

# Mortgage Classification using CatBoostclassifier

bokola

27 June, 2019

## Table of Contents

Executive Summary.....	1
Methodology .....	1
Description of Data.....	2
Data Recoding and Visualization.....	2
Data Recoding.....	2
Data Visualization.....	2
Feature Engineering.....	11
Model Building .....	11
Results and Conclusions.....	12

## Executive Summary

The objective was Predicting Mortgage Approvals from Government Data, taking into consideration the **demographics, location, property type, lender**, and other factors to predict whether a mortgage application was accepted or denied. It is a classification problem - A classifier is a machine learning model that separates the **label** into categories or **classes**. In other words, classification models are **supervised** machine learning models which predict a categorical label. We apply relevant skills to classify the label **accepted** - whether a mortgage application was accepted or declined using data obtained across the United States.

We used a CatBoostclassifier model with 1000 iterations, 8 cross-validations returning an accuracy of 0.72.

## Methodology

The section outlines data manipulation procedures, visualizations and model selection.

## Description of Data

Data consisted of 23 variables (21 possible features and 1 binary label(accepted) and 1 arbitrary identifier(row\_id). Both training and test sets spanned 500000 records

There were a number of variables in the train and test sets with missing values: applicant\_income - 39948, population - 22465, minority\_population\_pct - 22466, ffiecmedian\_family\_income - 22440, tract\_to\_msa\_md\_income\_pct - 22514, number\_of\_owner.occupied\_units - 22565, number\_of\_1\_to\_4\_family\_units - 22530. The missing entries were recoded to -999 for numeric variables.

## Data Recoding and Visualization

Variable recoding and visualization are key steps that ensure every variable is of the desired class/type. Visualization is key to understand the distribution and separation of values by the two levels of a binary label.

### Data Recoding

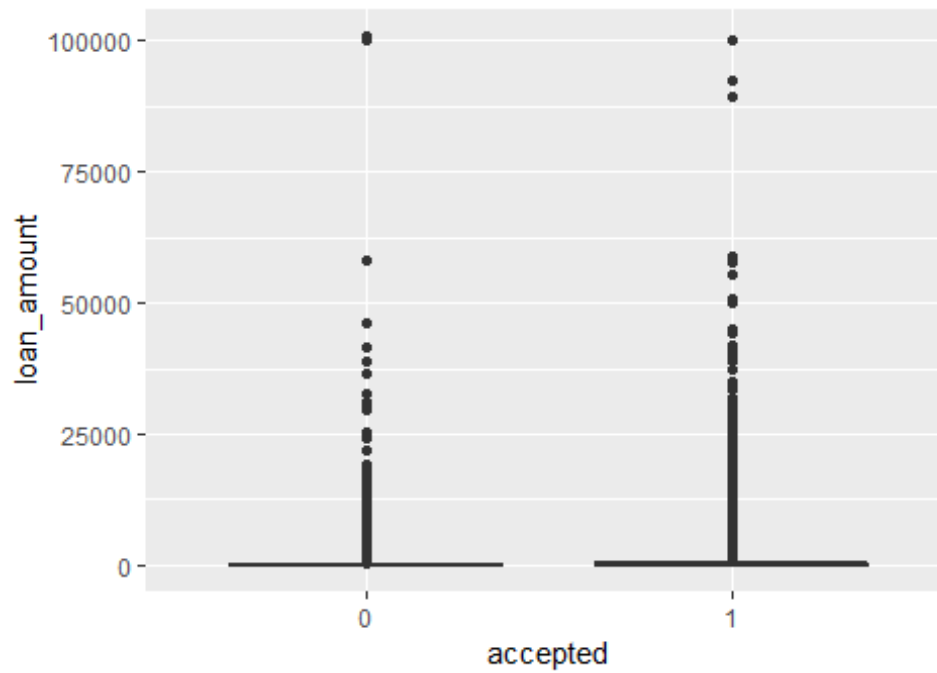
Our data had a number of categorical features (msa\_md, state\_code, county\_code, lender, loan\_type, property\_type, loan\_purpose, occupancy, preapproval, applicant\_ethnicity, applicant\_race, applicant\_sex) but which are read in as numeric. We began by recoding such variables into the desired character class and explicitly supplying the categories in both train and test sets. We then visualize the separation of values between different levels of the label - accepted(0,1).

