# Capstone Project : 3
## Credit Card Default Prediction

### By
### Mayank Mishra

# Points for Discussion

- **Understanding Problem Statement**
- **Data Set Information**
- **Feature Summary**
- **Approach Overview**
- **Data Preprocessing**
- **Exploratory Data Analysis**
- **Implementing Algorithm**
- **Challenges**
- **Conclusion**
- **Q&A**

# Understanding Problem Statement

- Topic – "Credit Card Default Prediction"
- Problem Statement -- Explore and Analyze the
  the prediction whether a customer will default on his/her
  Credit Card.

- Financial threats are displaying a trend about the risk of
  commercial banks as the improvement in the financial industry has arisen. In this way
  one of the biggest threats faced by them is the risk prediction of credit card clients.

- The goal of this project is to predict the credit card defaulter comprising variables
  like AGE, SEX, EDUCATION, LIMIT BAL, MARRIAGE and others.

# Data Set Information

- This dataset contains 30000 observations and 23 features of six months.
- There are nine categorical features in our dataset.
- This dataset is from the city of Taiwan and doesn't have any null or duplicate values.

```
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 25 columns):
 #   Column                      Non-Null Count   Dtype
---  ------                      --------------   -----
 0   ID                          30000 non-null   int64
 1   LIMIT_BAL                   30000 non-null   int64
 2   SEX                         30000 non-null   int64
 3   EDUCATION                   30000 non-null   int64
 4   MARRIAGE                    30000 non-null   int64
 5   AGE                         30000 non-null   int64
 6   PAY_0                       30000 non-null   int64
 7   PAY_2                       30000 non-null   int64
 8   PAY_3                       30000 non-null   int64
 9   PAY_4                       30000 non-null   int64
 10  PAY_5                       30000 non-null   int64
 11  PAY_6                       30000 non-null   int64
 12  BILL_AMT1                   30000 non-null   int64
 13  BILL_AMT2                   30000 non-null   int64
 14  BILL_AMT3                   30000 non-null   int64
 15  BILL_AMT4                   30000 non-null   int64
 16  BILL_AMT5                   30000 non-null   int64
 17  BILL_AMT6                   30000 non-null   int64
 18  PAY_AMT1                    30000 non-null   int64
 19  PAY_AMT2                    30000 non-null   int64
 20  PAY_AMT3                    30000 non-null   int64
 21  PAY_AMT4                    30000 non-null   int64
 22  PAY_AMT5                    30000 non-null   int64
 23  PAY_AMT6                    30000 non-null   int64
 24  default payment next month  30000 non-null   int64
```

# 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 Se
- Y: Default payment

# Approach Overview

**AI**

Data Cleaning

- Data Exploration

Modelling

- Find information on documented columns values.

- Clean data for further analysis.

- Analyze the data with EDA

- Logistic Regression
- XGBoost
- 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 decreases.

# Count of credit card basis of eduation
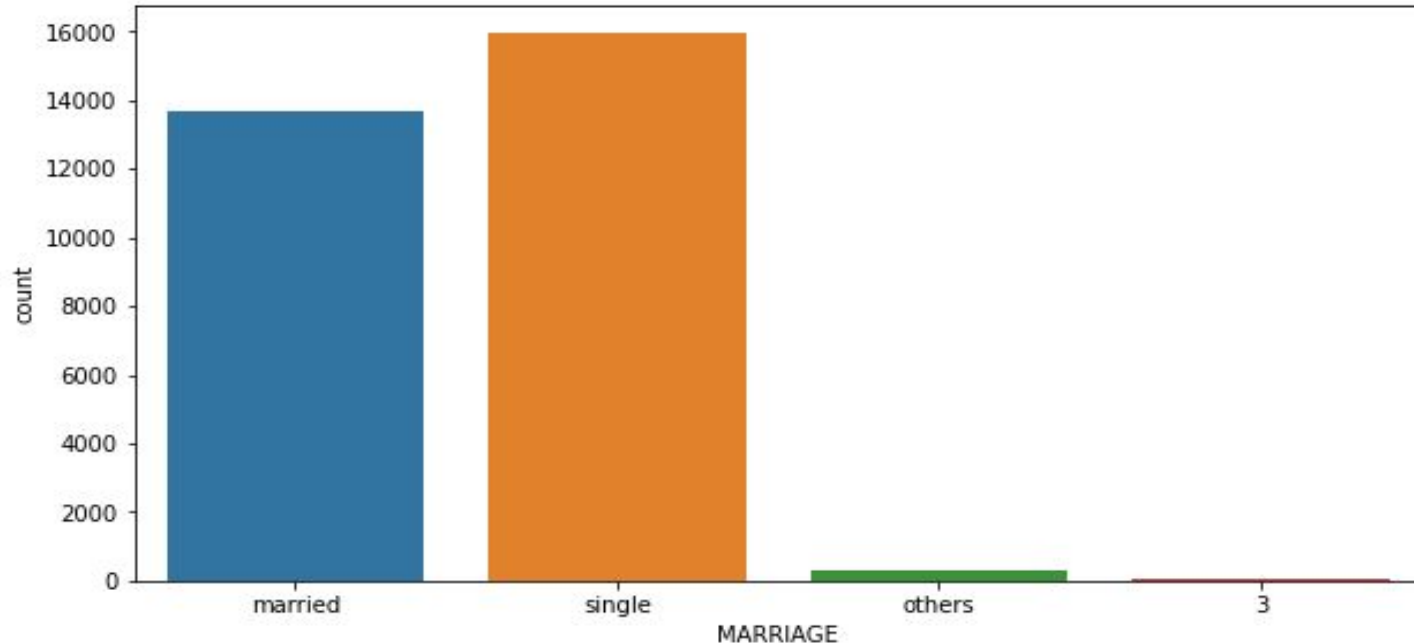
- Most number of credit card are used by university students.

