Recall score
- Model evaluation

AI

# Logistic Regression

**AI**
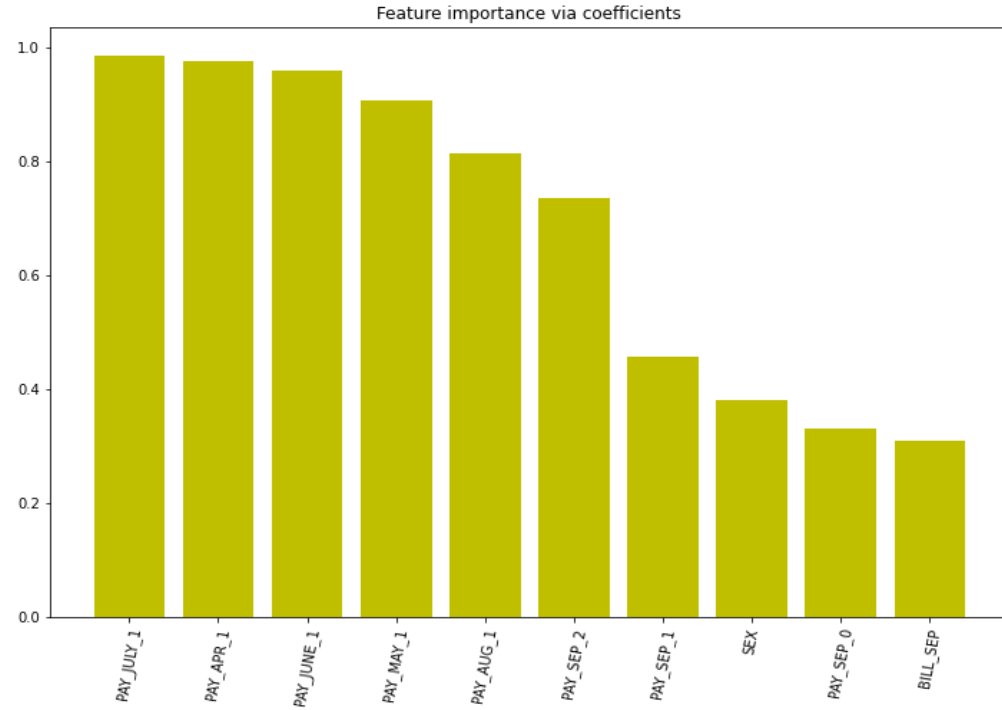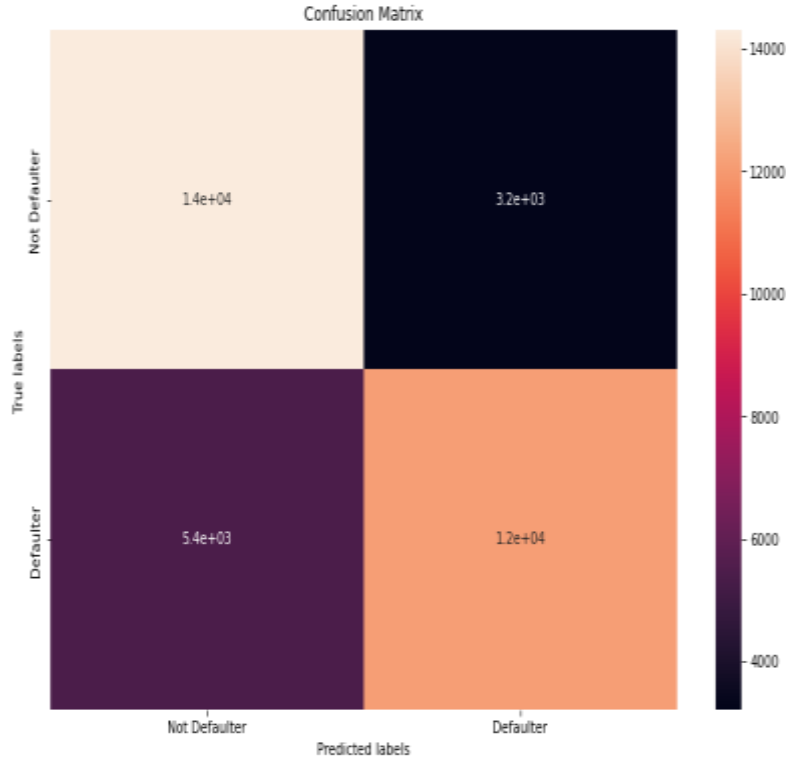
Parameters

C = 0.01

 Penalty = L2

We have implemented logistic regression
and we getting f1_sore approx 73.4%.
As we have imbalanced dataset,
F1 score is better parameter.



