# Credit Card Default Prediction Using Data Science & Machine Learning

Cody Dailey

Ishaan Dave

Yang Ge

Vipul Shinde

# Summary

Introduction

Exploratory Data Analysis

Classification

- Logistic Regression
- Random Forest Classifier

Predictions

Conclusions

# Introduction

- Credit risk plays a major role in the banking industry business. Banks main activities involve granting loan, credit card, investment, mortgage, and others.

- Credit card has been one of the most booming financial services by banks over the past years. However, with the growing number of credit card users, banks have been facing an escalating credit card default rate.

- As such data science and machine learning can provide solutions to tackle the current phenomenon and management credit risks.

- Thus, Logistic Regression and Random Forest Classifier models are used to predict if the applicants for the credit card will default or not soon based on the given input data.
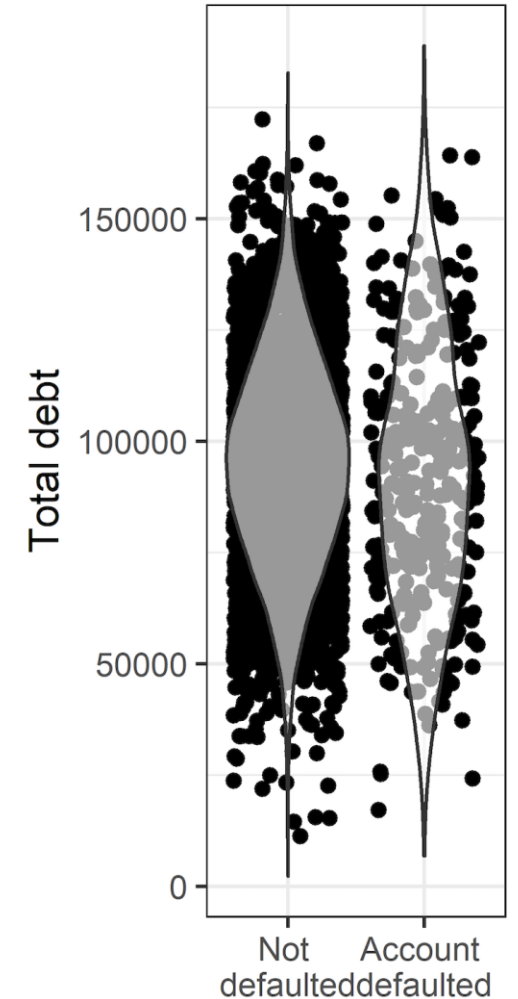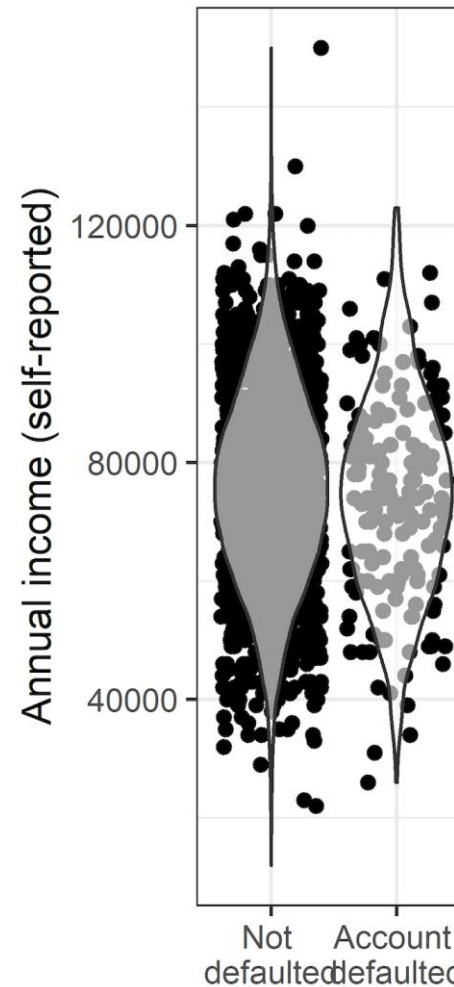
# Exploratory Data Analysis

