

---

# HomeCredit ScoreCard Model

---

Ridwan Akmal

Home Credit Indonesia Data Scientist Virtual  
Internship Program Rakamin Academy



## Problem Research

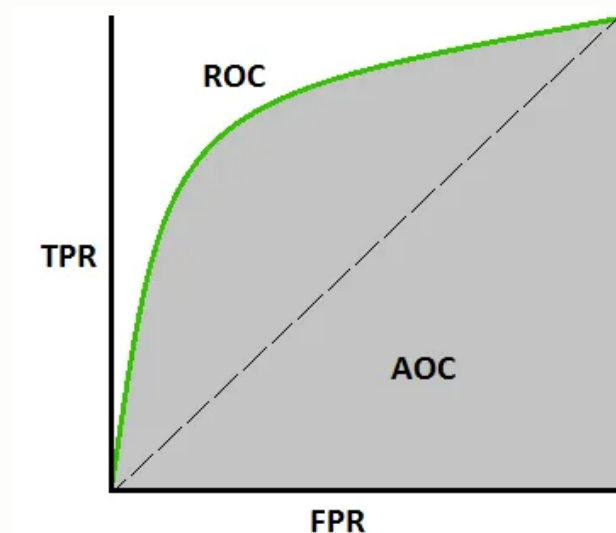
Home Credit Indonesia as a company provides financing services in shopping. This multipurpose financing makes it very easy to shop from online or offline. Here Home Credit seeks to classify clients who have difficulty paying loans and those who have no problems. In this case Home Credit uses a variety of alternative data including telecommunication and transactional information.

**HOME  
CREDIT**

## Main Goals

Here Home Credit want to minimize the number of client who have a problem with repay the loan

