

Credit Risk Analysis - LendingClub Loan data

Sreejaya¹, Vuthy¹, & Suman¹

¹ City University of New York (CUNY)

Author Note

This report is submitted as part of the course project for 'DATA621 - Business Analytics and Data Mining' at CUNY SPS.

Abstract

LendingClub is an online lending platform for loans. Borrowers apply for a loan online, and if accepted, the loan gets listed in the marketplace. As an investor/lender, you can browse the loans in the marketplace, and choose to invest in individual loans at your discretion - basically purchase *notes* backed by payments based on loans. In this project, we will attempt to analyse and predict the if a loan is risky loan or not so that we can avoid investment in high-risk notes.

Keywords: delinquent, dti (debt-to-income ratio), credit utilization ratio, mortgage accounts, open credit lines

PREDICTIVE MODELING WITH LOGISTIC REGRESSION, AND RANDOM FOREST

Credit Risk Analysis - LendingClub Loan data

Introduction

LendingClub's peer to peer model can not only yield better interest rates over the traditional banking counterparts, treasury bonds and other such financial instruments, it also has advantages like lower overhead costs, lower cost of capital etc. Based on each loan application and credit report, every loan is assigned a grade ranging from A to G, and sub-grade 1 to 5, with a corresponding interest rate. The higher the interest rate, the riskier the grade.

The LendingClub's datasets contains comprehensive list of features (about 115), we will shortlist few features to employ to train our model for predictions. We will build different models, and test those against the validation dataset and select a better performing model.

Literature Review

We have reviewed few of the previous projects/analysis/kaggle competitions done in this area ("GiveMeSomeCredit," 2011) (Jitendra Nath Pandey, 2011) (Liang, n.d.) (Shunpo Chang, 2015). Majority of these were applying machine learning to improve the loan default prediction. Some of these determined that Random Forest appeared to be better performing, and others logistic regression would be better. However, real-world records often behave differently from curated data in competitions like kaggle, so, we will try to apply different regression techniques including the logistic regression, random forest on the real loan data during the years 2012-13 to continue search for a better predictive model. We will also include additional features like dti, credit utilization etc. from the LendingClub dataset, to see if they influence our target (credit risk) variable.

Methodology

Data Exploration

The dataset we used for this analysis is from publicly available data from LendingClub.com, which is basically an online market place that connects borrowers and investors. As part of the data exploration, we will perform *Exploratory data analysis* to better understand the relationships in the given data including correlations, feature distributions and basic summary statistics. We will also identify the outliers, missing data and any look for invalid data.

Data Preparation

As part of the data preparation, we will try to fix the data issues we noticed in our exploratory analysis, which involves treating the outliers, missing data, invalid data etc. We will also identify the data classifications, and create dummy variables wherever required, and will convert the categorical data to numeric data which would help us later in model building.

Model Development

Our primary objective here is to *identify risky loans*, which is binary outcome, so we will be using *logistics regression, naive bayes, and random forest* modeling techniques to examine and predict the credit risk using a training subset (80%) of the original full data.

Model Validation

A validation subset (20%) was used to test how well our candidate models predict the target variable. AUC and the model Accuracy will be used to select a better predicting model.

Experimentation and Results

Data Exploration and Preparation

Feature Selection: There are 111 fields and 188K observations. However, not all fields are intuitively useful for our model building, such as the loan ID, member ID, last payment month etc., so we will be removing such fields. We will also be removing features with majority of NAs (80% NAs). In order to label the dataset, we will classify the loans that defaulted, charged off, or were late on payments as negative cases, and those that were fully paid or current was classified as positive loans.

Outliers: Look for outliers, and remove or transform the outliers that might affect the model. We are going to consider the loans with annual income of less than 250K for the model building.