data 1, auc=0.8268197281020138

- The accuracy on test data is 0.751.
- The precision on test data is 0.687.
- The recall on test data is 0.788.
- The f1 score on test data is 0.734.
- The roc score on test data is 0.755.

# Logistic Regression(Cont.)



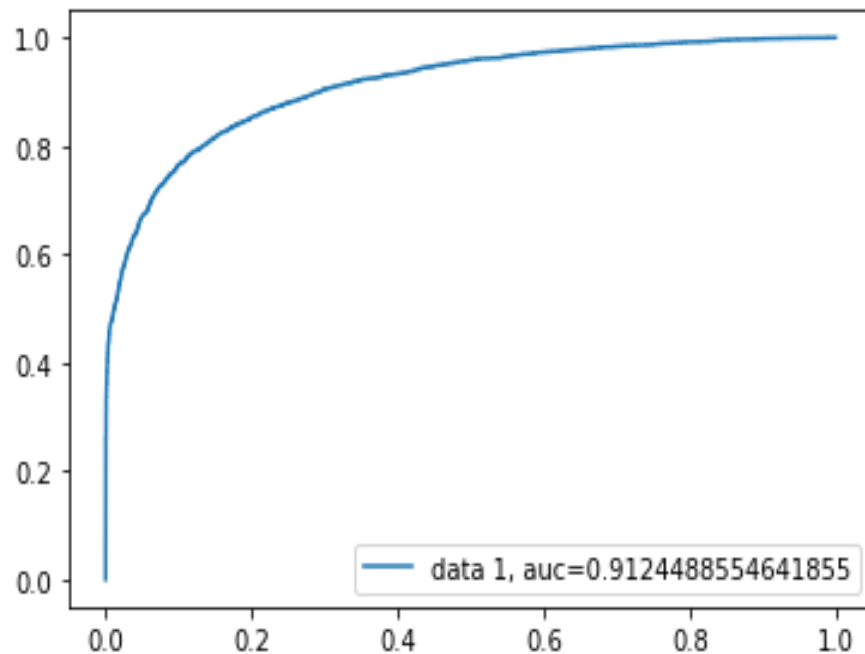- [[14305 3218]
- [ 5415 12108]]
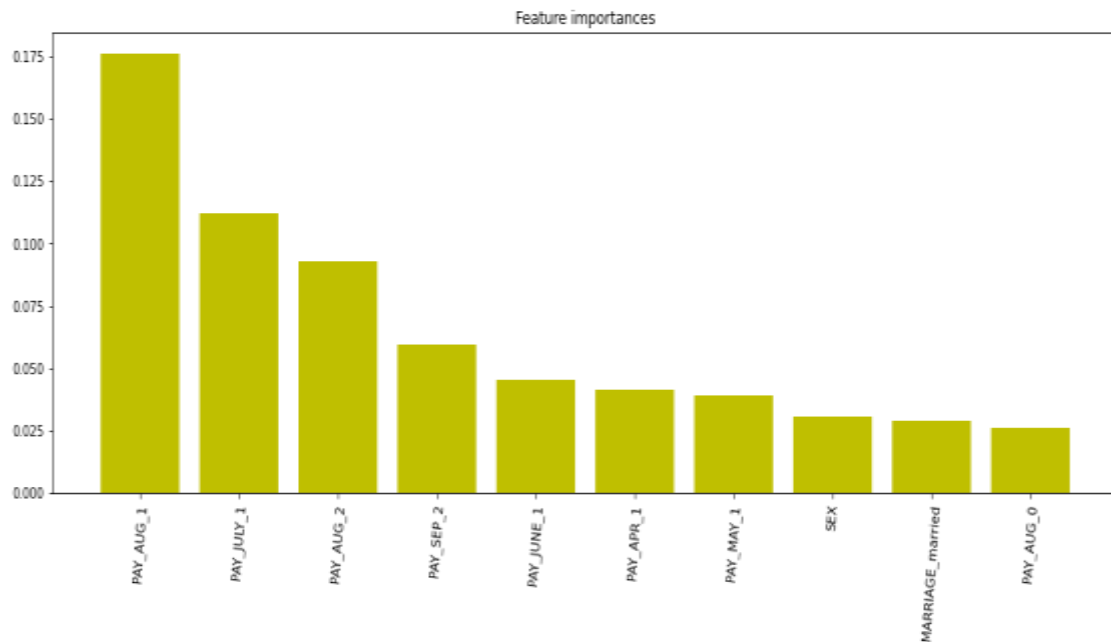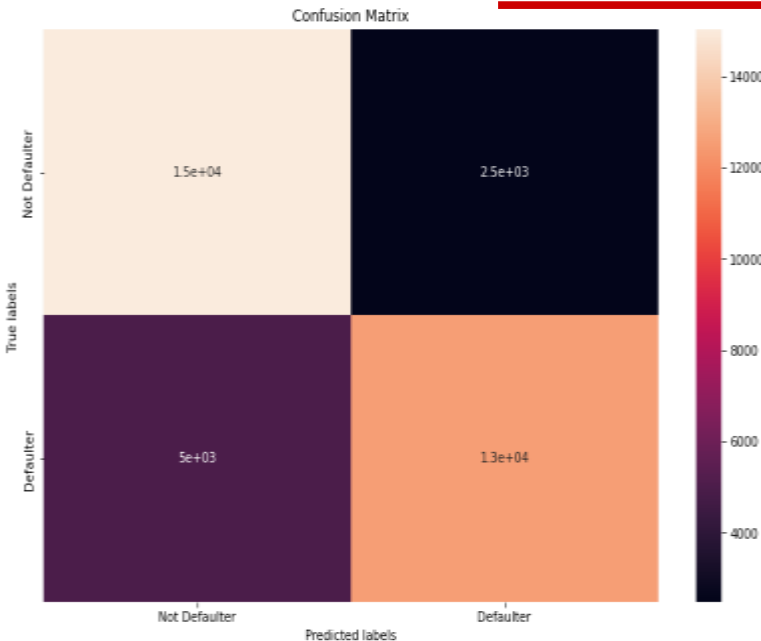
# XG Boost Classifier

