



# “Good Heart, Bad Credit Score”

Using Loan Applicant Personal Information as Additional Data Features to Better and More Fairly Evaluate Loan Applications

By Nebyou Zewde, Robert Ross, and Rowan Mockler

## Problem Definition

### Motivation

Outdated risk assessment models utilized by financial institutions continue to leave typically the most systematically disadvantaged unable to benefit from financial instruments that lead to increased ability to exercise financial freedom.

### Problem

In addition to FICO and DTI (traditional metrics used to determine loan terms), leverage loan applicant personal information to more fairly predict suitable loan interest rates.

## Challenges

### Memory constraints:

With our dataset with being comprised with over 2 million loan applications each with 330 features, running our model required cloud assistance and intense memory allocation

**Feature selection:** Although our data seemed rich with diverse features, we did determine that many were somewhat redundant. We would ideally leverage several other feature selection methods to improve selection

**Handling missing values:** With the data stemming from human input, of course several features were often represented as Null in many inputs, making it necessary to remove features where a significant portion of inputs lacked a value for it.

## Future Work

### Improve performance in the future by:

- Use model on a dataset of rejected applicants to determine potential bias
- Tuning hyperparameters
- Use the gradient boosting machine from the LightGBM library to refine features further.
- Identify zero and low importance features for removal

## Refs & Ack.

[1] NathanGeorge. (2019, April 10). All Lending Club loan data. Retrieved November 28, 2019, from <https://www.kaggle.com/wordsforthewise/lending-club>.

We'd like to thank the incredible CS221 teaching staff for a quarter full of learning.

## Data

### Source:

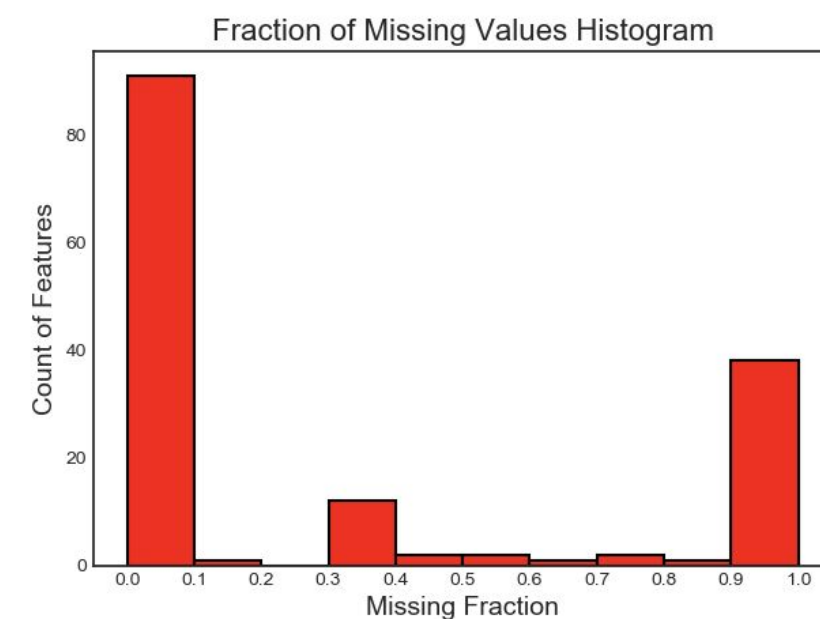
Lending Club dataset, 2 million+ loan applications (2007 - 2018), 330 features per loan.

### Feature Selection:

#### Method1:

Features with a high percentage of missing values

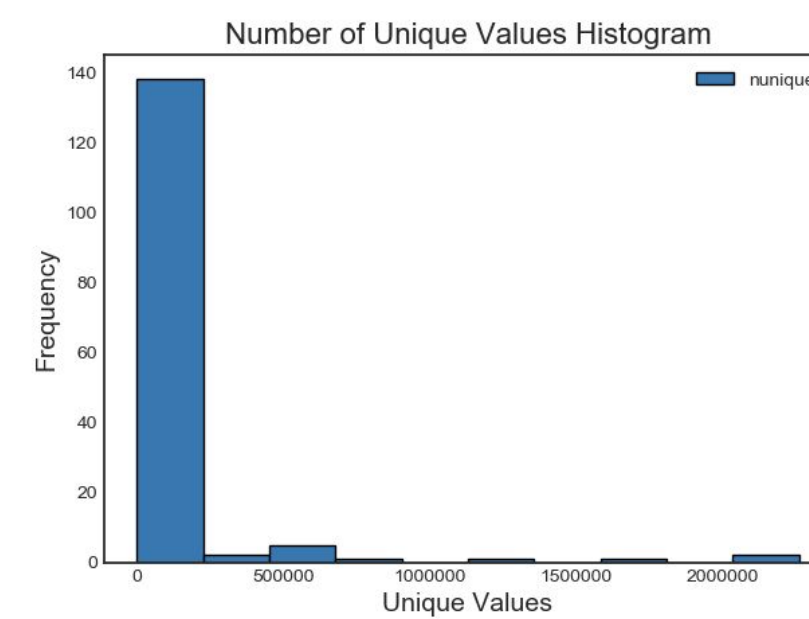
- Find features with a fraction of missing values above threshold of 0.6
- Identified 42 features from the dataset that have more than 60% of their values missing.
- Removed from dataset