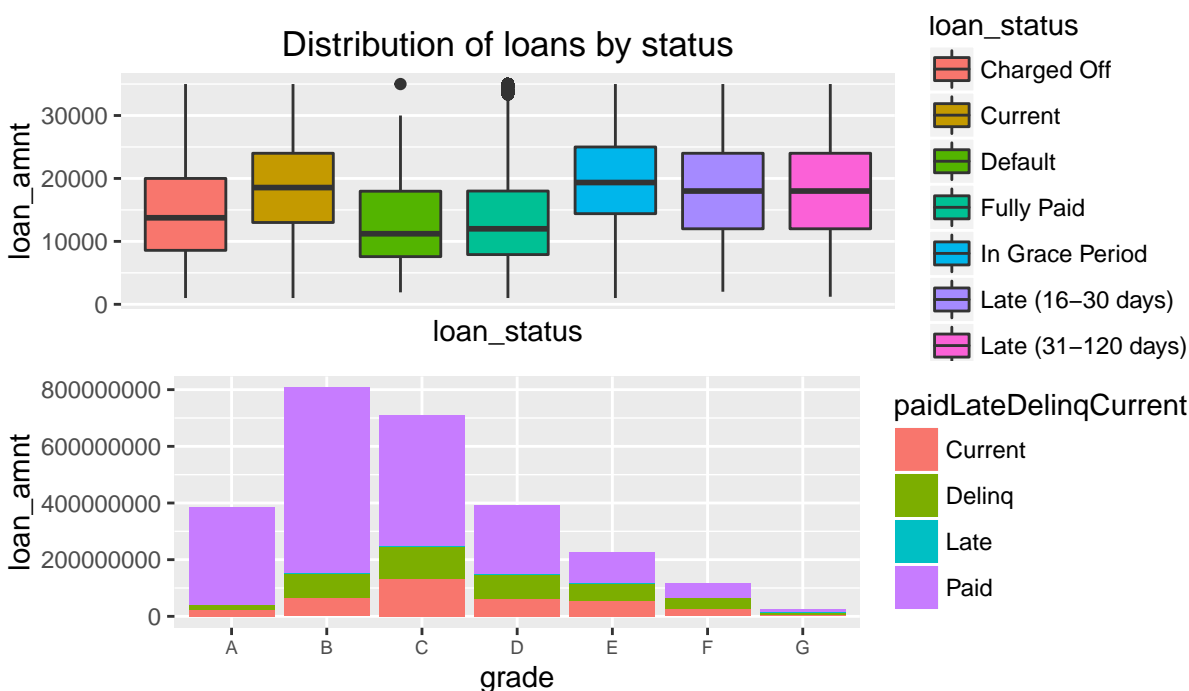


Figure 1. Distribution of loans by status and grade

PREDICTIVE MODELING WITH LOGISTIC REGRESSION, AND RANDOM FOREST

Tidy Data: Convert the date fields like issue date to proper date type, and remove % sign for interest rates, dti and convert those into numeric values. And consider matured loans only. [issue date + term months < today], and remove variables where more than 25% of the observations are missing values.

Factorize the loan status levels with proper ordering (“Charged Off”, “Default”, “Late (31-120 days)”, “Late (16-30 days)”, “In Grace Period”, “Current”, “Fully Paid”). Also loans issued by Lending Club fall into three categories of verification: “income verified,” “income source verified,” and “not verified.” The “home ownership” is another factor variable provided by the borrower during registration Or obtained from the credit report. The values are: RENT, OWN, MORTGAGE, OTHER.

Create a factor label *bad_loan* which indicates a loan is bad if it is delinquent, or late consider as negative, otherwise positive. This would be our response variable which indicates a loan is bad or not.

Shortlist Numeric variables: Identify the numeric variable, and compare how those distributions plot against the good, bad label, and pick a few variables with differences in the bad and good populations.(yhat, 2013)

Lets visualize the correlation graph of these numeric variable:

From our numerical feature list, it appears like the *total_pymnt*, and *total_rec_prncp* are having high correlation with the *bad_loan* classification.

