

Lending Club Dataset

Judah Drellich

Target Variable

- ❖ Separate the features (X) from the target variable (y).
 - ❖ `loan_status`: supposed to be the target variable.
 - ❖ Too many of the features have information that was taken after the loans originated
 - ❖ `next_pymnt_d` records the date that the next loan payment is due
 - ❖ Shift the time of the question to the time that the dataset was taken.
 - ❖ Find the value of the remaining loans
 - ❖ Target variable is the fraction of the loan amount that the borrower will be able to pay back

Data Wrangling

- All the features must have only numerical values
 - Some columns are numerical but need to be cleaned (Term: ‘ 36 months’)
 - Others are Categorical such as the type of employment that the borrower has.
- Imputation of values for this dataset is the most challenging aspect of this project.
 - Assumed that the event didn’t happen if the data was missing and assigned values accordingly.
 - For counting variables such as the number of accounts, I filled in with a zero
 - For Months since an event happened, I filled in a large number so that it the models would act like it never happened.

Data Wrangling

Multicollinearity

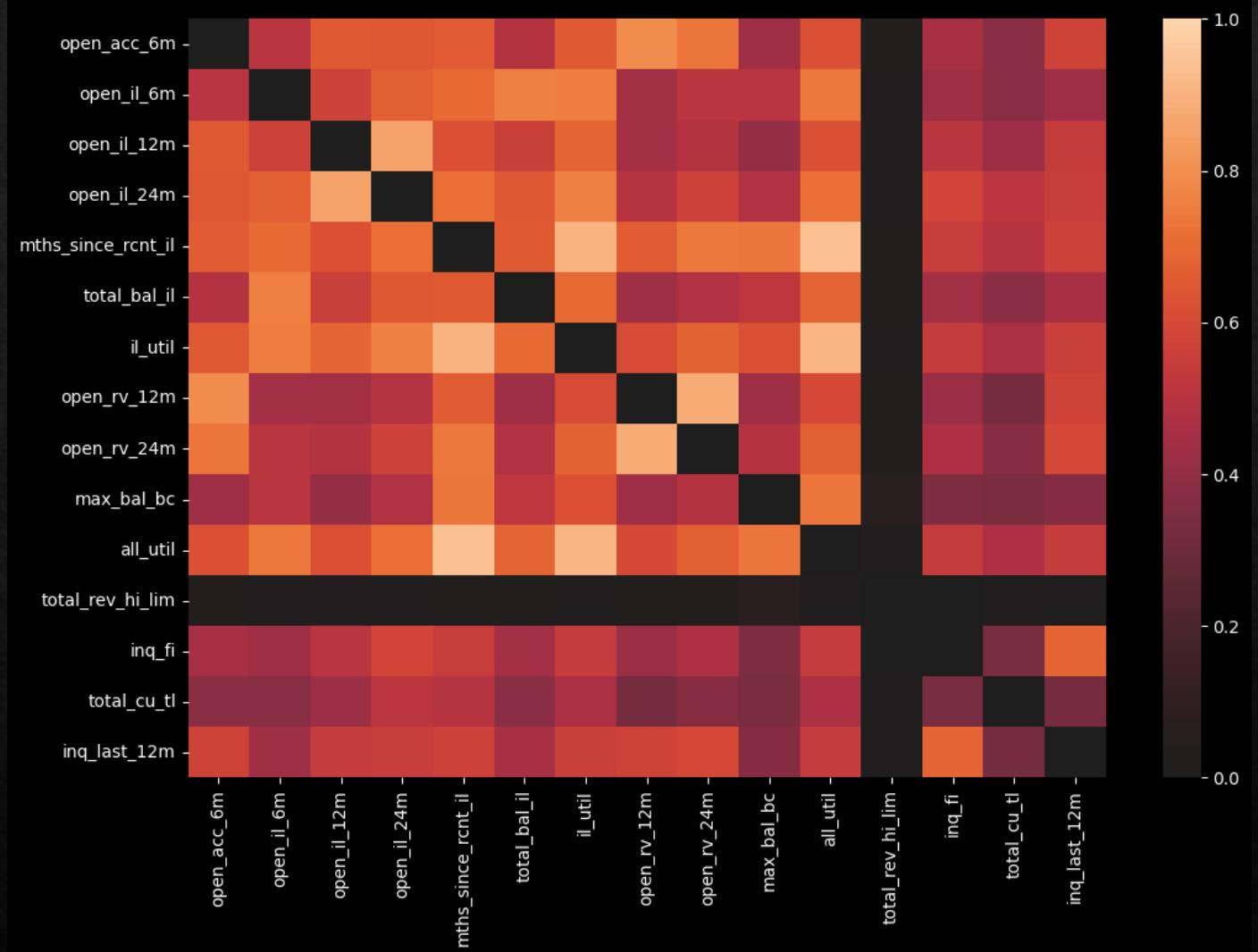
Initial heatmap for the correlations between the features



Data Wrangling

Multicollinearity

Focusing on the bright square of the previous heatmap



Data Wrangling

- Lasso Regularization is a technique for variable selection that uses linear regression to evaluate the effect that each features has on a target variable.

$$Loss(\beta_1, \dots, \beta_n) = SSD + \alpha \sum_{i=1}^n |\beta_i|$$

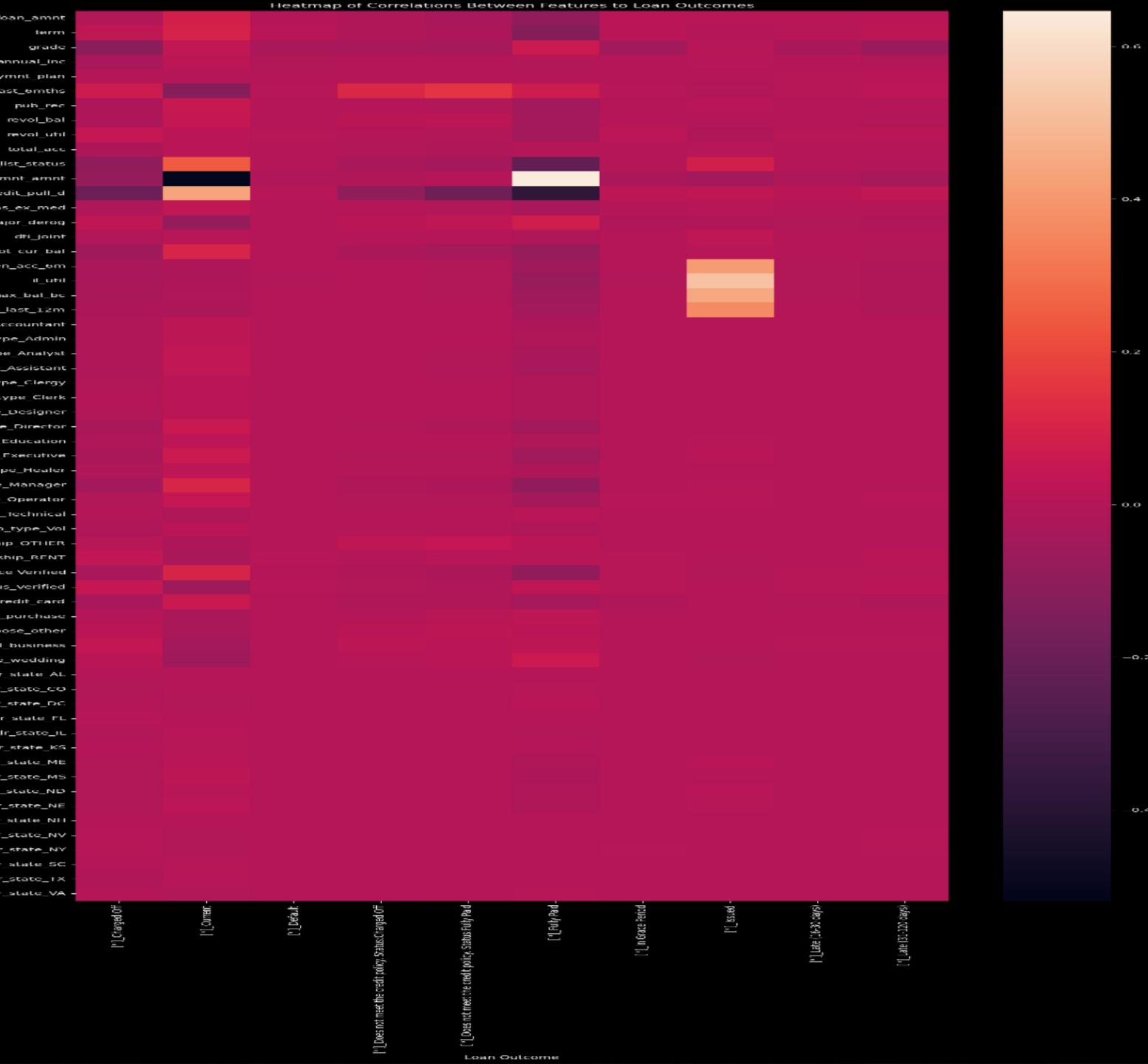
- I used Lasso Regularization to create two datasets
 - Small: 12 features, binary intrusion or not
 - Big: 61 features, multi class model, determines intrusion types