- **Parameters**

max_depth : 7

XG Boost is an implementation of gradient boosted decision trees
 designed for speed and performance that is dominative competitive machine learning.

Accuracy on test data after hyperparameter tuning is 0.83.
Precision on test data after hyperparameter tuning is 0.795.
recall on test data after hyperparameter tuning is 0.86.
f1 score on test data after hyperparameter tuning is 0.82.
roc score on test data after hyperparameter tuning is 0.83.
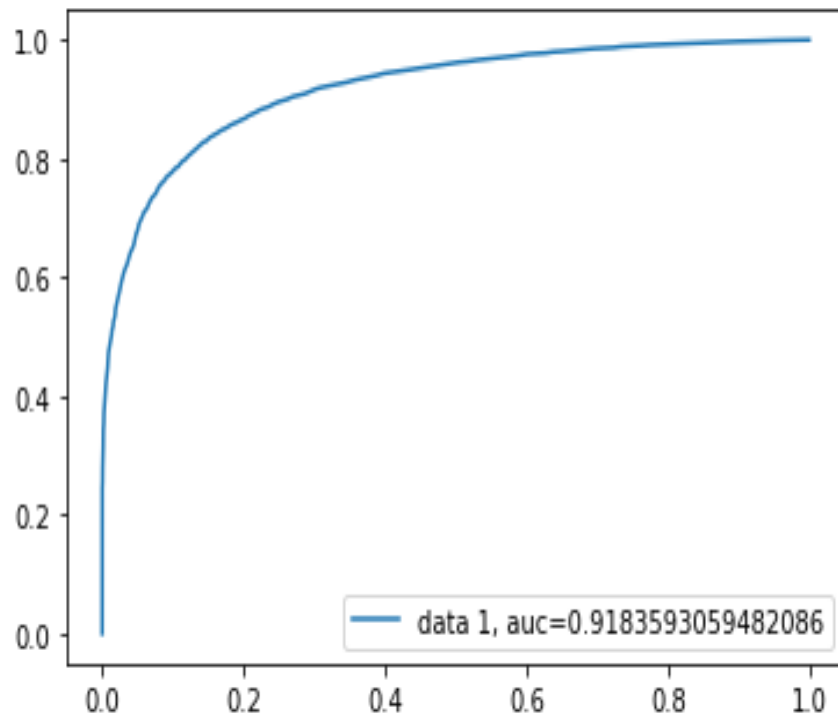
# XGBoost Classifier(Cont.)