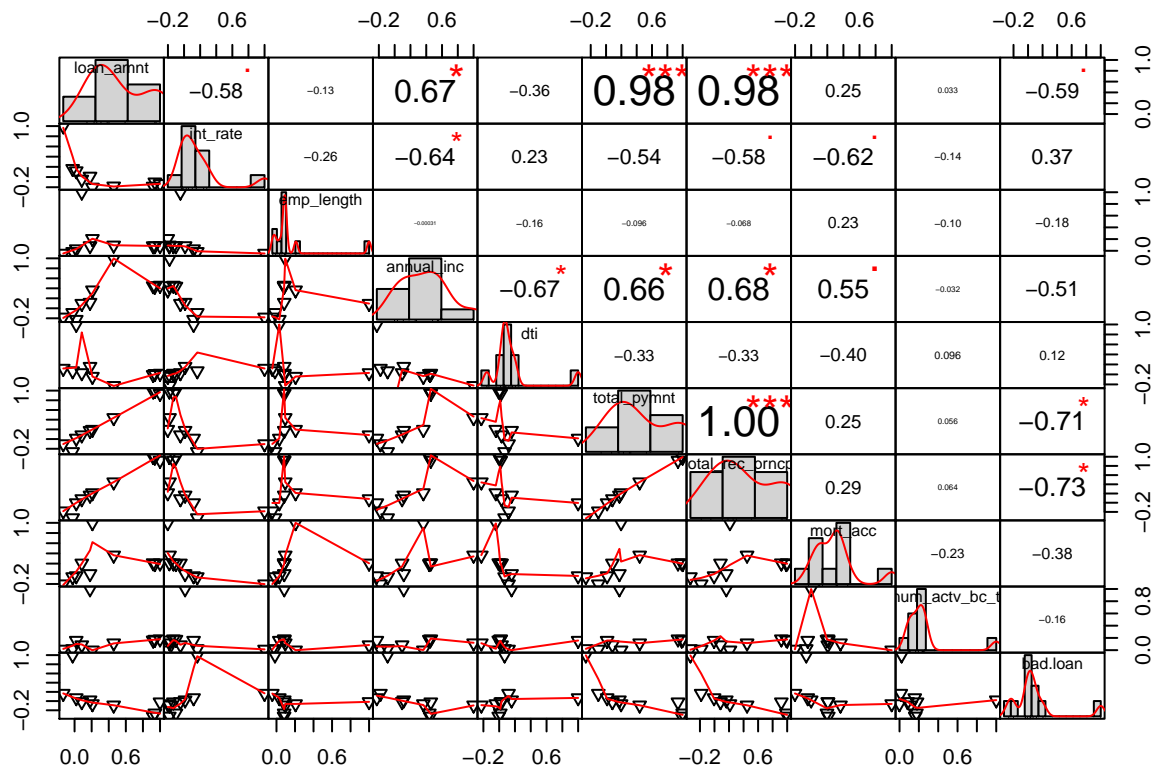


Figure 2. Correlations

Dummy variables: Since R's glm function will take care of the dummy variables, we just need to make sure that the categorical variables are factorized.

Split the dataset into training and test: We will randomly split our dataset into training (80%) and test (20%).

Number of observations in *training* dataset is 114190

Number of observations in *test* dataset is 28548

PREDICTIVE MODELING WITH LOGISTIC REGRESSION, AND RANDOM FOREST

Model Development

Model 1 - Logistic regression considering all predictors: Noticed high multicollinearity in the predictors, so removing the high VIF variables - `loan_amnt`, `int_rate`, `grade`, `total_pymnt` and `total_rec_prncp` and try to fit again.

Model 2 - Logistic regression removing high multicollinear predictors:

.rownames	Estimate	Std..Error	z.value	Pr...z..
(Intercept)	5.2784	0.8105	6.5123	0.0000
<code>emp_length</code>	0.0044	0.0027	1.6328	0.1025
<code>home_ownership</code> NONE	0.3928	0.4994	0.7866	0.4315
<code>home_ownership</code> OTHER	0.0658	0.5442	0.1209	0.9038
<code>home_ownership</code> OWN	0.0772	0.0342	2.2549	0.0241
<code>home_ownership</code> RENT	0.1914	0.0224	8.5245	0.0000
<code>annual_inc</code>	0.0000	0.0000	-21.9819	0.0000
<code>verification_status</code> Source Verified	0.1420	0.0243	5.8476	0.0000
<code>verification_status</code> Verified	0.1324	0.0211	6.2676	0.0000
<code>issue_d</code>	-0.0005	0.0001	-8.9192	0.0000
<code>dti</code>	0.0171	0.0012	13.6957	0.0000
<code>mort_acc</code>	-0.0330	0.0058	-5.6788	0.0000
<code>num_actv_bc_tl</code>	0.0239	0.0043	5.6010	0.0000

Model 3 - Random Forest, considering all predictors : With 500 trees, the random forest shows the below Gini index - a measure how each variable contribute to the homogeneity of the nodes and leaves in the resulting random forest. The *principal amount received to date*, *loan amount*, *payments received to date for total amount funded*, *debt to income ratio*, *interest rate*, and *annual income* are some of the high important variables here. Its also interesting to see that *issue_d* is also shown as one of the important variables.

