

Credit Risk Analysis - LendingClub Loan data

DATA 621: Business Analytics and Data Mining

Sreejaya, Vuthy and Suman

December 12, 2016

1 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 if a loan is risky or not so that we can avoid investment in high-risk notes.

2 Keywords:

delinquent, dti (debt-to-income ratio), credit utilization ratio, logistic regression, random forest

3 Literature Review

We have reviewed few of the previous projects/analysis/kaggle competitions done in this area (Kaggle, 2011) (Pandey and Srinivasan, 2011) (Liang, 2011) (Chang, Simon, and Genki, 2011). 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 and random forest on the real loan data during the years 2012-13 to 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.

The questions for logistic regression are, (1) what is the best way to describe the relationship between predictor and outcome, and (2) Is the relationship statistically significant - does knowing something about the predictor variable enable a better prediction of the outcome than knowing nothing about the predictor ?. The logistic regression equation takes the form:

$$\log odds = \log\left(\frac{p}{1-p}\right) = b_0 + b_1 * (predictor\ value)$$

where b_1 is the amount of increase in the log odds of the outcome given by an increase in one unit of the predictor, and b_0 is the intercept of the model. So, the logistic regression does *not* assume a linear relationship between the dependent and independent variables (but between the logit of the outcome and the predictor values). There are several sophisticated alternatives to logistic regression techniques available, such as decision trees, random forests, neural networks etc. (Reed and Wu, 2013)

Random forest is a nonparametric machine learning strategy that can be used for building a risk prediction model in survival analysis. In survival settings, the predictor is an *ensemble* formed by combining the results of many survival trees. (Mogenssen, Ishwaran, and Gerds, 2012)

4 Methodology

4.1 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.

4.2 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.

4.3 Model Development

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

4.4 Model Validation

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

5 Experimentation and Results

5.1 Data Exploration and Preparation

5.1.1 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.

5.1.2 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.

5.1.3 Tidy Data:

Convert the date fields such as issue date to a proper date format, and remove % sign for interest rates, dti and convert those into numeric values. Furthermore, we will consider only matured loans. [issue date + term months < today], and remove variables where more than 25% of the observations are missing values.