### Data Visualization

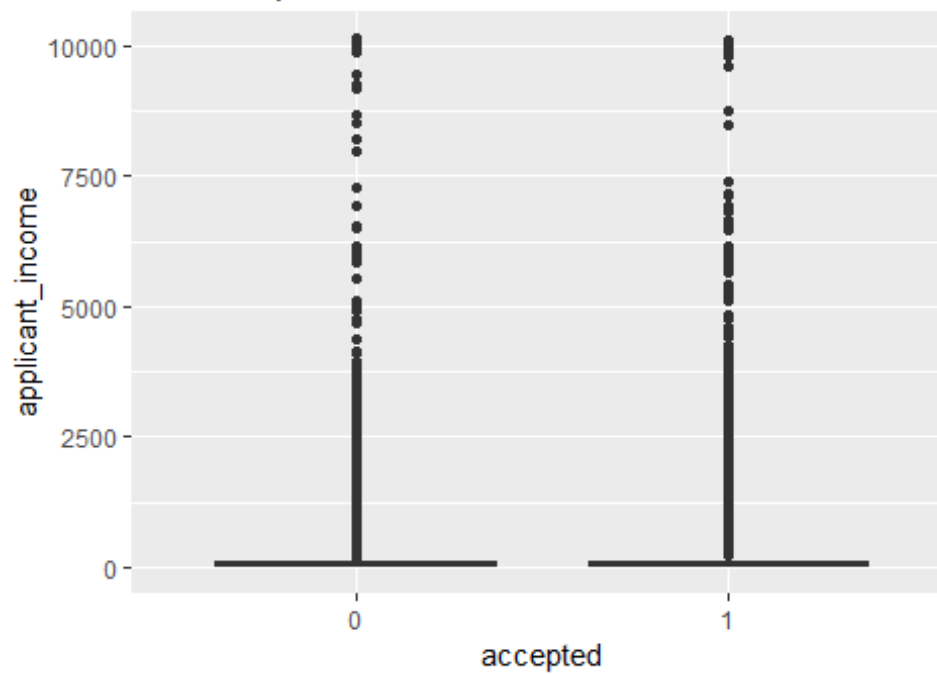
#### 1. Class separation by numeric variables

We explored the distinction of the distribution of values between the two levels of the label through boxplots.

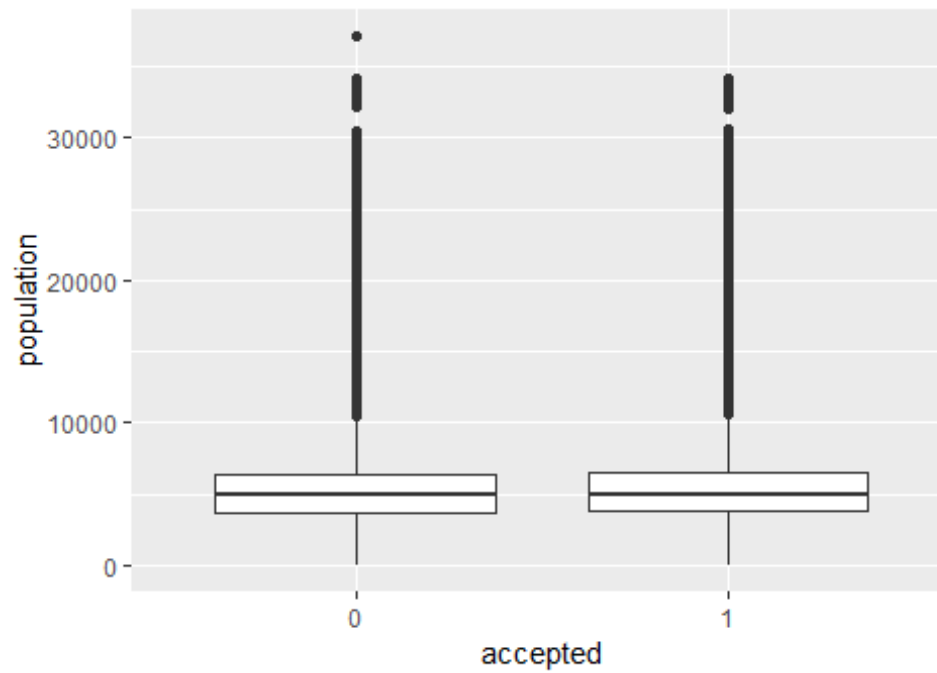
Box plot of loan\_amount  
vs. accepted



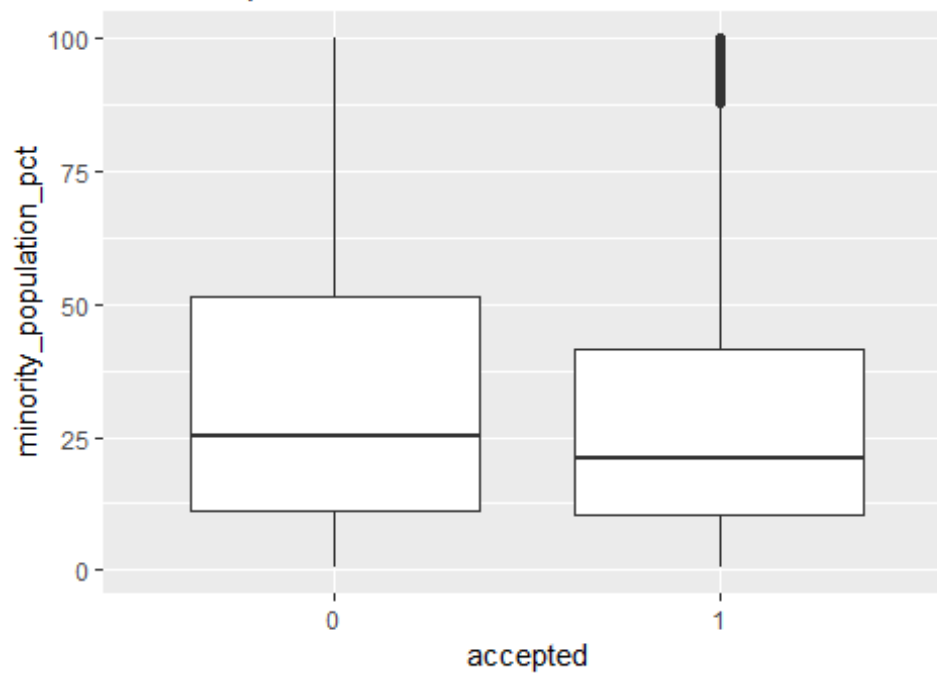
Box plot of applicant\_income  
vs. accepted