# Credit card defaulter on basis of Education

- High school and university students has high default risk.

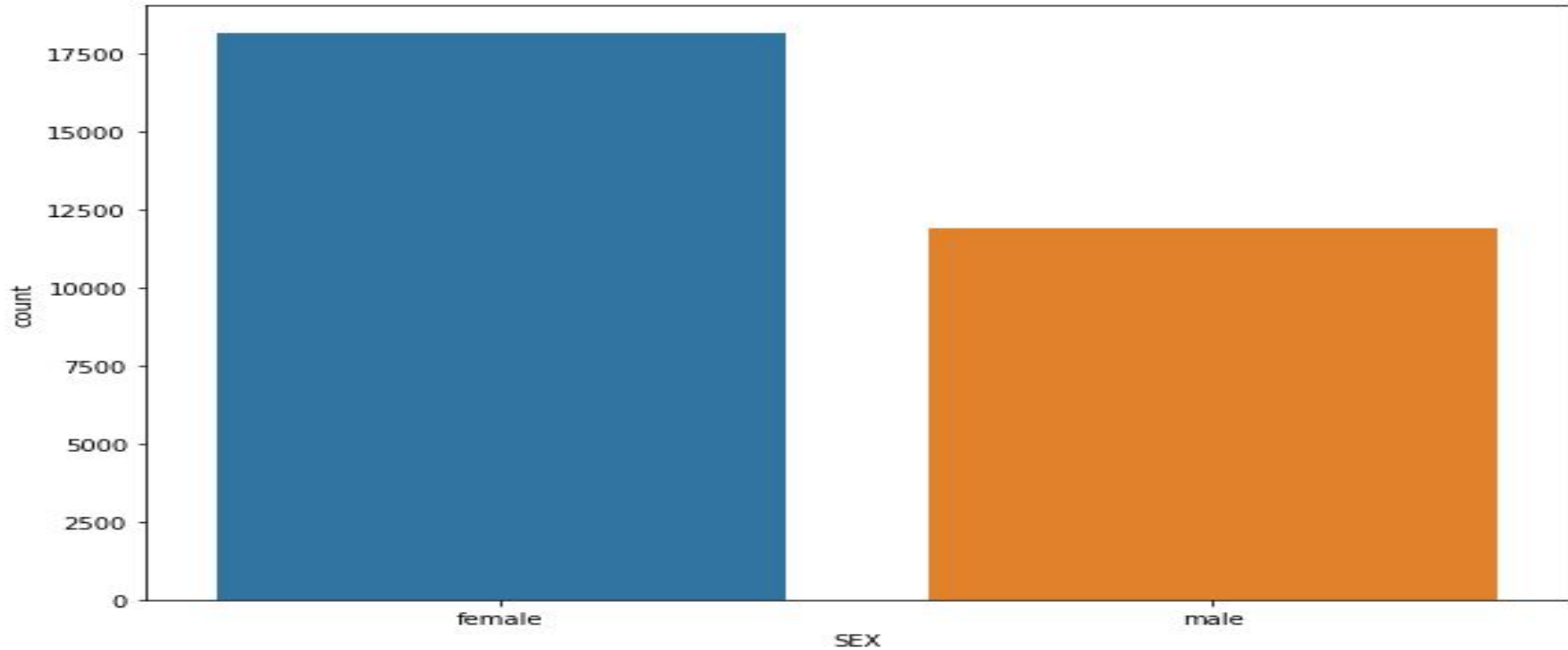# Count of credit card on basis of marital status

- Most number credit card are used by the people who are unmarried.