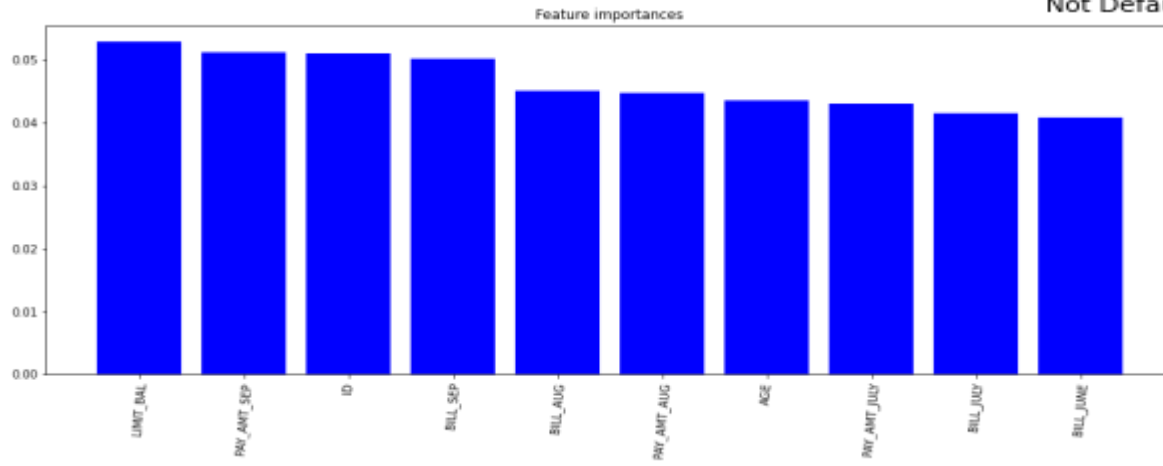[[15037 2486]
[ 4963 12560]]

# Random Forest

**Parameters**

Max_depth = 20

N_estimators = 200
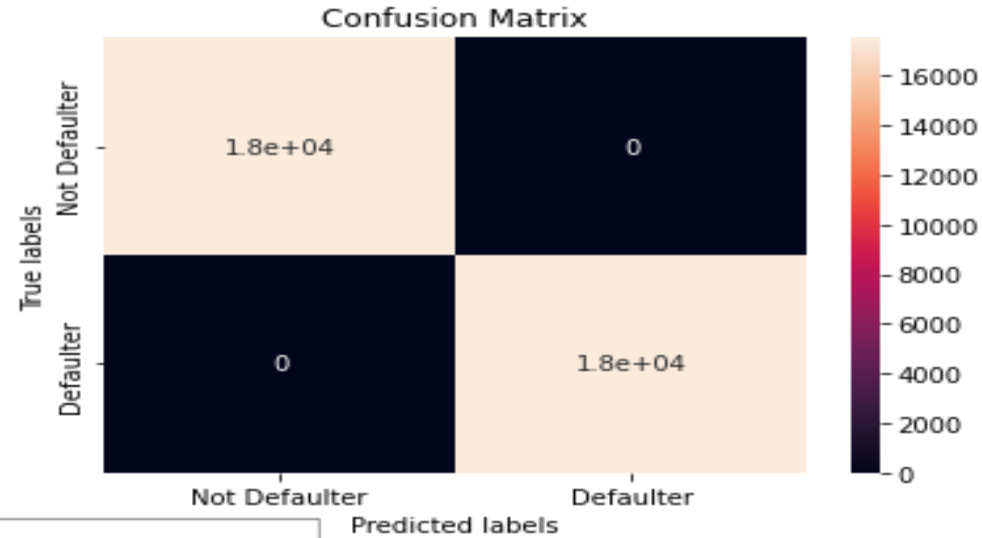
- The accuracy after hyperparameter tuning is 0.84.