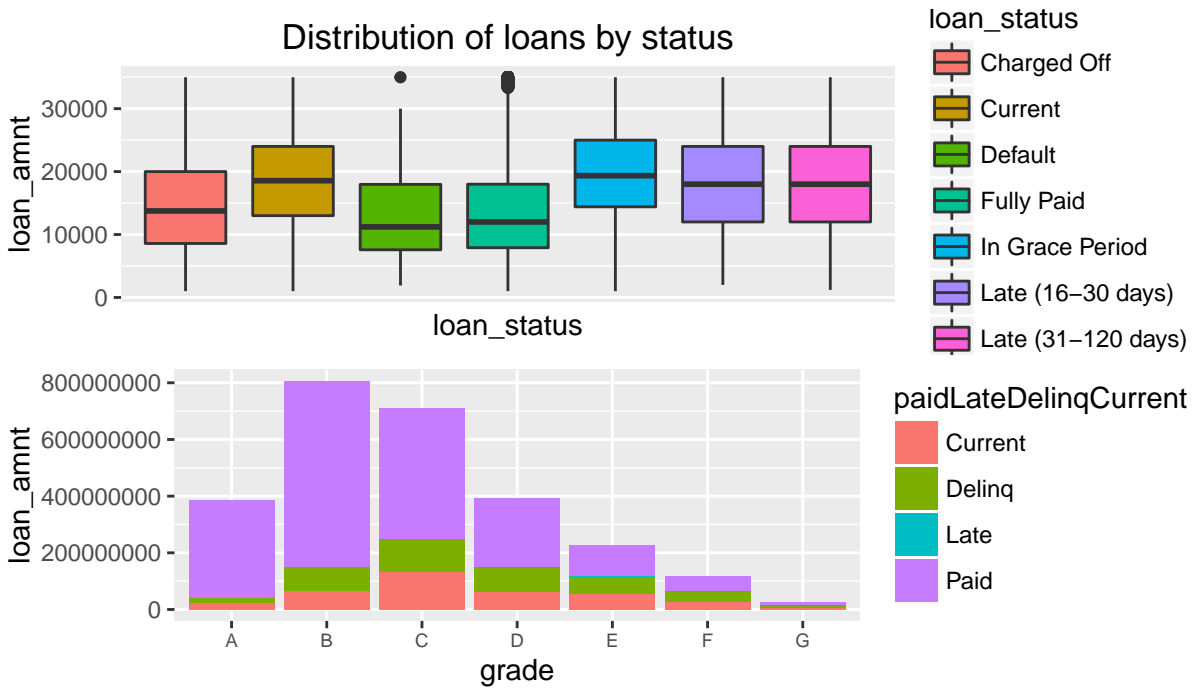


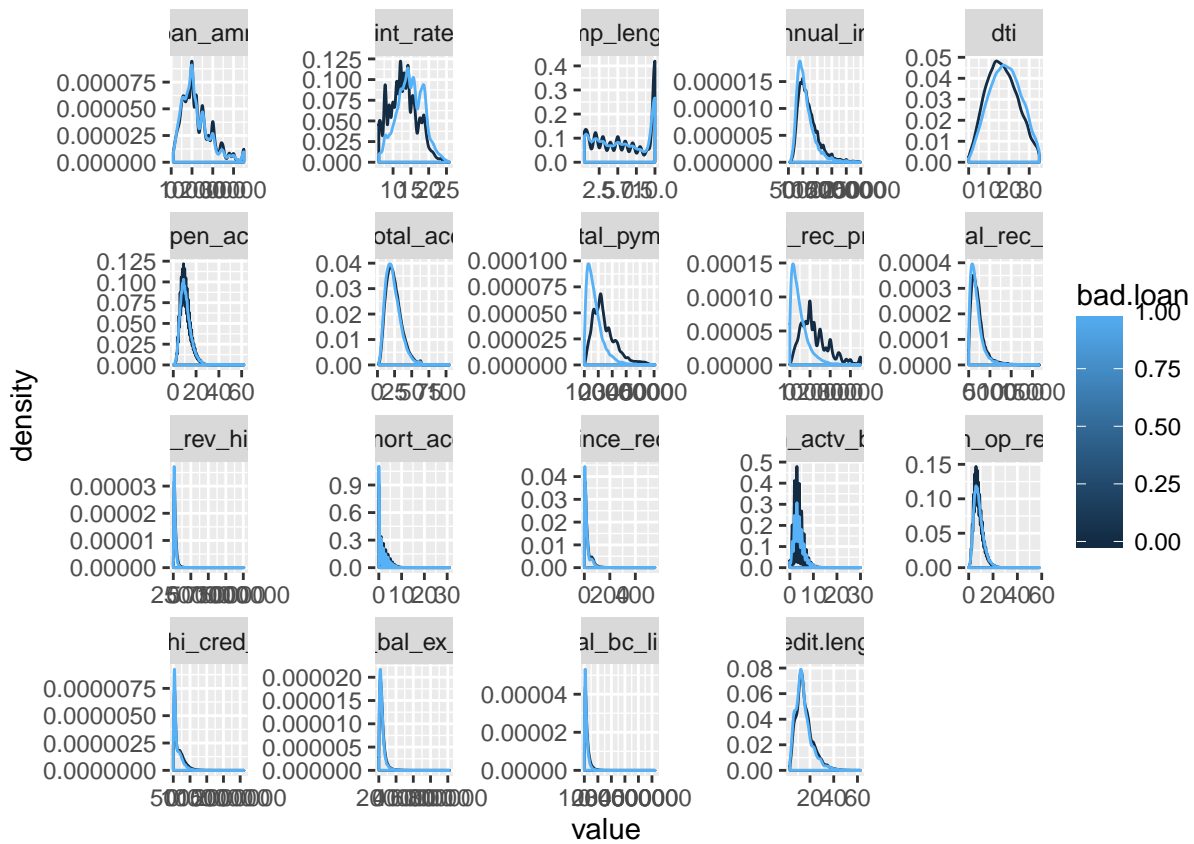
Figure 1: Distribution of loans by status and grade

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.

5.1.4 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]



Lets visualize the correlation graph of these numeric variable:

From our numerical feature list, it appears like the `total_pymnt`, and `total_rec_prncpl` has high correlation with the `bad_loan` classification.