Model 4 - Random Forest, removing high multicollinear predictors : Lets generate the *random forest* model by considering low VIF predictors only, and review the

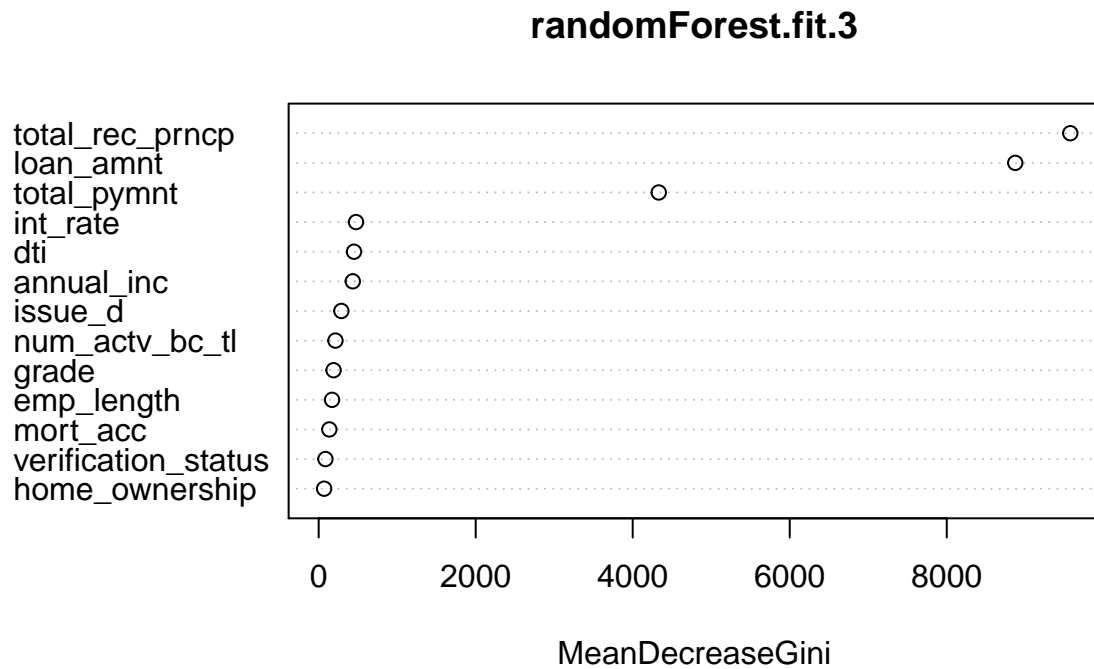


Figure 3. Random Forest - Variable Importance Plot

important features:

Model 5 - Logistic regression by considering the features from “random forest model”: We will consider the top 6 important features resulted from the *random forest* and build the logistic regression model.

.rownames	Estimate	Std..Error	z.value	Pr...z..
(Intercept)	5.3463	0.8100	6.6005	0.0000
dti	0.0184	0.0012	15.1492	0.0000
annual_inc	0.0000	0.0000	-22.1280	0.0000
issue_d	-0.0005	0.0001	-8.7677	0.0000
num_actv_bc_tl	0.0241	0.0043	5.6462	0.0000
emp_length	0.0017	0.0027	0.6410	0.5215
mort_acc	-0.0521	0.0052	-9.9972	0.0000