- The precision after hyperparameter tuning is 0.81.

- The recall after hyperparameter tuning is 0.86.

- The f1 score after hyperparameter tuning is 0.83.

- The roc score after hyperparameter tuning is 0.84.



data 1, auc=0.9183593059482086

# Random Forest(Cont.)

[[17523 0]
[ 0 17523]]



Confusion Matrix



Feature importances

# Support Vector Classifier

- **Parameters**
- C = 5
- Kernel = rbf

- A support vector machine (SVM) is a type of deep learning algorithm that performs supervised learning for classification or regression of data groups. In AI and machine learning, supervised learning systems provide both input and desired output data, which are labeled for classification.
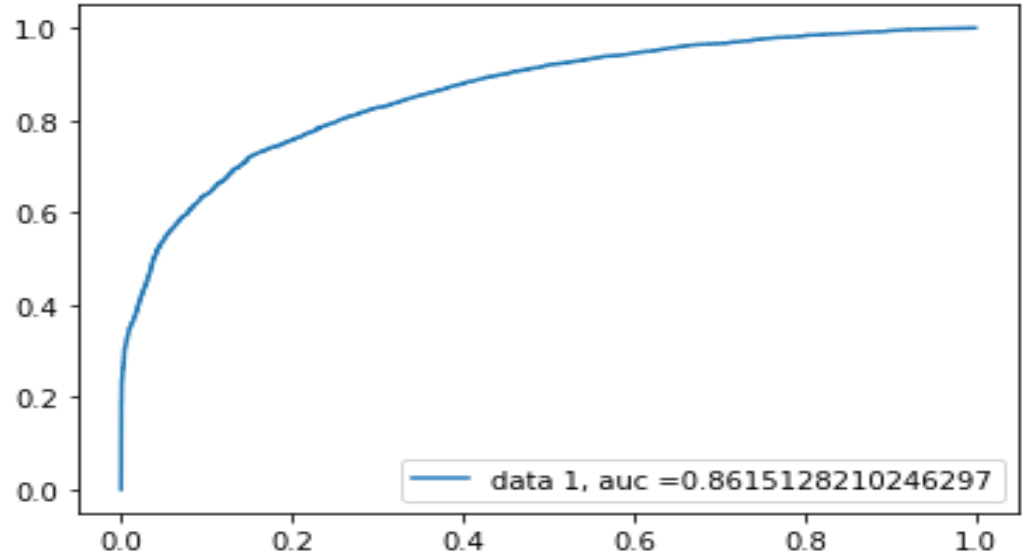
The accuracy score on test data is 0.78.

