
Home Credit Default Risk

Group 03

Zhengyuan Ding , Chenqin Yang

Jiayao Liu, Ziyu Lei, Zian Chen



Agenda

- Background
- Dataset Introduction
- Problem Statement
- Methodology
- Performance Improvements
- Schedule

Background

Home Credit tries to provide loans to underprivileged people who do not have sufficient credit histories. They used a variety of variables including telco and transactional information to predict the repayment ability of their clients.

The logo for Home Credit, featuring the words "HOME" and "CREDIT" in a bold, red, sans-serif font, stacked vertically. The letter "O" in "HOME" is stylized with a white circle inside.

This project will use the statistical and machine learning methods to help Home Credit fully unlock the potential of the data it collects from the clients.

Problem Statement:



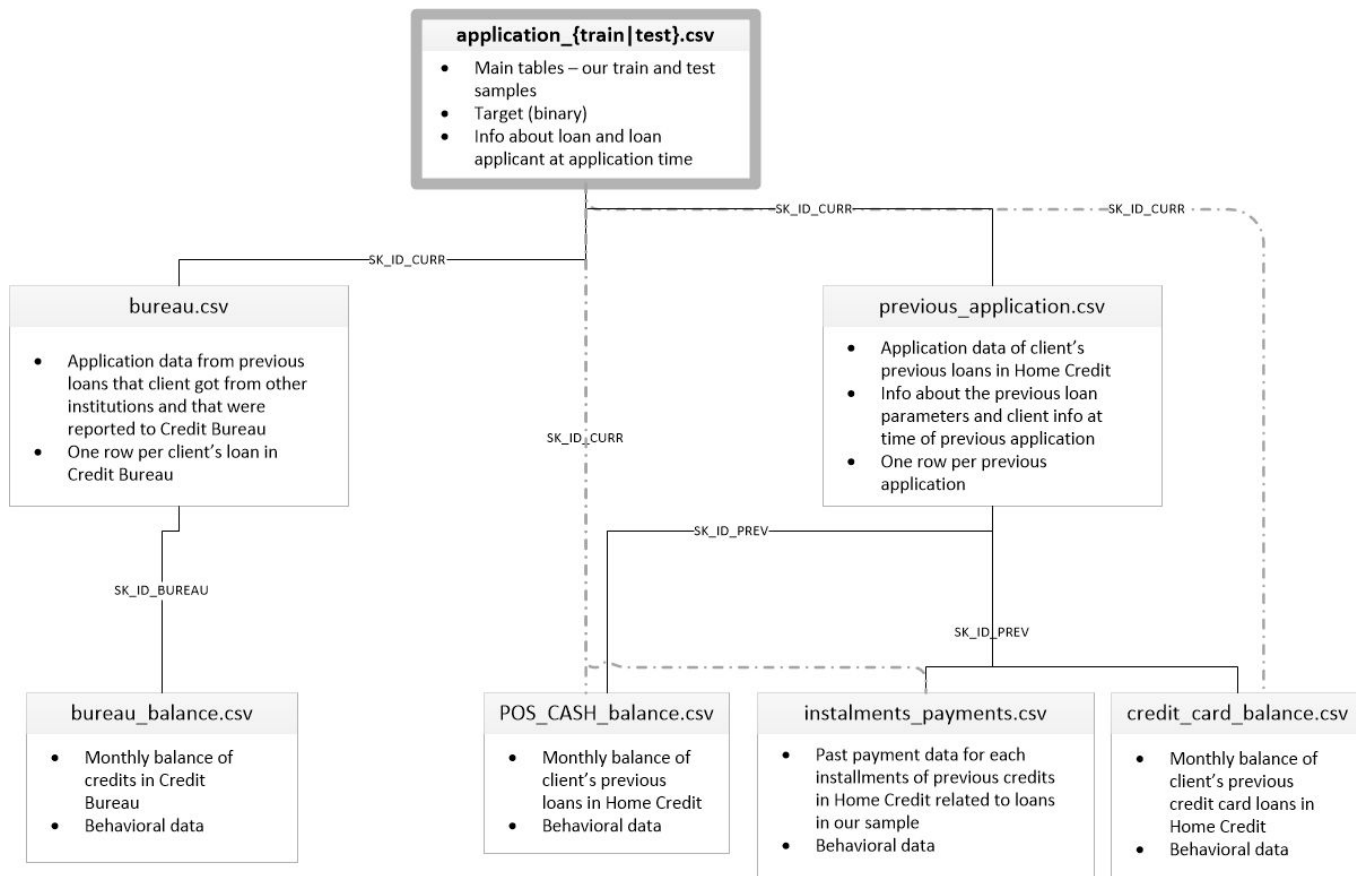
Can you predict how capable each applicant is of repaying a loan?

