

Modeling and Predictive Analysis of Credit Card Approval Outcome

ANLY 512 - Final Project 2020

Group 11
Team members:

Feiyang Jiang

Yanou Yang

Xuejun Zhang

Zihao Zhou

Directed By: Prof. Purna Gamage

Introduction

Nowadays, people use credit cards for almost everything, for example, online shopping, paying for electric bills, purchasing concert tickets and so on. Credit card issuers attract applicants by providing appealing benefits such as opening bonus and cash back by using their credit cards. However, creditors also suffer from huge losses if their customers fail to pay off the debts.

Credit risk is a huge challenge to creditors. Currently, the U.S. consumer debt was \$13.7 trillion in Q1 2019, which amounts to approximately \$40.0 trillion in the whole year of 2019. Also in 2019, around 60% of U.S. cardholders between 18 and 34 carried large balances and continually paid late fees, and among them, 8% are in serious delinquency. Moreover, 7 million Americans were at least 90 days behind their car loan payments("Reduce Credit Risk | Brighterion AI | A Mastercard Company", 2020). Creditors suffer from huge losses and need to identify the risks before they occur. This serious situation motivates our team to come up with a statistical model to assist the creditors to identify potential risky credit card applicants.

In this project, our team wanted to explore the personal background information of the applicants to further investigate the relationships between these predictors and the risk level of applicants in order to conduct predictive analysis. We will apply feature engineering to the predictors, which cover applicants' background information such as annual income, number of children, education level, in order to predict the level of risk of credit card applicants. This project will not only help the bank mitigate risk of

approving unqualified applicants, but also assist other financial sectors such as personal loans and mortgage to improve their premium and borrowing models.

This project has four parts. The first part is to define the response variable. According to the Basel Accord in 2007 (Thomas, 2009), users who overdue more than 90 days are considered unable to repay their debt; therefore, in order to control the risk under 3%, we choose to use 60 days overdue as the cutoff point for the users (“high risk” if ≥ 60 , “low risk” < 60). The second part is to clean the predictors since the predictors include categorical variable, binary variable, and continuous variable. The third part will be applying feature engineering, such as LASSO Regression, Principal Component Analysis, Forward and Backward selection, Classification Tree, Random Forest, and Bagging to determine the statistically significant predictors. The last part will be model evaluation. We will split the dataset into testing and training data. Among all models mentioned above, we will compare the test accuracy, recall, precision scores to select the best-fitted model for the dataset.

Analysis

i. Data Collection and Preprocessing

For the data collection of this project, the data we used is from Kaggle website. ("Credit Card Approval Prediction", 2020) The data contains two csv files. The first file has 18 columns/variables and 438558 rows that includes the applicants' application record: ID, Gender, ownership of car, ownership of realty, count of children, income,

income type, education type, family status, housing type, birth days, employed days, whether applicants filed the mobile phone number, whether applicants filed the work phone number, whether applicants filed the phone number, whether applicants filed their emails, occupation, and count of family members. The second file has 3 columns/variables and 1048576 rows that includes applicants' credit record:ID, month balance, and risk status. Month balance means the month of the extracted data is the starting point, backwards. For example, in month balance, "0" is the current month, "-1" is the previous month, and so on. Status explains the applicants' payments overdue situation and it has the following rules: "0" means 1-29 days past due; "1" means 30-59 days past due; "2" means 60-89 days overdue; "3" means 90-119 days overdue; "4" means 120-149 days overdue; "5" means over 150 days overdue which means the bad debt; "C" means paid off that month and "X" means no loan for the month.

To preprocess the two dataset, we first used the largest absolute value in the month balance column from the second data file to find the applicant credit history in month and then created a new column named Begin_Month. Then, we joined two files on the applicants' ID. Rows with missing values are dropped. Next, the target column named risk is created by determining whether the applicant has payment overdue. If a applicant' status shows the applicant has no loan, all loans paid off, and less than 60 days overdue, he will be marked as lower risk applicant which is "0" in the risk column; if a applicant' status shows the applicant has payments overdue 60 days, he will be marked as higher risk applicant. Then, since begin_month and risk columns are created, the month balance and status columns are dropped. Birth days and employed days in

days unit are transformed to ages and employed years in years unit. All the categorical values are converted to numerical values, such as code_gender column in male and female are transformed to “1” as female and “0” as male. In addition, some numerical columns’ values are distributed too widely, so bins are created to replace the original columns. For example, the income column varies from \$27000 to \$1575000, so there are three bins in \$0 - \$40k, \$40k - \$100k, \$100k - \$2000k corresponding to “1”, “2”, “3” created. Finally, it is found that the dataset is very imbalanced, because there are only 2957 lower risk applicants and 24, 404 higher risk applicants. The synthetic minority oversampling technique abbreviated as SMOTE is used to oversampled and undersampled the lower risk applicants rows and higher risk applicants rows. In the end, the dataset the team used has 31 columns/variables and 20700 rows.

ii. Methodology

SMOTE

As mentioned in the data preprocessing section, SMOTE is used to undersample the lower risk applicants and oversample the higher risk applicants. In SMOTE, undersampling the majority class is intuitive, and oversampling the minority class uses k-nearest neighbors to control the way the new examples are created. ("SMOTE function | R Documentation", 2020) Precisely, in the dataset, for each originally existing lower risk applicant class example X, new examples will be created. These new examples will

be generated using the information from the k nearest neighbors of each lower risk applicant. In this project, k of 5 is used.

K-means Clustering

Clustering seeks a partition of the data into distinct groups so that the observations within each group are quite similar to each other. K-means Clustering seeks to partition the observations into a pre-specified number of clusters.

The idea behind K-means clustering is that a good clustering is one for which the within-cluster variation is as small as possible.

Details of K-means clustering algorithm(James, Witten, Hastie and Tibshirani, n.d.):

Algorithm: K-Means

1. Randomly assign a number, from 1 to K , to each of the observations. These serve as initial cluster assignments for the observations.

2. Iterate until the cluster assignments stop changing:

2.1 For each of the K clusters, compute the cluster centroid. The K th cluster centroid is the vector of the p feature means for the observations in the k th cluster.

2.2 Assign each observation to the cluster whose centroid is closest (where closest is defined using Euclidean distance).

K-means clustering assumes spherical shapes of different clusters of data and works badly at non-spherical shapes of clusters. But K-means provide a good sketch to better understand the inner relationship of the dataset.

PCA

PCA produces a low-dimensional representation of a dataset. It finds a sequence of linear combinations of the variable that have maximal variance, and are mutually uncorrelated. Through it, we can directly decrease the number of feature variables, thereby narrowing down the important features and saving on computations.

PCA mainly has three steps:

1. Compute the covariance matrix of the data
2. Compute the eigenvalues and eigenvectors of the covariance matrix
3. Use the eigenvalues and eigenvectors to select only the most important principle components and transform your data into uncorrelated linear combinations to reduce dimensionality.

PCA not only can do dimension reduction, but also can provide good visualization and understanding of the dataset. To this dataset, PCA primarily works as this job.

Forward Selection and Backward Selection

Forward selection is a type of stepwise regression which has 4 steps:

1. Start with no variables in the model
2. Test the addition of each variable using a chosen model fit criterion
3. Add the variable (if any) whose inclusion gives the most statistically significant improvement of the fit
4. Repeat this process until none improves the model to a statistically significant extent.

Backward selection is a type of stepwise regression which has 4 steps:

1. Involve starting with all candidate variables
2. Testing the deletion of each variable using a chosen model fit criterion
3. Deleting the variable (if any) whose loss gives the most statistically insignificant deterioration of the model fit
4. Repeating this process until no further variables can be deleted without a statistically insignificant loss of fit.

In this case, forward selection and backward selection are used to select the best subset of predictors to simplify our model and prevent overfitting. C_p , BIC and adjusted R^2 are used as the criteria of these two procedures.

LASSO

Lasso is a regression analysis method that improves the prediction accuracy and interpretability of regression models by altering the model fitting process to select only a subset of the provided covariates.

Lasso uses the l1 penalty to force the sum of the absolute value of the regression coefficients to be less than a fixed value, which shrinks certain coefficients to zero, so we can choose a simpler model that does not include those coefficients. Before using lasso, we scale the variables.

Logistic Regression

Logistic regression is a kind of generalized linear model (GLM) that uses a logistic function as its link function to model a variable with multinomial distribution. Here the predictors selected by LASSO are used to build the model.

Generalized Additive Model

Generalized additive model (GAM) is a generalized linear model in which the linear predictor depends linearly on unknown smooth functions of some predictor variables, and interest focuses on inference about these smooth functions. Also, the predictors selected by LASSO are used to build the model. Compared with the logistic regression model above, all the predictors other than the binary predictors `FLAG_WORK_PHONE` and `FLAG_PHONE` are replaced by smoothing splines.

Random Forest

Random forest is one of the tree based models which can eliminate the basic tree models' problem in high variance, especially when the dataset's ratio of n (number of observations) by p (number of variables/dimension) is large. Because the dataset has many categorical values, the random forest classification model is chosen. The random forest classification model performs following steps below (James, Witten, Hastie and Tibshirani, n.d.):

Algorithm: Random Forest in Classification

1. Use bootstrap method to take B repeated samples from the training dataset
 2. Create B decision trees using the B bootstrapped samples from the last step, numbers of variables randomly sampled as candidates at each split and numbers of sampled trees need to be determined in this step. Notice here random forest is not using all of the dataset.
 3. For the test dataset, record the class predicted by each of the B decision trees and take the majority vote: the overall prediction which is the most commonly occurring majority class among the B prediction will be the final random forest model.
-

Bagging is a special case of random forest when the number of variables randomly sampled at each split (m) in trees equals the number of variables (p) in the dataset. Normally, in random forest classification model, half of number of variables

($p/2$) and square root of number of variables (\sqrt{p}) are chosen as the number of variables randomly sampled at each split to create trees.