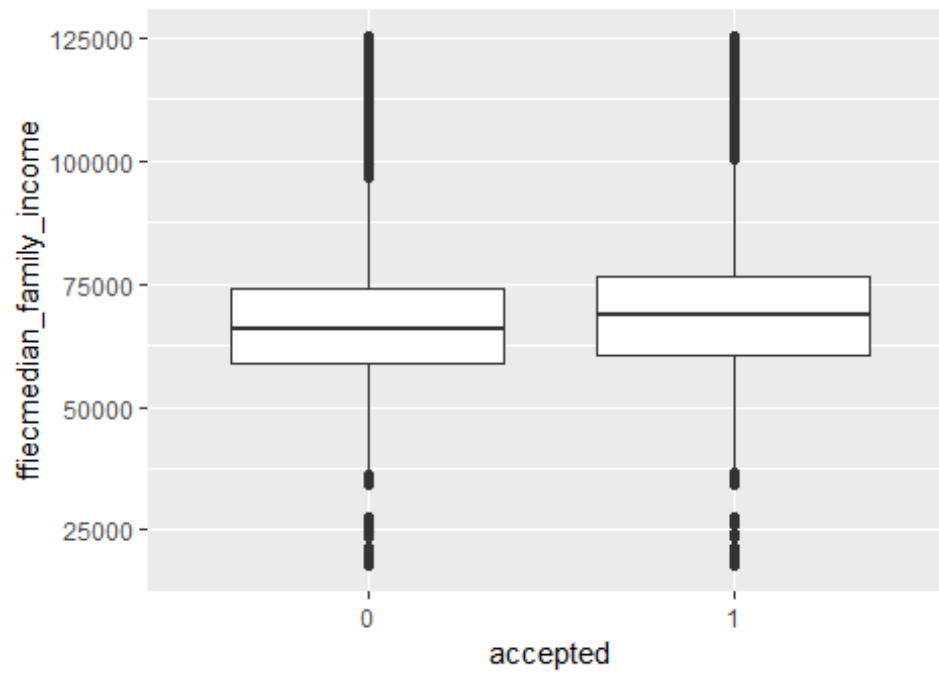
Box plot of population  
vs. accepted



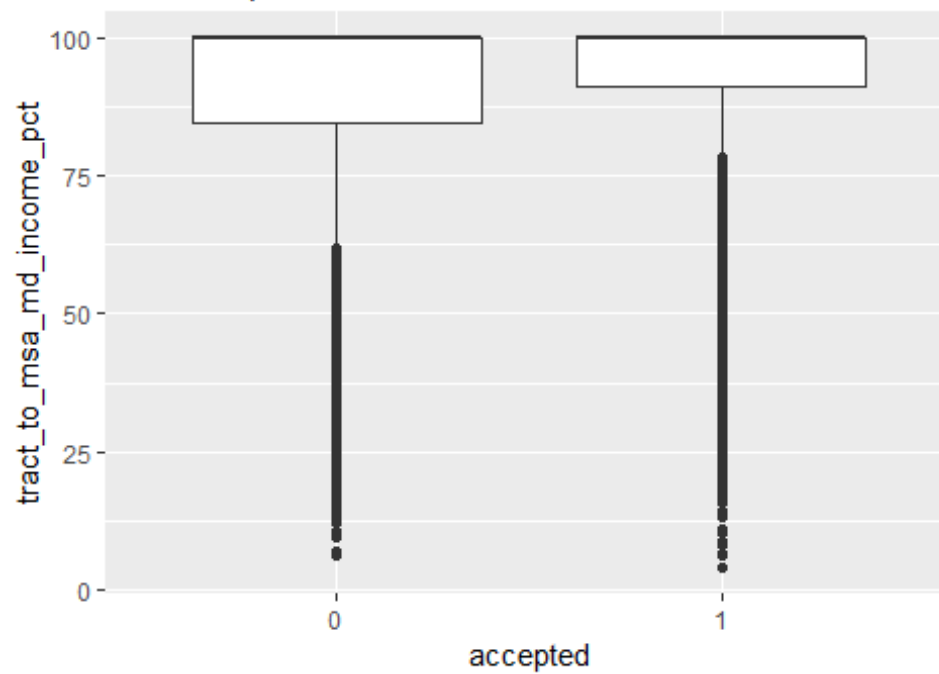
Box plot of minority\_population\_pct  
vs. accepted