The precision score on test data is 0.71.
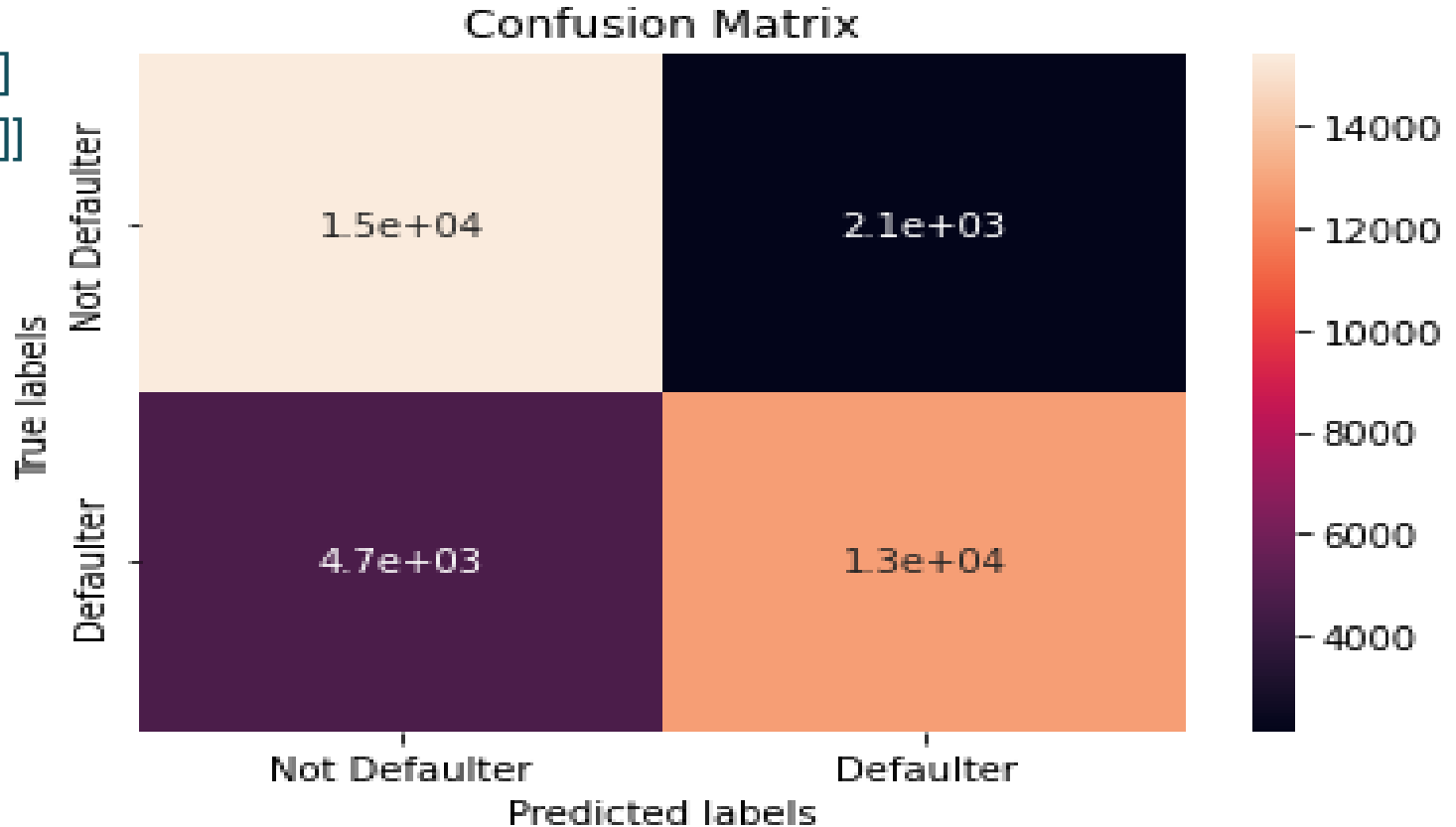
The recall on test data is 0.82.

The f1 score on test data is 0.76.
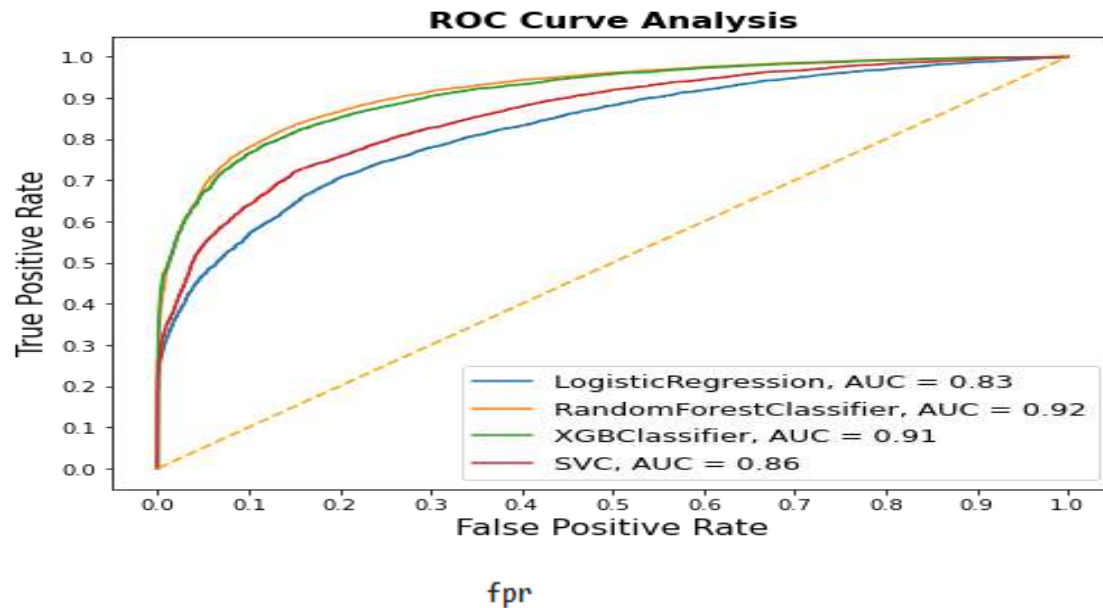
The roc score on tes data is 0.78.



