# Loan Approval Prediction

This is a Binary Classification problem. We must correctly classify whether a loan will be approved or denied based on the given features.

## Processed Data

Modified Data had a shape of (514,17) along with:
- Number of Integer-Categorical Columns        = 7
- Number of String-Categorical Columns         = 6
- Number of String-Boolean Columns             = 1
- Number of Numeric-Boolean Columns            = 1
- Number of ID Columns                         = 1

The Features are grouped in Numerical and Categorical variable:

| Numerical Features | <ul><li>LoanPayoffPeriodInMonths</li><li>RequestedAmount</li><li>InterestRate</li><li>YearsAtCurrentEmployer</li><li>YearsInCurrentResidence</li><li>Age</li><li>NumberOfDependantsIncludingSelf</li><li>CurrentOpenLoanApplications</li></ul> |
|---|---|
| Categorical Features | <ul><li>LoanReason</li><li>Co-Applicant</li><li>RentOrOwnHome</li><li>TypeOfCurrentEmployment</li><li>CheckingAccountBalance</li><li>DebtsPaid</li><li>SavingsAccountBalance</li></ul> |
| Target | <ul><li>WasTheLoanApproved</li></ul> |

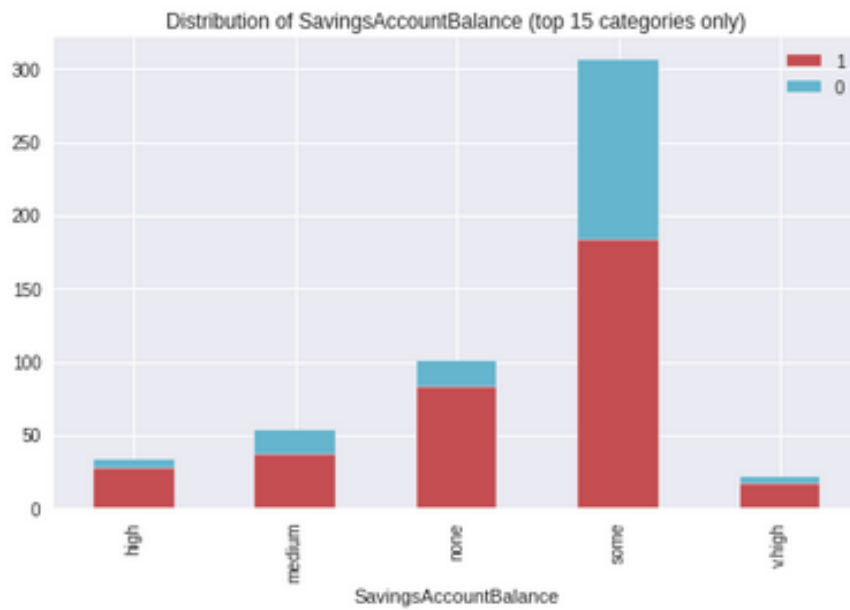The given data set was imbalanced with 'Loan Approved' being the majority class

# Univariate and Bivariate Analysis

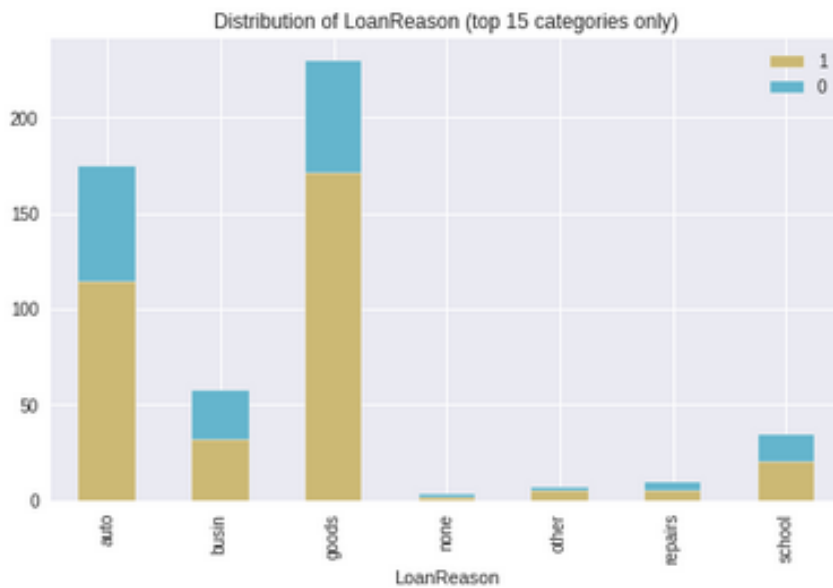| Variable | Type | Comments |
|----------|------|----------|
| WasTheLoanApproved | Dependent variable | 67% of people had their loans approved |
| RequestedAmount | Independent variable | Mean value of 4000, values lie in the range of (1024-18400) |
| InterestRate | Independent variable | Mode value of 2, values lie in the range of range (0-4) |
| Age | Independent variable | Mean value of 36 and values lie in the range of (19 to 75) |