Results

i. EDA

Exploratory data analysis is performed to get the general pattern of some variables. Figure 1 shows the education type of the applicants distributed in our dataset. Education types “1”, “2”, “3”, “4”, “5” stand for “Incomplete Secondary”, “Secondary/secondary special”, “Incomplete higher”, “Higher education” and “Academic degree”. Figure 1 explains that most of the applicants have secondary/secondary special and higher education. Only a few applicants’ education level is in incomplete secondary and academic degrees.

Figure 2 shows the applicants’ age distribution, most of applicants are between 29-50 years, and fewer applicants are under 29 years.

Figure 3 shows the applicants’ occupation distribution, most of applicants’ work as laborers, managers, drivers, sales staff etc. The occupation types distributions look generally balanced.

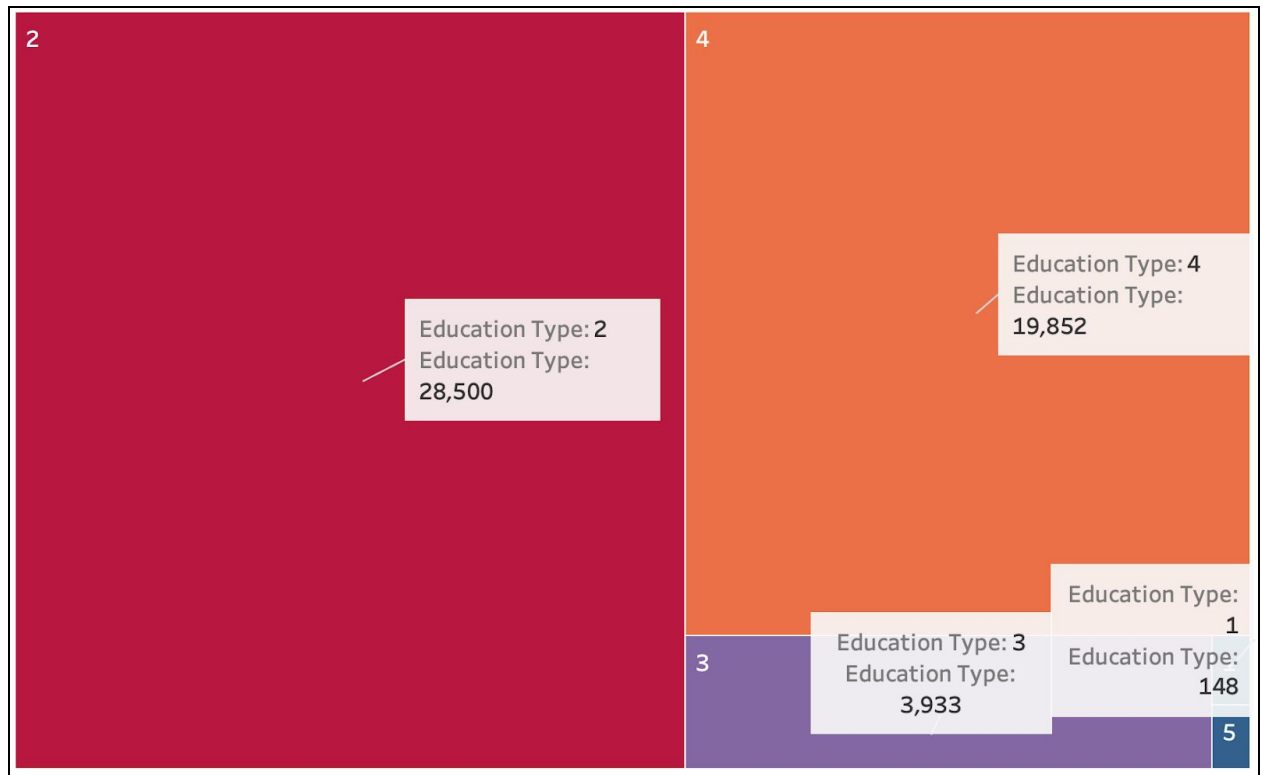


Figure 1: Applicants' Education Types Distribution

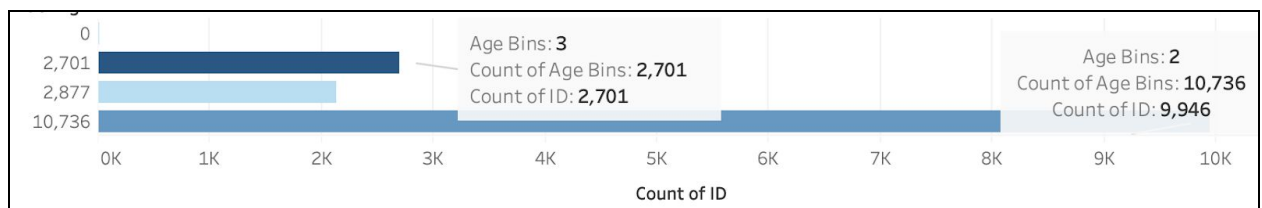


Figure 2: Applicants' Ages Distribution

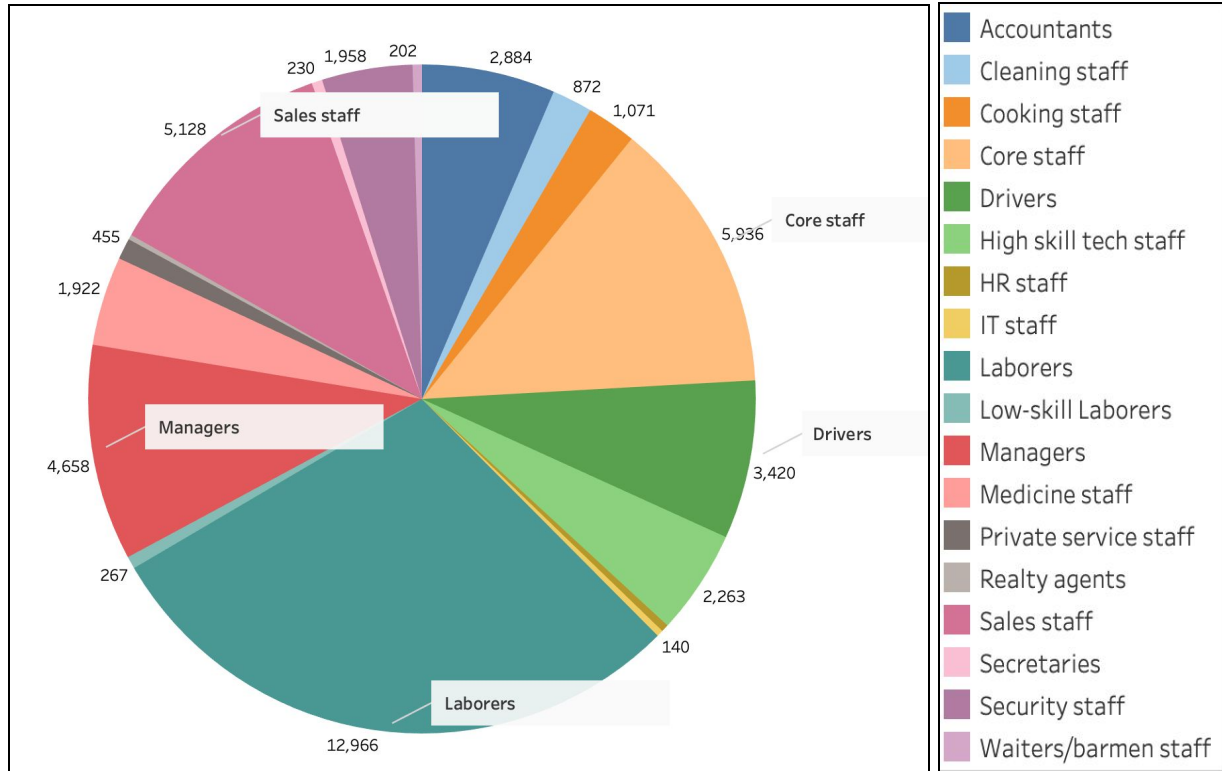


Figure 3: Applicants' Occupations Distribution

ii. Unsupervised Learning

K-means Clustering

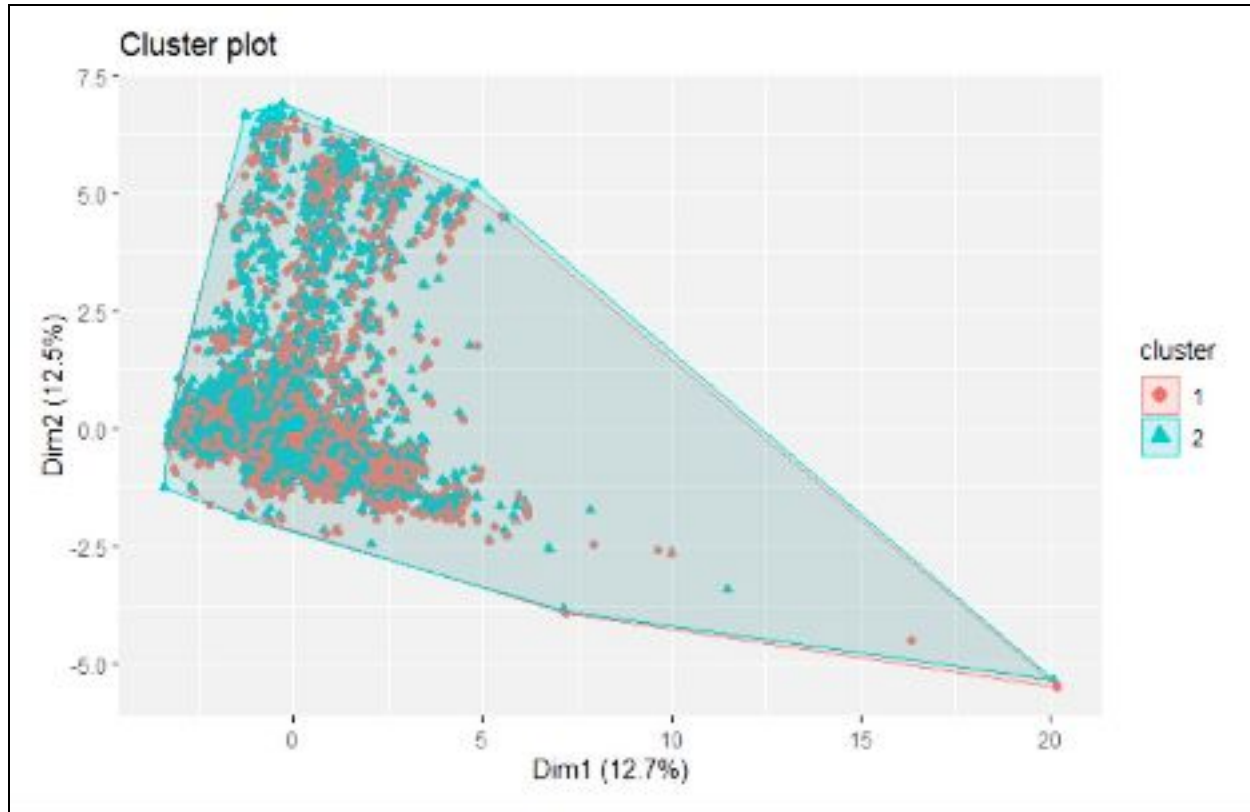


Figure 4: K-means Clustering in 2D

K-means Clustering results are plotted in the above graph, red point is cluster 1 and blue point is cluster 2. Points within cluster 1 are the most similar to each other, which means this kind of applicant has the same features, so as cluster 2.