Lasso
Regularization

Feature
Selection

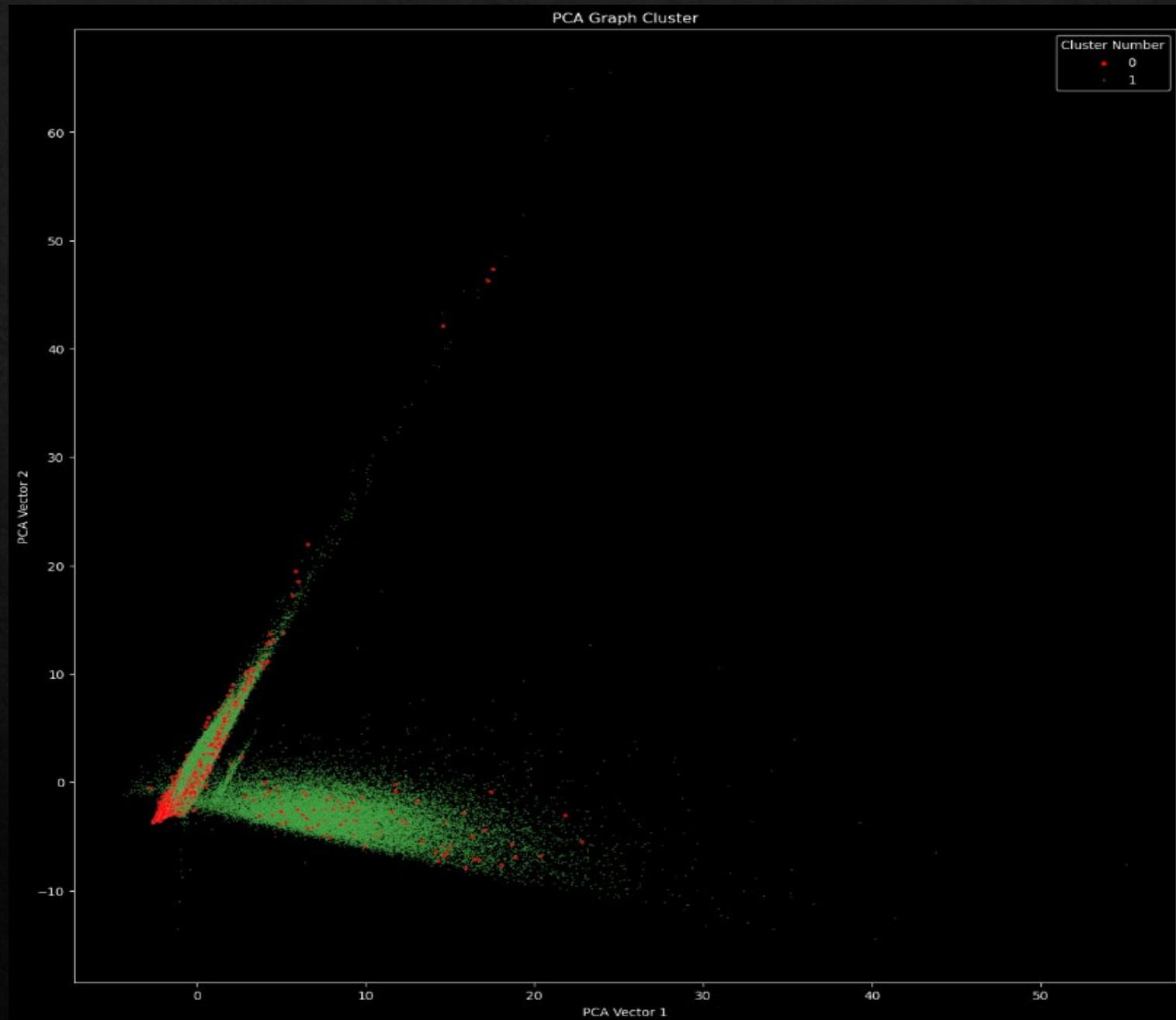
EDA

Correlation Heatmap
between features and
loan statuses

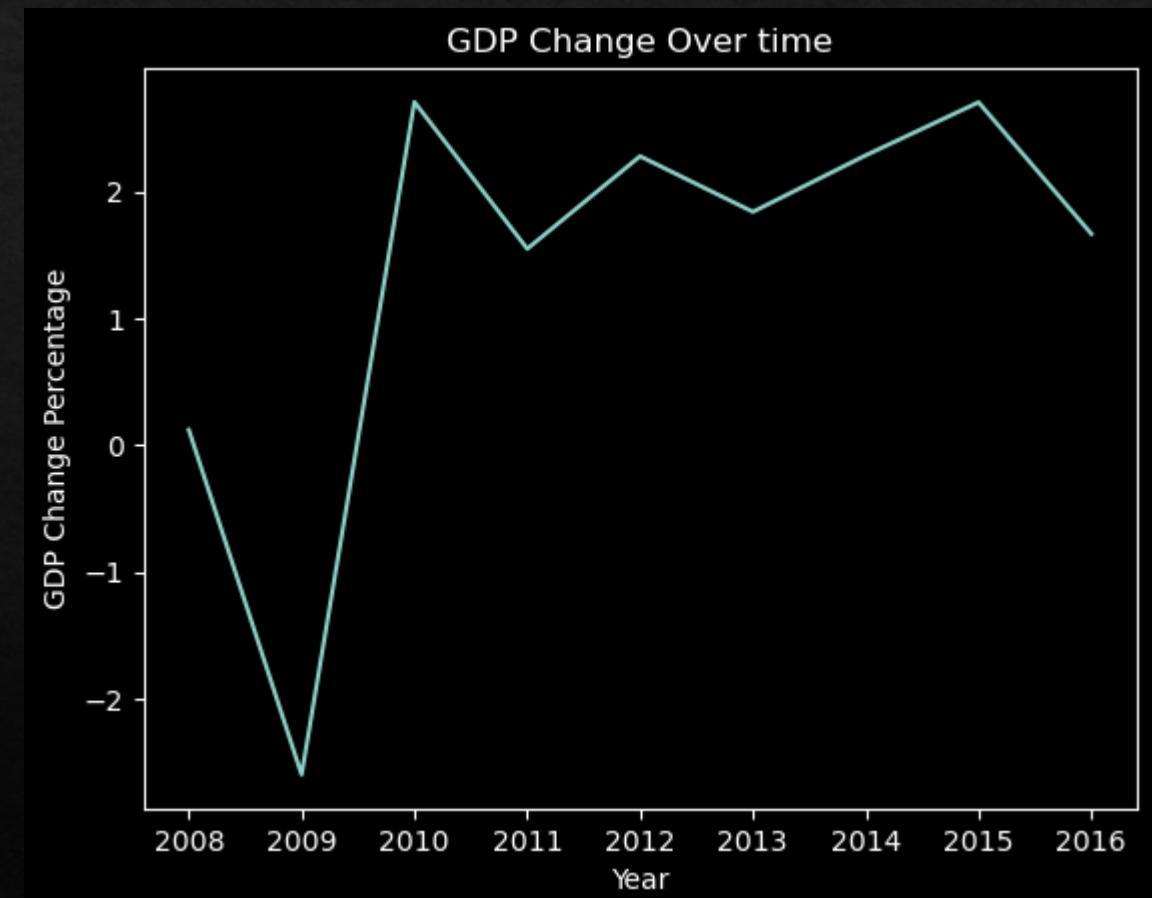
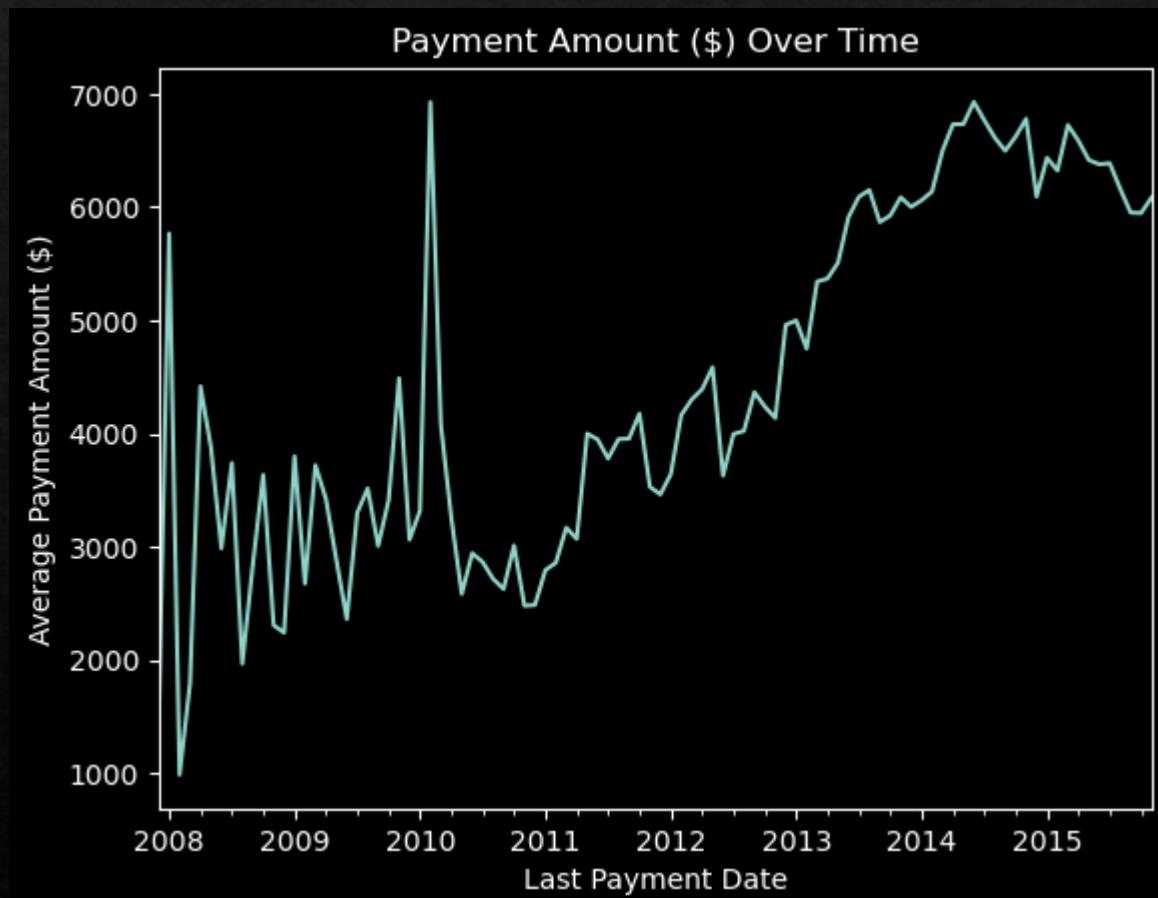