# Credit card defaulter on basis of marital status

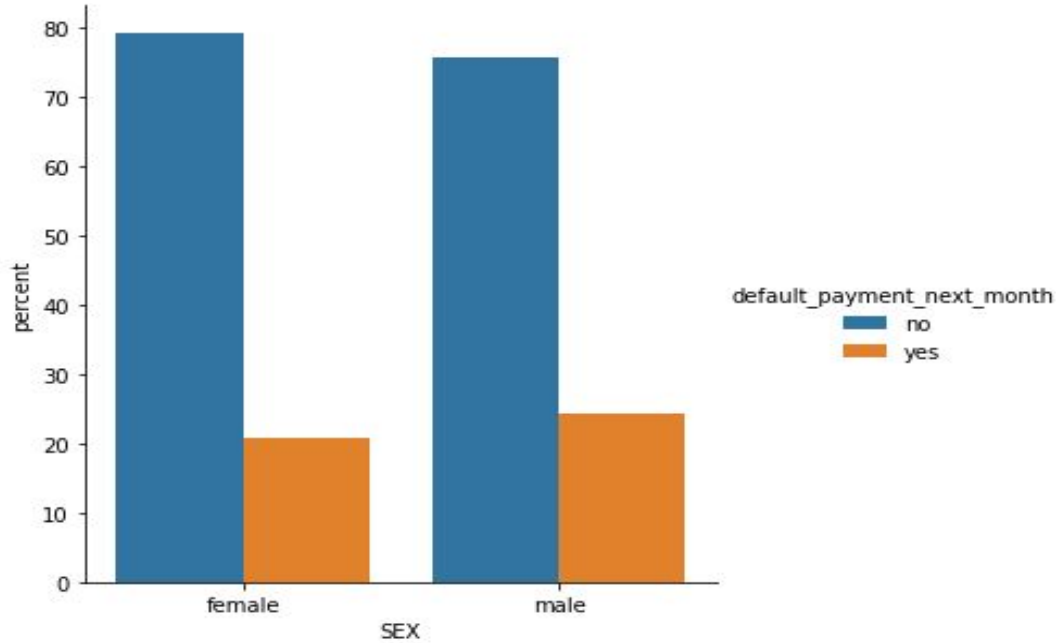- No significant correlation of default risk and marital status.

# Count of credit card in basis of Gender



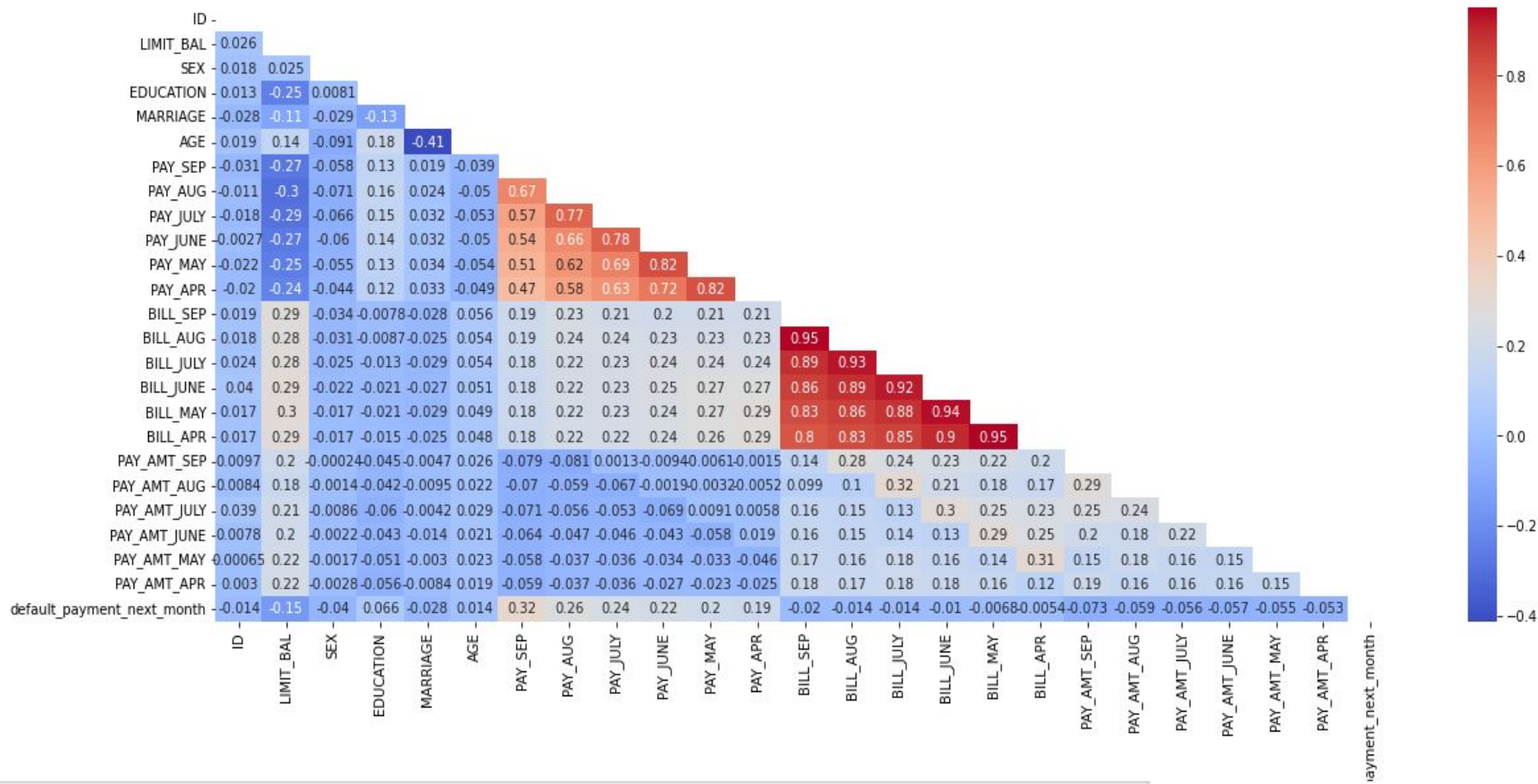- No of female credit card holder are more as compared to male.