We applied K-means Clustering with $K = 2$. From the plot above, it is observed that two clusters are overlapped with each other and hard to be separated by clustering. It means applicants even in different clusters still have many same features, which

prevents K-means to segment the applicant from non-risky to risky. Besides, the two clusters are not-circle shaped either.

In sum, after applying the K-means Clustering method, the result shows that we need supervised learning to determine if a customer is risky.

PCA

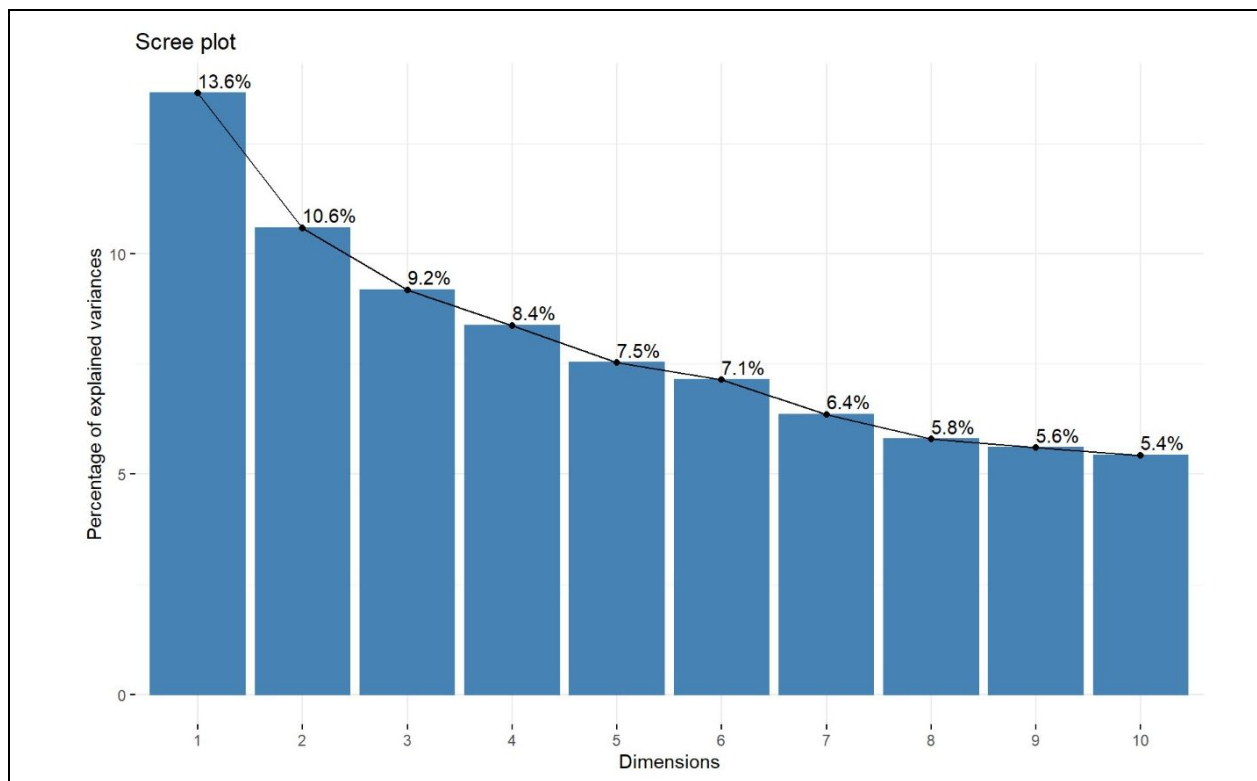


Figure 5: Percentage of Explained Variables vs. Dimensions

The above graph is the histogram of the percentage of explained variables and dimensions. From the graph, the first dimension only explains 13.6% variance. The histogram shows that even the first two principles can only explain 20% variance of the predictors, which would tell the overall movement in the original dataset. In other

words, if we use principal components to determine the risk, then it would include at least 5 principal components, we are unable to reduce the original dimensions.

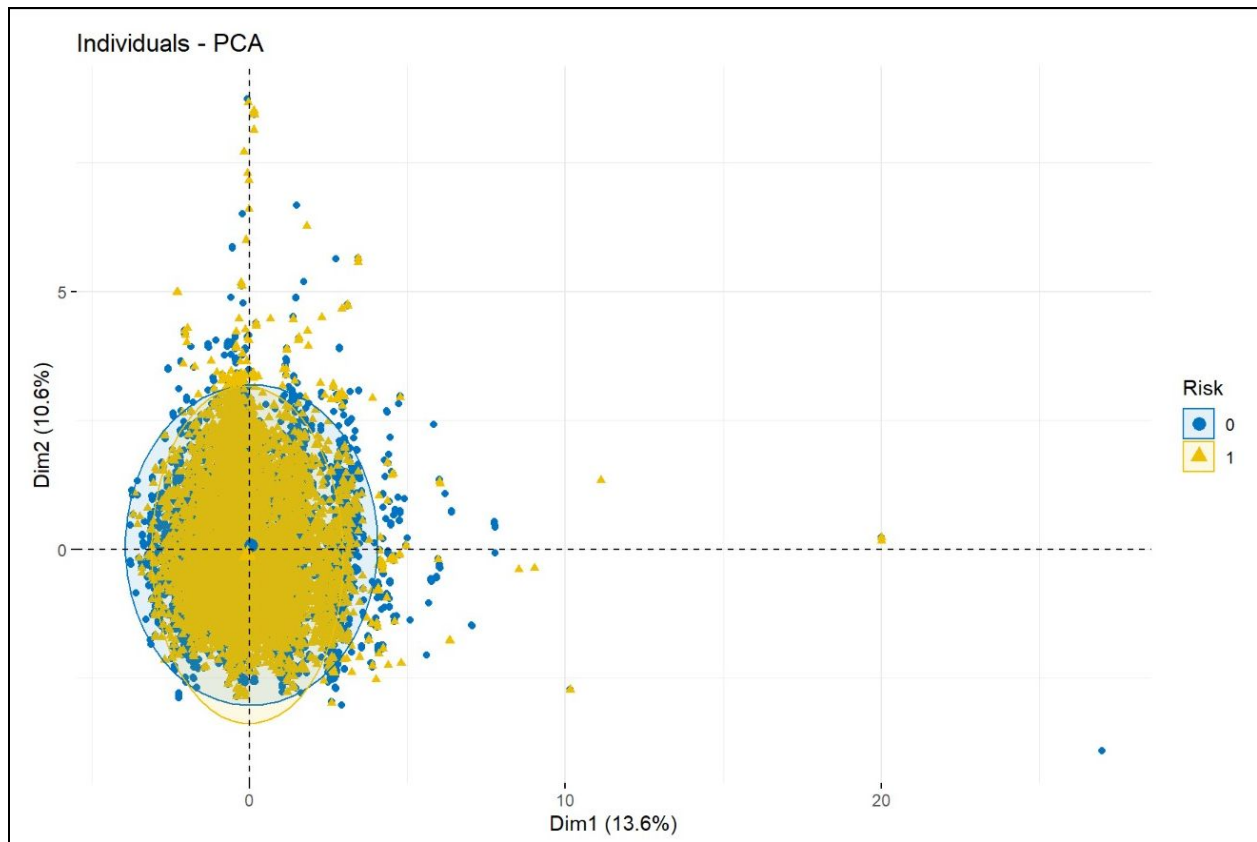


Figure 6: Individuals embedding on 2D

The above graph shows individuals embedding into 2D after PCA. From the plot, yellow points and blue points are closely clustered and overlapped.

By embedding the first two principal components into 2D plot, risk 0 and risk 1 points are overlapped as well, indicating whether being risky or not is determined by multiple principal components.