EDA

*PCA graph of the features
with bad outcomes colored
in red and good outcomes
colored in green*



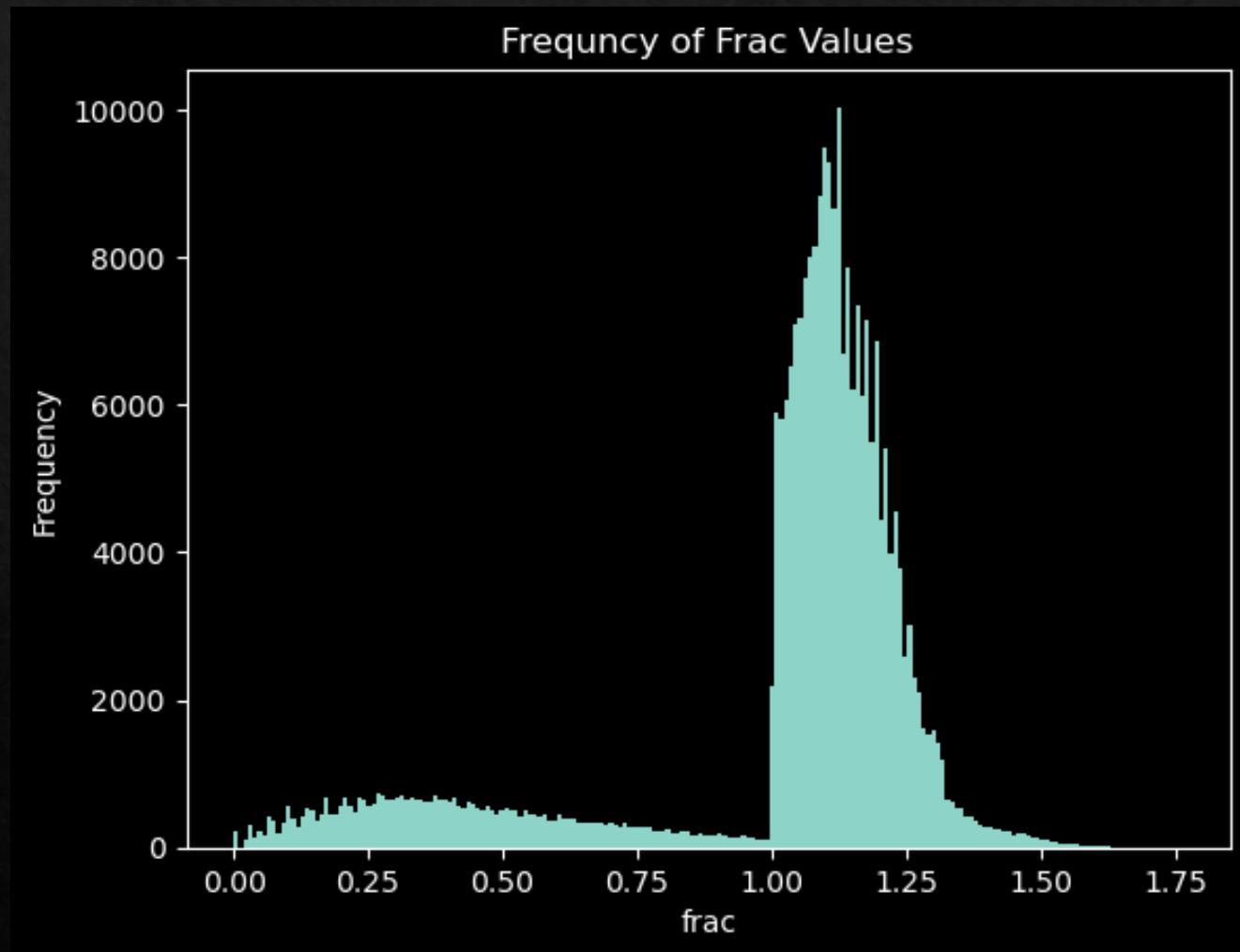
External Factors



Preprocessing

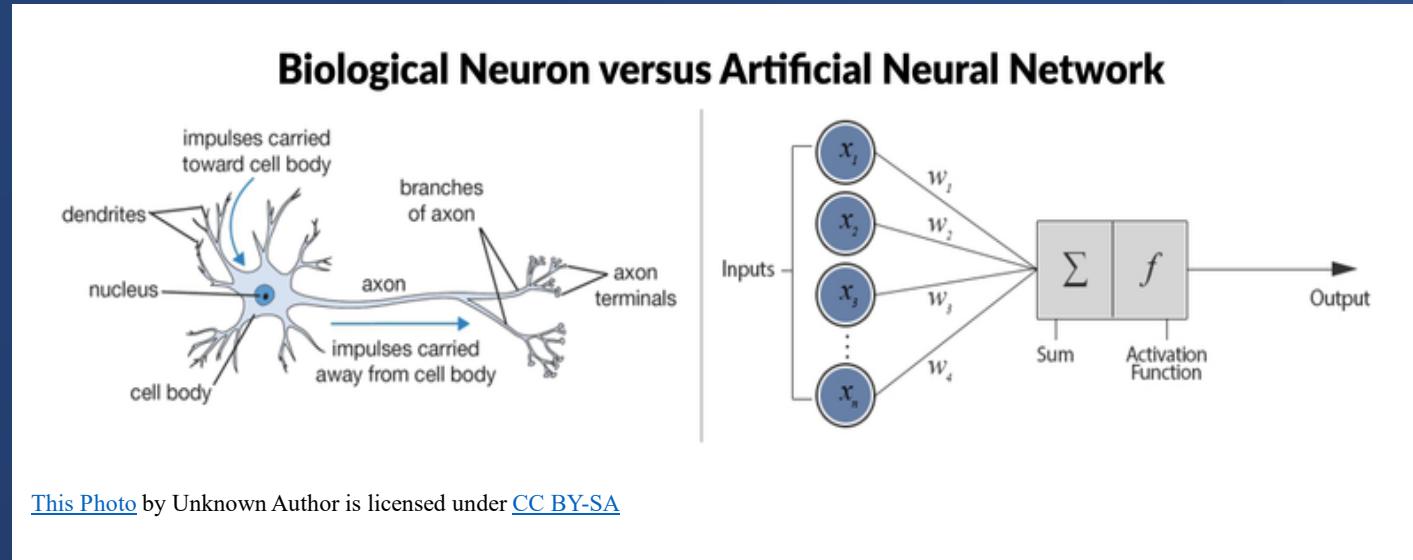
- ❖ $Expected\ Payment = Principal * (1 + Interest\ Rate) * Years$
 - ❖ Where *Interest Rate* is in decimal form
- ❖ $Actual\ Payment = Total\ Payment + Total\ Late\ Fees + Recoveries$
- ❖ $frac = \frac{Actual\ Payment}{Expected\ Payment}$
- ❖ $Expected\ Payment = Principal * (1 + Interest\ Rate) * Years - Actual\ Payment$