**3 datasets**

- Training – used to create model

- Validation – qualify performance

- Test – evaluate final model performance

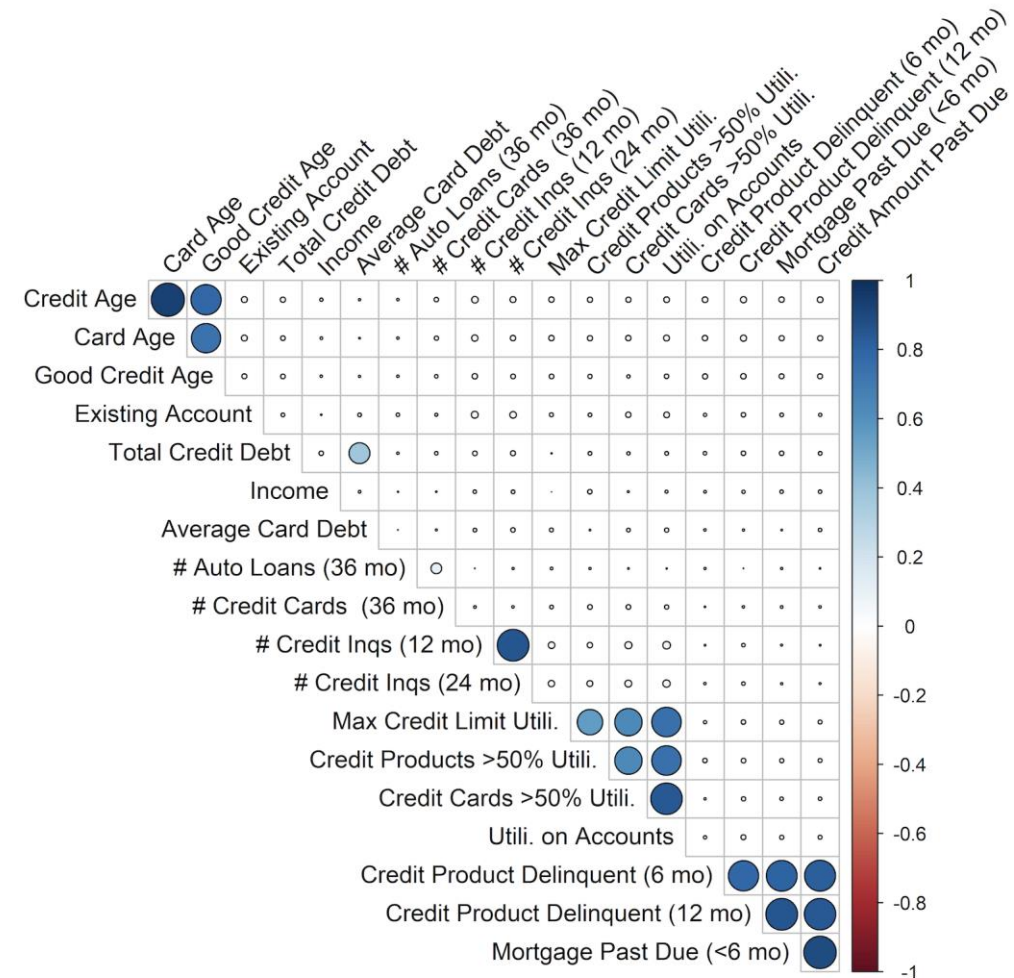# Exploratory Data Analysis

**Univariate Statistics**

- Mean +/- SD and normality for continuous

- N (%) for categorical

- Median [IQR] for discrete counts

- N missing for each variable

- Compared between outcome groups (Did / did not default)