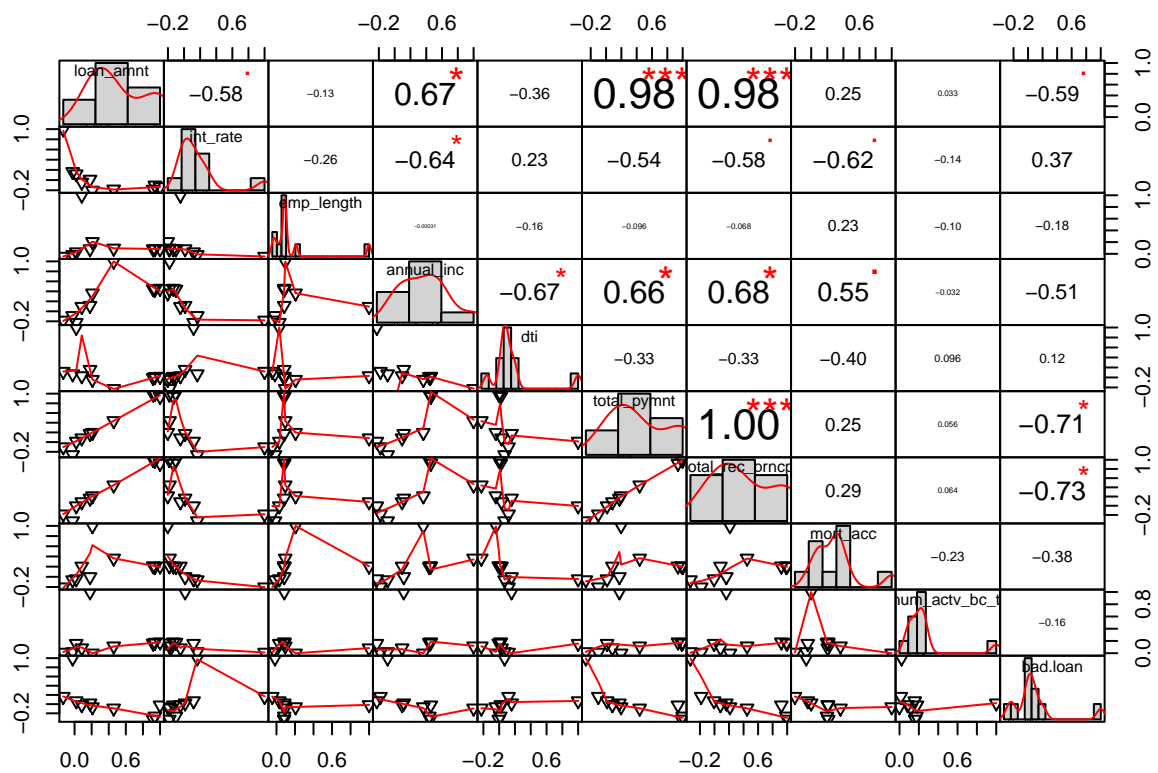


Figure 2: Correlations

5.1.5 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.

5.1.6 Split the dataset into training and test:

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

```
## Observations: 114,190
## Variables: 14
## $ loan_amnt      <int> 14400, 15000, 10000, 15000, 35000, 5600, 4...
## $ int_rate       <dbl> 12.12, 11.14, 10.74, 7.90, 14.30, 18.75, 2...
## $ grade          <ord> B, B, B, A, C, D, E, B, C, B, A, C, B, D, ...
## $ emp_length     <dbl> 8.0, 6.0, 2.0, 2.0, 2.0, 6.0, 2.0, 0.5, 10...
## $ home_ownership <fctr> RENT, MORTGAGE, RENT, MORTGAGE, MORTGAGE,...
## $ annual_inc     <dbl> 45000, 78000, 170000, 104000, 170000, 6200...
## $ verification_status <fctr> Source Verified, Verified, Source Verifie...
## $ issue_d        <date> 2012-09-01, 2012-12-01, 2012-03-01, 2013-...
## $ dti            <dbl> 22.32, 22.72, 4.59, 19.54, 9.00, 21.43, 3....
## $ total_pymnt    <dbl> 17247.966, 17194.223, 9498.950, 16299.200,...
## $ total_rec_prncp <dbl> 14400.0, 15000.0, 7474.0, 15000.0, 35000.0...
## $ mort_acc       <dbl> 0, 2, 0, 6, 2, 1, 0, 0, 5, 5, 1, 0, 2, 2, ...
## $ num_actv_bc_tl <dbl> 4, 4, 0, 2, 3, 3, 3, 1, 8, 2, 4, 6, 11, 4,...
## $ bad_loan       <fctr> 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,...
```