EDA - frac





Modelling



$$Output = f(\sum_1^n w_i x_i) + b$$

Neural Networks

Diagram of the similarity between biological neuron and a node in a neural network

Formula for a node in a neural network

Metrics

- ◆ MAE

$$\diamond MAE = \frac{\sum_{i=1}^n |error_i|}{n}$$

- ◆ MSE

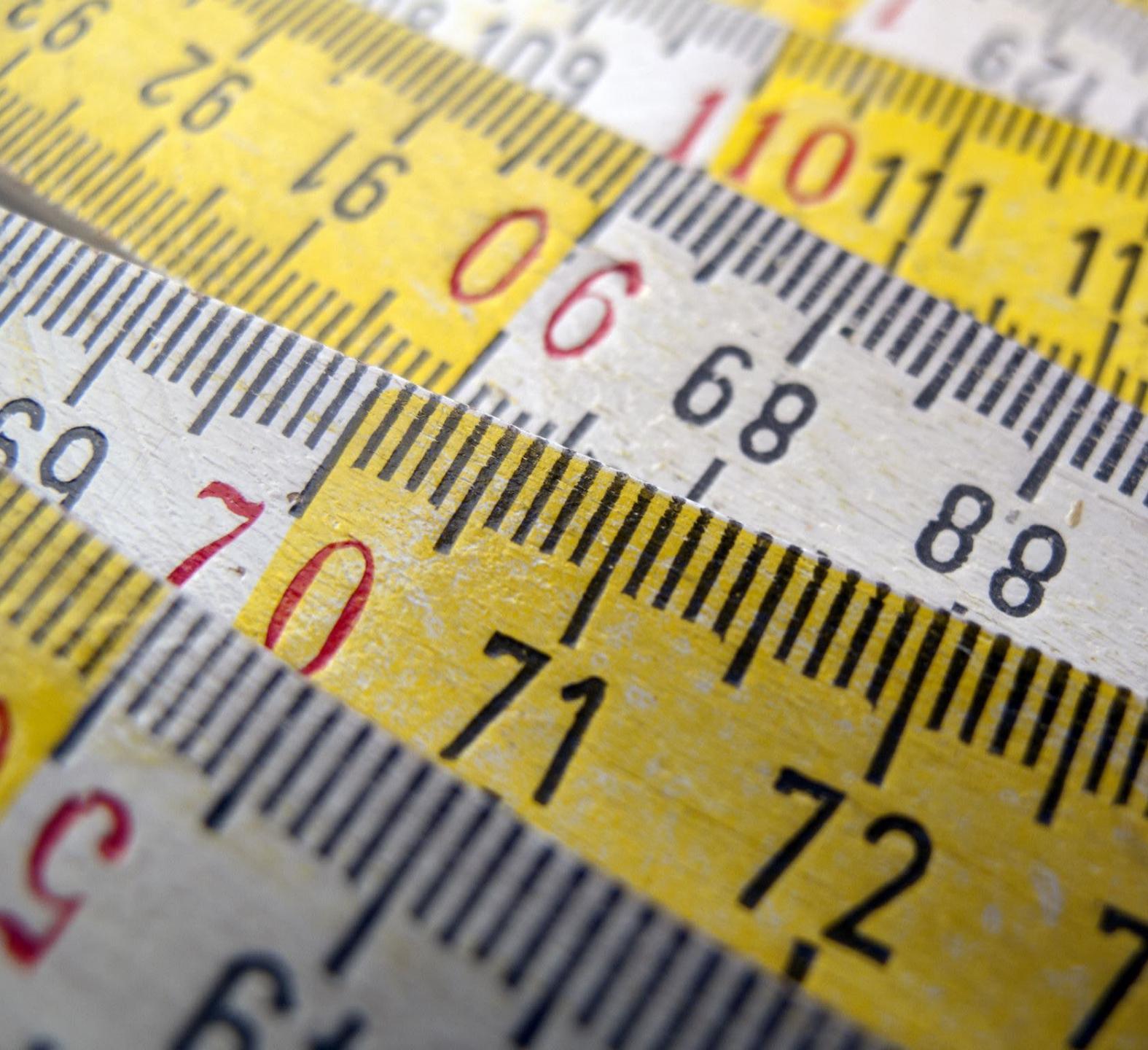
$$\diamond MSE = \frac{\sum_{i=1}^n error_i^2}{n}$$

- ◆ RMSE

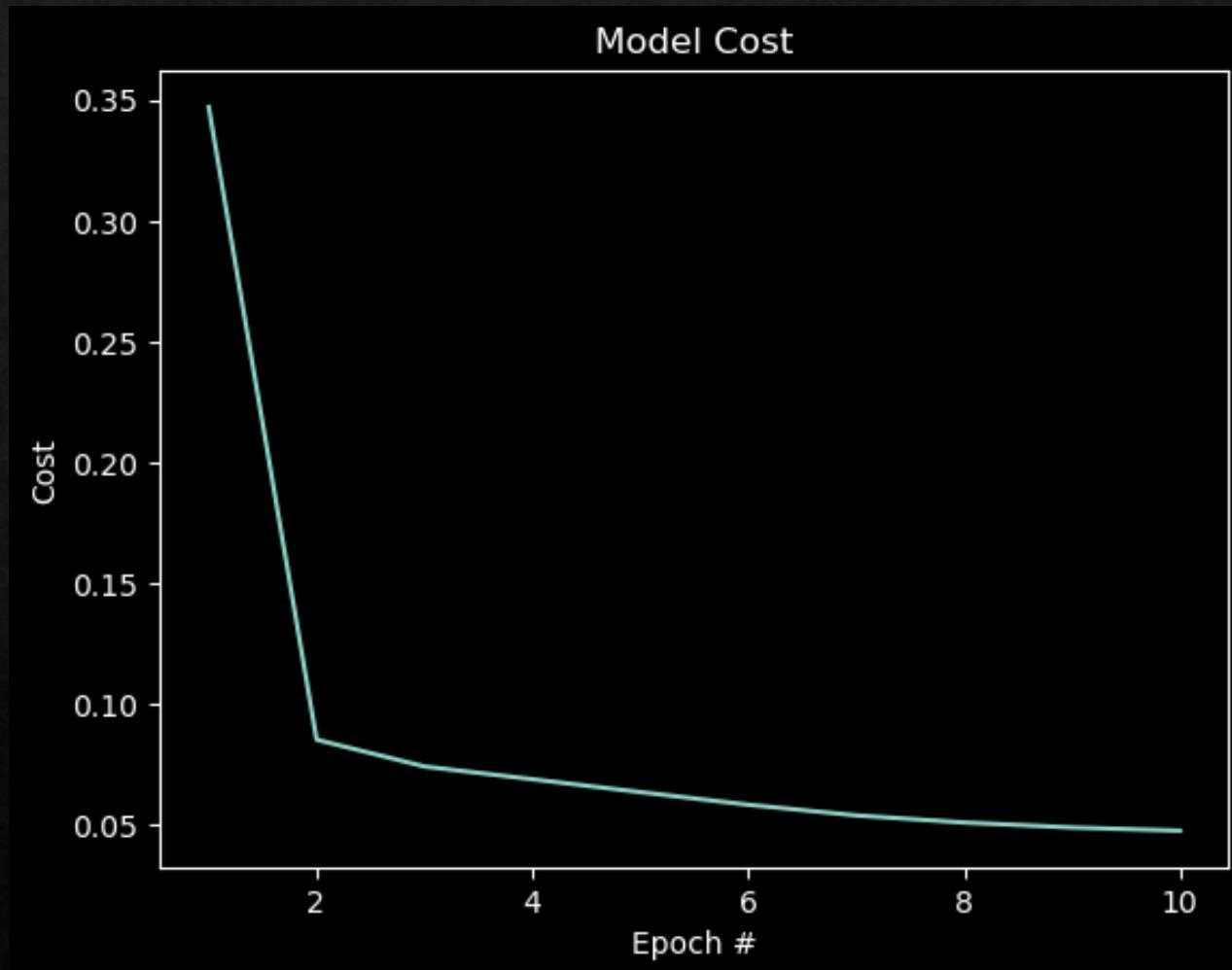
$$\diamond RMSE = \sqrt[2]{\frac{\sum_{i=1}^n error_i^2}{n}}$$