Dataset

- Application.csv
 - Main table
 - One row represents one loan in our data sample.
- Bureau.csv
 - All client's previous credits provided by other financial institutions that were reported to Credit Bureau (for clients who have a loan in our sample)
- Bureau_balance.csv
 - Monthly balances of previous credits in Credit Bureau
- POS_CASH_balance.csv
 - Monthly balance snapshots of previous POS (point of sales) and cash loans that the applicant had with Home Credit

Dataset (Cont'd)

- Credit_card_balance.csv
 - Monthly balance snapshots of previous credit cards that the applicant has with Home Credit
- Previous_application.csv
 - All previous applications for Home Credit loans of clients who have loans in our sample
- Installments_payments.csv
 - Repayment history for the previously disbursed credits in Home Credit related to the loans in our sample



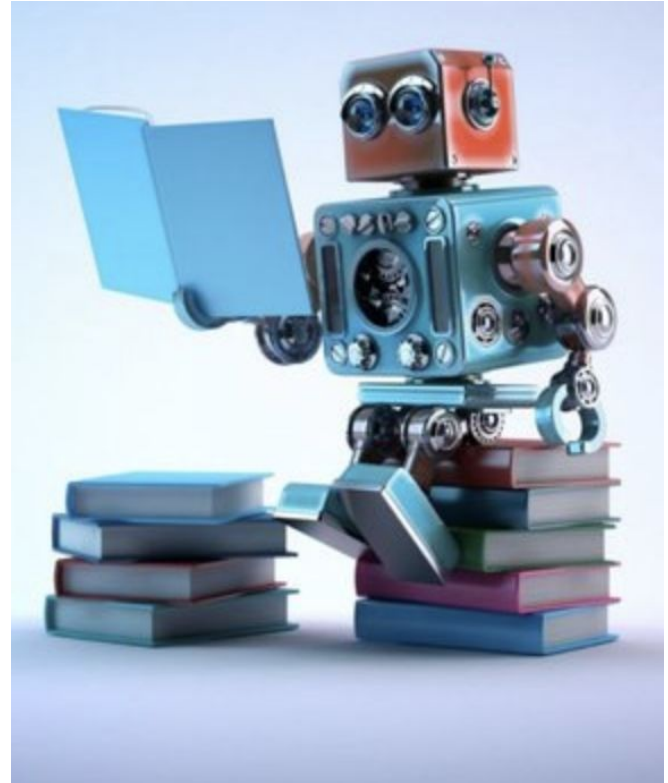
Methodology

Binary Classification Models

- Baseline model: Naive Bayes
- Logistic Regression
- SVM
- Random Forest
- Gradient Boosting Method

Evaluation Metrics

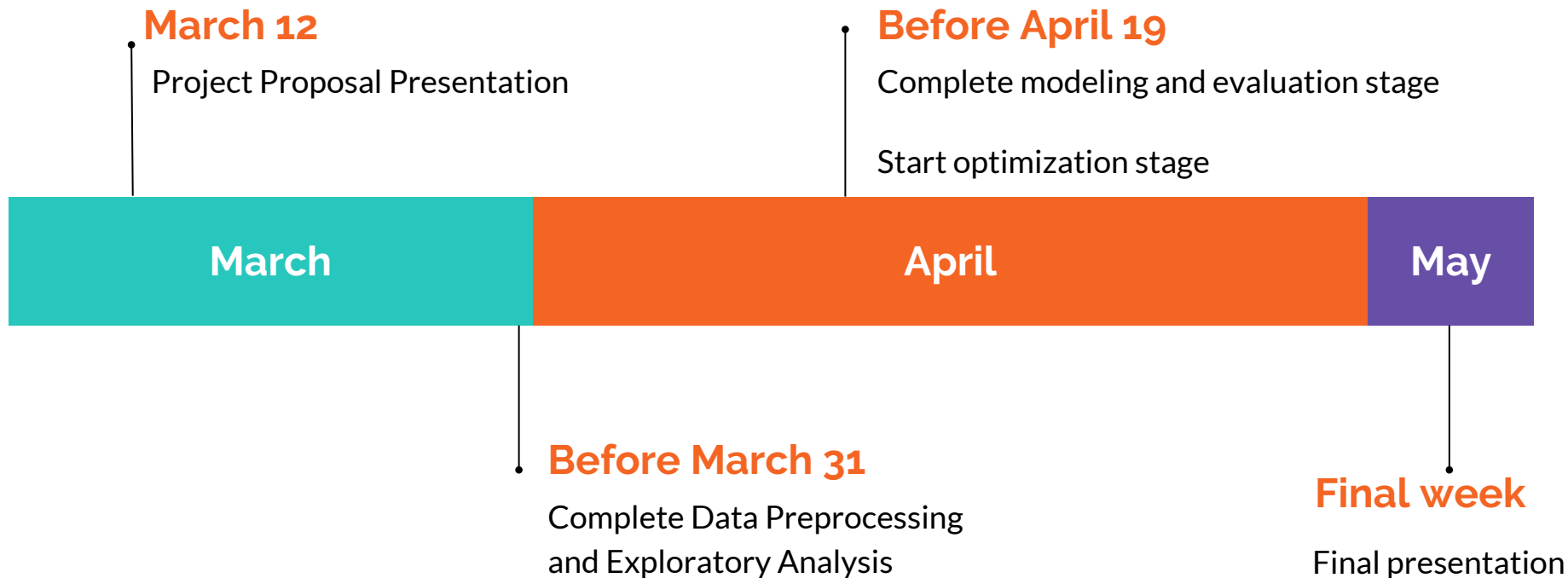
- AUC



Possible Performance Improvements

- Cython – Operations between columns (groupby count, max, min, mean)
- Check code sanity by following Python performance tips
 - Itertools, function call overhead, built-in tools.....
- Do Python performance tuning by using cProfile and line profiler
- Numba: JIT Compiling, Vectorize decorator
- PySpark

Schedule





Thank you!