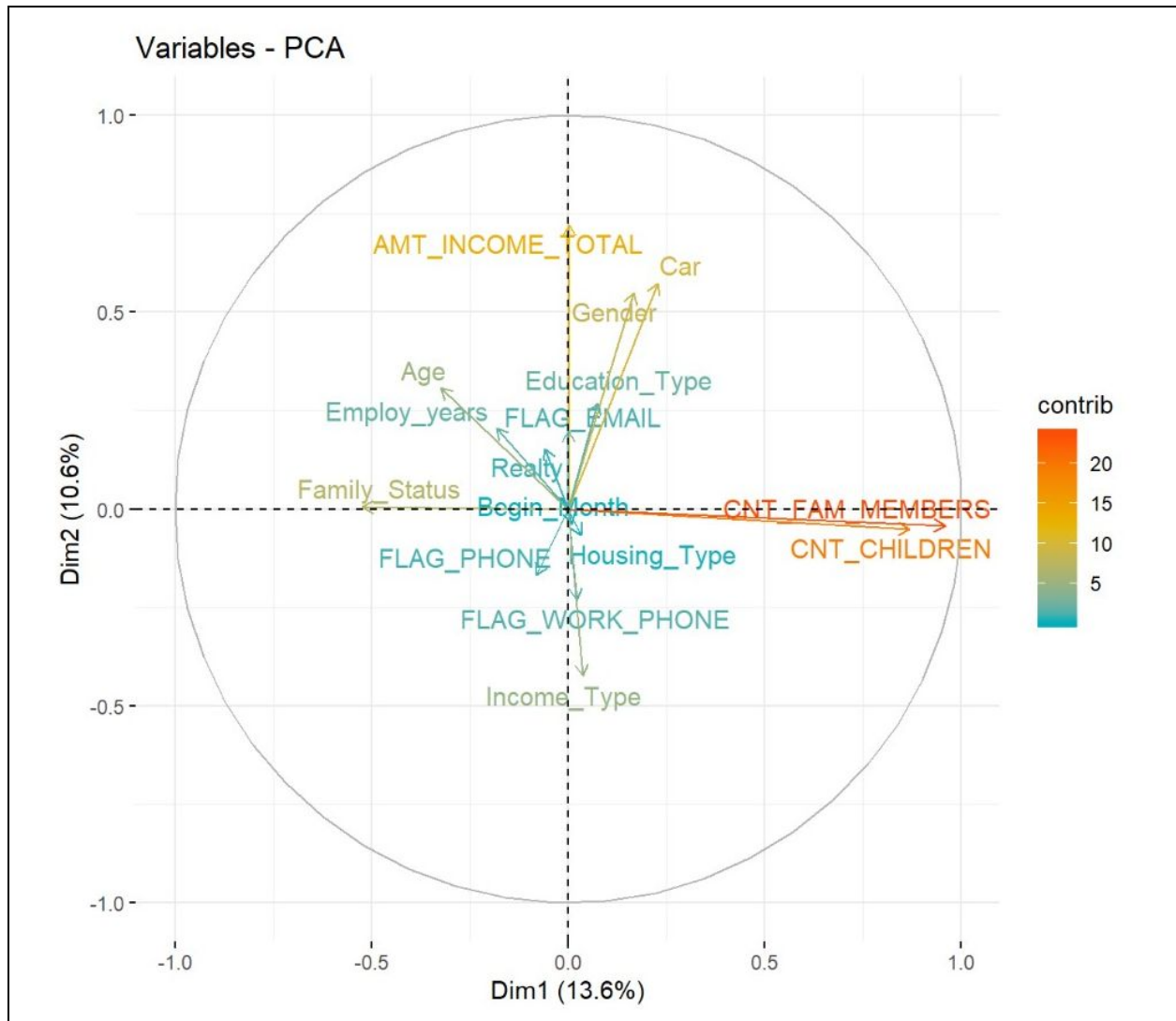


Figure 7: Variables embedding on 2D

This graph shows which variable contributes the most to explain the variance. It's obvious to find most of the variables only contribute less than 10% to the total variance.

PCA results show that variables like the number of children, the number of family members, and AMT account balance explain the most variance, which can be paid more attention to.

iii. Feature Selection

Forward Selection and Backward Selection

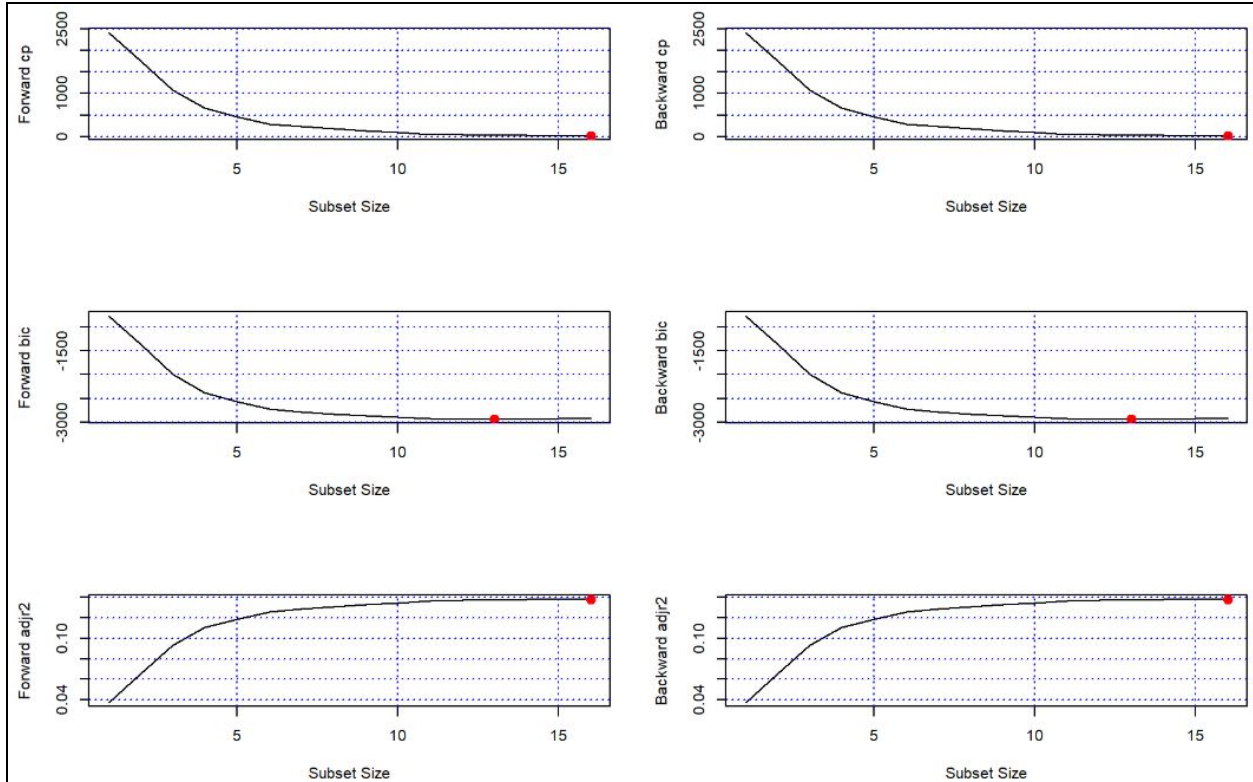


Figure 8: The forward and backward selection

The C_p and adjusted R^2 criteria suggest that we should choose a subset that includes all 16 predictors, while the BIC criteria suggests that we should choose a subset with size 14. All the criteria suggest a large subset to reach optimum, which is consistent with the finding of unsupervised learning. However, the plot shows that when the subset size is 6, the criteria is already good.

LASSO Regression

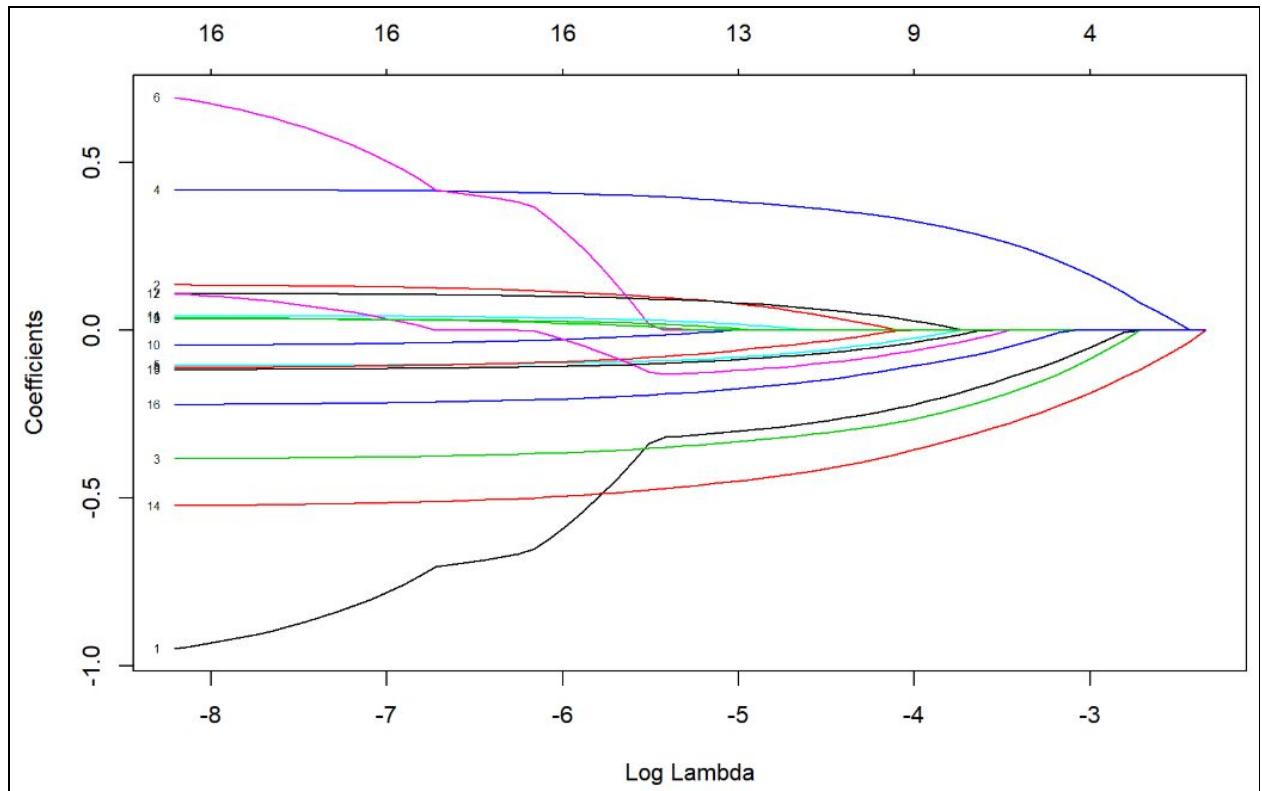


Figure 9: The coefficient selected by LASSO

At $\log(\lambda) = -3.5$, the LASSO selects 6 variables “Count of Children”, “Filling in Work Phone Number”, “Filling in Phone Number”, “Ages”, “Education Type” and “Family Status” as important features to determine whether the credit cards applicant has high risk, the selection is the same with forward and backward selection.