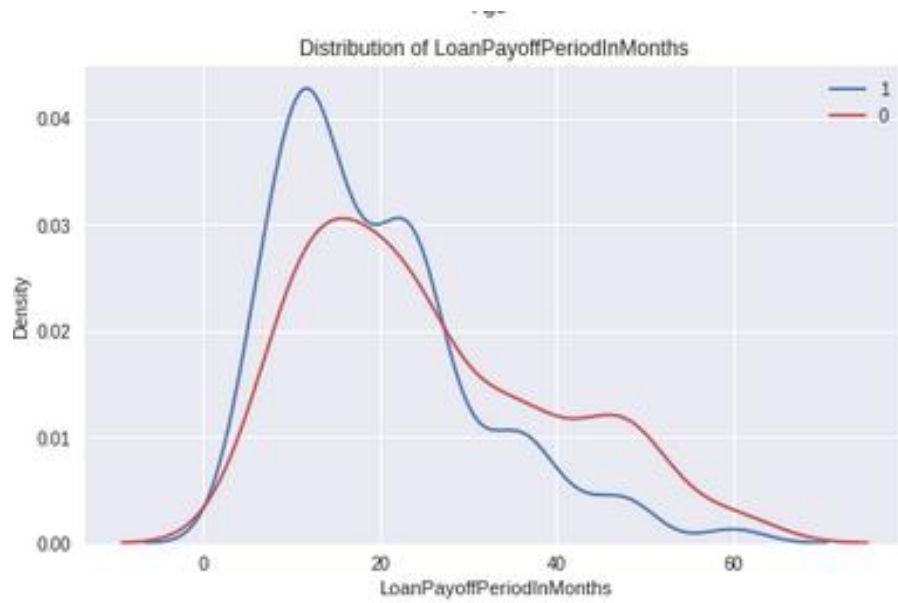Analysis of 'CheckingAccountBalance' with Target Variable: 'WasTheLoanApproved'