Number of observations in *training* dataset is 114190

Number of observations in *test* dataset is 28548

5.2 Model Development

5.2.1 Model 1 - Logistic regression considering all predictors:

	GVIF	Df	GVIF ^{1/(2*Df)}
loan_amnt	210.211433	1	14.498670
int_rate	13.789034	1	3.713359
grade	13.158463	6	1.239559
emp_length	1.103136	1	1.050303
home_ownership	1.473040	4	1.049607
annual_inc	1.451972	1	1.204978
verification_status	1.160992	2	1.038024
issue_d	1.732473	1	1.316234
dti	1.127885	1	1.062019
total_pymnt	270.667965	1	16.451990
total_rec_prncp	330.533383	1	18.180577
mort_acc	1.518237	1	1.232168
num_actv_bc_tl	1.442362	1	1.200984

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.

5.2.2 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
emp_length	0.0044	0.0027	1.6328	0.1025
home_ownershipNONE	0.3928	0.4994	0.7866	0.4315
home_ownershipOTHER	0.0658	0.5442	0.1209	0.9038
home_ownershipOWN	0.0772	0.0342	2.2549	0.0241
home_ownershipRENT	0.1914	0.0224	8.5245	0.0000
annual_inc	0.0000	0.0000	-21.9819	0.0000
verification_statusSource Verified	0.1420	0.0243	5.8476	0.0000
verification_statusVerified	0.1324	0.0211	6.2676	0.0000
issue_d	-0.0005	0.0001	-8.9192	0.0000
dti	0.0171	0.0012	13.6957	0.0000
mort_acc	-0.0330	0.0058	-5.6788	0.0000
num_actv_bc_tl	0.0239	0.0043	5.6010	0.0000

5.2.3 Model 3 - Random Forest, considering all predictors :

Since the decision trees could be susceptible to noise in the training data, lets directly go with random forest, which would randomly select observations and variables, and build several trees and takes majority vote from all these trees for our prediction. 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.

```
##
## Call:
## randomForest(formula = bad_loan ~ ., data = loans.ss.train, do.trace = F,          type = "class", importan
##               Type of random forest: classification
##               Number of trees: 500
## No. of variables tried at each split: 3
##
##               OOB estimate of  error rate: 0.9%
## Confusion matrix:
##      0      1 class.error
## 0 99679      24 0.0002407149
## 1  1009 13478 0.0696486505
```

Here, Out Of Bag error estimate can be taken as estimate of performance on unseen data.

5.2.4 Model 4 - Random Forest, removing high multicollinear predictors :

Lets generate the *random forest* model by considering low VIF predictors only and review the important features:

5.2.5 Model 5 - Logistic regression considering the features from “random forest model”:

We thought, even if we can not use random forest as primary tool for some reason, we can still use it to select important variables from our data. 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

.rownames	Estimate	Std..Error	z.value	Pr...z..
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