iv. Supervised Learning

Logistic Regression

Coefficients:					
	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	2.69080	0.11138	24.159	<2e-16	***
CNT_CHILDREN	-0.45058	0.02602	-17.317	<2e-16	***
FLAG_WORK_PHONE	-0.99548	0.04571	-21.776	<2e-16	***
FLAG_PHONE	0.98741	0.03993	24.727	<2e-16	***
Family_Status	-0.19087	0.02142	-8.911	<2e-16	***
Age	-0.04979	0.00191	-26.064	<2e-16	***
Education_Type	-0.22425	0.01977	-11.340	<2e-16	***

Figure 10: Coefficients for logistic regression

The p-values are all less than the significance level 0.05, suggesting that the risk is significantly correlated with all the predictors selected.

The coefficient shows that:

The applicants that are currently married have higher risks of arrear.

The applicants with more children have lower risks of arrear.

The applicants that provide work phone numbers have lower risks of arrear.

The applicants that provide phone numbers have higher risks of arrear.

The older applicants have lower risks of arrear.

The applicants with higher education level have lower risks of arrear.

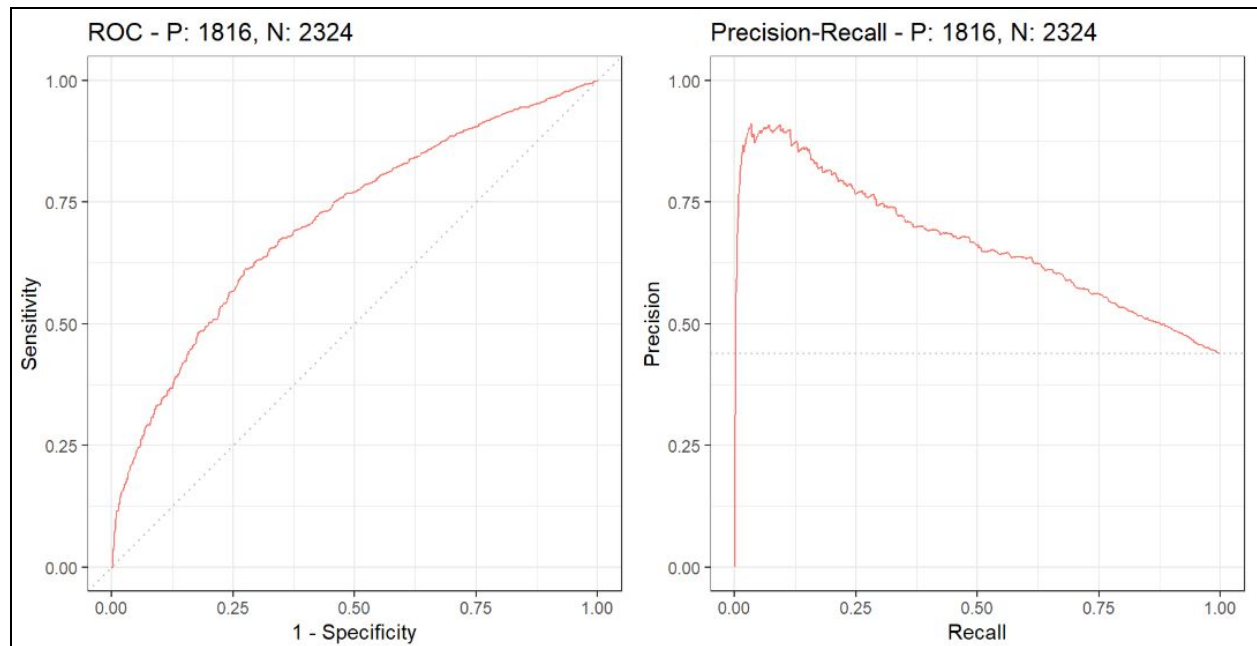


Figure 11: The ROC and precision-recall curve for Logistic Regression

The ROC and precision-recall curve show a tradeoff between true positive rate and false positive rate and between precision and recall. The precision score here is the proportion of the applicants that truly have risks among the predicted applicants with risks, and recall is the proportion of applicants with risks that are correctly identified by the model among the applicants that truly have risks. In this case, our goal is to identify the applicants with risks, therefore, we are more concerned about the recall in the precision-recall tradeoff.

Under a 0.4 threshold, the recall is 69.93%, and the precision is 57.89%.

Generalized Additive Model with Smoothing Spline

Anova for Parametric Effects						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
s(CNT_CHILDREN)	1	213.3	213.29	211.727	< 2.2e-16	***
FLAG_WORK_PHONE	1	186.0	185.97	184.604	< 2.2e-16	***
FLAG_PHONE	1	536.9	536.88	532.947	< 2.2e-16	***
s(Family_Status)	1	47.9	47.92	47.572	5.494e-12	***
s(Age)	1	560.7	560.67	556.568	< 2.2e-16	***
s(Education_Type)	1	128.2	128.24	127.303	< 2.2e-16	***
Residuals	16540	16662.0	1.01			

Figure 11: The ANOVA table of Generalized Additive Model with Smoothing Spline

The p-values of ANOVA are all less than the significance level 0.05, suggesting that all the predictors selected have significant influence on the risk.

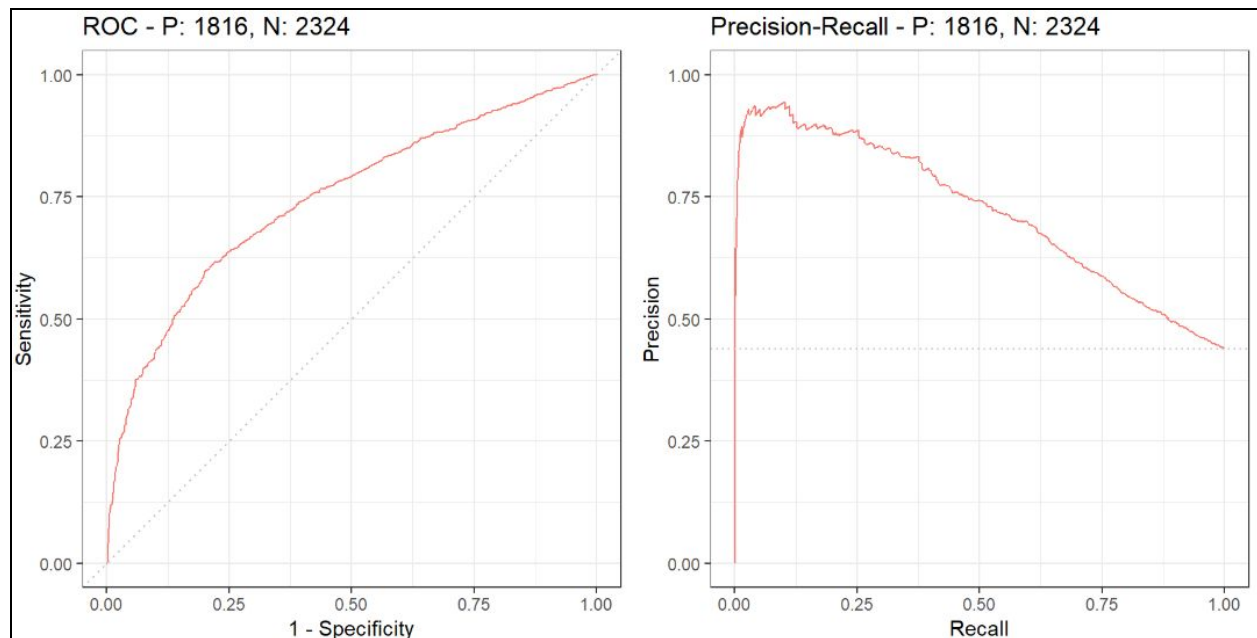


Figure 12: The ROC and precision-recall curve for GAM model with Smoothing Splines

The above plots are the ROC and precision-recall curve, which shows that recall and precision of the generalized additive model is a little higher. For a 0.4 threshold, the recall is 71.20%, and the precision is 60.51%, which is a little better. The logistic regression failed to give a higher recall probably because there are still some binary categorical predictors in this model, although they can be transformed into dummy variables, the performance of the model will be influenced.

