

Predicting loan default:

Combining detailed customer data with machine learning models to better predict default

The problem

Loan defaults represent a large risks for banks

If they can better predict who will default, they can better price their loans

Can they use better data to help them optimize their operations?

US consumer loan default rate

4.62%



2.91%

2.01%

2.29%

2009

2012

2015

2018

Data Sources and Tools

Loan Data provided by Home Credit Bank.

Competition - the winning prize of \$70K!

Can I get close to the winning score?!

kaggle

HOME
CREDIT

aws

PostgreSQL



Application data:
Static data for all applications
Info about loan and loan applicant

Credit Bureau:
client's previous credits provided by
other financial institutions

Credit Bureau balance:
Monthly balances of previous
credits in Credit Bureau

Previous Applications:
All previous applications for Home
Credit loans of clients who have
loans in these data set

Cash loan balance:
Monthly balance
snapshots of previous
POS (point of sales) and
cash loans

Credit card balance:
Monthly balance snapshots
of previous credit cards
that the applicant has with
Home Credit

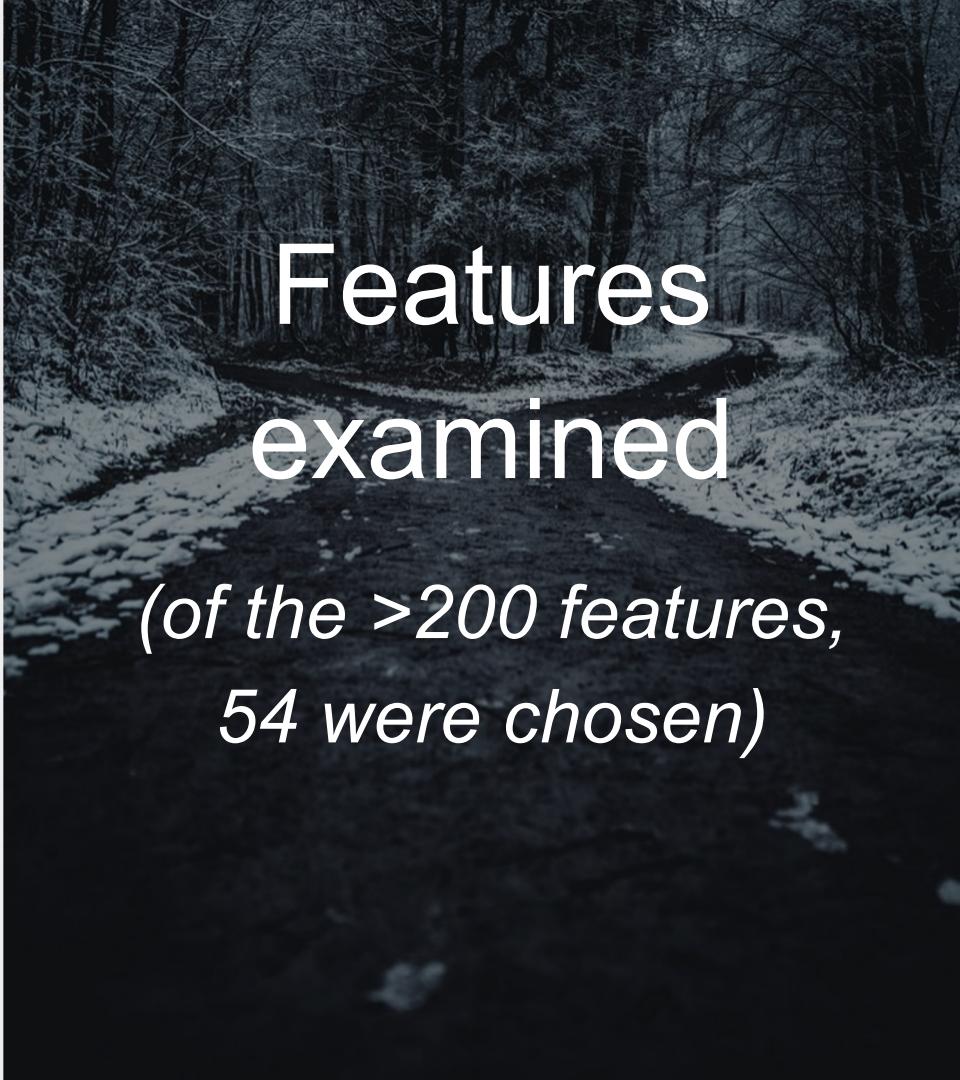
**Installments
payments:**
Repayment history for the
previously disbursed
credits in Home Credit

Some of the selected features include:

- Client's monthly income;
- Monthly payments;
- Whether client owns a real estate;
- Whether client could be reached via phone;
- Education level;
- Payment history, if any.