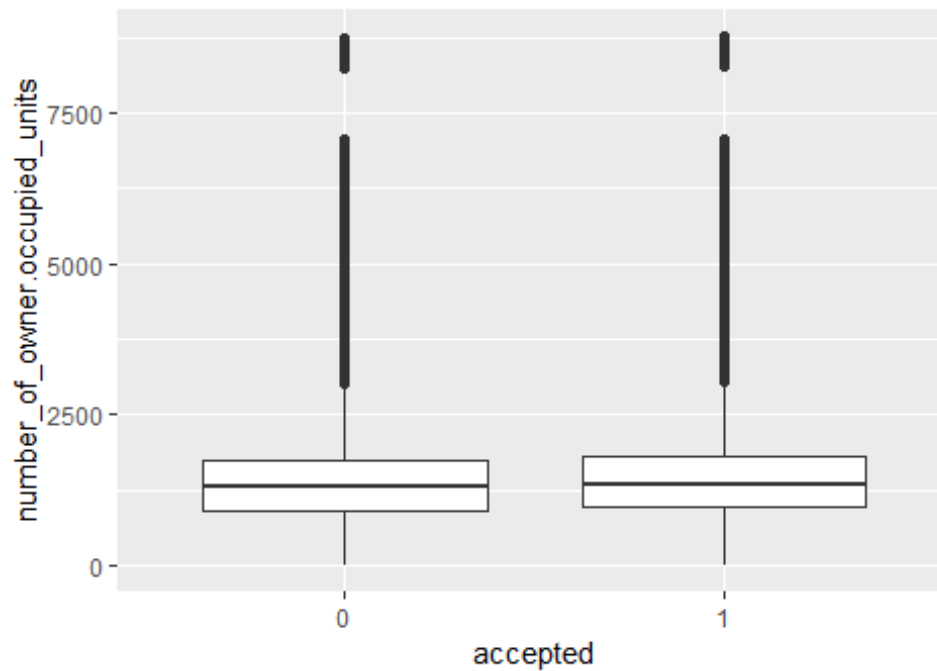
Box plot of ffilecmedian\_family\_income  
vs. accepted



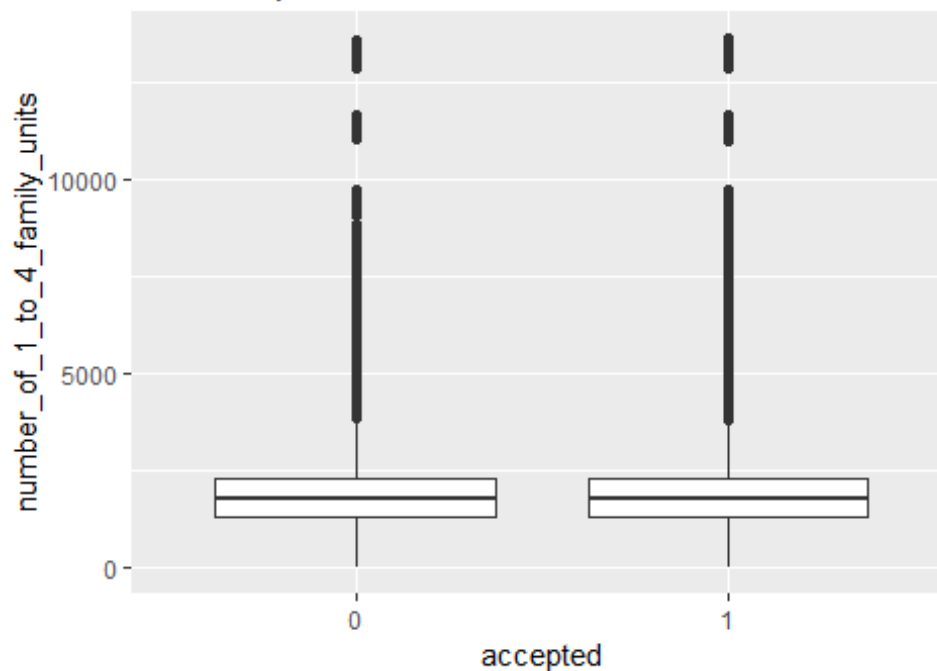
Box plot of tract\_to\_msa\_md\_income\_pct  
vs. accepted