# Credit card defaulter on basis of Gender



- Males credit card holder has more default as compared to females
- Males defaulter are around 25% and females defaulter are around 20%
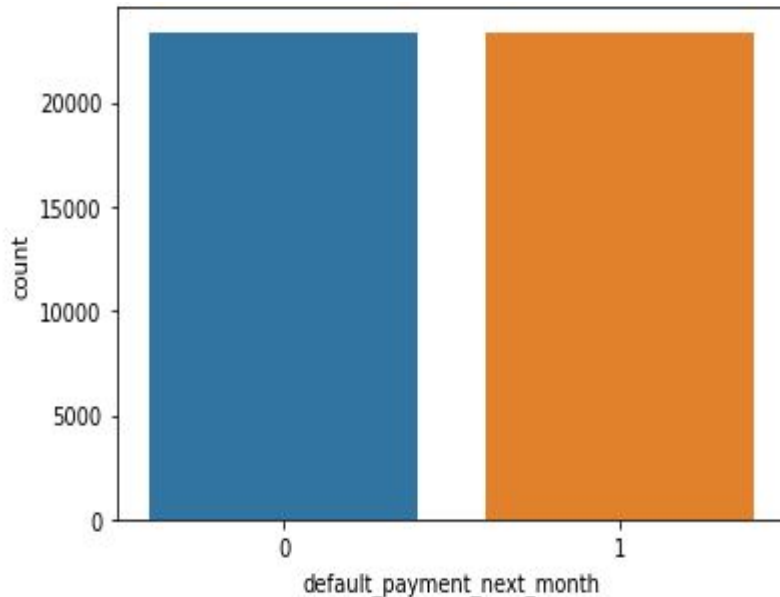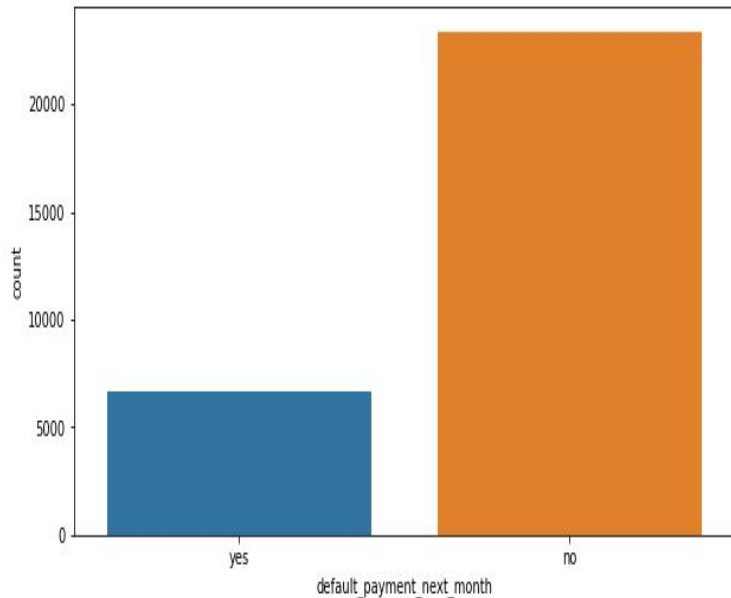
# Correlation Matrix

# SMOTE(Synthetic Minority Oversampling Technique)

- We can see we have imbalanced dataset.

- Defaulters are less than non defaulters.

- We have solve the imbalance by SMOTE.

# Confusion Matrix

- Confusion matrix is a performance measurement for machine learning classification problem where output can be two or more classes.
- It is a table with 4 different combinations of predicted and actual values.

Actual Values

|                  | Positive (1) | Negative (0) |
|------------------|--------------|--------------|
| Positive (1)     | TP           | FP           |
| Negative (0)     | FN           | TN           |

Predicted Values

# **<u>Modeling Overview</u>**

❑ Supervised learning / Binary classification

❑ Imbalance data with 78% non defaulters and 22% defaulters.

❑ <u>Models Used</u>

❑ Logistic Regression

❑ XGBoost CLF

❑ Random Forest CLF

❑ Support Vector Classifier (SVC)

# Modelling Approach

**AI**

Data Preprocessing

Data Fitting & Tuning

Model Evaluation

- Feature Selection
- Feature engineering
- Train Test split
- SMOTE Oversampling

- Start with default parameter
- Hyperparameter Tuning
- Measure RUC AUC on training data