Some created features include:

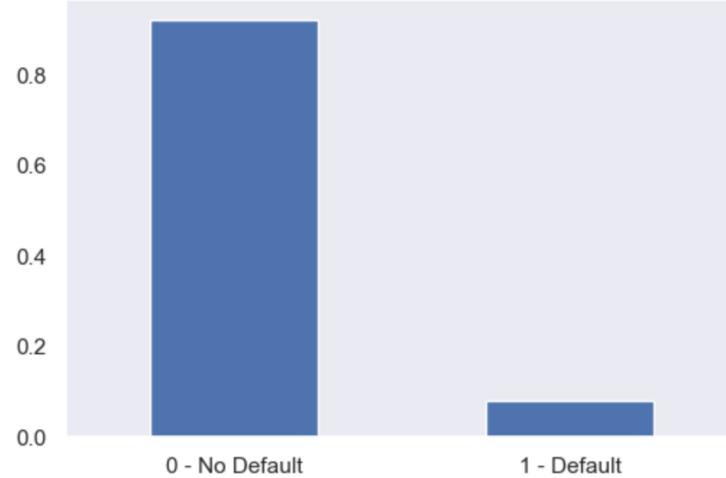
- Income to Debt ratio;
- Average number of applications that client had previously;
- Rejection rate if client applied for the loan before;
- Average number of days that client had past due



Features examined

*(of the >200 features,
54 were chosen)*

Target value distribution

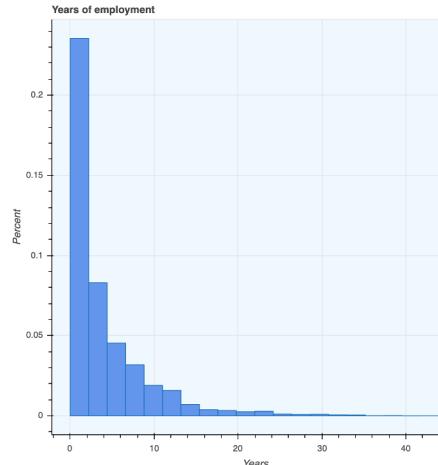
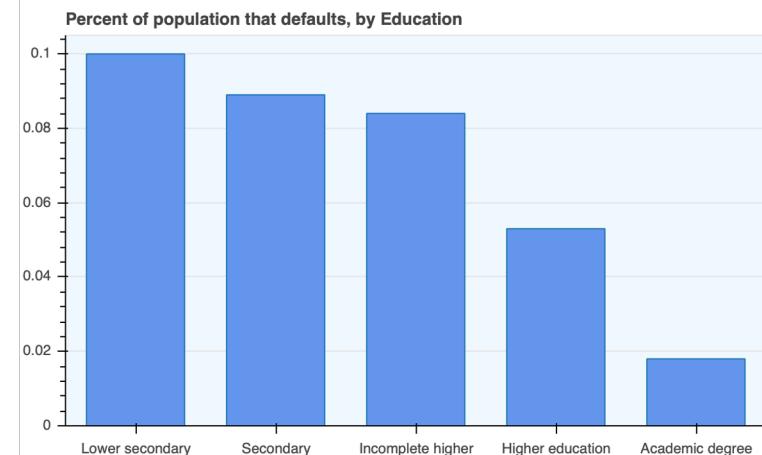