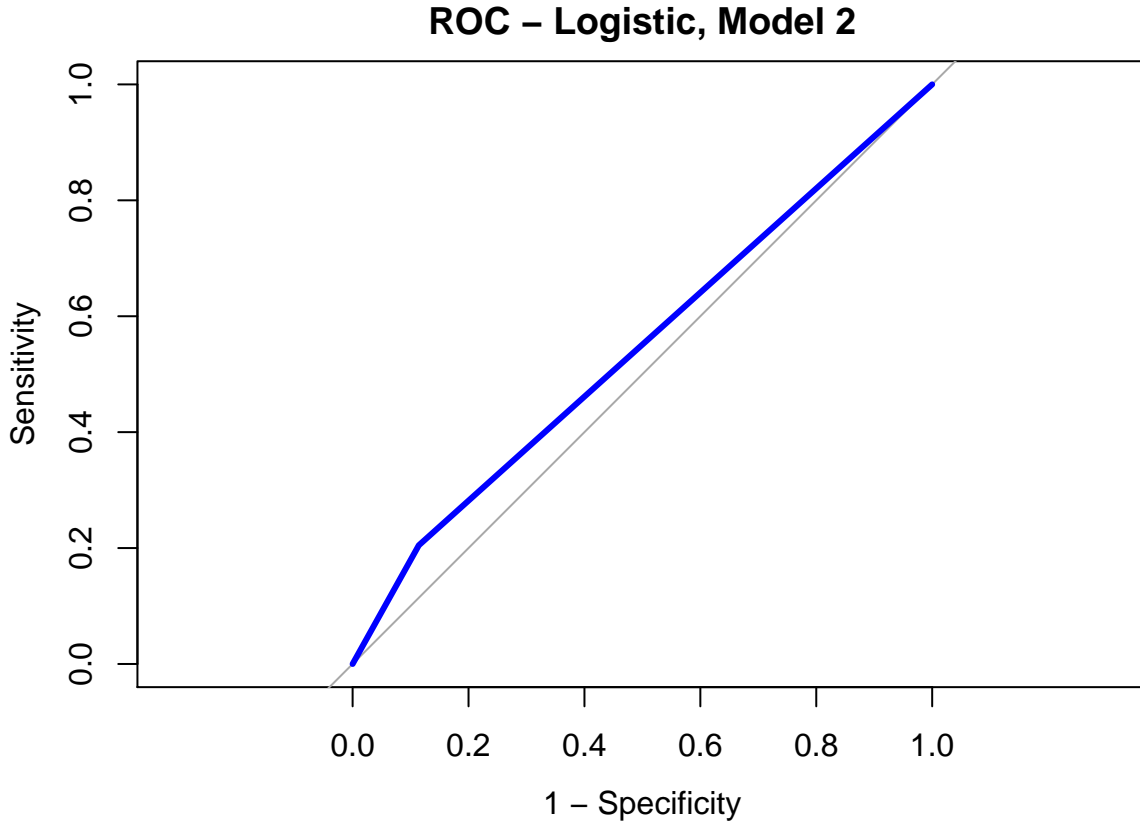
5.3 Model Validation

Now that our models are trained, 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.test, plot = FALSE,
##
## Data: as.numeric(bad.loan) in 25842 controls (factor(predicted) 0) < 2706 cases (factor(predicted) 1).
## Area under the curve: 0.5453
## 95% CI: 0.5374-0.5531 (DeLong)
```

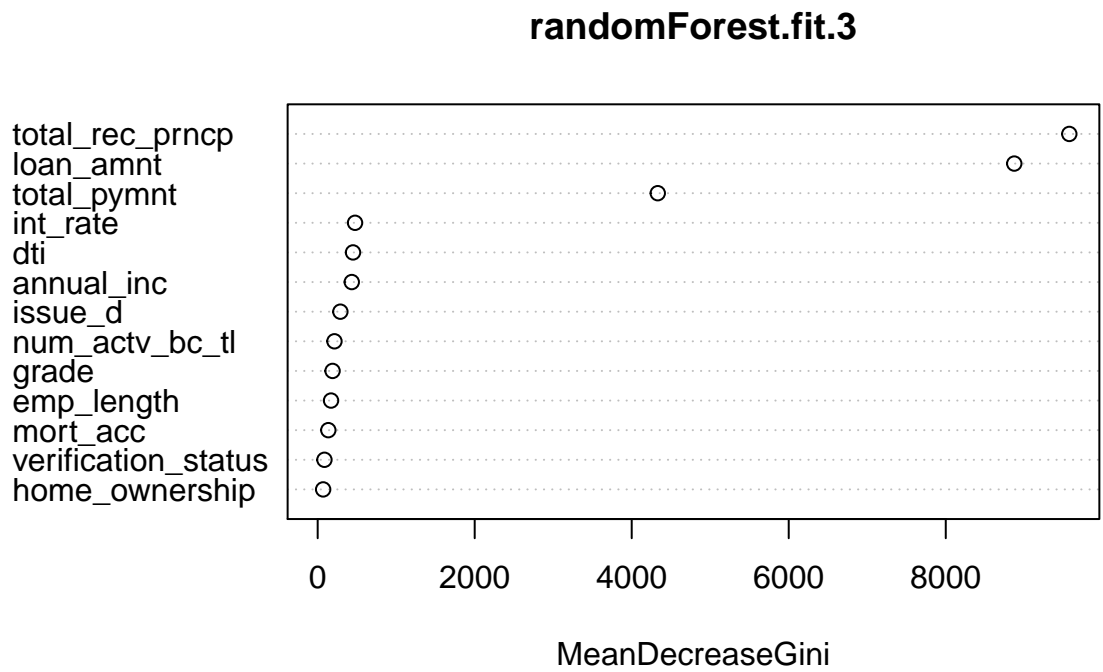



Figure 3: Random Forest - Variable Importance Plot