Random Forest

Three random forests models are performed to the dataset. 17 variables are used as predictors. The response variable is the risk level of the applicants. Predictors include count of children, count of family members, ages, genders and so on. As shown in Table 1, the bagging model is performed when the number of variables randomly sampled at each split (Mtry) equals to the number of predictors 17. The train recall, precision and accuracy are 83.51%, 70.22% and 81.40% respectively. As the Mtry is changed to 17/2, the random forest model has train recall, precision and accuracy of 85.69%, 70.39% and 82.34% respectively. As the Mtry is changed to $\sqrt{17}$, 4 variables randomly sampled at each split in sample trees and the random forest model has train recall, precision and accuracy of 93.43%, 68.40% and 84.49% respectively. In Table2, The test recall, precision and accuracy are 85.67%, 71.75% and 83.09% respectively. As the Mtry is changed to 17/2, the random forest model has test recall, precision and accuracy of 85.61%, 72.25% and 83.00% respectively. As the Mtry is changed to $\sqrt{17}$, 4 variables randomly sampled at each split in sample trees and the random forest model has test

recall, precision and accuracy of 94.48%, 70.04% and 84.93% respectively. Test data recall, precision and accuracy are all higher than train data recall, precision and accuracy, so there's no evidence of overfitting. Model with the highest recall is selected which is highlighted in Table 1 and Table 2. Table 3 shows the optimal model the team chose in r, and 40,000 of trees sampled, to ensure that every input row gets predicted at least a few times. Recall in 94.28% means that about 94% of applicants with actual risk can be captured by model. Precision in 70.93% means that about 70% of applicants with true high risk situations under the higher risks predicted class. Accuracy in 84.93% means about 85% of applicants are predicted correctly in the model.

Table 1: Different Trees Models Train Results by Setting the Same Ntree = 500

Model	Mtry	Train Recall	Train Precision	Train Accuracy
Bagging	17	83.51%	70.22%	81.40%
Random Forest	17/2	85.69%	70.39%	82.34%
Random Forest	sqrt(17)	93.43%	68.40%	84.49%

Table 2: Different Trees Models Test Results by Setting the Same Ntree = 500

Model	Mtry	Test Recall	Test Precision	Test Accuracy
Bagging	17	85.67%	71.75%	83.09%

Random Forest	17/2	85.68%	72.25%	83.00%
Random Forest	sqrt(17)	94.48%	70.04%	84.93%

Table 3: The Final Optimal Random Forest Model Settings and Result

Model	Ntree	Mtry	Test Recall	Test Precision	Test Accuracy
Random Forest	40,000	sqrt(17)	94.28%	70.93%	84.95%

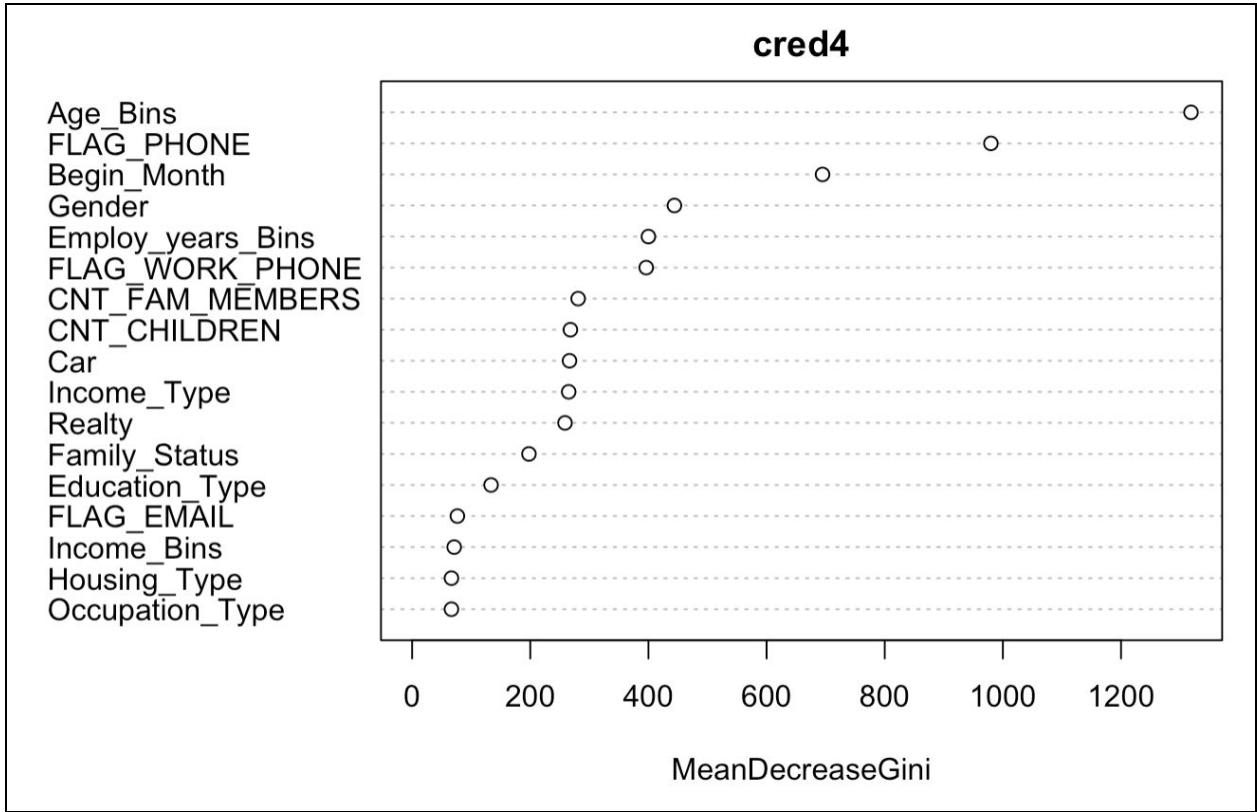


Figure 13: Feature Importance Plot of the Final Optimal Random Forest Model

From the feature importance plot of the final optimal random forest model, applicants' ages, we noticed whether filing the phone numbers, credit history in months, genders and employed years are the top important variables in classifying whether the applicant has high risk or not.

Conclusions

In conclusion, the statistically significant features we selected are *education level*, *whether filled in work phone number*, *family status*, *number of children*, and *employment history*, and *age*. We have four major findings from the above modeling and analysis. First of all, we would like to advise creditors to be cautious with applicants with higher education levels. This is a surprising finding because we would consider people with higher education levels to earn more; therefore they are more likely to repay their debts. However, according to the article, “*1 million people default on their student loans each year*” by CNBC, “Within four years after leaving school, nearly a quarter of the borrowers had defaulted”, and this article also found that those people are in financial distress (Nova, 2020). Although people with higher education may have the ability to repay their debts, the student loans they have when they pursue higher education make them more vulnerable to financial burdens.

Secondly, we would like to advise creditors to be cautious with applicants with short employment history. Similar to the reasons stated in the previous finding that applicants with short employment history may suffer from problems with low income,

short credit history, unpaid student loans (Nova, 2020). Therefore, they have a higher probability of defaulting compared to people who have accumulated wealth for several years.

Thirdly, applicants who provide work phone numbers are found to be low-risk. This finding is intuitive because if the applicants provide work phone numbers, then it indicates that they have a stable job and consistent income every month; therefore, they are more likely to fulfill their debt responsibility.

Lastly, we would like to advise creditors to be cautious with applicants with no or few children. This is a surprising finding as well because we would think families with more children would have more financial burdens and less likely to be able to repay their debts. According to the article, “*Families With Young Kids Are More Likely to Live in Poverty*”, it states that families with few children tend to be young parents(Ryan, 2020). Due to the lack of workplace support and income to young parents and lack of experience in raising young children, they are more likely to be in poverty(Ryan, 2020). Therefore, they are less likely to repay debts on time.

In all, we hope this project would help creditors and financial sectors such as mortgage and personal loans to mitigate risk of approving unqualified candidates by analyzing the background information they provide in the application process.

Discussion

We faced several challenges in completing this project. First of all, the original dataset is highly unbalanced that only 10% of applicants are low-risk and nearly 90% of

applicants are high-risk. When we first fit the logistic model to the data, the accuracy is as high as 87%, but both recall and precision scores are very low. Therefore, we had to resample the data to achieve better recall and precision scores. The second difficulty we faced was to choose the number of predictors. Although the forward and the backward selection chose 14-17 predictors, we decided to build a simpler model with 6 predictors because the Cp scores were not improved much after adding predictors. Interpreting the model is also challenging because some of the findings are counterintuitive. We researched through a good amount of related works and academic papers to understand why our results align with the reality.

The next steps for this project are to reassess the risk level of applicants by incorporating more factors such as demographic background, whether the applicants have car loans, how many credit cards do applicants already have etc... We will also be cautious in overfitting the model because adding more predictors are not necessarily better in predicting the outcomes. We will continue to research more about creditors' current risk models because having risky customers can also be profitable for creditors since they will pay late fees and penalties, so we are curious how creditors distinguish the applicants who default and applicants who are risky but profitable to creditors.

References

James, G., Witten, D., Hastie, T. and Tibshirani, R., n.d. *An Introduction To Statistical Learning*.

Reduce Credit Risk | Brighterion AI | A Mastercard Company. (2020). Retrieved 29 April 2020, from https://brighterion.com/reduce-credit-risk-with-ai/?gclid=CjwKCAjw4pT1BRBUEiwAm5QuRyL1HWOW16l8gDFwNFxyXfZl3evYmB-Chxo3fml7buNudCbRtoSWJBoCO88QAvD_BwE

Nova, A. (2020). *More than 1 million people default on their student loans each year*. Retrieved 29 April 2020, from <https://www.cnbc.com/2018/08/13/twenty-two-percent-of-student-loan-borrowers-fall-into-default.html>

Ryan, L. (2020). *Families With Young Kids Are More Likely to Live in Poverty*. Retrieved 29 April 2020, from <https://www.thecut.com/2016/12/families-with-young-children-have-lower-incomes-report.html>

Credit Card Approval Prediction. (2020). Retrieved 29 April 2020, from https://www.kaggle.com/rikdifos/credit-card-approval-prediction#application_record.csv