Data is imbalanced

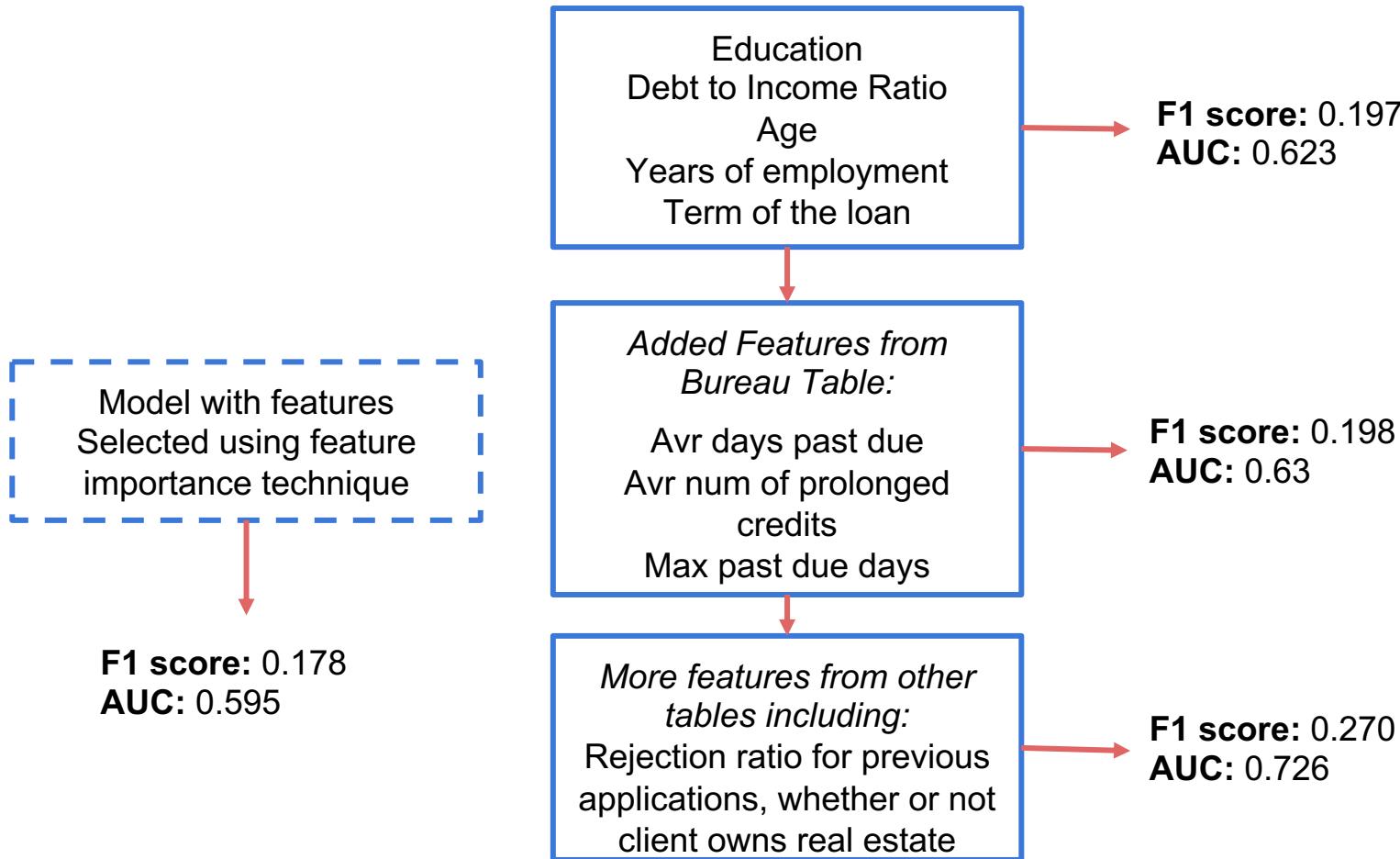
Metrics for analysis:

F1 score and AUC

Will be using soft classification

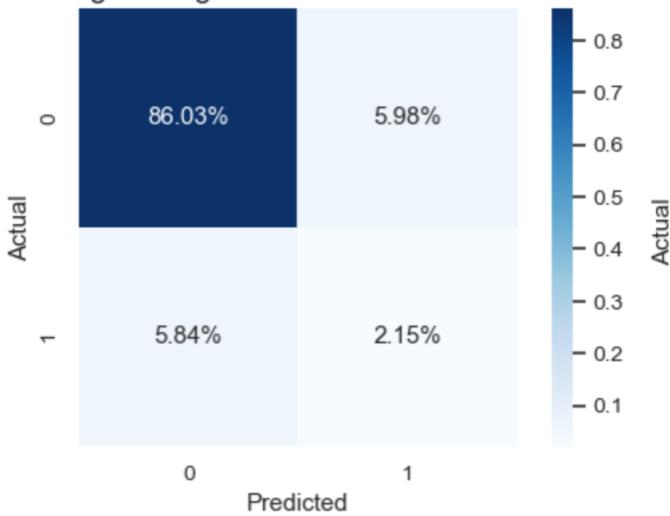


Base model: Logistic regression

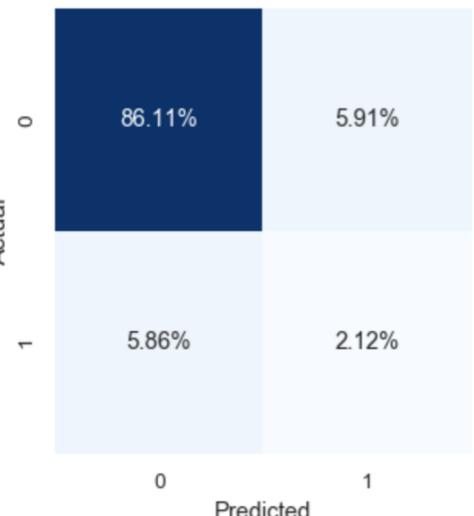


	F1 Score	AUC
Logistic regression	0.270	0.726
Random Forest	0.286	0.737
XGBoost	0.286	0.74