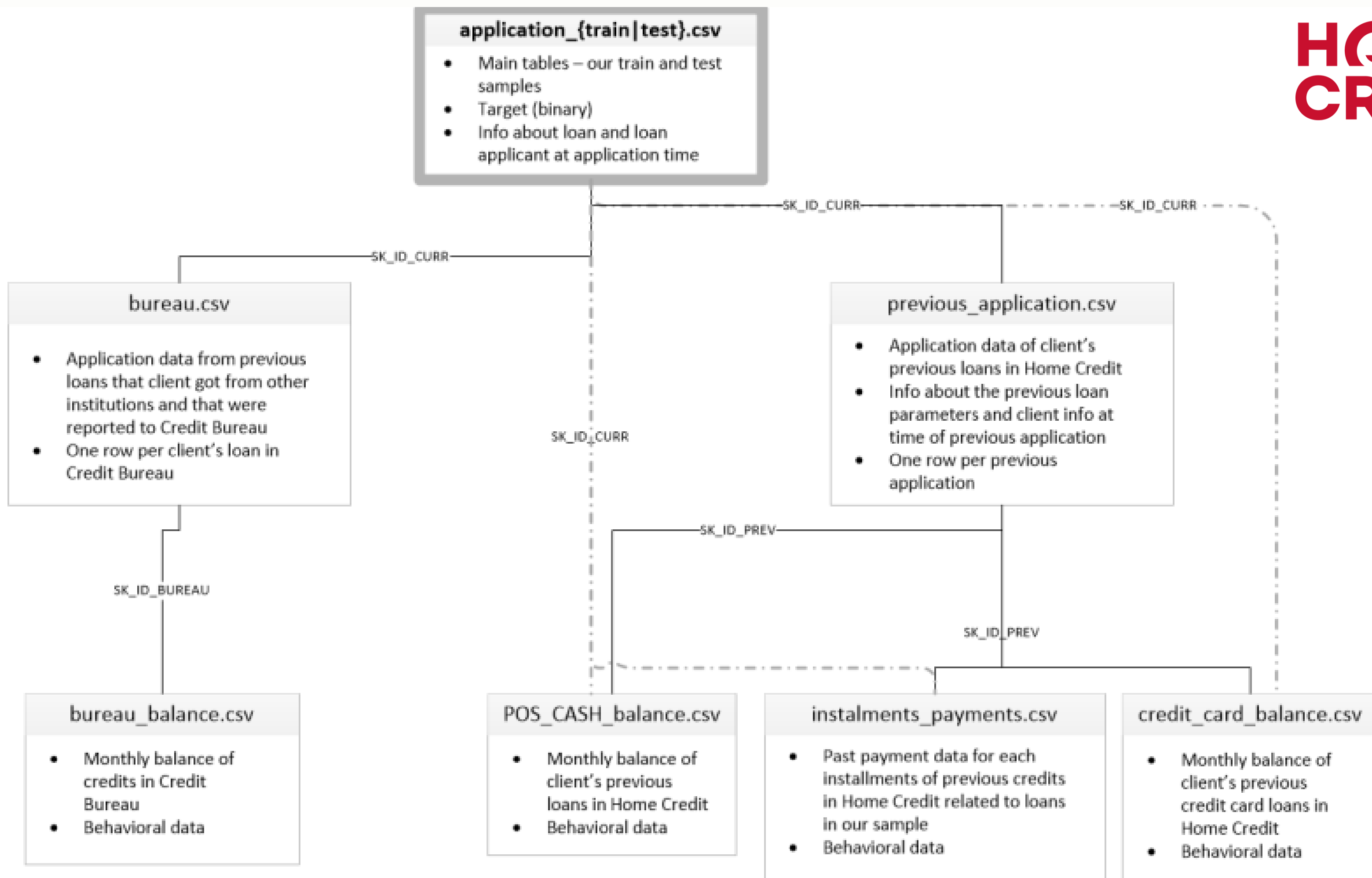
## Model Evaluation



Area under ROC curve

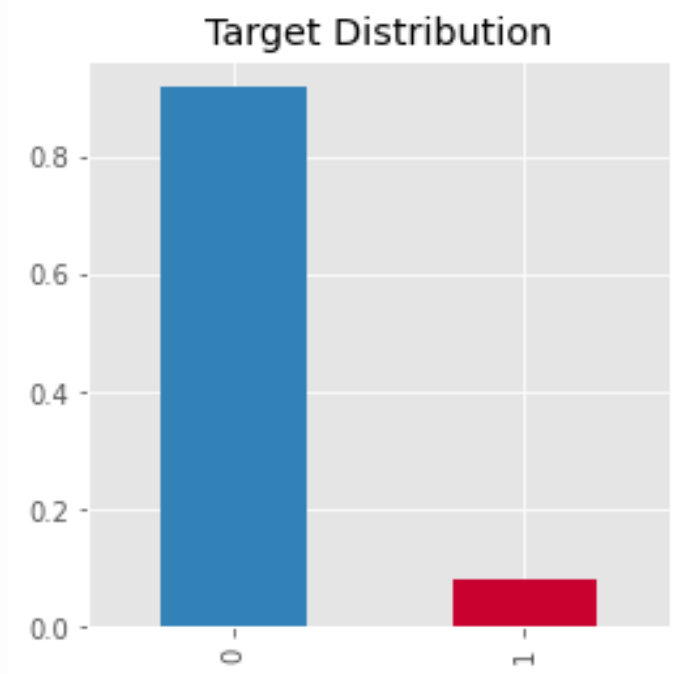


# Dataset & Description



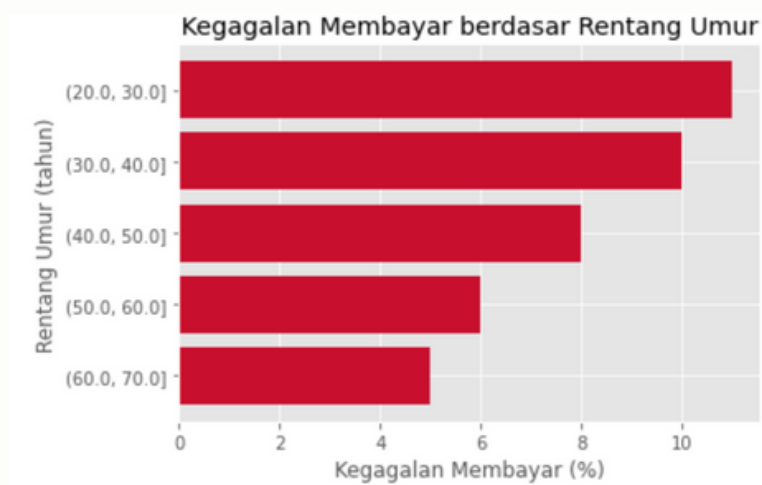
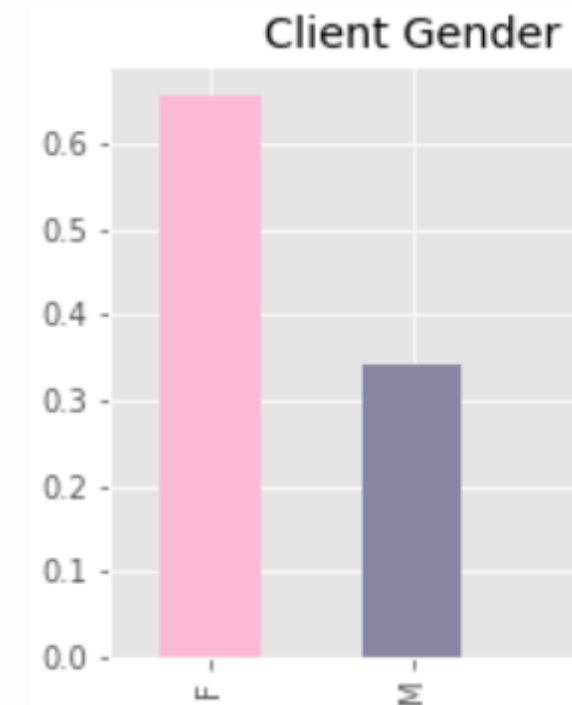