data 1, auc =0.8615128210246297

# Support Vector Classifier(Cont.)

[[15386 2137]
[ 4731 12792]]



Confusion Matrix

# Plotting ROC AUC for all the models

**AI**



ROC Curve Analysis

| Classifiers | fpr | tpr | auc |
|---|---|---|---|
| LogisticRegression | [0.0, 0.0, 0.0, 0.0001712035610340695, 0.00017... | [0.0, 0.0001712035610340695, 0.077726416709467... | 0.826820 |
| RandomForestClassifier | [0.0, 0.0, 0.0, 0.0, 0.0001712035610340695, 0.... | [0.0, 0.033042287279575415, 0.0544427324088341... | 0.918359 |
| XGBClassifier | [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, ... | [0.0, 0.0001712035610340695, 0.002568053415511... | 0.912449 |
| SVC | [0.0, 0.0, 0.0, 0.0001712035610340695, 0.00017... | [0.0, 0.0001712035610340695, 0.184557438794726... | 0.861513 |

# Feature Importance for recommended model

- "LIMT_BAL","BILL_SEP" AND "PAY_AMT_SEP" are the most recent 2 months' payment status and they are the strongest predictors of future payment default risk.

# Challenges

- Data Cleaning

- Data mining

- Feature Engineering

- Feature Selection

- Model optimization

- Hyperparameter Tuning

- Deciding the flow of presentation

# Overall Conclusion

Random Forest model and XG Boost model both has same recall, so if the business cares recall the most than both of this model are best candidate. If the balance of recall and precision is most important metric than Random Forest is the ideal model. Random Forest has recall and precision both higher than the other model applied. Hence, I would recommend Random Forest for this dataset.

|   | Classifier | Train Accuracy | Test Accuracy | Precision Score | Recall Score | F1 Score |
|---|---|---|---|---|---|---|
| 0 | Logistic Regression | 0.753667 | 0.751498 | 0.687211 | 0.788605 | 0.734425 |
| 1 | Xgboost CLF | 0.916139 | 0.833077 | 0.795069 | 0.860478 | 0.826482 |
| 2 | Random Forest CLF | 1.000000 | 0.841722 | 0.812703 | 0.862777 | 0.836992 |
| 3 | Support Vector CLF | 0.804029 | 0.782657 | 0.711693 | 0.829409 | 0.766055 |

# Overall Conclusion

There were not huge gap but female clients tended to default the most.
- Labels of the data were imbalanced and had a significant difference.
- Gradient boost gave the highest accuracy of 82% on test dataset.
- Repayment in the month of September tended to be most important feature for our machine learning model.
- The best accuracy is obtained for the Random forest and XGBoost classifier.
- Data categorical variables had minority classes which were added to their closest majority class.
- From above table we can see that XGBoost Classifier having Recall = 86%, F1-score = 82%, and ROC Score=
 83% and Random forest Classifier having Recall =86%, F1-score = 82% and ROC Score = 84%.
- XGBoost Classifier and Random Forest Classifier are giving us the best Recall, F1-score, and ROC Score among other algorithms. We can conclude that these two algorithms are the best to predict whether the credit card is default or not default according to our analysis on this dataset.

THANK YOU