# Exploratory Data Analysis

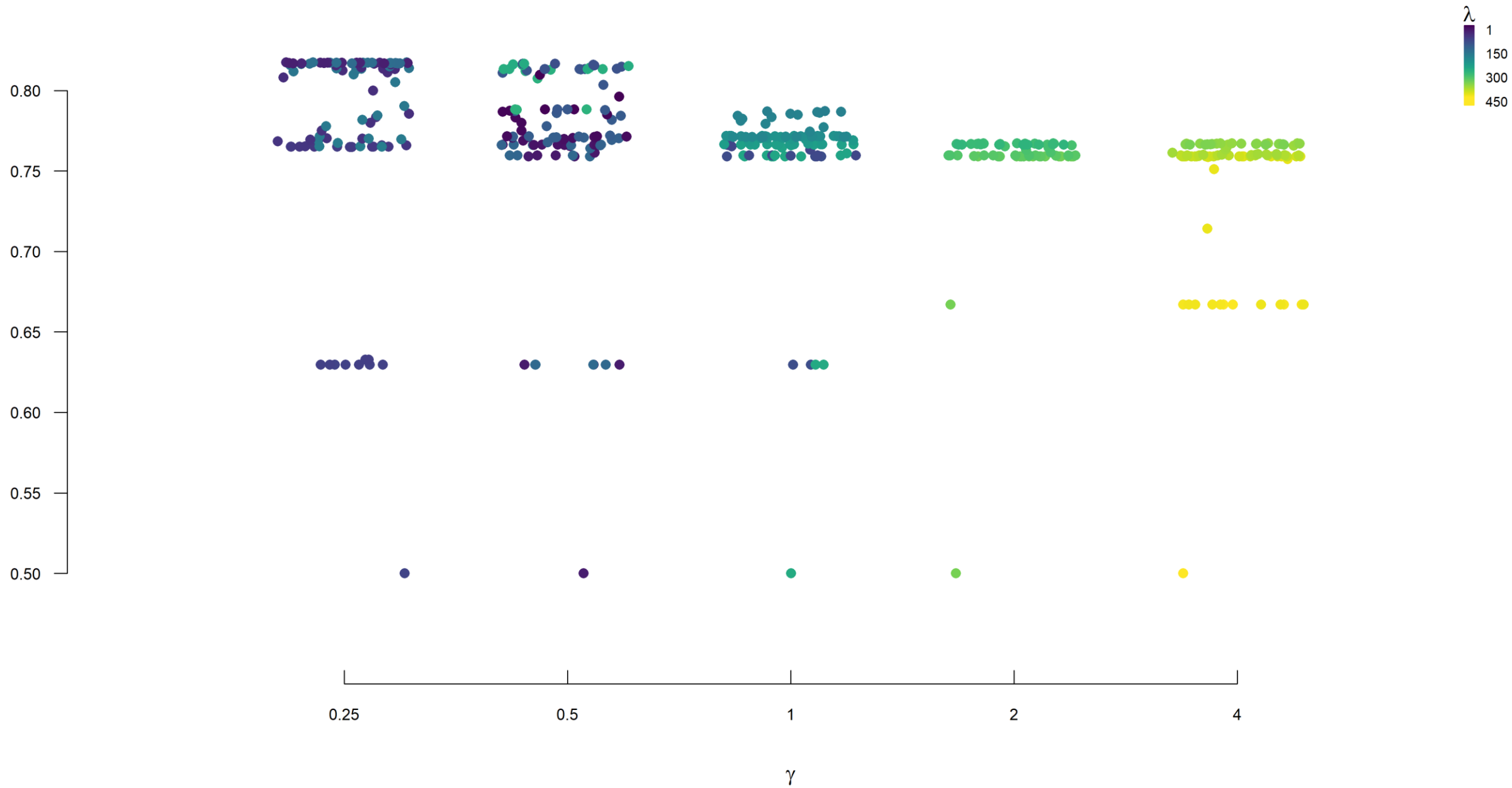**Correlations between variables**

- 20 potential predictors

- Suspected multicollinearity