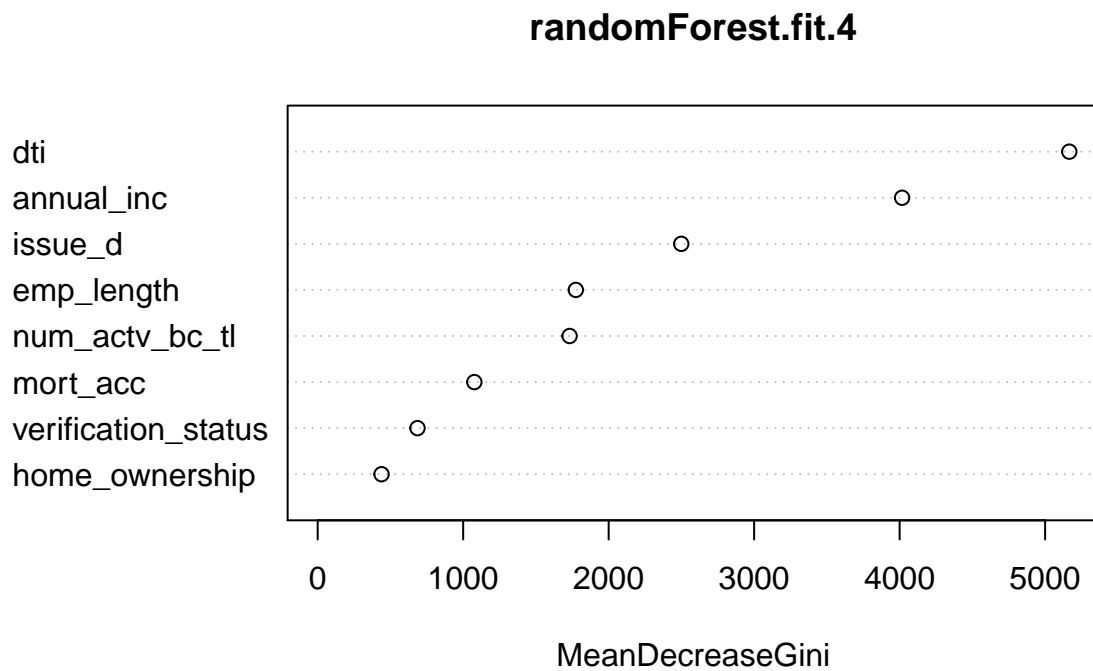


Figure 4. Random Forest - Variable Importance Plot

Model Validation

Now we have got the trained models, lets apply those to the test dataset and derive the performance measures.

$$Sensitivity = \frac{TP}{TP + FN}$$

$$Specificity = \frac{TN}{TN + FP}$$

$$Accuracy = \frac{TP + TN}{(TP + FN) + (FP + TN)}$$

##

Call:

roc.formula(formula = factor(predicted) ~ as.numeric(bad_loan), data = loans.ss.t