Box plot of number\_of\_owner.occupied\_units  
vs. accepted



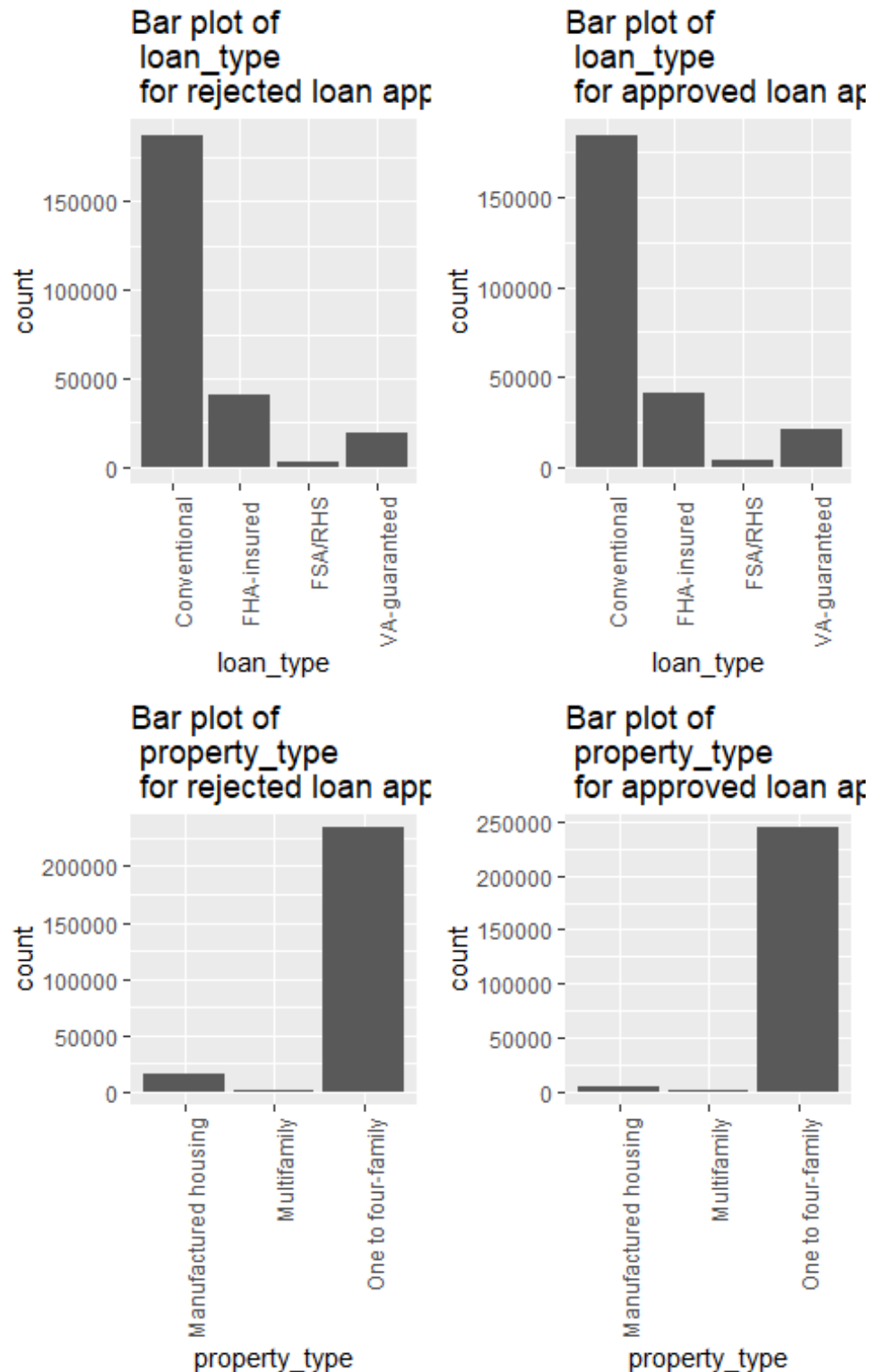
Box plot of number\_of\_1\_to\_4\_family\_units  
vs. accepted



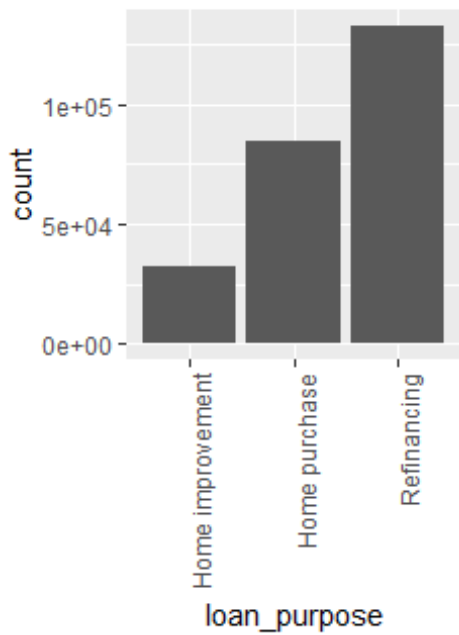
From the plots, three variables - loan\_amount, minority\_population\_pct and tract\_to\_msa\_md\_income\_pct show clear separation of values between the two levels. Next we explore this separation by categorical variables.

## 1. Class separation by categorical variables

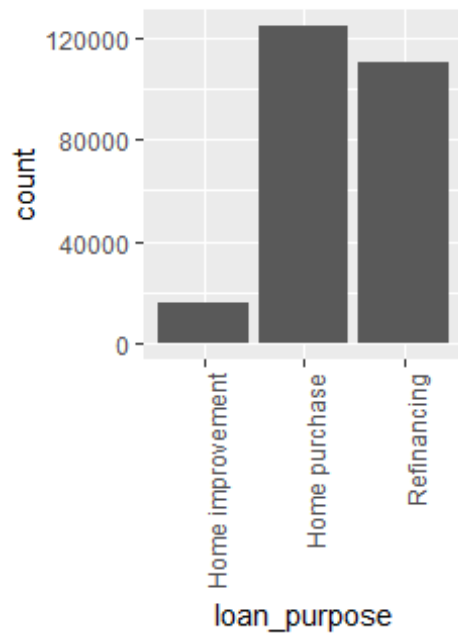
We explore the distinction of the distribution of values between the two levels of the label through bar plots.