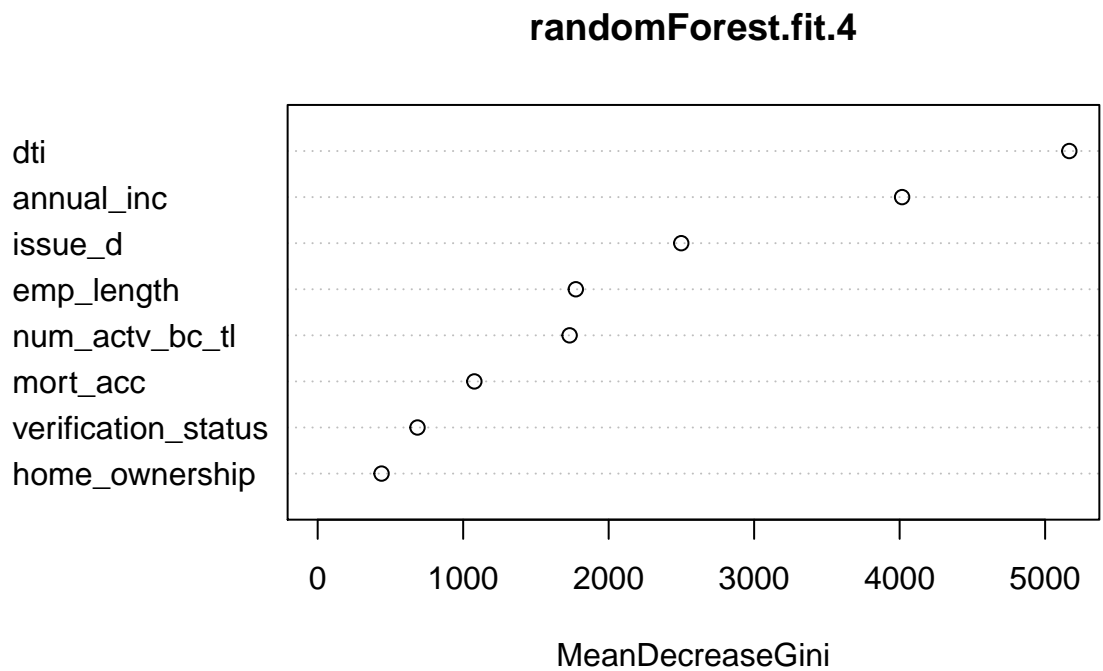
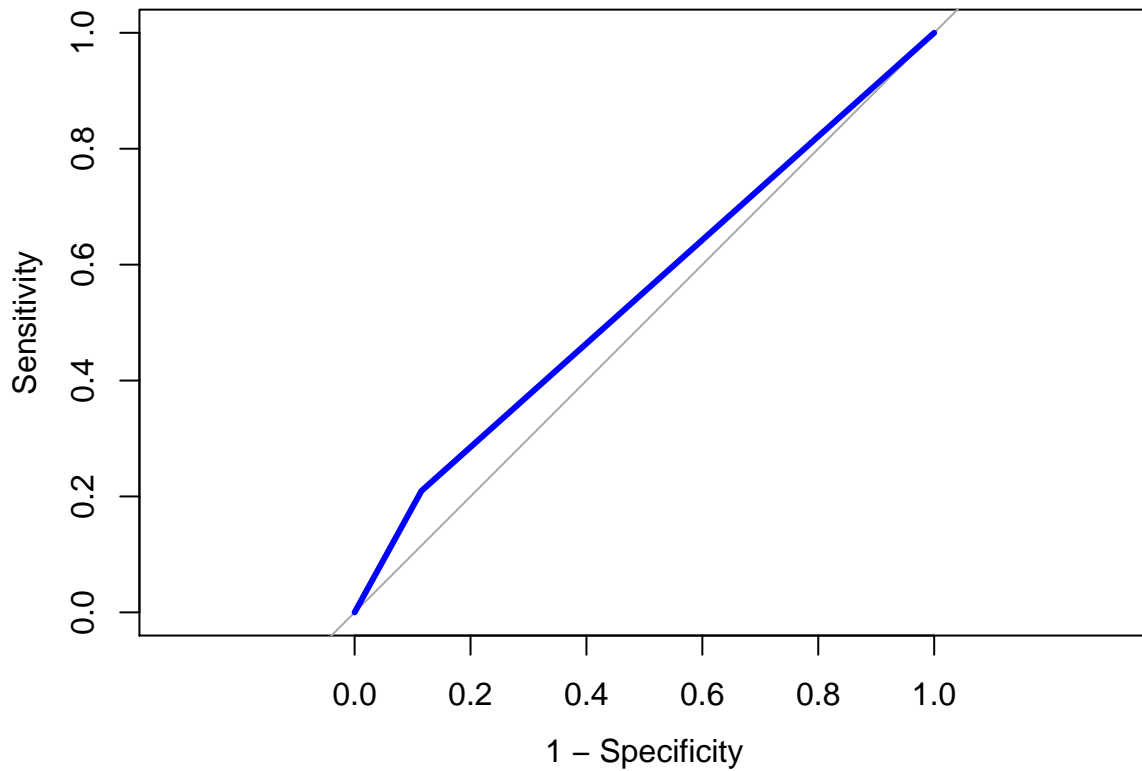


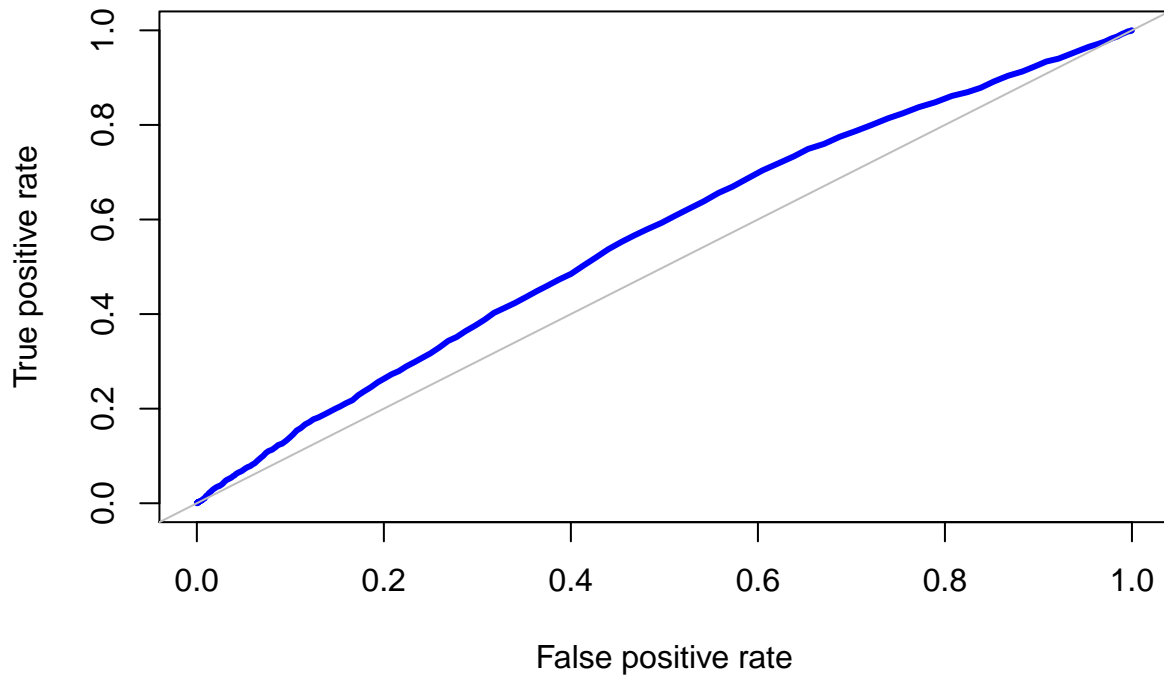
Figure 4: Random Forest - Variable Importance Plot

ROC – Logistic, Model 3

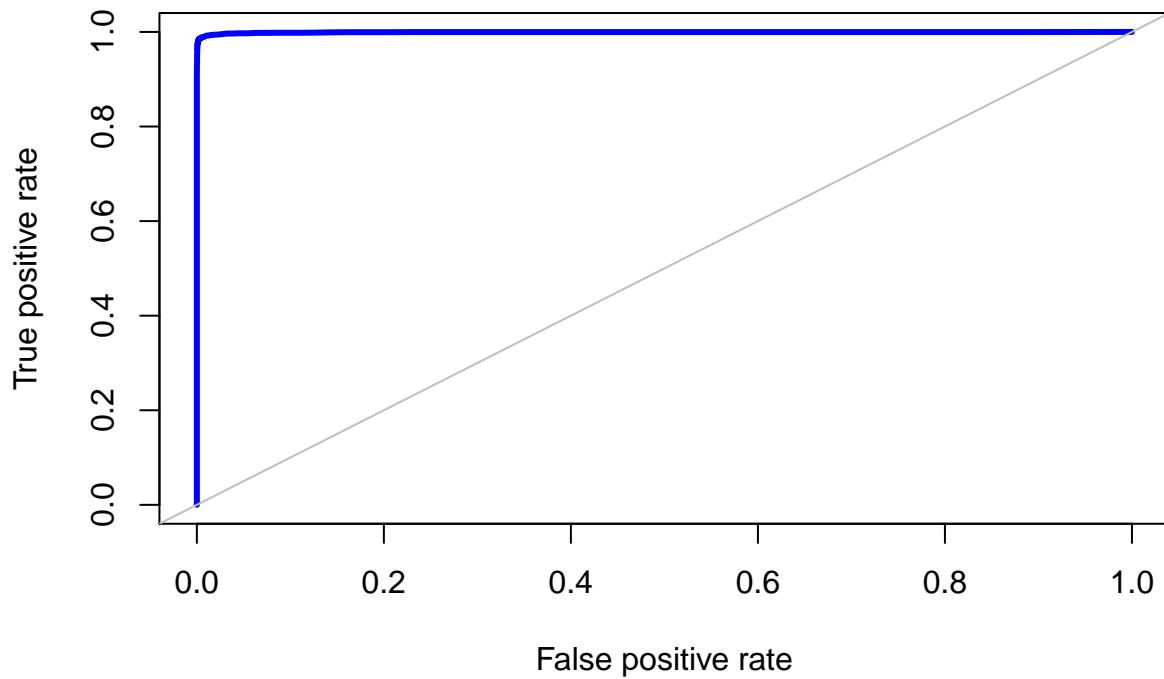