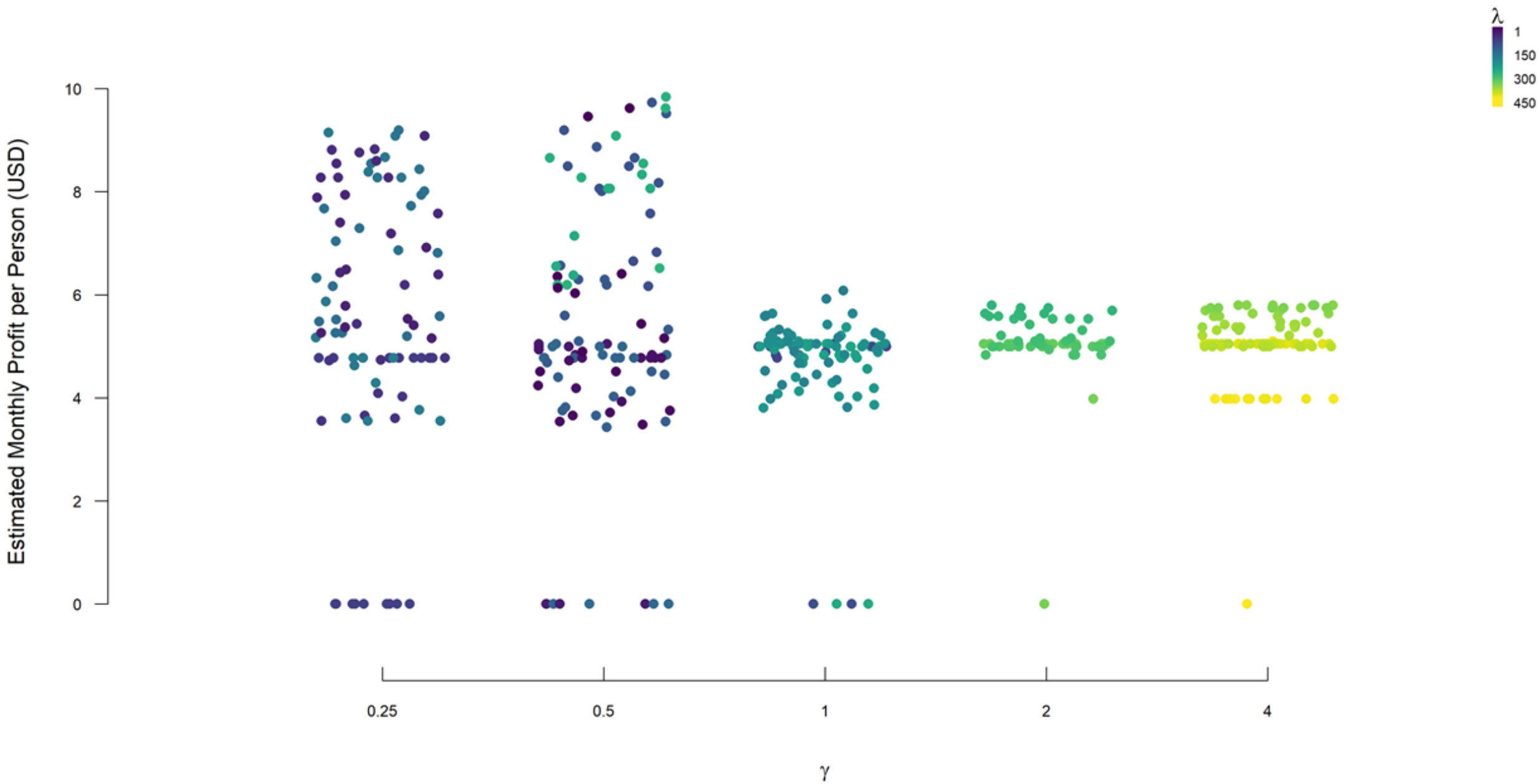- Correlation matrix to examine what variables may vary together
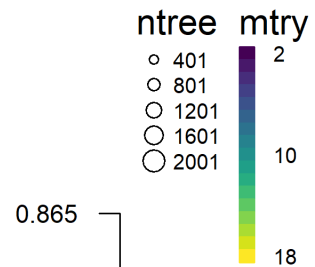
- 