- Model testing
- Precision score
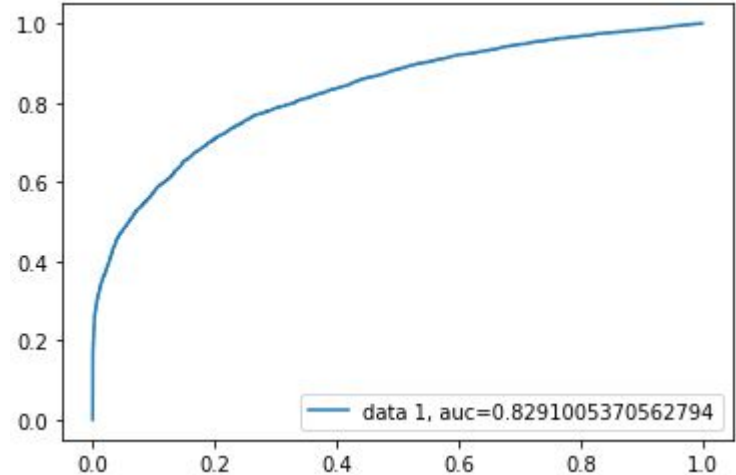- Recall score
- Model Evaluation
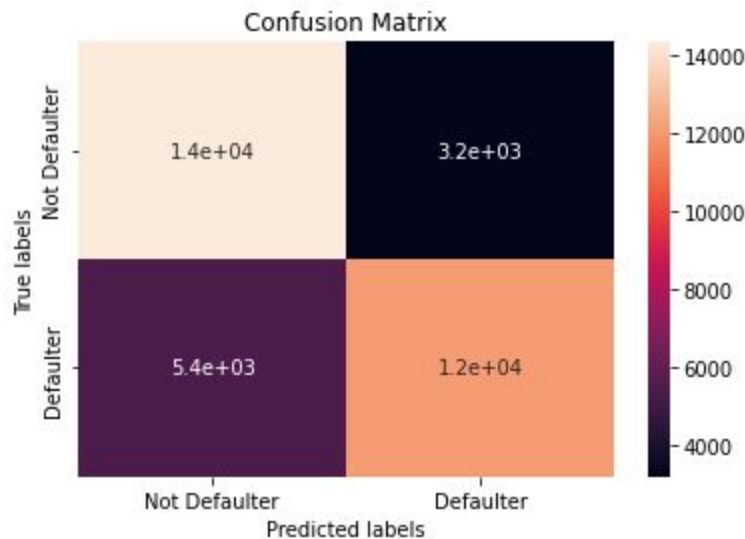
# Logistic Regression

- **Parameters**

☐ C = 0.01
☐ Penalty = L2

- We have implemented logistic regression.
- we getting f1_sore approx. 73%. As we have imbalanced dataset, F1- score is better parameter
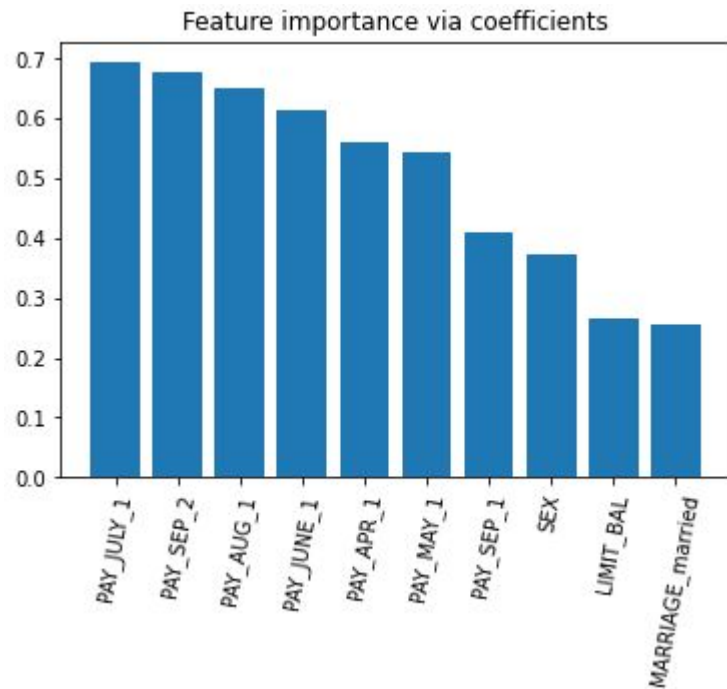


```
The accuracy on test data is 0.7536380756719739
The precision on test data is 0.6935456257490156
The recall on test data is 0.7882856586884608
The f1 score on test data is 0.7378870673952641
The roc score on test data is 0.7573554229557236
```

# Logistic Regression(Cont.)

Confusion Matrix



Feature importance via coefficients
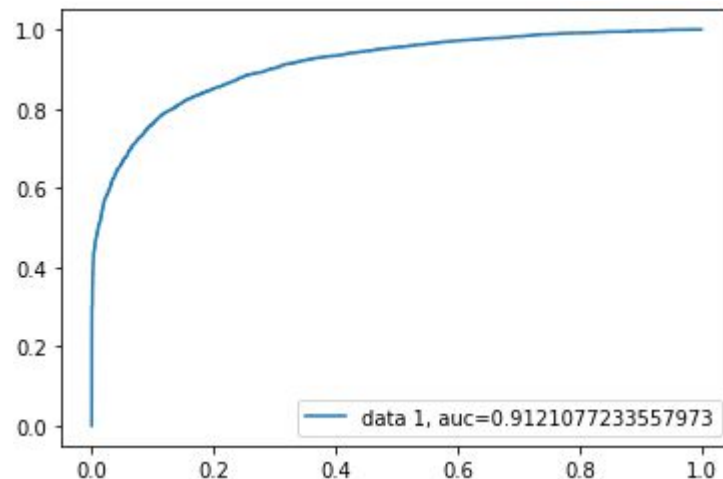
```
[[14329   3194]
 [ 5450 12073]]
```