- ◆ Percent Error

$$\diamond Percent\ Error = \frac{\sum_{i=1}^n \frac{|error_i|}{pred_i}}{n}$$



Training



Testing (Neural Network)

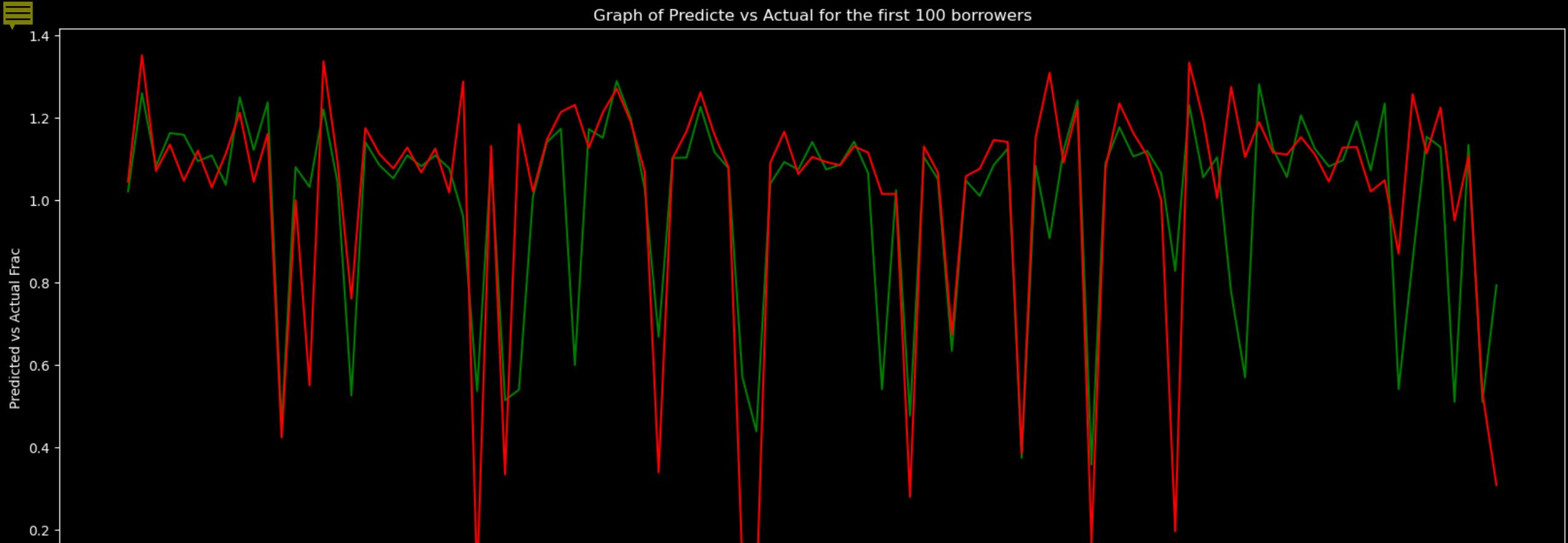
MSE = 0.0490,

RMSE = 0.2213,

MAE = 0.1386,

Mean Percent Error of 14.6275%.

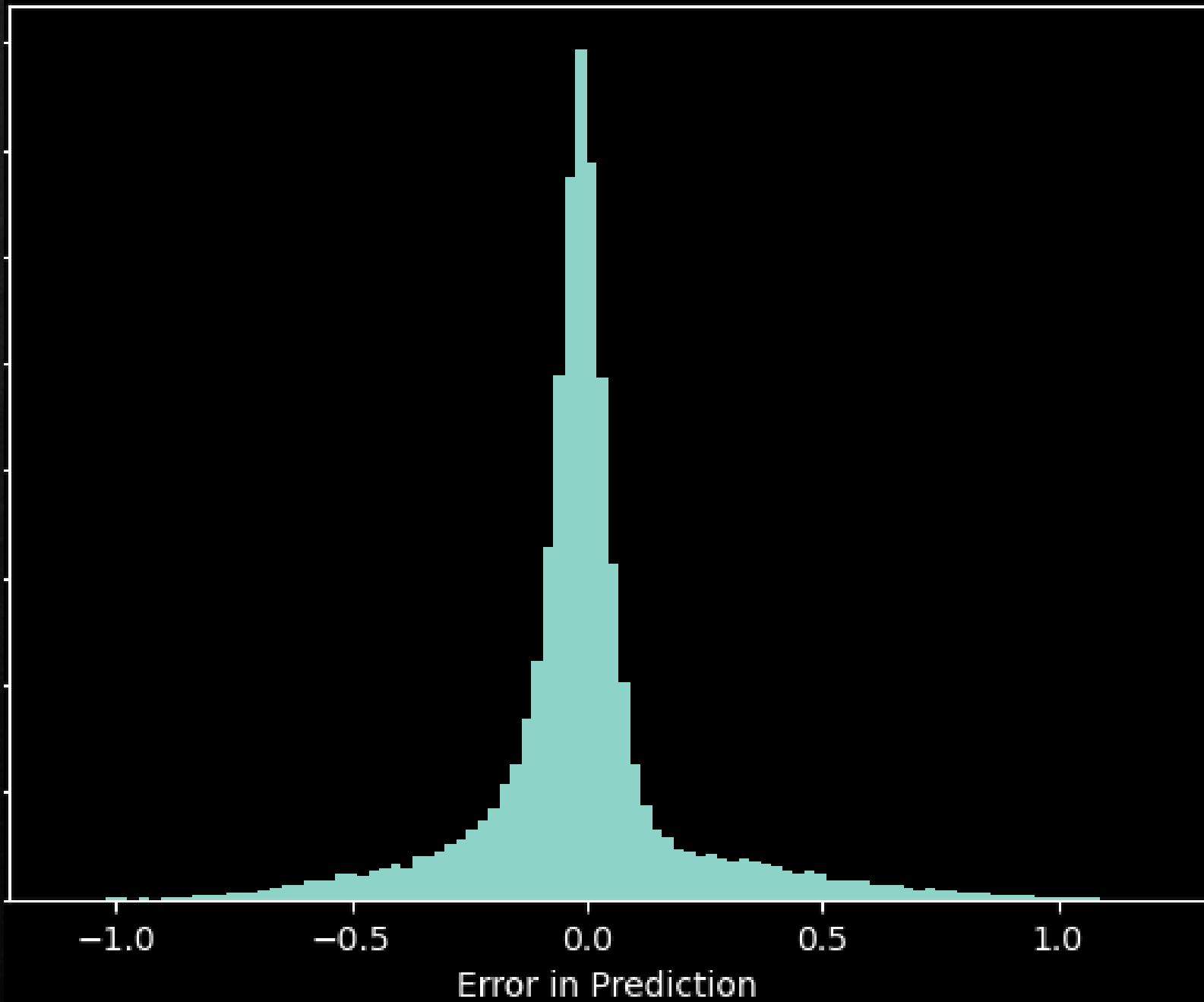
Graph of Predicted vs Actual for the first 100 borrowers



Testing NN cont.