##

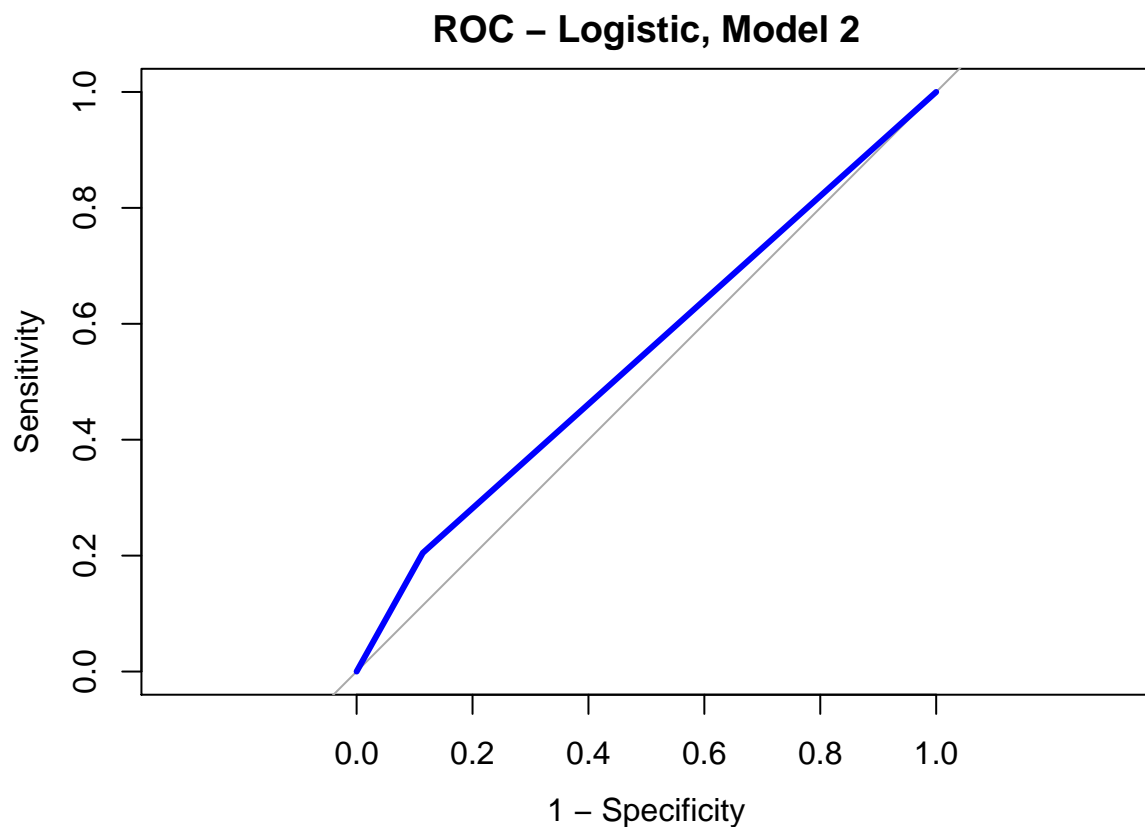


Figure 5

```
## Data: as.numeric(bad.loan) in 25842 controls (factor(predicted) 0) < 2706 cases (fact
## Area under the curve: 0.5453
## 95% CI: 0.5374-0.5531 (DeLong)

##
## Call:
## roc.formula(formula = factor(predicted) ~ as.numeric(bad.loan),      data = loans.ss.t
##
## Data: as.numeric(bad.loan) in 26271 controls (factor(predicted) 0) < 2277 cases (fact
## Area under the curve: 0.5471
## 95% CI: 0.5385-0.5557 (DeLong)
```

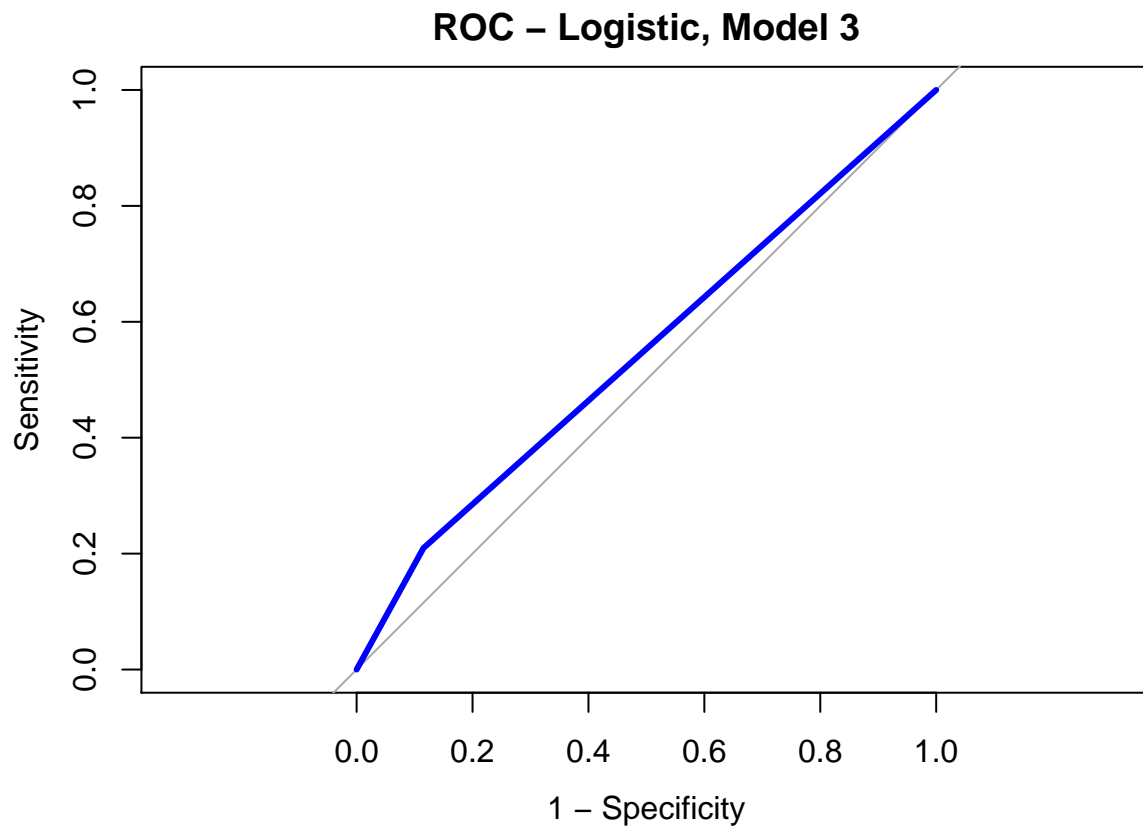


Figure 6

Lets compare the selected logistic and random forest models:

Model	Accuracy	AUC
Logistic, Model 2	0.8212	0.5453
Logistic, Model 5	0.8309	0.5471
Random Forest, Model 4	0.8772	0.5636
Random Forest, Model 3	0.9915	0.9993

From the above, it is evident that *Random Forest, Model 3* performed well here. It has got the AUC [$P(\text{predicted TRUE}|\text{actual TRUE})$ Vs $P(\text{FALSE}|\text{FALSE})$], and the accuracy [$P(\text{TRUE}|\text{TRUE}).P(\text{actual TRUE}) + P(\text{FALSE}|\text{FALSE}).P(\text{actual FALSE})$], are both higher compared to the other models.

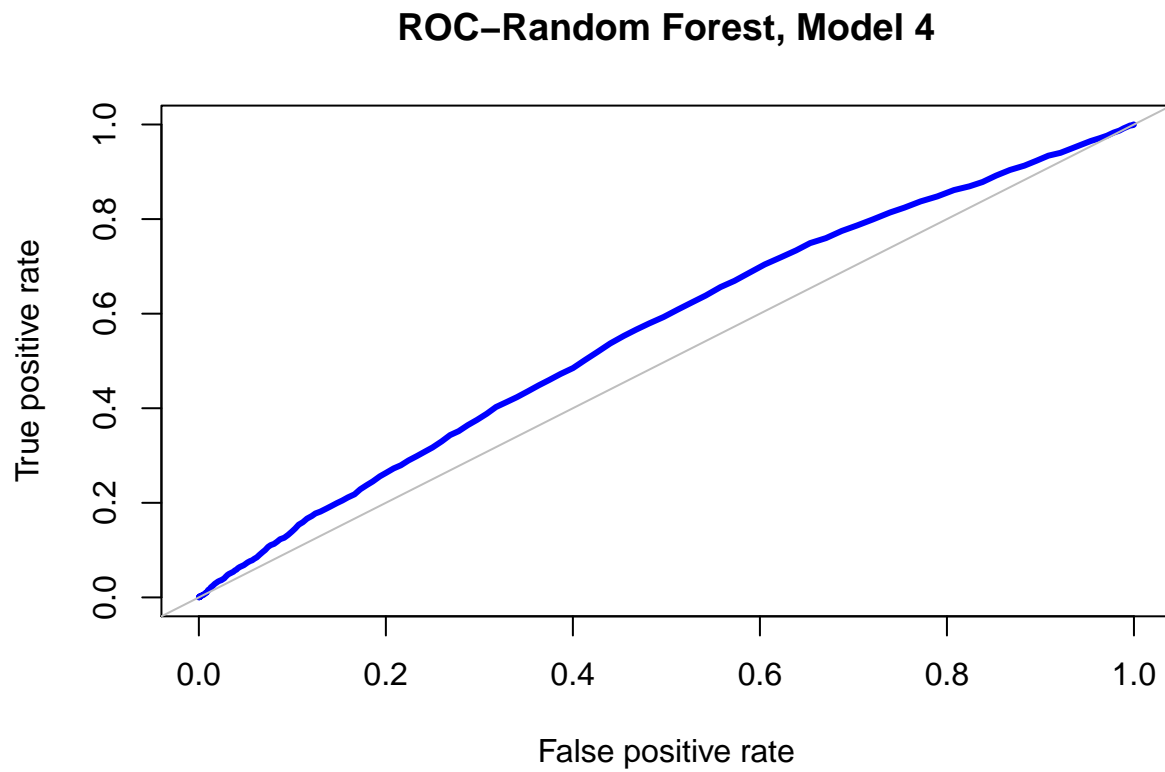


Figure 7

Conclusion

Depending on the features we have chosen for this study, the *random forest* model provided the best outcome out of the few we analysed. So, the limitation here is our selection of the predictors (14 out of 111) and possibly the overall dataset. The logistic regression models interestingly showed high accuracy and low AUC, this could be due the cutoff point we have chosen for the accuracy, and AUC here indicates that the probability of classifying a loan is risky is roughly around 55%.

