



# **A Machine Learning Approach to Credit Card Default Prediction**

**BY ALISHBA TAHIR**

# **1. Problem Statement**

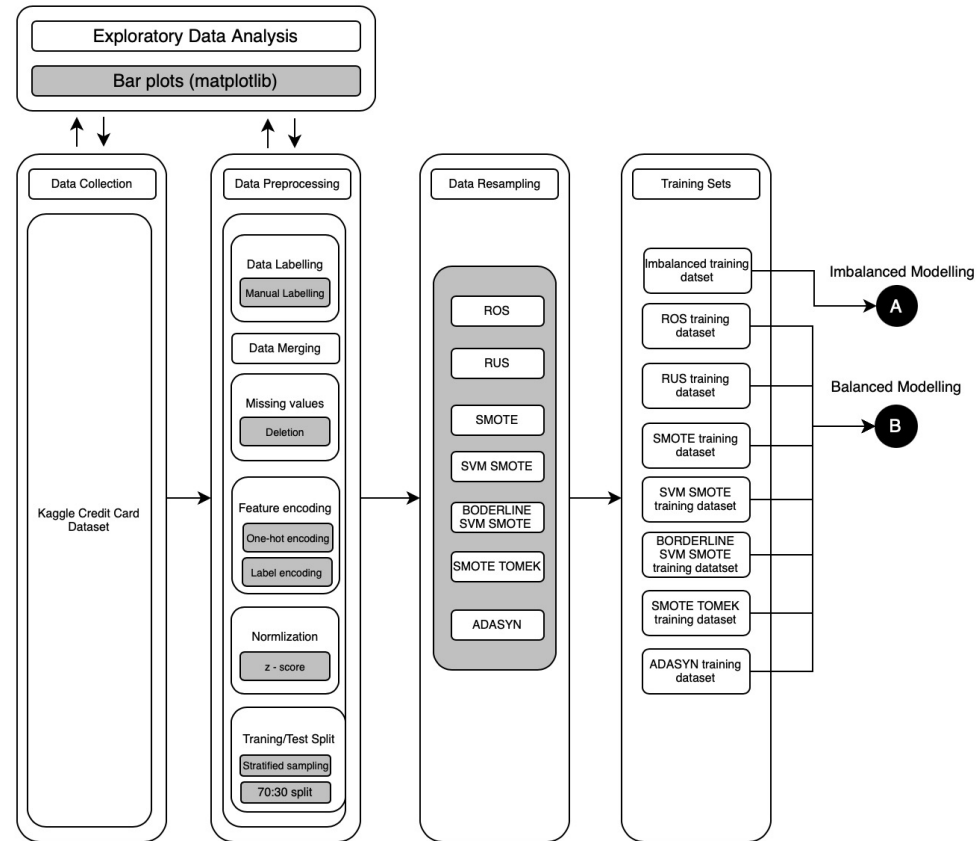
- The surge in credit card and bank loans has come with an increase in loan defaults, costing banks billions annually. To navigate this rising risk, accurate and efficient credit assessment has become more crucial than ever.