#### Method 2:

Features with a single unique value

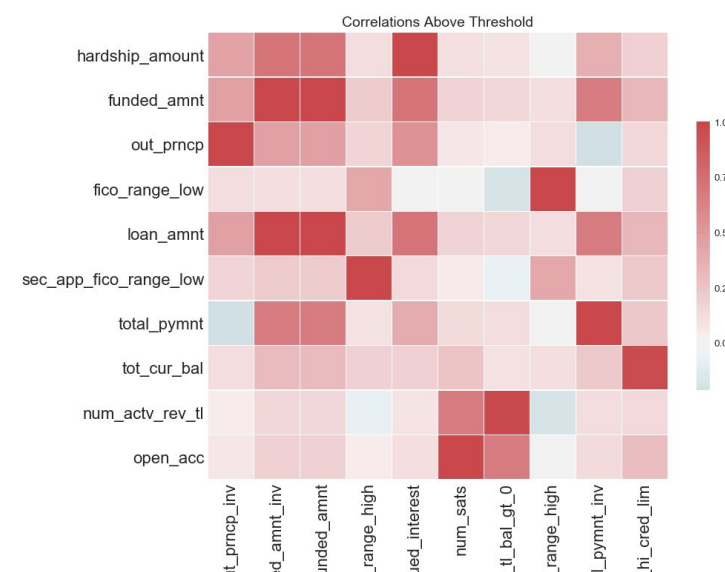
- Find columns that have a single unique value.
- A feature with only one unique value has zero variance therefore not useful in ML.
- Found and removed 4 features



#### Method 3:

Collinear (highly correlated) features

- Found collinear features based on a specified correlation coefficient value.
- For correlated features, identified one of the features for removal
- Found and removed 10 features with a correlation magnitude greater than 0.97.



## Method



Loan classification

Predict outcome of loan (default vs fully paid)

Interest rate prediction

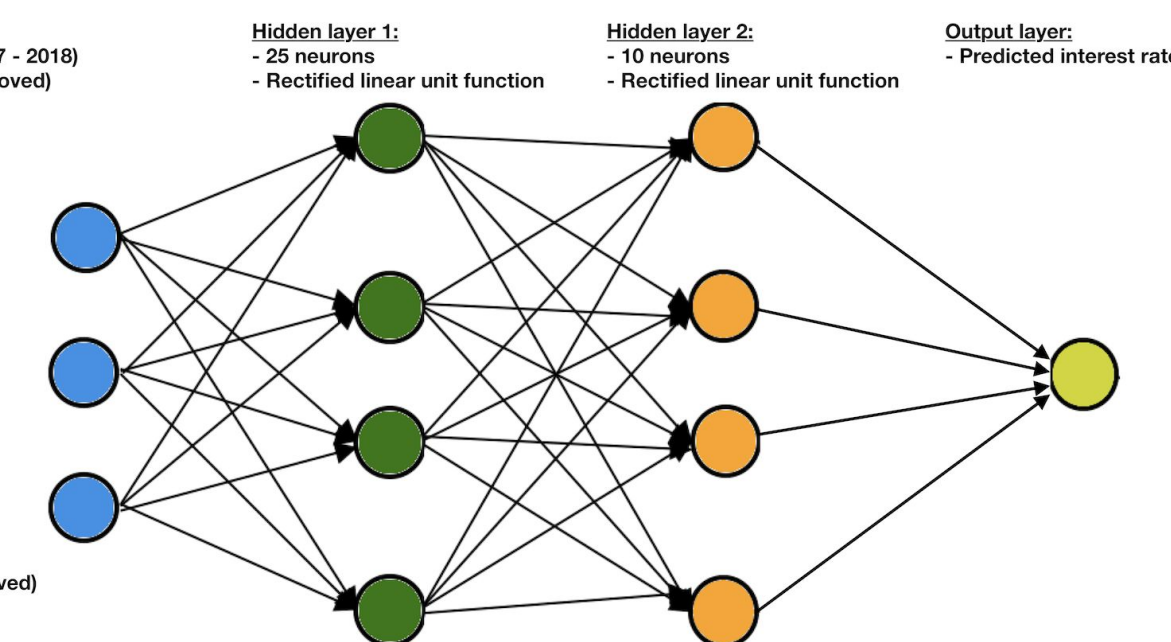
Predict appropriate interest rate for loans predicted to be fully paid

1.1) General strategy

Input layer:  
- 2,000,000 loans (2007 - 2018)  
- 330 features (51 removed)

id  
member\_id  
loan\_amnt  
funded\_amnt  
funded\_amnt\_inv  
term  
int\_rate  
installment  
grade  
sub\_grade  
emp\_title  
emp\_length  
home\_ownership  
annual\_inc  
...