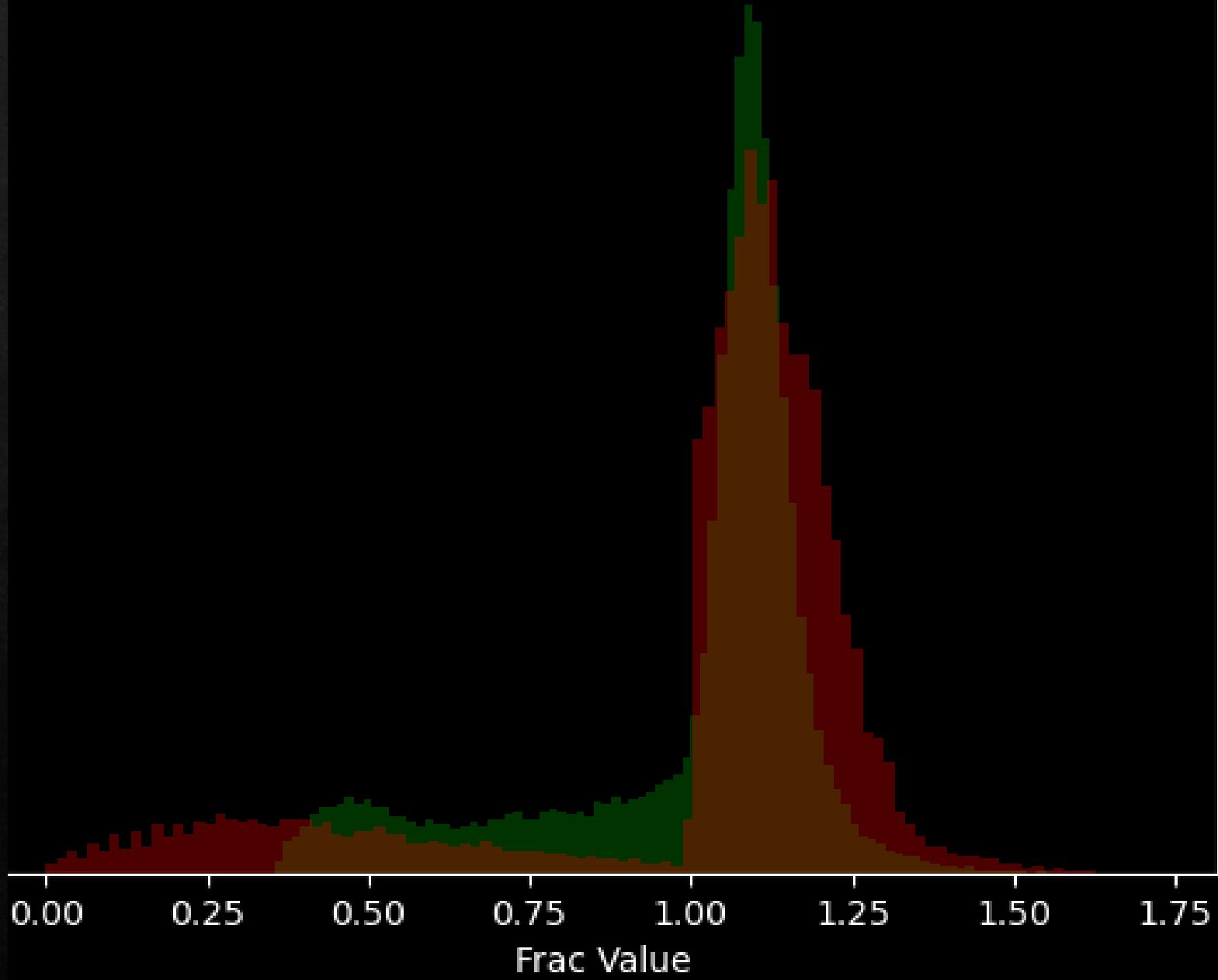
*Testing NN
cont.*

Neural Net Distribution of Error



*Testing NN
cont.*

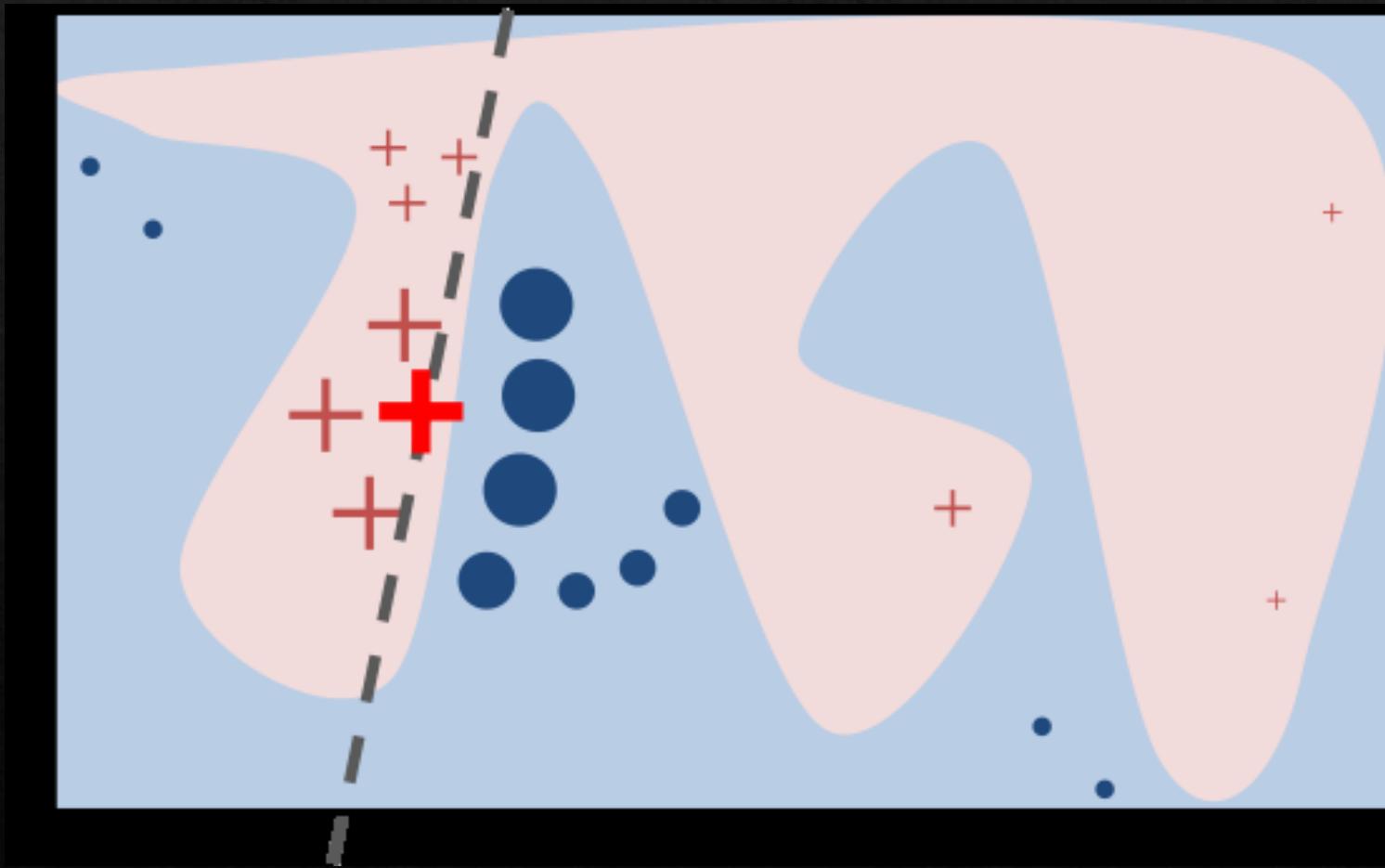
Histogram of the Neural Net Predictions vs. Y-Test



LIME

- ❖ LIME stands for Locally Interpretable Model Explanation and gives simpler explanations to how complex machine learning models arrive at a prediction
- ❖ $\xi(x) = \underset{g \in G}{\operatorname{argmin}}(\mathcal{L}(f, g, \pi_x) + \Omega(g))$
- ❖ I created a function that could take in the number of any loan and create a LIME explanation for that loan

LIME cont.

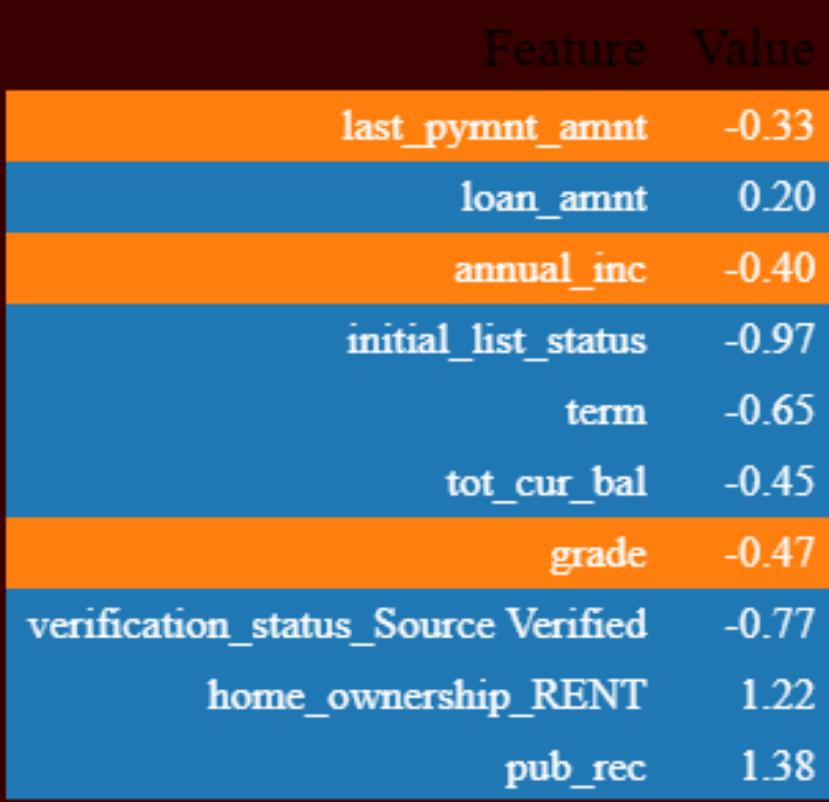
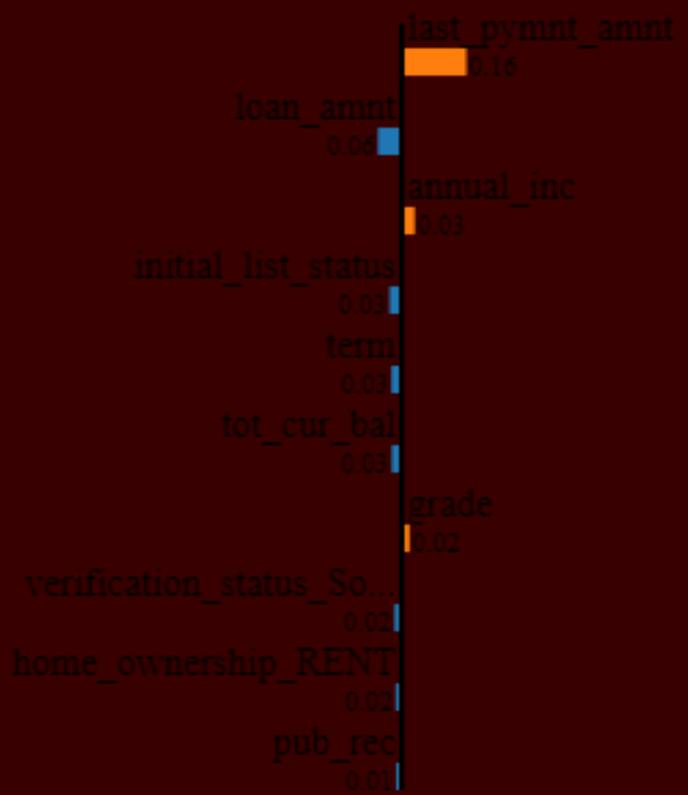