# Exploratory Data Analysis

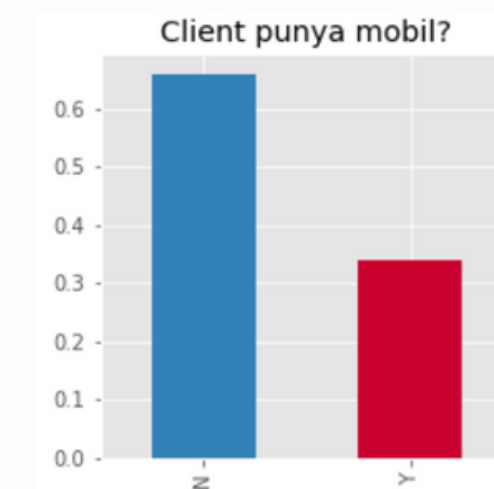


The visualization shows that there is an imbalance in the data where more clients have no problems paying off loans, then will do sampling before modeling

Here we also get insight where clients are mostly dominated by women and clients who dont have children



From here we can see that a group in a range youngest old & group range shortest year employed have a problem which is failure to repay the loan

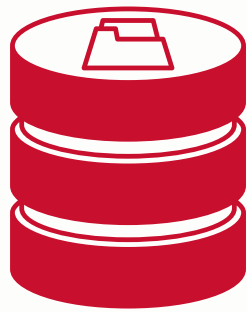


Most of the clients owned realty but dont have a car



# Data Preprocessing

## RAW DATA



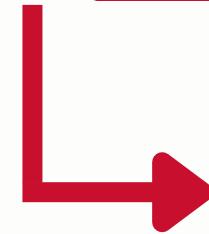
### Data Cleaning



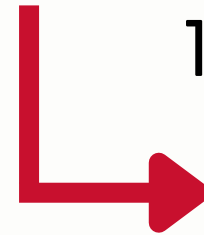
### Feature Engineering



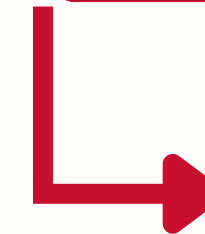
### Model Building



1. Replace XNA values with NaN
2. Detect & Handling Missing values
3. Detect redundant data



1. Convert day birth to years and get client age
2. Calculating number of documents, IAP, & EITC
3. Scalling the numerical features
4. Encoding the categorical features
5. Drop features with  $VIF > 10$



1. Oversampling smote with ratio 2:1
2. Build model with various algorithm like Logistic Regression, Adaboost, & KNN
3. Model evaluation, compare all model which one is the best then all