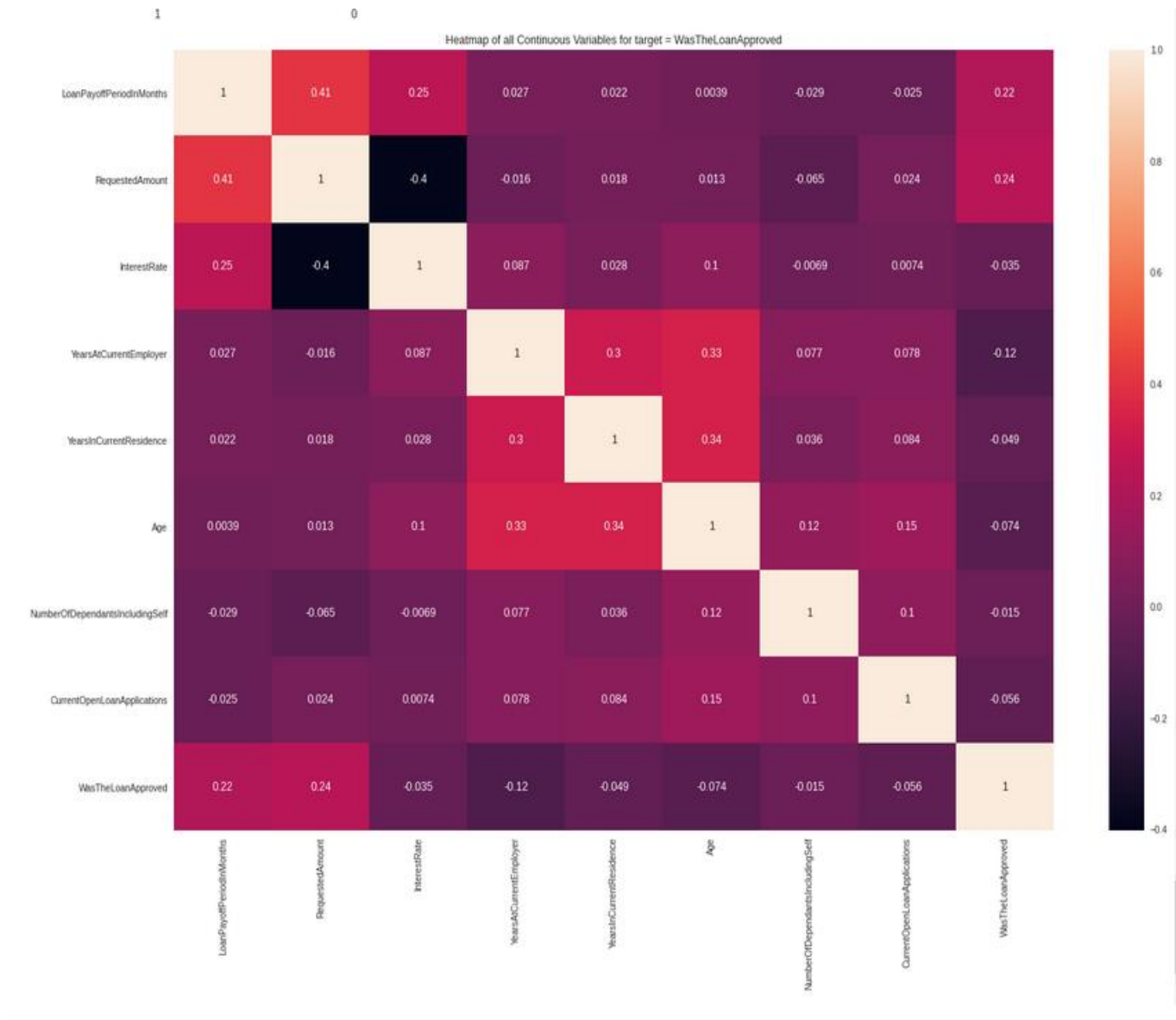Analysis of 'SavingsAccountBalance' with Target Variable: 'WasTheLoanApproved'



Distribution of SavingsAccountBalance (top 15 categories only)

Analysis of 'LoanReason' with Target Variable : 'WasTheLoanApproved'



Distribution of LoanReason (top 15 categories only)

Analysis of 'LoanPayoffPeriodInMonths' with Target Variable : 'WasTheLoanApproved'



Distribution of LoanPayoffPeriodInMonths

# Correaltion Matrix between different features



From the Correaltion Matrix we can infer that:

- Loan pay off periods in months is correlated with Requested Amount
- Loan pay off periods in months is correlated with Interstate
- Years at Employer is correlated with Years in Current residence
- Years at Employer is correlated with Age
- Years in Current residence is correlated with Age

# Model Building and Feature Engineering

From Feature selection algorithm the *important features* were:

- CheckingAccountBalance
- LoanPayoffPeriodInMonths
- RequestedAmount
- SavingsAccountBalance
- CurrentOpenLoanApplication
- Age
- YearsAtCurrentEmployer
- LoanReason

For model validation, I used accuracy, precision, and recall.

The best base line accuracy was for XgBoost which had the highest accuracy of close to 69%

After up-sampling and scaling, the ensemble voting model of "XGB","RF","DT","ADB","GB" showed an accuracy of 80%

## Submission Files:

1) Data.csv
2) Model_experimentation.ipynb
3) Data_explore.ipynb
4) Preprocessing.ipynb

## Assumptions:

I used inner join of the .tsv files to avoid data imputation problem. The data had 515 rows for training/testing.