Predicted value

0.45
min) 0.66
(max) 1.85

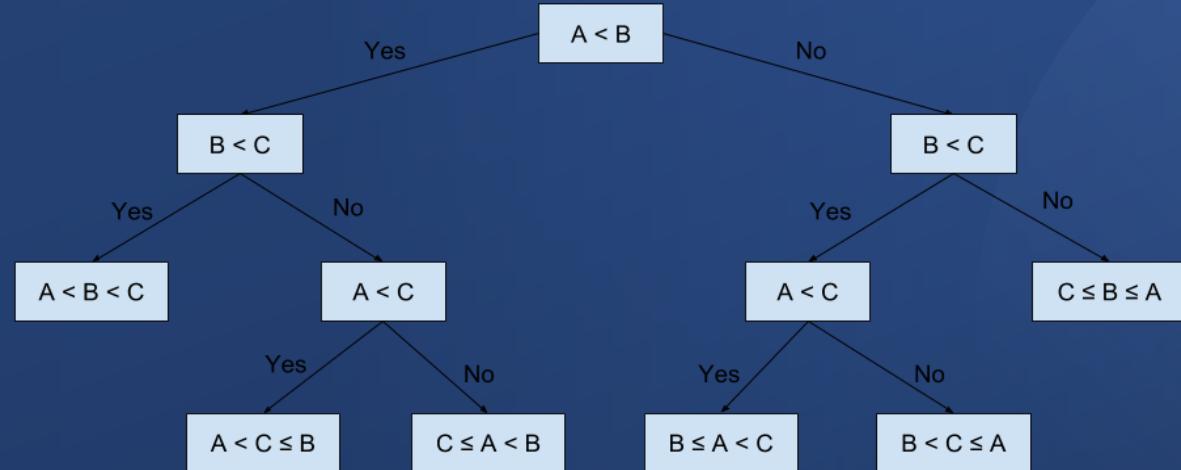
negative

positive



LIME cont.

Modelling



$$Gini = 1 - \sum_{i=1}^n p_i^2$$

Random Forest

Model of a single decision tree in a random forest

Purity measurement formula

Testing (Random Forest)

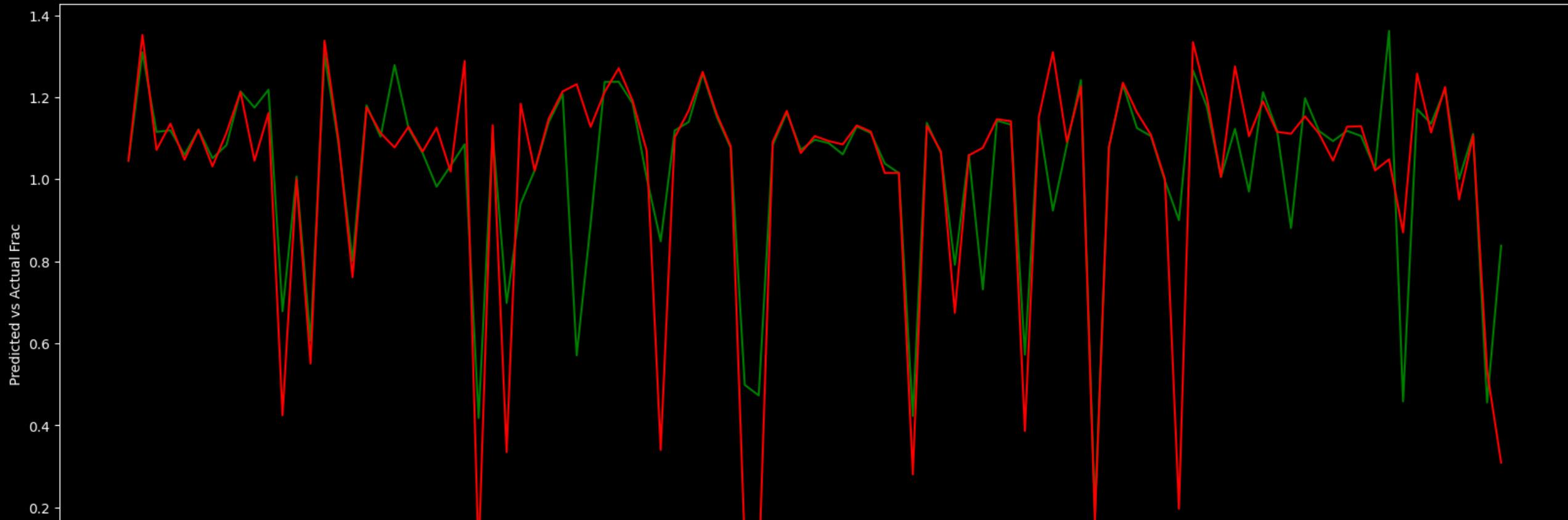
MSE = 0.0309,

RMSE = 0.1758,

MAE = 0.0872

Mean Percent Error of 13.6849%.

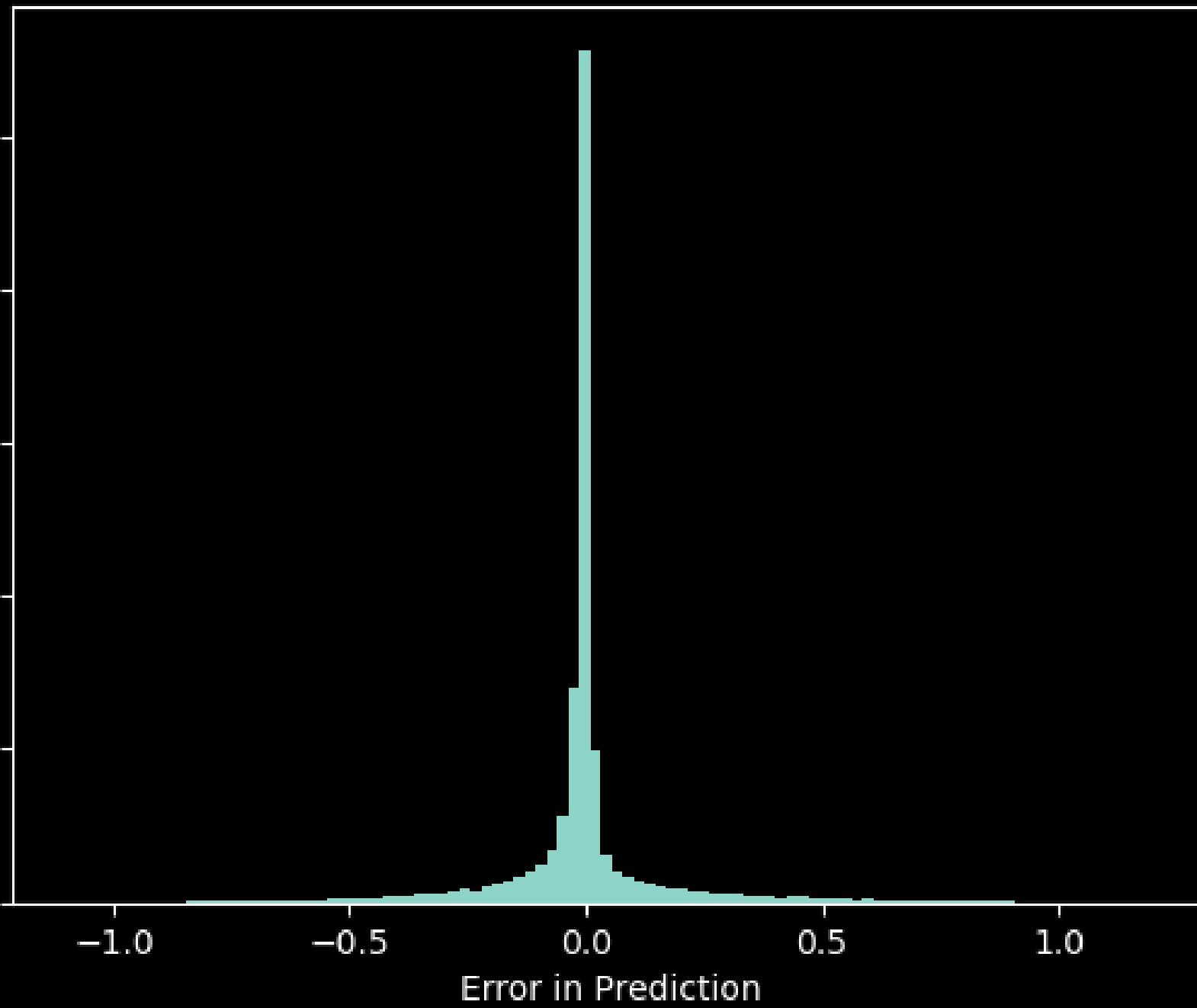
Graph of Predicted vs Actual for the first 100 borrowers