1.2) Feed forward neural network model used for predicting appropriate interests rates



### Overview:

We use a two prong approach of first predicting whether or not a given loan would default or be fully paid; as an extension of the model, if a given loan was predicted to be fully paid, we predicted the respective suitable interest rate for that loan.

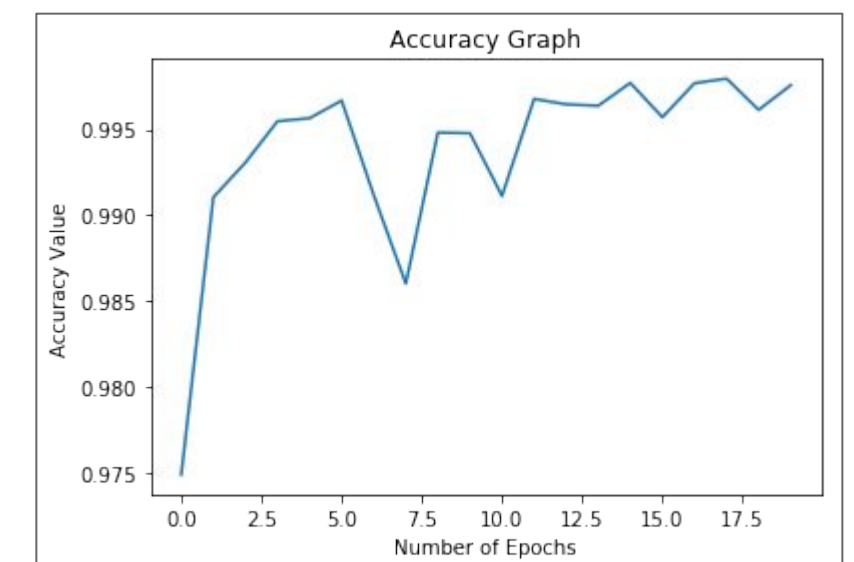
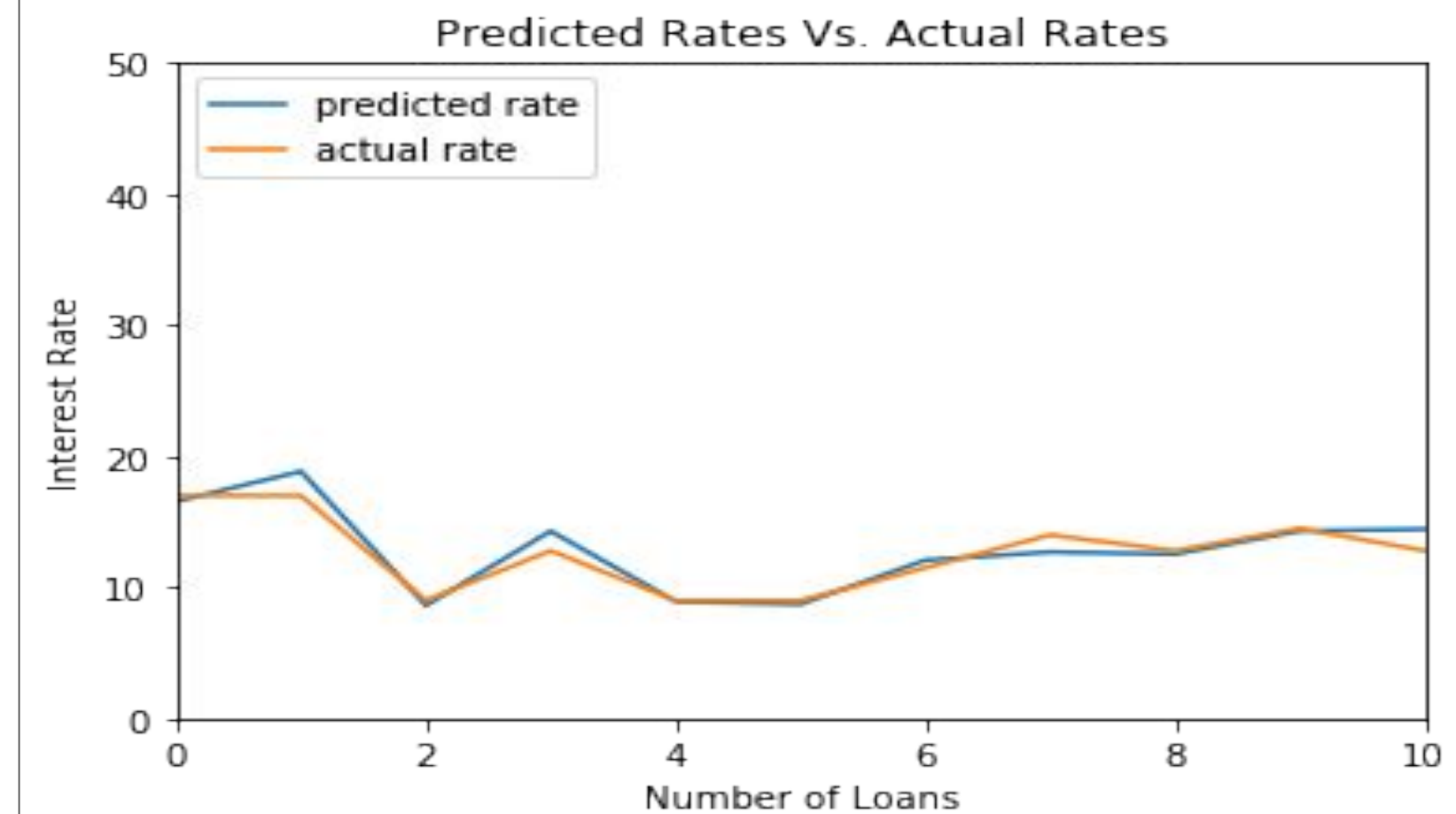
### Step 1: Classification of a given loan