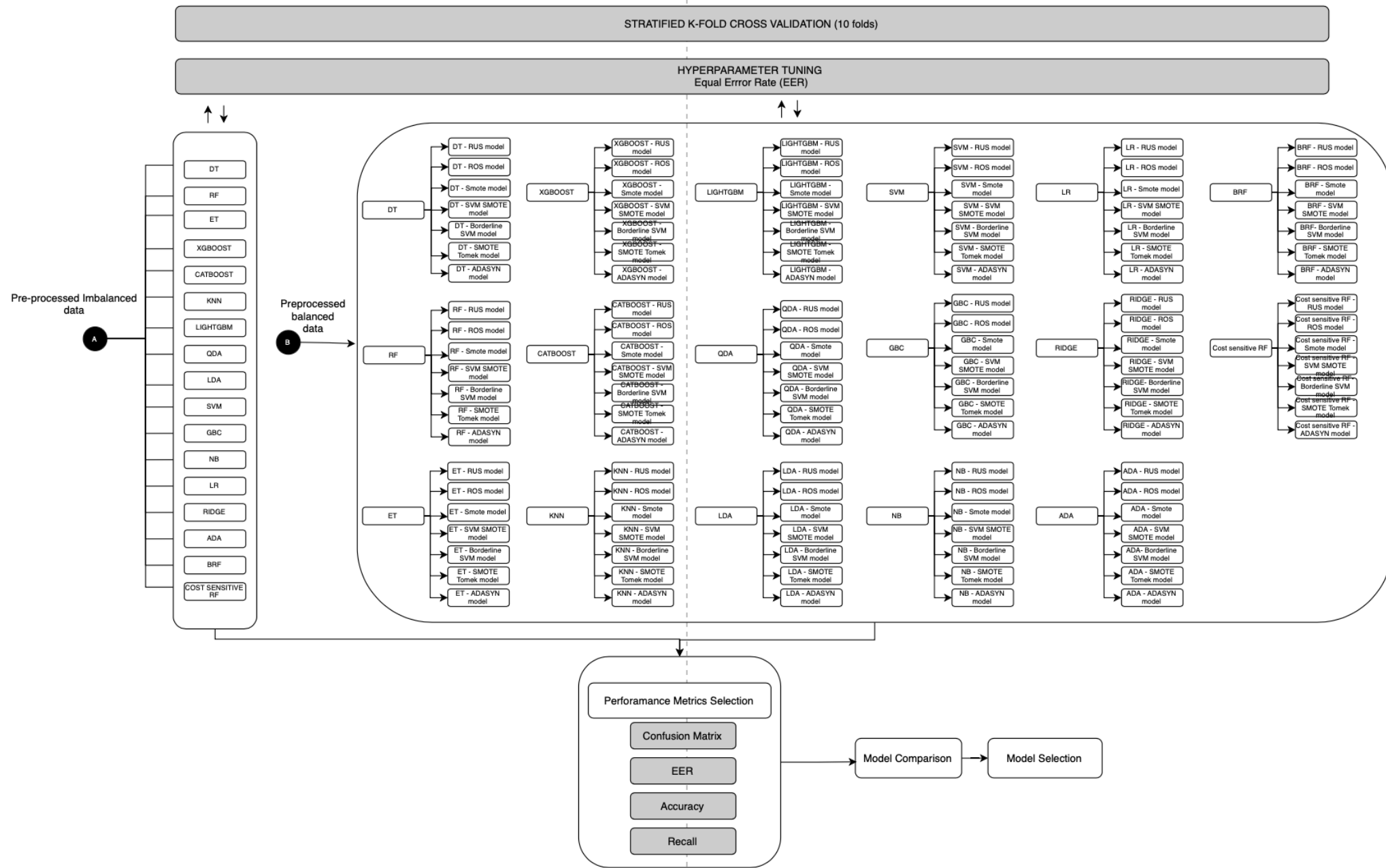
# **2. Project Aim**

- This project aims to identify machine learning models that can effectively predict credit card and loan defaults.