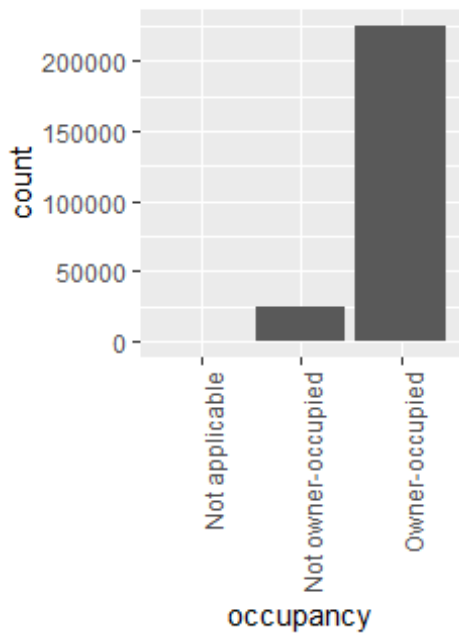
Bar plot of  
loan\_purpose  
for rejected loan appl



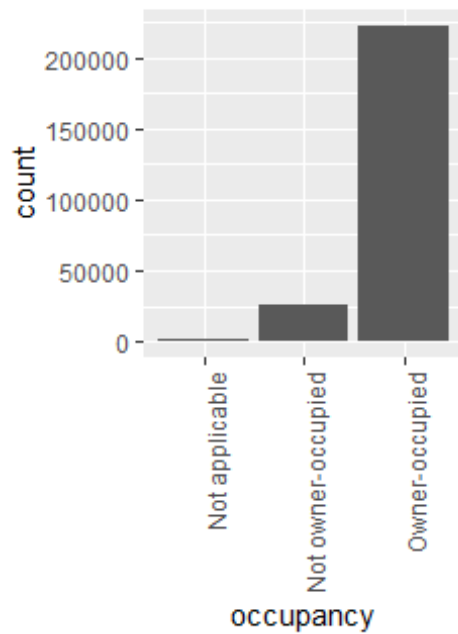
Bar plot of  
loan\_purpose  
for approved loan ap