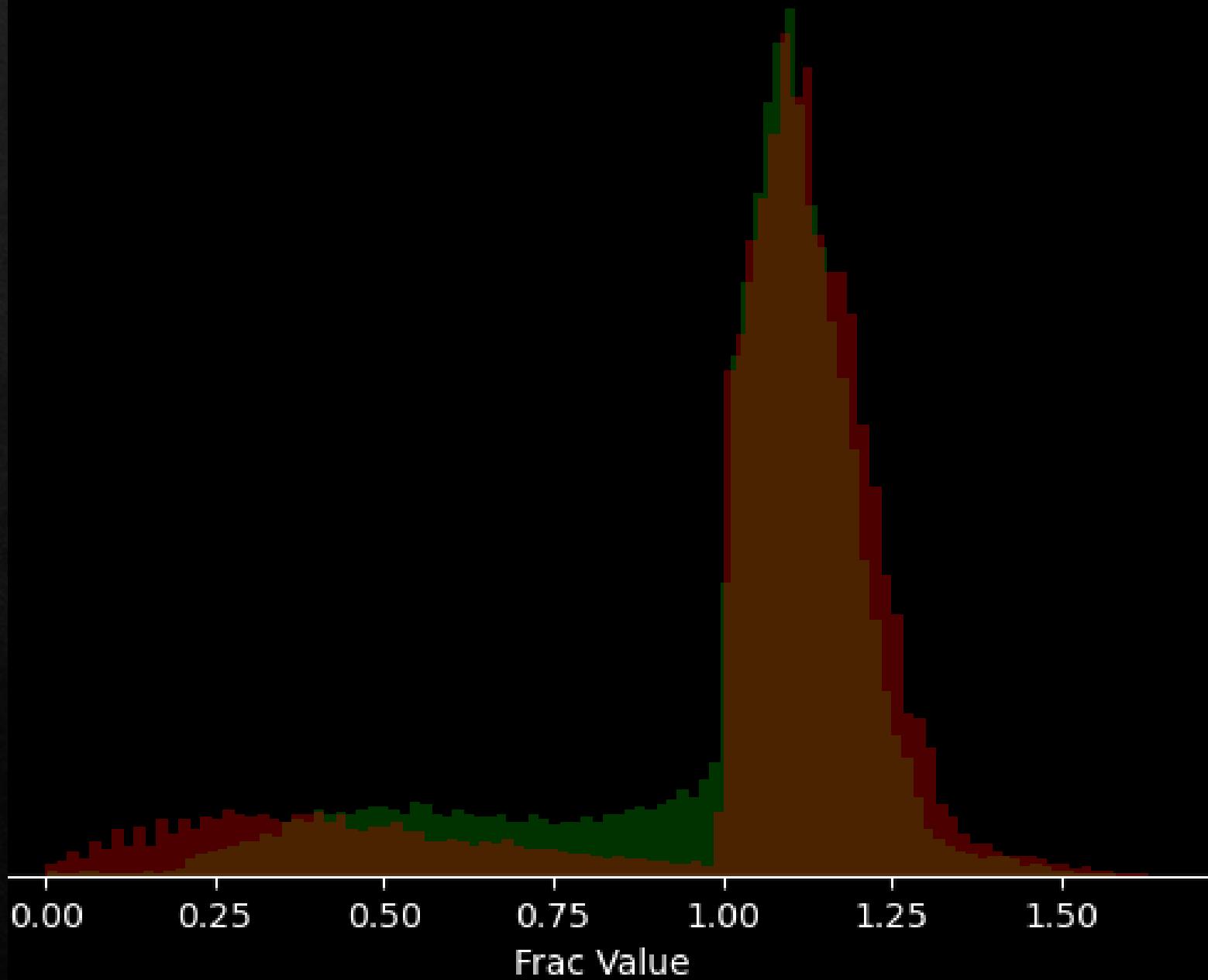
Testing (RF) cont.

Testing (RF) cont.

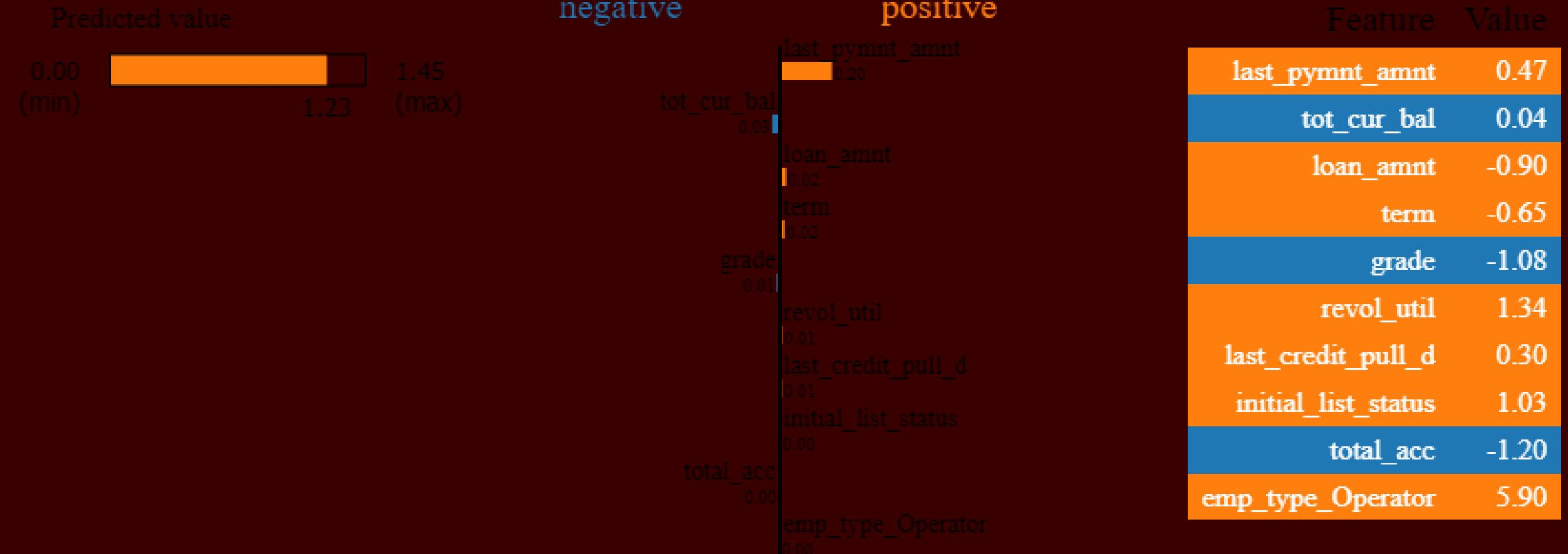
Random Forest Distribution of Error



Histogram of the Random Forest Predictions vs. Y-Test



Testing (RF) cont.



LIME (RF)

Conclusion

- ❖ The most important methods from this paper are:
 - ❖ Random Forest Modelling
 - ❖ Locally Interpretable Modelling Explanation (LIME)
- ❖ With most of the dataset comprised of current loans, it is impossible to evaluate the accuracy and precision of supervised learning methods.