# 3. Methodology

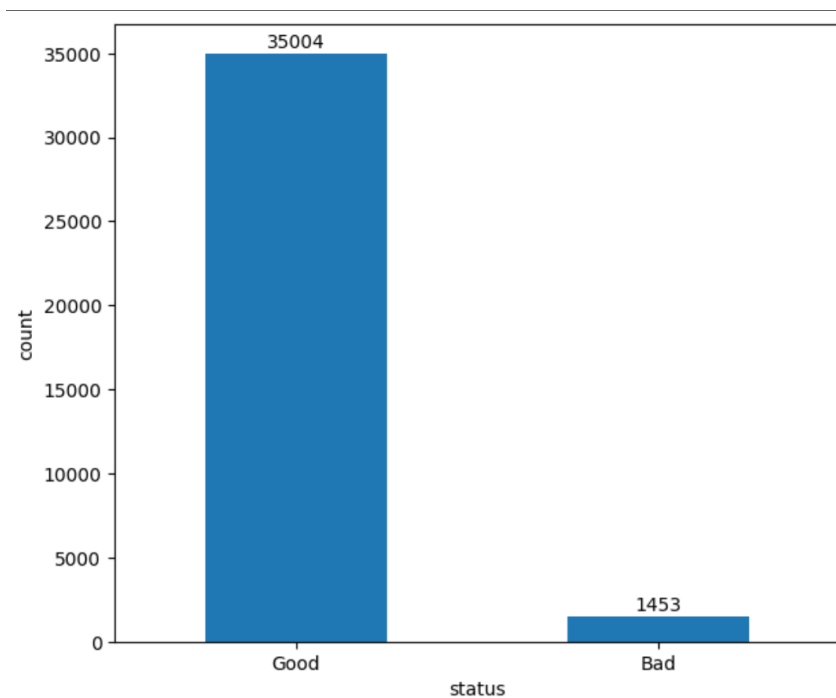
## DATA PREPROCESSING



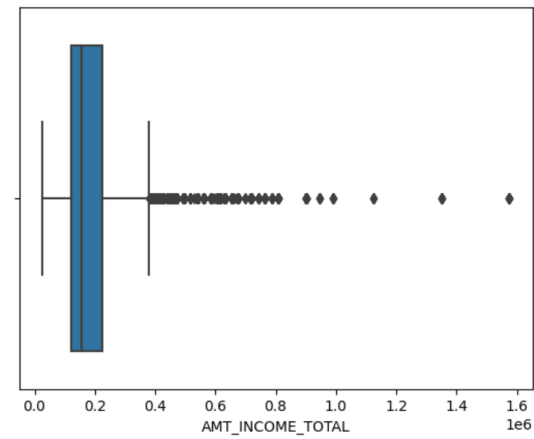
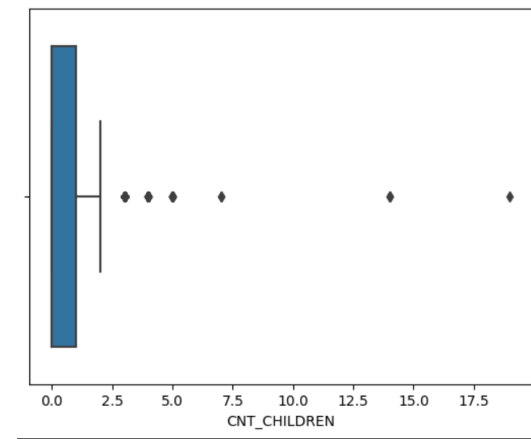
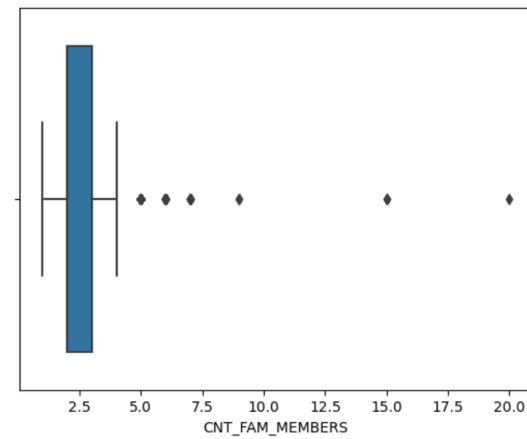
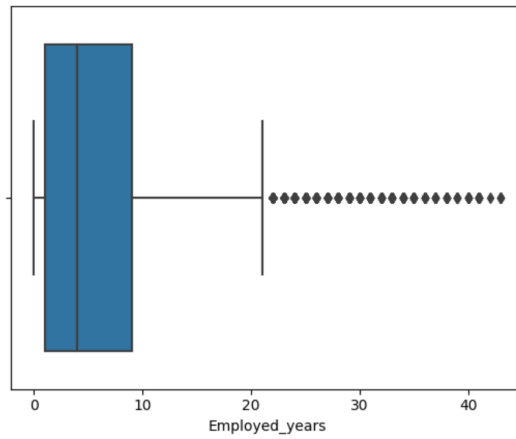