SMOTE function | R Documentation. (2020). Retrieved 29 April 2020, from <https://www.rdocumentation.org/packages/DMwR/versions/0.4.1/topics/SMOTE>

project

YanouY

2020/4/23

```
data = read.csv("CleanV2.csv")

#install.packages('kernlab')
#install.packages('e1071')
#install.packages('factoextra')

library(tidyverse) # data manipulation and visualization

## Warning: package 'tidyverse' was built under R version 3.6.3

## -- Attaching packages ----- tidyverse 1.3.0 --

## √ ggplot2 3.2.1    √ purrr  0.3.3
## √ tibble  3.0.1    √ dplyr  0.8.4
## √ tidyr   1.0.2    √ stringr 1.4.0
## √ readr   1.3.1    √ forcats 0.5.0

## Warning: package 'tibble' was built under R version 3.6.3

## Warning: package 'tidyr' was built under R version 3.6.3

## Warning: package 'readr' was built under R version 3.6.3

## Warning: package 'forcats' was built under R version 3.6.3

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

library(gridExtra) # plot arrangement

## Warning: package 'gridExtra' was built under R version 3.6.3

##
## Attaching package: 'gridExtra'

## The following object is masked from 'package:dplyr':
##
##   combine

library(ggplot2)
library(kernlab) # SVM methodology

##
## Attaching package: 'kernlab'

## The following object is masked from 'package:purrr':
##
##   cross

## The following object is masked from 'package:ggplot2':
##
##   alpha

library(e1071) # SVM methodology

## Warning: package 'e1071' was built under R version 3.6.3

library(cluster) # clustering algorithms

## Warning: package 'cluster' was built under R version 3.6.3

library(factoextra) # clustering algorithms & visualization

## Warning: package 'factoextra' was built under R version 3.6.3

## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

K-means

```
use <- data[, c('CNT_CHILDREN', 'FLAG_WORK_PHONE', 'FLAG_PHONE', 'FLAG_EMAIL', 'CNT_FAM_MEMBERS', 'Begin_Month', 'Gender', 'Car', 'Realty', 'Income_Type', 'Family_Status', 'Housing_Type', 'Occupation_Type', 'Income_Bins', 'Age_Bins', 'Employ_years_Bins', 'Education_Type', 'Income_Type')]
y <- data[, 'Risk']

smp_size <- floor(0.5 * nrow(use))
sample_ <- sample(1:length(y), size = smp_size)
#use_ <- use[sample_, c('Income_Type', 'Family_Status', 'Begin_Month', 'Occupation_Type', 'Housing_Type')]
use_ <- use[sample_, ]
distance <- get_dist(t(use_))
fviz_dist(distance, gradient = list(low = "#00AFBB", mid = "white", high = "#FC4E07"))

km.out=kmeans(use,2,nstart=15)
cluster_ <- km.out$cluster - 1

# Precision
precision <- length(cluster_[cluster_*y == 1])/(length(cluster_[cluster_-y == 1])+length(cluster_[cluster_*y == 1]))
precision

## [1] 0.4552146

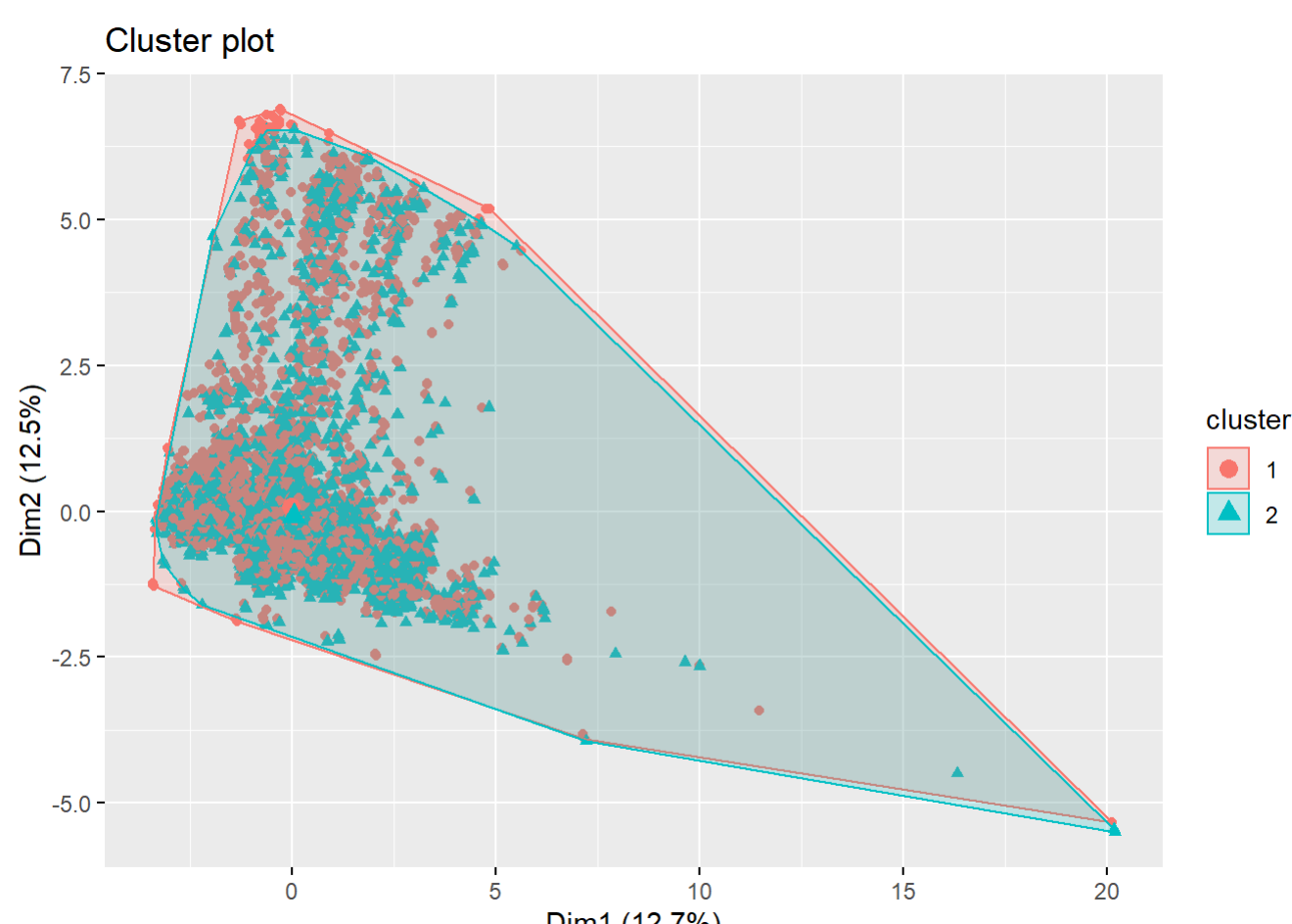
# Recall
recall <- length(cluster_[cluster_*y == 1])/(length(cluster_[cluster_-y == 1])+length(cluster_[cluster_*y == -1]))
recall

## [1] 0.8355852

km.out$centers

##      CNT_CHILDREN FLAG_WORK_PHONE FLAG_PHONE FLAG_EMAIL CNT_FAM_MEMBERS
## 1      0.4328588      0.2289701      0.3584082      0.09861483      2.228272
## 2      0.4355312      0.2322480      0.3619631      0.09120983      2.267432
## Begin_Month Gender      Car Realty Income_Type Family_Status Housing_Type
## 1    15.17483 1.379498 1.411644 1.680733      3.648549      2.306795      2.315027
## 2    40.18138 1.328254 1.406175 1.643242      3.793629      2.243644      2.259298
## Occupation_Type Income_Bins Age_Bins Employ_years_Bins Education_Type
## 1      2.315027      2.985183 1.911349      1.623254      4.161175
## 2      2.259298      2.911223 1.911060      1.626005      4.136469

fviz_cluster(km.out, geom = "point", data = use)
```



H-clust

```
d <- dist(use, method = "manhattan")
hc1 <- hclust(d, method = "complete")
plot(hc1, cex = 0.6, hang = -1)

Cluster Dendrogram

Height
100
80
60
40
20
0

d
hclust(*, "complete")

sub_grp <- cutree(hc1, k = 2)
table(sub_grp)

## sub_grp
##      1      2
## 15447  5252

result = sub_grp-1

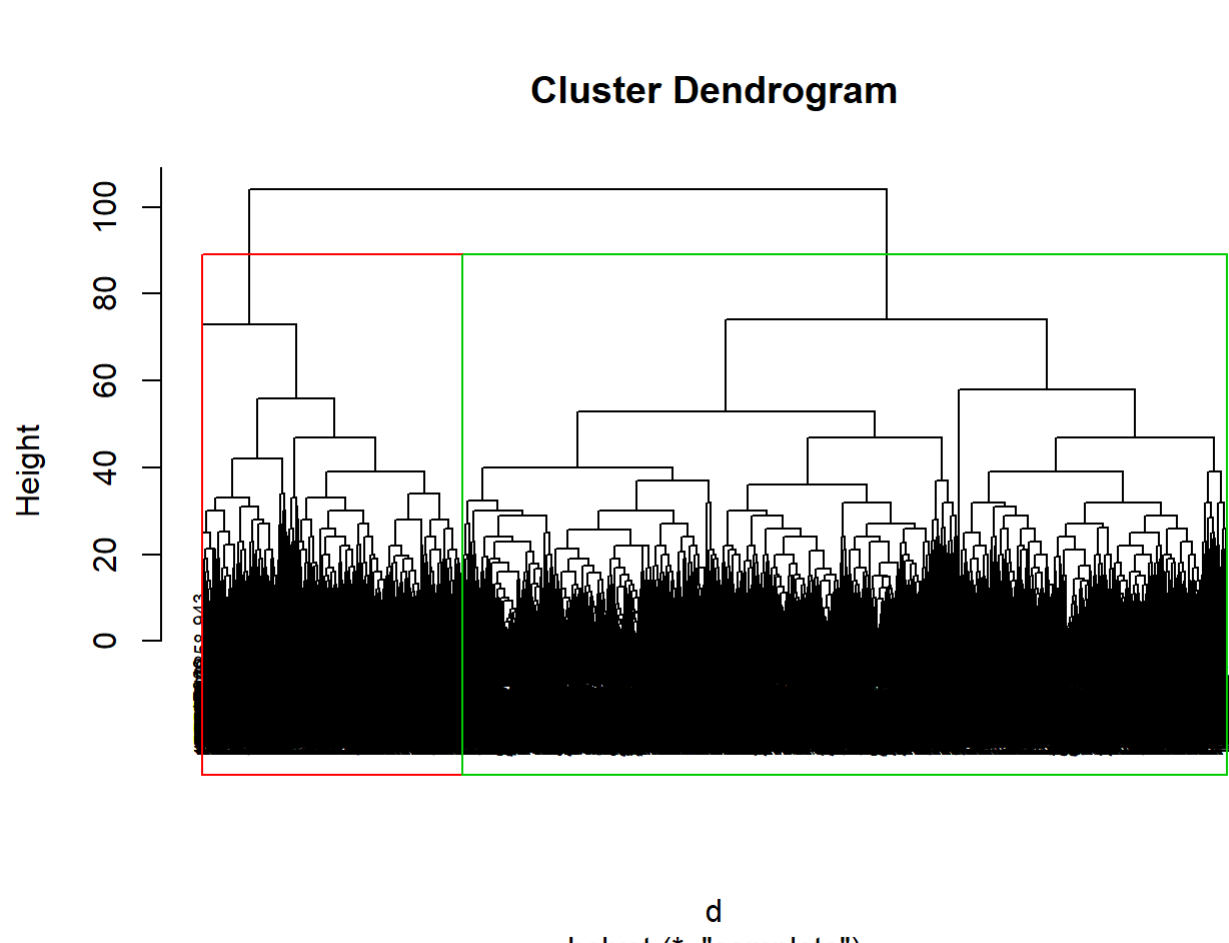
# Precision
precision <- length(result[result*y == 1])/(length(result[result-y == 1])+length(result[result*y == 1]))
precision