Future work - we will try including more predictors, and possibly loading data before the year 2012. We will also plan on including the *Naïve Bayes* classification analysis to the model suite.

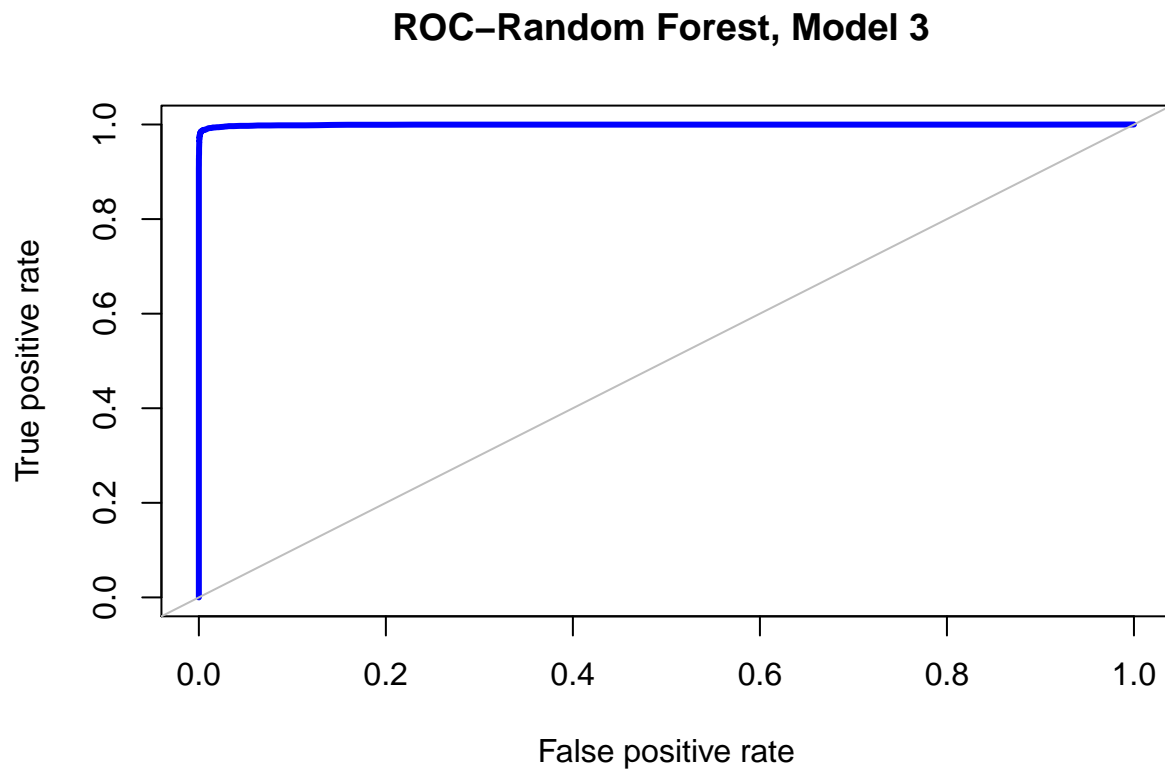


Figure 8

References

- GiveMeSomeCredit. (2011). Retrieved from <https://www.kaggle.com/c/GiveMeSomeCredit>
- Jitendra Nath Pandey, M. S. (2011). Predicting probability of loan default. Project Report. Retrieved from <http://cs229.stanford.edu/proj2011/PandeySrinivasan-PredictingProbabilityOfLoanDefault.pdf>
- Liang, J. (n.d.). Predicting borrowers chance of defaulting on credit loans. Kaggle Submission. Retrieved from <http://cs229.stanford.edu/proj2011/JunjieLiang-PredictingBorrowersChanceOfDefaultingOnCreditLoans.pdf>
- Shunpo Chang, G., Simon. (2015). Predicting default risk of lending club loans. Project Report. Retrieved from http://cs229.stanford.edu/proj2015/199_report.pdf
- yhat. (2013). Machine learning for predicting bad loans. Retrieved from

<http://blog.yhat.com/posts/machine-learning-for-predicting-bad-loans.html>