# XGBoost Classifier

- **Parameters**

  ❑ max_depth = 10
  ❑ min_child_weight =6

- The XGBoost model for classification is called XGBClassifier. We can create and fit it to our training dataset. Models are fit using the scikit-learn API and the model.fit() function.



```
The accuracy on test data after hyperparameter tuning is  0.8340181475774696
The Precision on test data after hyperparameter tuning is  0.7952405410032529
The recall on test data after hyperparameter tuning is  0.8621009651076467
The f1 score on test data after hyperparameter tuning is  0.8273221123875679
The roc score on test data after hyperparameter tuning is  0.8360393608506138
```
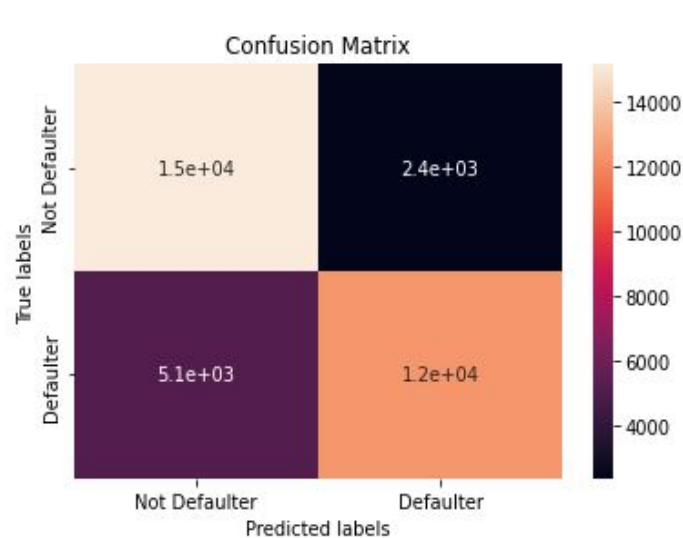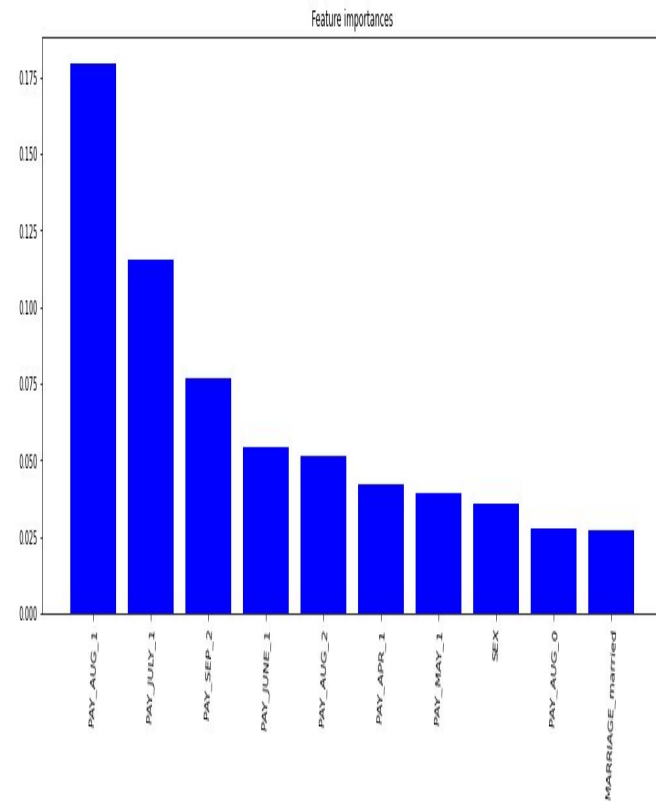
# XGBoost Classifier(Cont.)

**AI**

## Confusion Matrix



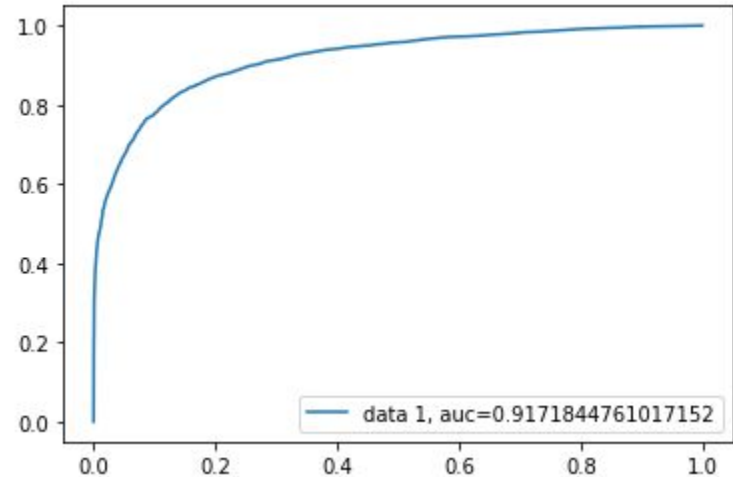|  | importance_xgb |
|---|---|
| PAY_AUG_1 | 0.179251 |
| PAY_JULY_1 | 0.115708 |
| PAY_SEP_2 | 0.076841 |
| PAY_JUNE_1 | 0.054122 |
| PAY_AUG_2 | 0.051261 |
| PAY_APR_1 | 0.042170 |
| PAY_MAY_1 | 0.039029 |
| SEX | 0.035905 |
| PAY_AUG_0 | 0.027767 |
| MARRIAGE_married | 0.027361 |