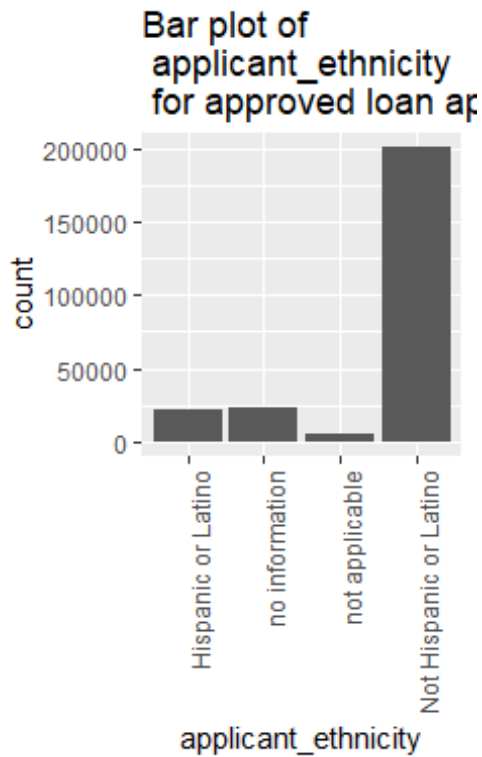
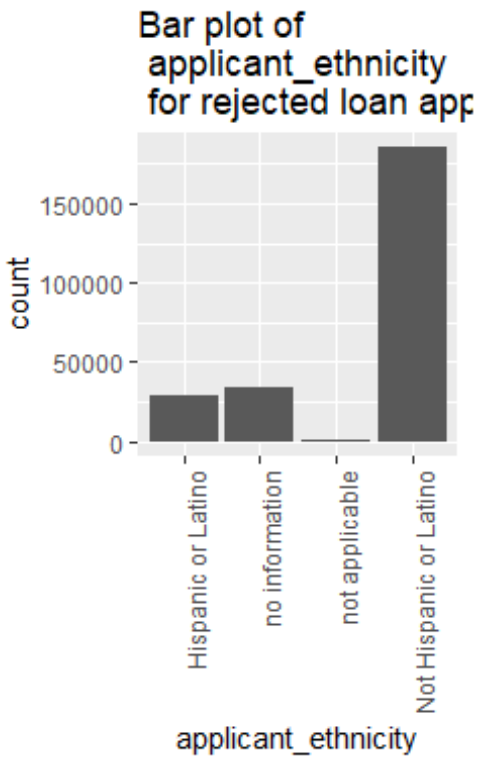
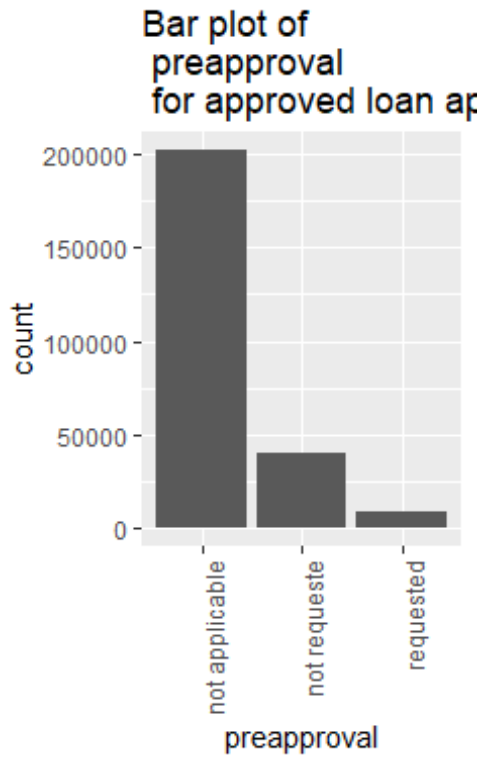
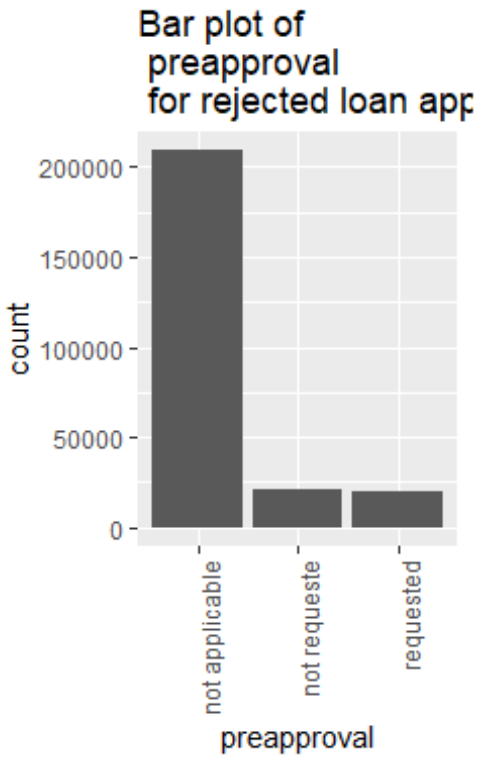
Bar plot of  
occupancy  
for rejected loan appl



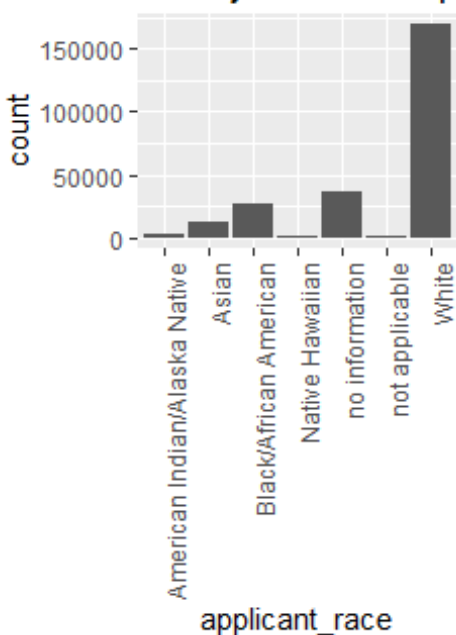
Bar plot of  
occupancy  
for approved loan ap



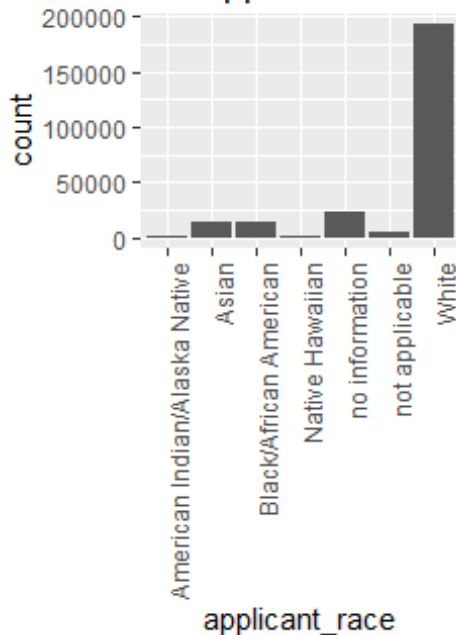




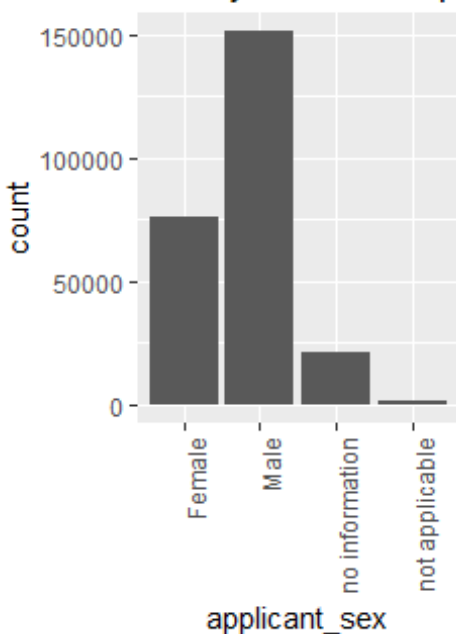
Bar plot of  
applicant\_race  
for rejected loan app