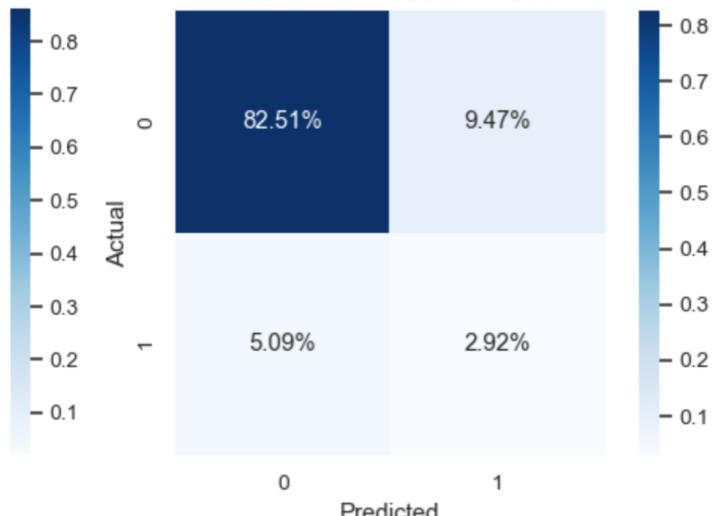
Logistic regression confusion matrix



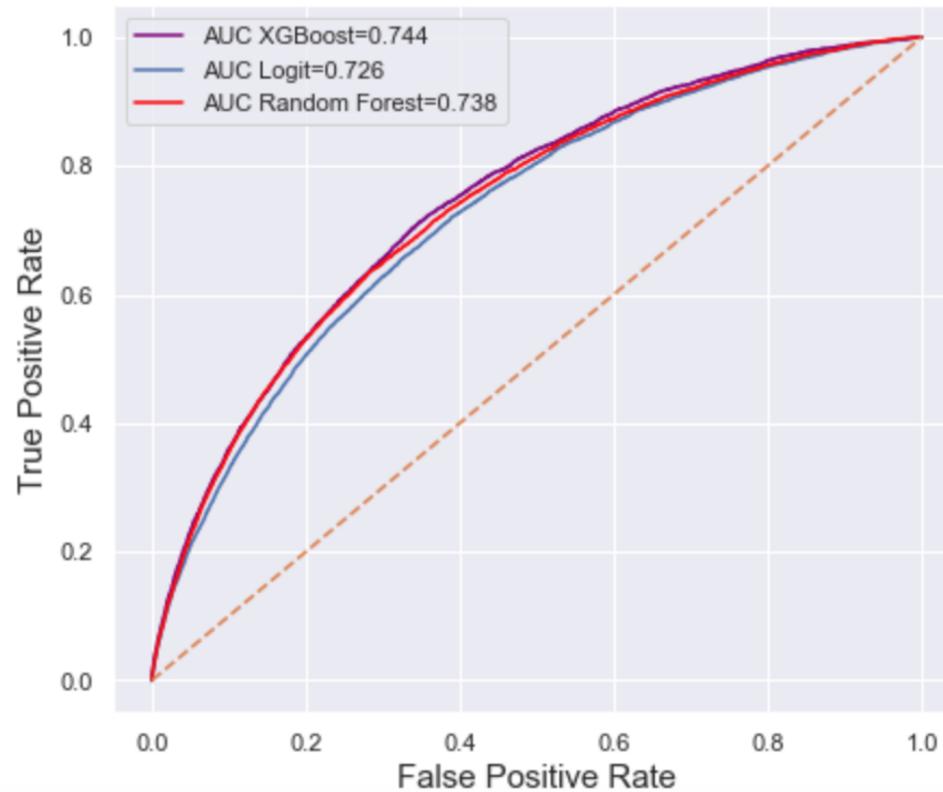
Random Forest confusion matrix



XGBoost confusion matrix



ROC curve



A photograph of a clear glass filled with a variety of coins, including pennies, nickels, dimes, quarters, and some foreign currency. A small green sprout with three leaves grows from a thin wooden stick standing upright in the center of the glass. The background is a dark, textured surface.

Results and improvements:

1. The winning score for the competition was 0.8. My model had a score of 0.74, which is a good start but also leaves room for improvement
2. To improve my model I would go again through my data and use more feature engineering.