Feature importances

```
[[15145  2378]
 [ 5124 12399]]
```
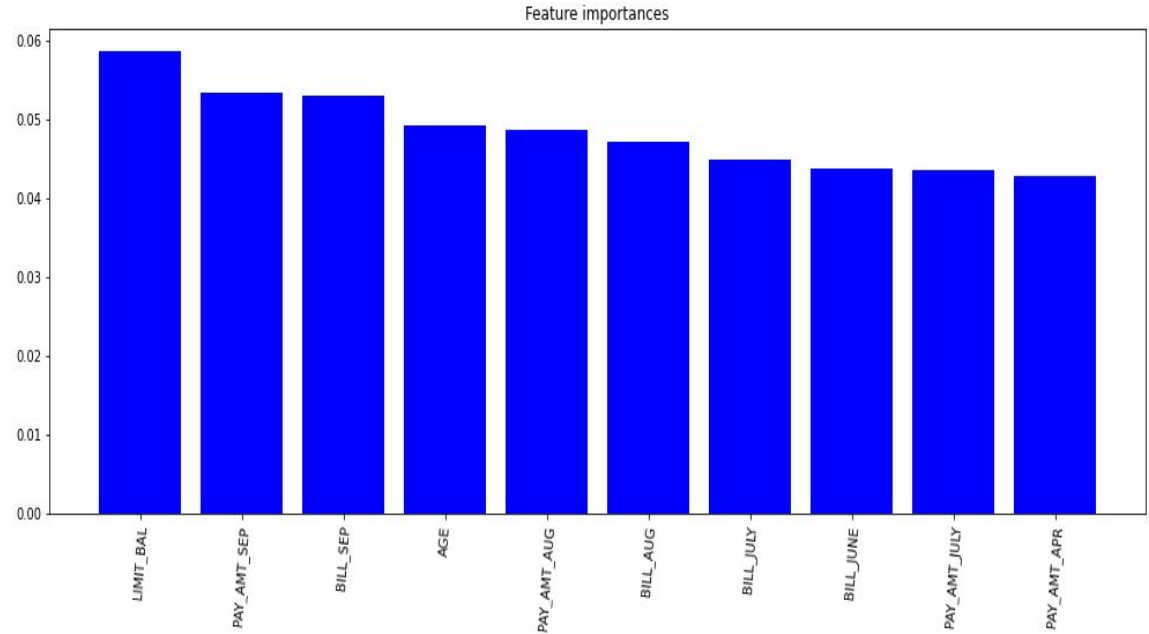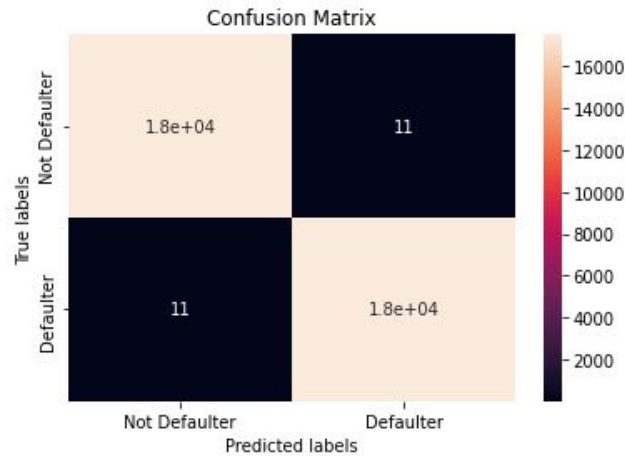
# Random Forest

- **Parameters**

  ⬜ Max_depth = 40
  ⬜ N_estimators = 250



```
The accuracy after hyperparameter tuning is 0.8436055469953775
The precision after hyperparameter tuning is  0.8164697825714775
The recall after hyperparameter tuning is  0.8633236784938451
The f1 score after hyperparameter tuning is  0.8392432908051034
The roc score after hyperparameter tuning is  0.8446205920887542
```