# 4. Challenges

- 4.1. Imbalanced Data



## 4.2. Outliers



## 4.3. Data Merging

application.csv

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 438557 entries, 0 to 438556
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ID                    438557 non-null  int64
1   CODE_GENDER          438557 non-null  object
2   FLAG_OWN_CAR         438557 non-null  object
3   FLAG_OWN_REALTY      438557 non-null  object
4   CNT_CHILDREN         438557 non-null  int64
5   AMT_INCOME_TOTAL     438557 non-null  float64
6   NAME_INCOME_TYPE     438557 non-null  object
7   NAME_EDUCATION_TYPE  438557 non-null  object
8   NAME_FAMILY_STATUS   438557 non-null  object
9   NAME_HOUSING_TYPE    438557 non-null  object
10  DAYS_BIRTH           438557 non-null  int64
11  DAYS_EMPLOYED        438557 non-null  int64
12  FLAG_MOBIL           438557 non-null  int64
13  FLAG_WORK_PHONE      438557 non-null  int64
14  FLAG_PHONE           438557 non-null  int64
15  FLAG_EMAIL           438557 non-null  int64
16  OCCUPATION_TYPE      304354 non-null  object
17  CNT_FAM_MEMBERS      438557 non-null  float64
dtypes: float64(2), int64(8), object(8)
```

credit.csv

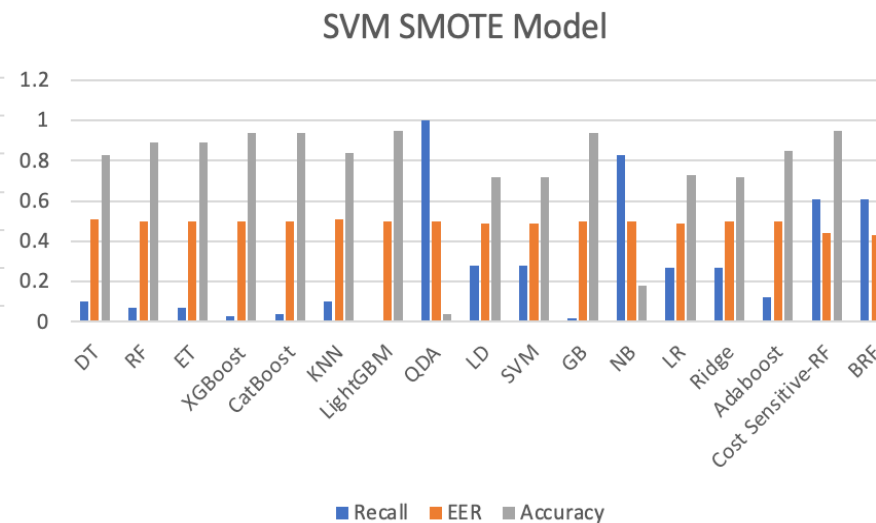
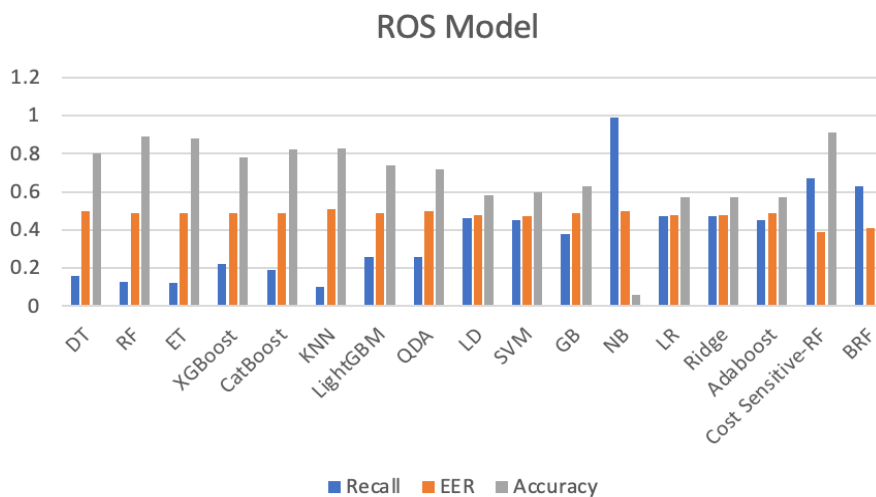
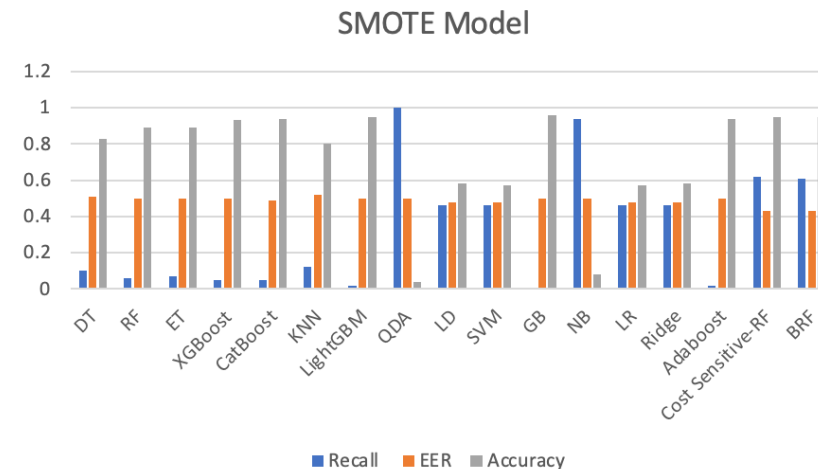
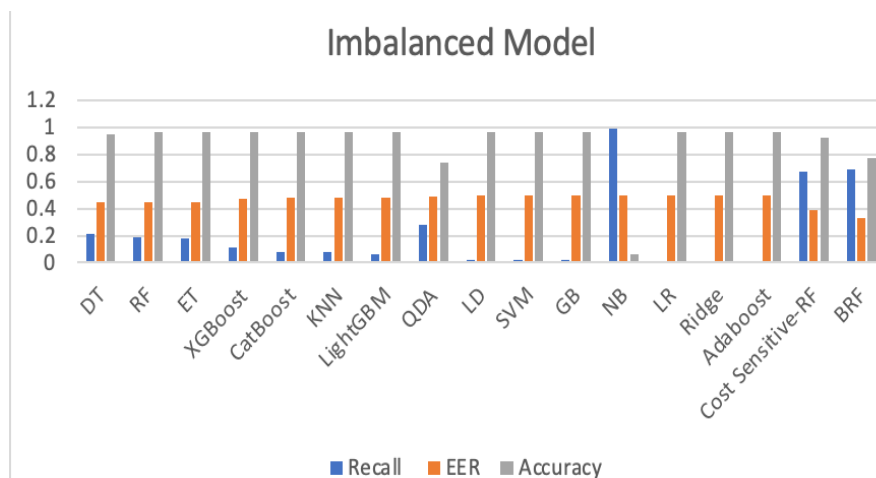
```
RangeIndex: 1048575 entries, 0 to 1048574
Data columns (total 3 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ID                    1048575 non-null  int64
1   MONTHS_BALANCE       1048575 non-null  int64
2   STATUS               1048575 non-null  object
```

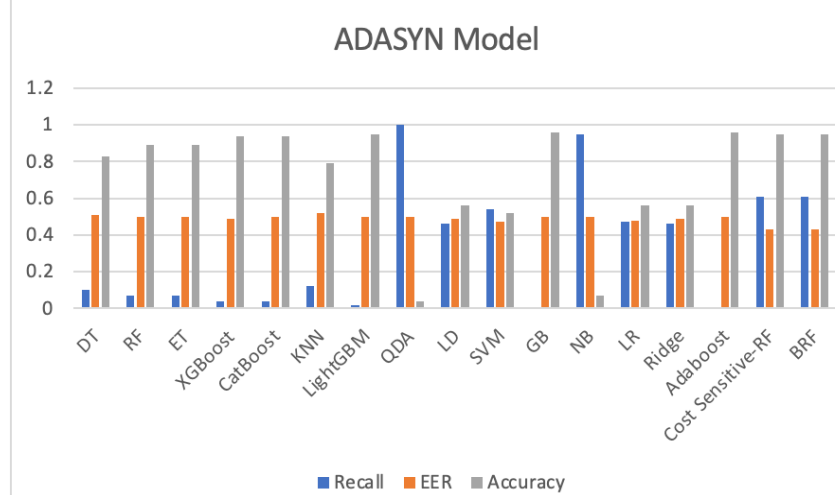
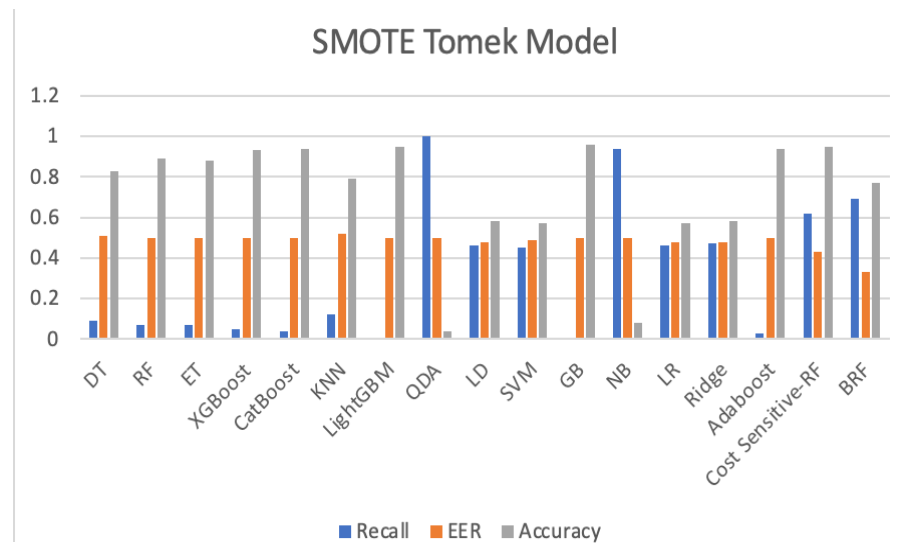
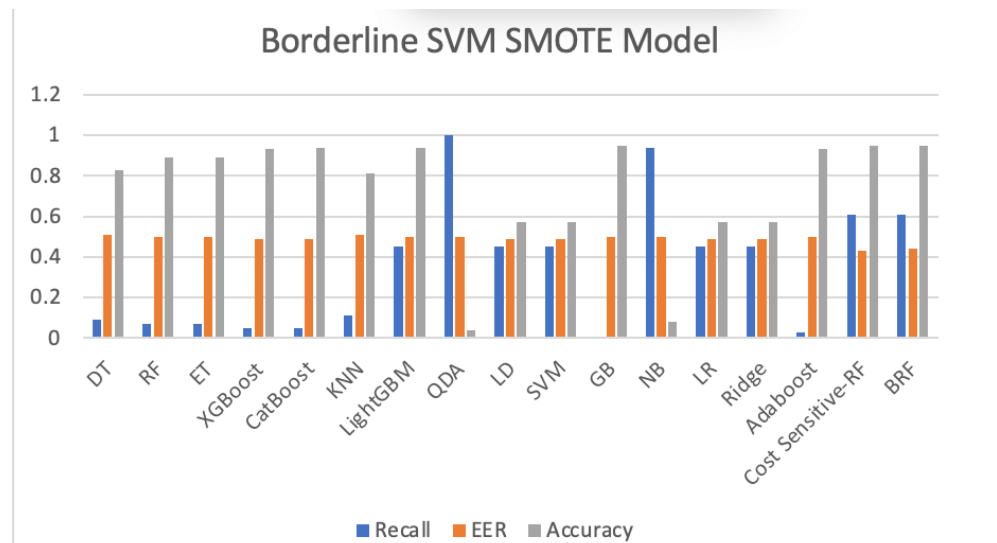