Used binary classification to predict whether a given loan would default or be fully paid

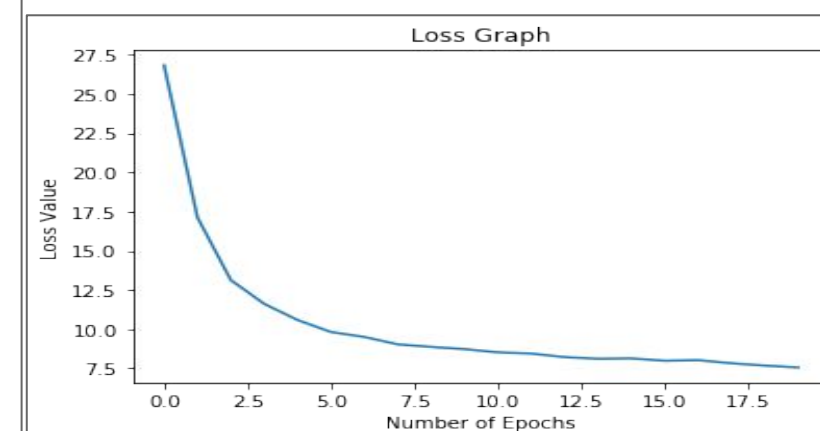
### Step 2: Predict loan interest rate

Utilized two fully connected dense layers with a rectified linear unit activation functions. This model was trained on the loans that were predicted fully paid by our classifier.

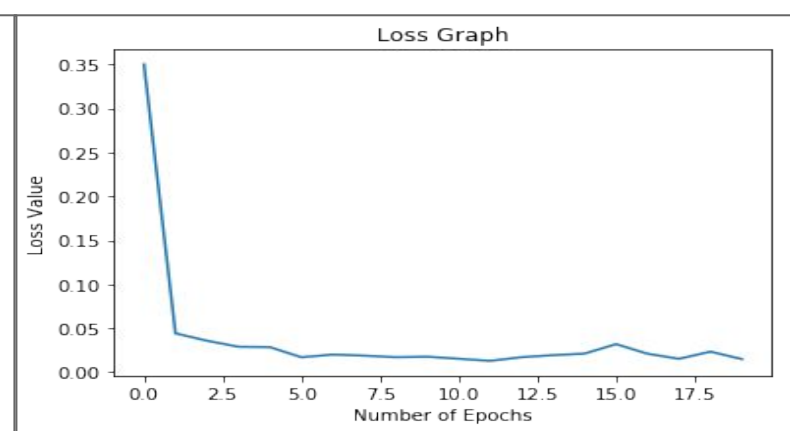
## Results



Regression Loss



Classification Loss



## Analysis

### Predicting loan status:

- We saw a general decrease of binary cross entropy losses
- Despite our use of dropout regularization and early stopping to counter the unbalanced nature of our dataset (lack of negative examples), the model did seem to suffer from overfitting
  - Possible solution would be to utilize the categorical features by implementing one-hot encoding.

### Predicting loan interest rate:

- We saw a general decrease of mean absolute percentage error losses.
- We had an average distance between our prediction and the true value of .08.