# Feature Selection

- We want to select covariates that are important in predicting whether or an individual will default or not and "penalize" (shrink) those that are not important

- Adaptive Least Absolute Shrinkage and Selection Operator (LASSO) Logistic Regression was performed to select covariates and predict binary outcome

- $logit(E(Y|\boldsymbol{X}) = \boldsymbol{\beta}^T(\boldsymbol{X}),$
  - Y is the binary outcome (whether or not an individual defaulted)
    - $\boldsymbol{X} = (1, X_1, \ldots, X_p)$ is the vector of covariate values
    - $\boldsymbol{\beta} = (\beta_0, \beta_1, \ldots, \beta_p)$ vector of regression parameters we want to estimate

- Estimation is an optimization problem in the form $L(\boldsymbol{\beta}) - \lambda * pen(\boldsymbol{\beta})$, where $pen(\boldsymbol{\beta})$ is the "penalty" and $\lambda$ is the "tuning parameter" (amount of shrinkage)

# Predictions and Comparisons

- Logistic Regression

| | | Truth | | |
|---|---|---|---|---|
| | | Defaulted | Did Not Default | |
| Prediction | Defaulted | 398 | 4,150 | 4,548 |
| | Did Not Default | 3 | 449 | 452 |
| | | 401 | 4,599 | 5,000 |

- Random Forest

| | | Truth | | |
|---|---|---|---|---|
| | | Defaulted | Did Not Default | |
| Prediction | Defaulted | 390 | 3,060 | 3,450 |
| | Did Not Default | 11 | 1,539 | 1,550 |
| | | 401 | 4,599 | 5,000 |

| Metric | LASSO | RF |
|---|---|---|
| Accuracy | 16.94% | 38.58% |
| Error | 83.06% | 61.42% |
| Profit | $31205.95 | $97091.97 |
| Profit.mp | $6.24 | $19.42 |
| PPV | 0.088 | 0.113 |
| NPV | 0.993 | 0.993 |
| Sensitivity | 0.993 | 0.973 |
| Specificity | 0.098 | 0.335 |
| Youden | 0.090 | 0.307 |
| d2c | 0.902 | 0.666 |

# Conclusions

- Rejection example based on Random Forest that shows how a change in a variable impacts model performance.

- As seen in figure to the right, those who already have an account with bank XYZ do not receive favorable treatment



| | Client Observations | Mean Accepted |
|---|---|---|
| Total Credit Debt | 69676.72 | 95654.06 |
| Average Monthly Credit Card Debt | 13010.56 | 13344.64 |
| Credit Age | 302 | 326.1709 |
| Credit Age (Good) | 164 | 162.6121 |
| Credit Card Age | 288 | 296.4327 |
| Non-mortgage Deliquencies (12 Months) | 0 | 0 |
| Non-mortgage Deliquencies (6 Months) | 0 | 0 |
| Mortgage Deliquencies (6 Months) | 0 | 0 |
| Total Credit Past Due | 0 | 0 |
| Credit Inquiries (12 Months) | 2 | 1 |
| Credit Card Inquiries (24 Months) | 4 | 2 |
| Opened Credit Cards (36 Months) | 0 | 0 |
| Opened Auto Loans (36 Months) | 0 | 0 |
| Credit Card Utilization (All) | 0.5721114 | 0.4230993 |
| Percentage Credit Products >50% Utilization | 0.5027835 | 0.4469897 |
| Credit Utilization (Highest Limit) | 0.6682015 | 0.4465359 |
| Percentage Credit Cards >50% Utilization | 0.5390094 | 0.407966 |
| Account @ Bank XYZ | No Account | No Account |
| Reported Income | 82000 | 76779.98 |
| Residence State | FL | GA |

# Thank you!