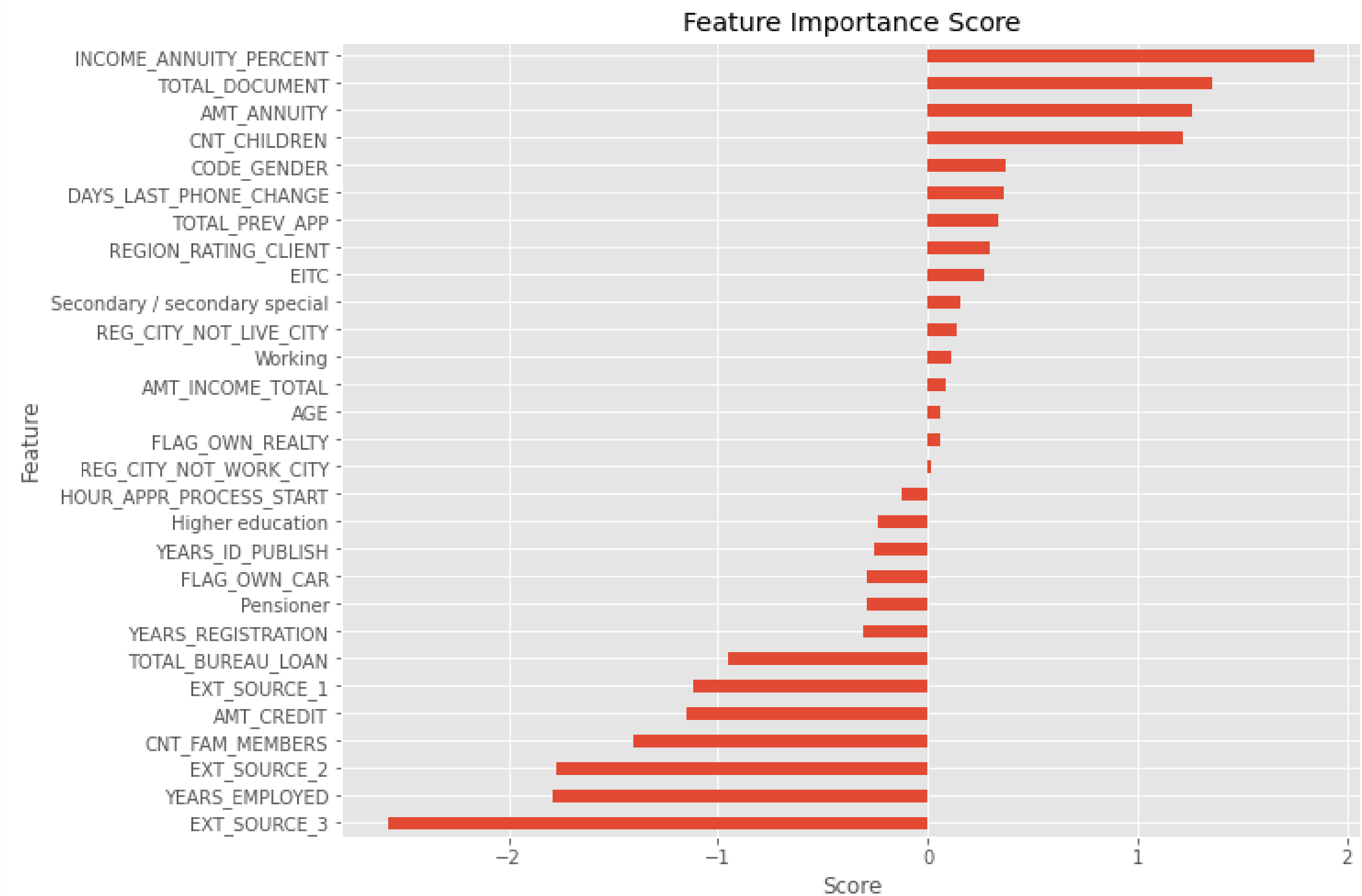
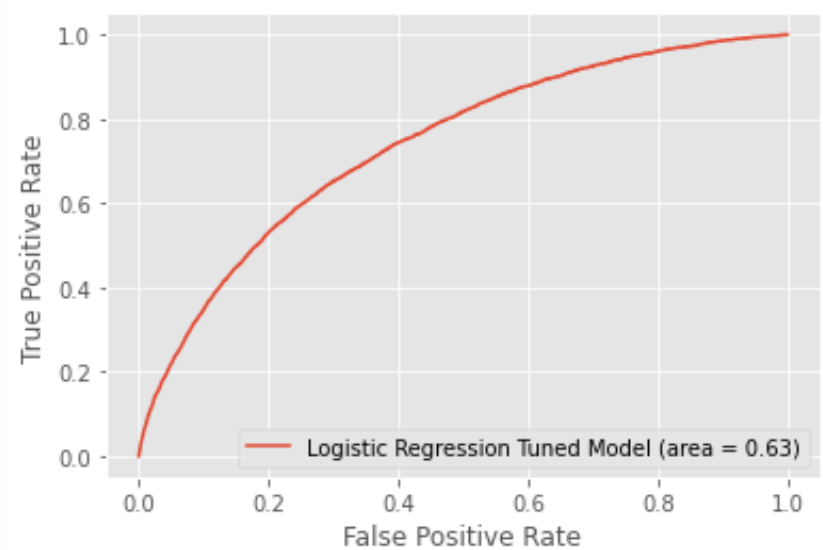
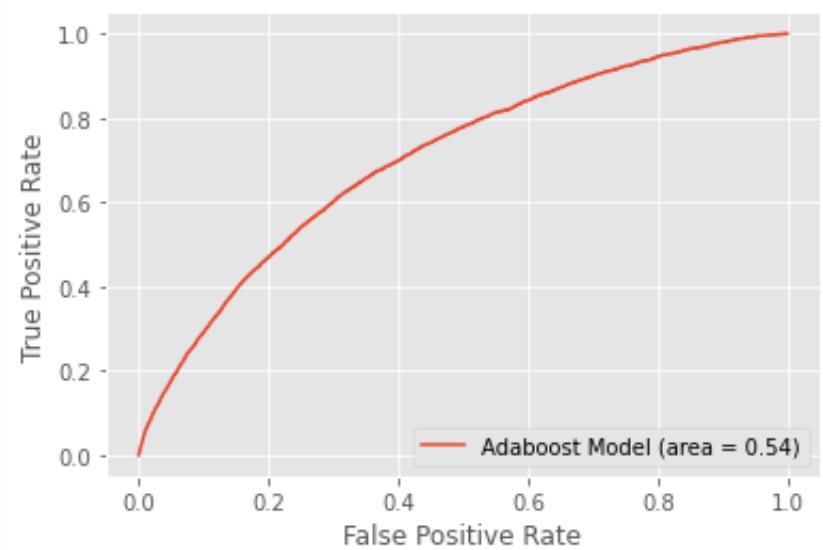
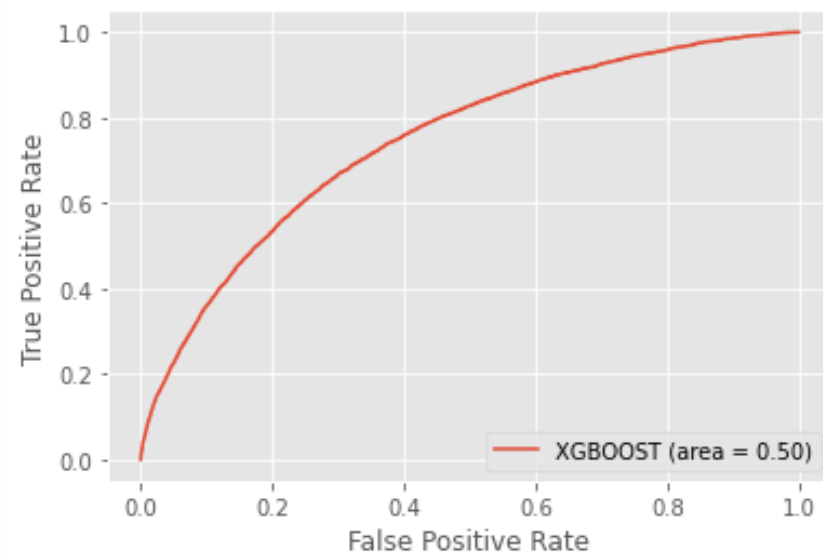


# Model Comparasion

	Training acc score	Testing acc score	Area under ROC
Logistic Regression	0.715	0.843	0.63
XGBOOST	0.936	0.919	0.5
Adaptive Boosting	0.830	0.907	0.54



As previously mentioned we make a area under ROC as a model eval then with the highest score it can be ascertained that the winner is logistics regression, but i feel there is an overfit model



We can see here that five most important features is INCOME\_ANNUITY\_PERCENT, TOTAL\_DOCUMENT, AMT\_ANNUITY, CNT\_CHILDREN, CODE\_GENDER



# Business Recommendation

1. From the visualizations we know that the feature NAME\_INCOME\_TYPE with values student and Business Man have no problems to repay the loan(100% application approved) and next one is feature OCCUPATION\_TYPE with values Managers, High skill tech staff, and Accountants so we can create campaign for Student, Business Man, Managers High skill tech staff, and Accountants so that more can be interested & apply for loan.
2. Next, we can see that clients with maternity leave and unemployed status have a difficult tendency to repay loans, where both reach a percentage of more than 30%, for that we can recommend other types of loan contracts that are suitable for both clients like that.



backenddeviwan@gmail.com

Get in touch  
and let's  
have a chat!



<https://github.com/RidwendDev>