**AI**

# Support Vector Classifier

- **Parameters**

  - **C = 5**
  - **Kernel = 'rbf'**



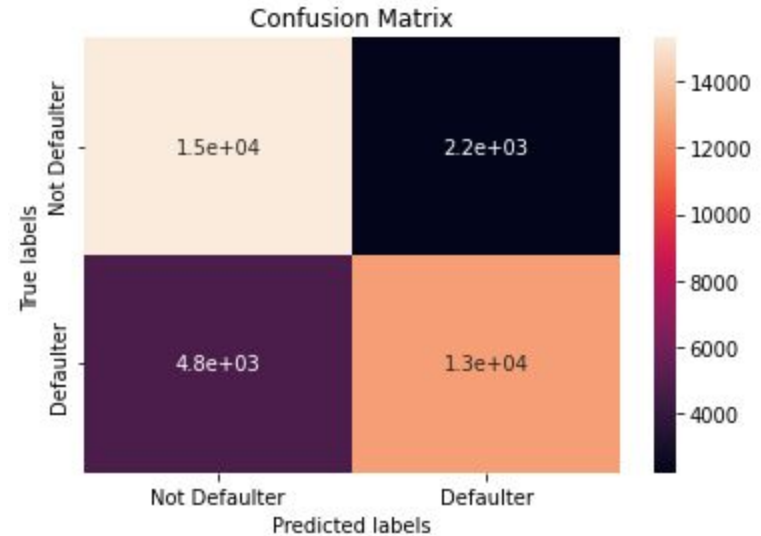data 1, auc =0.8616977126635619

The accuracy score on test data is  0.7817154596815614
The precision score on test data is  0.7151172744393084
The recall on test data is 0.8250049377839226
The f1 score on test data is  0.7661408657373442
The roc score on tes data is  0.7868037228578172

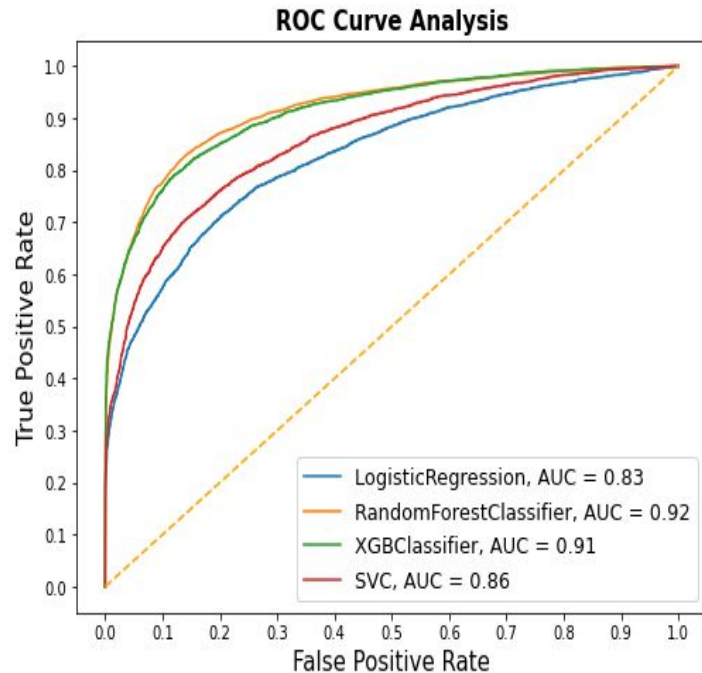# Support Vector Classifier(Cont.)

- **Confusion Matrix**

$$[[15313 \quad 2210]$$
$$[\ 4800 \quad 12723]]$$

# Plotting ROC AUC for all the models

- Curve Analysis of all model.

| Classifiers | fpr | tpr | auc |
|---|---|---|---|
| LogisticRegression | [0.0, 0.0, 0.0, 0.0001712035610340695, 0.00017... | [0.0, 0.0001712035610340695, 0.091936312275295... | 0.829101 |
| RandomForestClassifier | [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, ... | [0.0, 0.036466358500256806, 0.0369799691833590... | 0.917184 |
| XGBClassifier | [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, ... | [0.0, 0.0001712035610340695, 0.003081664098613... | 0.912108 |
| SVC | [0.0, 0.0, 0.0, 0.0001712035610340695, 0.00017... | [0.0, 0.0001712035610340695, 0.189008731381612... | 0.861698 |



**ROC Curve Analysis**

- LogisticRegression, AUC = 0.83
- RandomForestClassifier, AUC = 0.92
- XGBClassifier, AUC = 0.91
- SVC, AUC = 0.86

# Model Recommendation

- We recommend recall = 0.8, however, the threshold can be adjusted to reach higher recall.

# Feature Importance for recommended model

**AI**

- "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.



Feature Importance

# *Challenges*

- Data Cleaning
- Data mining
- Feature Engineering
- Feature Selection
- Model optimization
- Hyperparameter Tuning
- Deciding the flow of presentation

# Overall Conclusion

- Random Forest model and XGBoost 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 | 0.753353 | 0.753638 | 0.693546 | 0.788286 | 0.737887 |
| 1 | Xgboost CLF | 0.905781 | 0.834018 | 0.795241 | 0.862101 | 0.827322 |
| 2 | Random Forest CLF | 0.999372 | 0.843606 | 0.816470 | 0.863324 | 0.839243 |
| 3 | Support Vector CLF | 0.799977 | 0.781715 | 0.715117 | 0.825005 | 0.766141 |

# 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 the 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.

- In general, all models have comparable accuracy. Nevertheless, because the classes are imbalanced (the proportion of non-default credit cards is higher than default) this metric is misleading. Also, accuracy does not consider the rate of false positives (non-default credits cards that were predicted as default) and false negatives (default credit cards that were incorrectly predicted as non-default). Both cases have negative impact on the bank, since false positives leads to unsatisfied customers and false negatives leads to financial loss.

- 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 = 83% 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