```
##  
## Call:  
## roc.formula(formula = factor(predicted) ~ as.numeric(bad.loan),      data = loans.ss.test, plot = FALSE,  
##  
## Data: as.numeric(bad.loan) in 26271 controls (factor(predicted) 0) < 2277 cases (factor(predicted) 1).  
## Area under the curve: 0.5471  
## 95% CI: 0.5385-0.5557 (DeLong)
```

ROC–Random Forest, Model 4



ROC–Random Forest, Model 3



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. 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.

6 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 and multinomial logistic regression techniques to the model suite.

7 References

Chang, S., Simon and Genki. Predicting Default Risk of Lending Club Loans. 2011. URL: http://cs229.stanford.edu/proj2015/199_report.pdf.

Kaggle. GiveMeSomeCredit. 2011. URL: <https://www.kaggle.com/c/GiveMeSomeCredit>.

Liang, J. Predicting Borrowers Chance Of Defaulting On Credit Loans. Kaggle Submission. 2011. URL: <http://cs229.stanford.edu/proj2011/JunjieLiang-PredictingBorrowersChanceOfDefaultingOnCreditLoans.pdf>.

Mogensen, U., H. Ishwaran and T. Gerds. “Evaluating Random Forests for Survival Analysis”. In: Journal of Statistical Software (2012). URL: <https://www.jstatsoft.org/article/view/v050i11>.

Pandey, J. and M. Srinivasan. Predicting Probability of Loan Default. 2011. URL: <http://cs229.stanford.edu/proj2011/PandeySrinivasan-PredictingProbabilityOfLoanDefault.pdf>.

Reed, P. and Y. Wu. “Logistic regression for risk factor modeling in stuttering research”. In: Journal of Fluency Disorders (2013). URL: https://www.ucl.ac.uk/speech-research-group/publications/publications-files/reed_2013_logistic_regression.pdf.