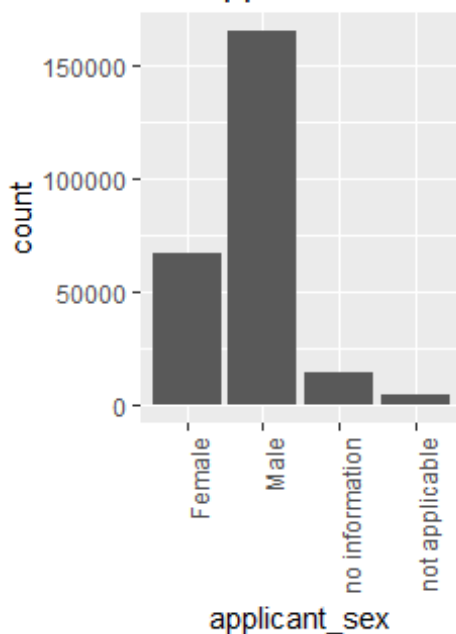
Bar plot of  
applicant\_race  
for approved loan app

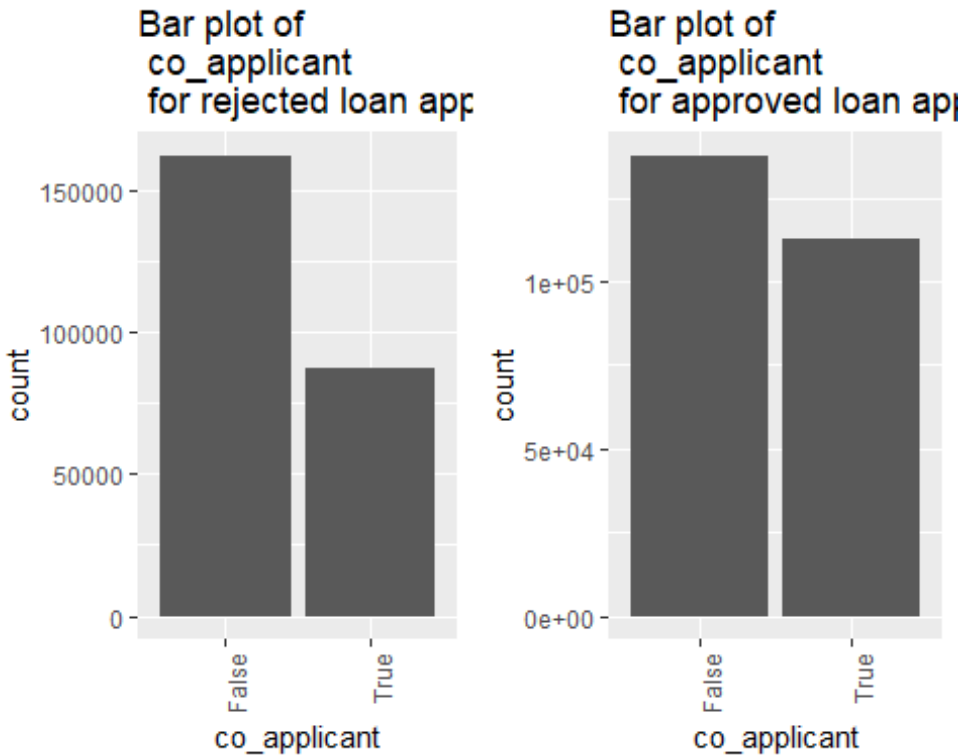


Bar plot of  
applicant\_sex  
for rejected loan app



Bar plot of  
applicant\_sex  
for approved loan app





We observe the following:

1. Some features such as property\_type, loan\_purpose, applicant\_ethnicity, applicant\_race and co\_applicant have an almost significant different distribution of categories between the label categories.
2. Others features such as preapproval, applicant\_sex, occupancy show small differences, but these differences are unlikely to be significant.

Next we do feature engineering in readiness for modelling.

## Feature Engineering

Data preparation is an important step. We recoded loan\_amount and applicant\_income to categorical strings which helped in distinguishing levels within the dataset and imputed missing numeric variables with means.

Feature engineering and modelling were done in python

## Model Building

Modelling was done with python, selecting the CatBoostclassifier library for the classification. As earlier indicated we settled on 8 cross-validations with 1000 iterations to minimize bias and variance of the model.

## Results and Conclusions

The model was able to correctly classify 72% of the accepted values with AUC of 0.8. There were discrepancies with the supplied test values indicating that our model was overfitting. However, it offers a good starting point to build on and improve.