## 4.4. Hyperparameter Tuning

- Navigating the vast parameter space and identifying the most optimal set of parameters to maximize model performance and avoid overfitting was a complex task.
- Fine-tuning ensemble models such as random forest using GridCVSearch proved to be highly resource-intensive, requiring a substantial 5 hours to optimize a single model within the constraints of a limited parameter set.



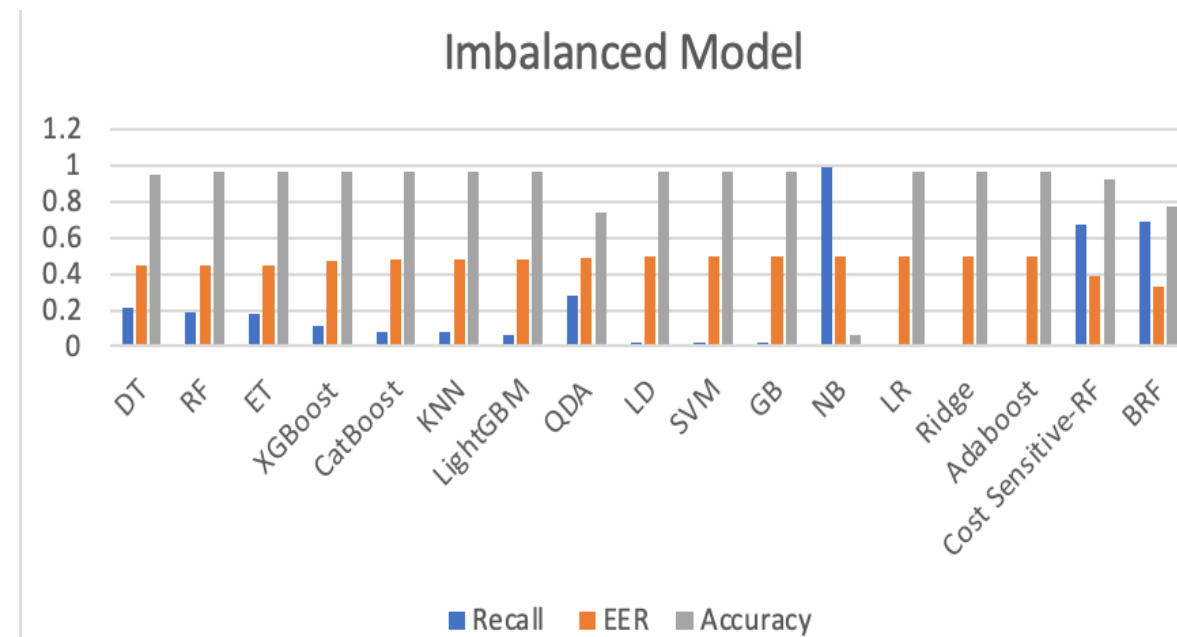
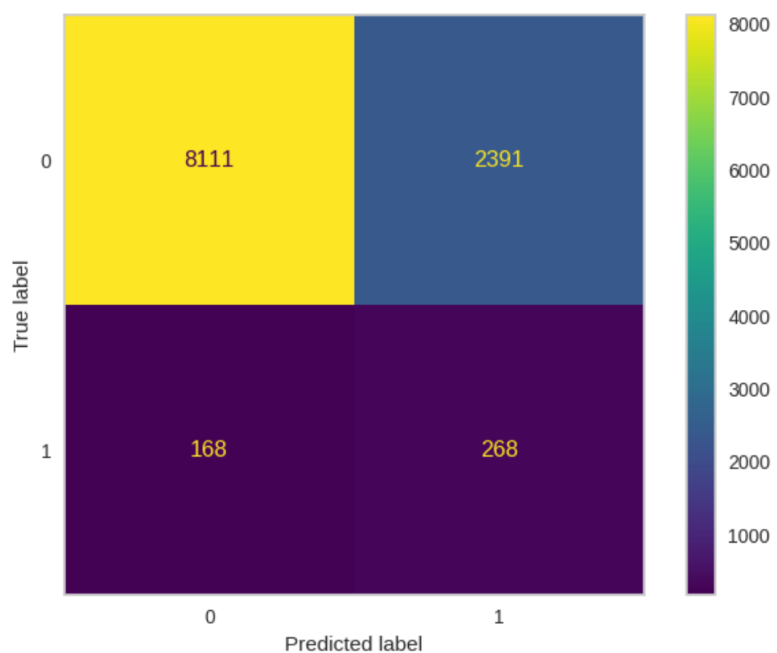
# Model Results





# Model Selection

- Balanced Random Forest



# Conclusion

- To find the optimal model to predict loan defaulters we built seventeen machine learning classification models Decision Trees, Random Forest, Light Gradient Boosting Machine, Linear Discriminant Analysis, Ridge Classifier, Logistic Regression, CatBoost, AdaBoost, Extreme Gradient Boosting, K-Nearest Neighbors, SVM with a Linear Kernel, Quadratic Discriminant Analysis, Extra Trees, Gradient Boosting, Naïve Bayes, Balanced Random Forest and Cost Sensitive Learning (Random Forest).
- To tackle the imbalance data problem we used different variations of datasets using seven different data sampling techniques namely SMOTE, SVM SMOTE, BORDERLINE SVM SMOTE, SMOTE TOMER, ADASYN, Random Undersampling, and Random Oversampling.
- Balanced Random Forest with imbalanced data and default parameters stands out with its well-rounded performance with an accuracy of 77% and lowest EER score of 0.33.

# Limitations

- We did not investigate any feature selection methods.
- We could not conduct a more extensive search for optimal hyperparameters for the selected models.
- We did not incorporate deep learning algorithms.

THANK YOU!