## [1] 0.3461538

# Recall
recall <- length(result[result*y == 1])/(length(result[result-y == 1])+length(result[result*y == -1]))
recall

## [1] 0.5294118

plot(hc1, cex = 0.6)
rect.hclust(hc1, k = 2, border = 2:5)
```



Final Project

```
library("glmnet")
library("factextra")
library("leaps")
library("gam")
library("precrec")
library("ModelMetrics")
```

PCA

```
df = read.csv("CleanV3.csv")
df = subset(df, select= c(ID, CODE_GENDER, FLAG_OW_N_CAR, FLAG_OW_N_REALTY, NAME_INCOME_TYPE, NAME_EDUCATION_TYPE, NAME_FAMILY_STATUS, NAME_HOUSING_TYPE, FLAG_MOBIL, OCCUPATION_TYPE, Occupation_Type, Income_Bins, Age_Bins, Emplo
y_years_Bins))
y = df$Risk
table(y)

## y
##      0      1
## 11828  8871

X.model = model.matrix(Risk ~ .-1, data = df)
head(X.model)

##      CNT_CHILDREN AMT_INCOME_TOTAL FLAG_WORK_PHONE FLAG_PHONE FLAG_EMAIL
## 1              0          279800              0              0              0
## 2              0           202500              0              0              0
## 3              0          157500              0              1              0
## 4              1          148500              0              0              0
## 5              0          180000              0              0              0
## 6              0          292500              0              0              0
##      CNT_FAM_MEMBERS Begin_Month Gender Car Realty Income_Type Family_Status
## 1              2          2          6      2      2      5      2
## 2              2          0          1      1      1      2      2
## 3              2          17         1      1      1      3      1
## 4              3          14         2      2      2      5      2
## 5              2          21         2      2      2      5      2
## 6              1          32         1      2      2      5      5
##      Housing_Type Age Employ_years Education_Type
## 1              2   39              11          2
## 2              2   56              10          2
## 3              3   50               8          4
## 4              2   41              14          2
## 5              2   59               7          2
## 6              2   33               4          4

X.std = scale(X.model)
head(X.std)

##      CNT_CHILDREN AMT_INCOME_TOTAL FLAG_WORK_PHONE FLAG_PHONE FLAG_EMAIL
## 1 -0.5858234      0.73366355 -0.5694531 -0.8004671 -0.3342278
## 2 -0.5858234      0.06006953 -0.5694531 -0.8004671 -0.3342278
## 3 -0.5858234      0.5694531 -0.5694531  1.4361611 -0.3342278
## 4  0.7571539      -0.47890558 -0.5694531 -0.8004671 -0.3342278
## 5 -0.5858234     -0.16446181 -0.5694531 -0.8004671 -0.3342278
## 6 -0.5858234      0.06006953 -0.5694531 -0.8004671 -0.3342278
##      CNT_FAM_MEMBERS Begin_Month Gender Car Realty Income_Type
## 1 -0.2862711 -1.4040977  1.4177210  1.2673667  0.7462296  0.7586914
## 2 -0.2862711 -1.0970025 -0.7774155 -0.8712062  0.7462296  0.7586914
## 3 -0.2862711 -0.6651030 -0.7774155 -0.8712062 -1.4773941 -0.4331775
## 4  0.8564825 -0.8666812  1.4177210  1.2673667  0.7462296  0.7586914
## 5 -0.2862711 -0.3965206  1.4177210  1.2673667  0.7462296  0.7586914
## 6 -1.4290247  0.3423033 -0.7774155  1.2673667  0.7462296  0.7586914
##      Family_Status Housing_Type Age Employ_years Education_Type
## 1 -0.3389826 -0.3123778 -0.81771928  0.64018424 -0.6276262
## 2 -0.3389826 -0.3123778  1.81126405  0.47217039 -0.6276262
## 3 -1.5763680  0.7515577  1.16574422  0.13614269  1.6751719
## 4 -0.3389826 -0.3123778  0.19746327  1.14422578 -0.6276262
## 5 -0.3389826 -0.3123778  2.13402517 -0.03107115 -0.6276262
## 6  3.3731735 -0.3123778 -0.66323091 -0.53591270  1.6751719
```

Feature Selection - LASSO

```
fit.lasso = glmnet(X.std, y, family = "binomial", alpha = 1)
plot(fit.lasso, xvar="lambda", label=TRUE)

##      16      16      16      13      9      4
##      0
##      4
##      8
##      12
##      16
##      20
##      24
##      28
##      32
##      36
##      40
##      44
##      48
##      52
##      56
##      60
##      64
##      68
##      72
##      76
##      80
##      84
##      88
##      92
##      96
##      100
##      104
##      108
##      112
##      116
##      120
##      124
##      128
##      132
##      136
##      140
##      144
##      148
##      152
##      156
##      160
##      164
##      168
##      172
##      176
##      180
##      184
##      188
##      192
##      196
##      200
##      204
##      208
##      212
##      216
##      220
##      224
##      228
##      232
##      236
##      240
##      244
##      248
##      252
##      256
##      260
##      264
##      268
##      272
##      276
##      280
##      284
##      288
##      292
##      296
##      300
##      304
##      308
##      312
##      316
##      320
##      324
##      328
##      332
##      336
##      340
##      344
##      348
##      352
##      356
##      360
##      364
##      368
##      372
##      376
##      380
##      384
##      388
##      392
##      396
##      400
##      404
##      408
##      412
##      416
##      420
##      424
##      428
##      432
##      436
##      440
##      444
##      448
##      452
##      456
##      460
##      464
##      468
##      472
##      476
##      480
##      484
##      488
##      492
##      496
##      500
##      504
##      508
##      512
##      516
##      520
##      524
##      528
##      532
##      536
##      540
##      544
##      548
##      552
##      556
##      560
##      564
##      568
##      572
##      576
##      580
##      584
##      588
##      592
##      596
##      600
##      604
##      608
##      612
##      616
##      620
##      624
##      628
##      632
##      636
##      640
##      644
##      648
##      652
##      656
##      660
##      664
##      668
##      672
##      676
##      680
##      684
##      688
##      692
##      696
##      700
##      704
##      708
##      712
##      716
##      720
##      724
##      728
##      732
##      736
##      740
##      744
##      748
##      752
##      756
##      760
##      764
##      768
##      772
##      776
##      780
##      784
##      788
##      792
##      796
##      800
##      804
##      808
##      812
##      816
##      820
##      824
##      828
##      832
##      836
##      840
##      844
##      848
##      852
##      856
##      860
##      864
##      868
##      872
##      876
##      880
##      884
##      888
##      892
##      896
##      900
##      904
##      908
##      912
##      916
##      920
##      924
##      928
##      932
##      936
##      940
##      944
##      948
##      952
##      956
##      960
##      964
##      968
##      972
##      976
##      980
##      984
##      988
##      992
##      996
##      1000
##      1004
##      1008
##      1012
##      1016
##      1020
##      1024
##      1028
##      1032
##      1036
##      1040
##      1044
##      1048
##      1052
##      1056
##      1060
##      1064
##      1068
##      1072
##      1076
##      1080
##      1084
##      1088
##      1092
##      1096
##      1100
##      1104
##      1108
##      1112
##      1116
##      1120
##      1124
##      1128
##      1132
##      1136
##      1140
##      1144
##      1148
##      1152
##      1156
##      1160
##      1164
##      1168
##      1172
##      1176
##      1180
##      1184
##      1188
##      1192
##      1196
##      1200
##      1204
##      1208
##      1212
##      1216
##      1220
##      1224
##      1228
##      1232
##      1236
##      1240
##      1244
##      1248
##      1252
##      1256
##      1260
##      1264
##      1268
##      1272
##      1276
##      1280
##      1284
##      1288
##      1292
##      1296
##      1300
##      1304
##      1308
##      1312
##      1316
##      1320
##      1324
##      1328
##      1332
##      1336
##      1340
##      1344
##      1348
##      1352
##      1356
##      1360
##      1364
##      1368
##      1372
##      1376
##      1380
##      1384
##      1388
##      1392
##      1396
##      1400
##      1404
##      1408
##      1412
##      1416
##      1420
##      1424
##      1428
##      1432
##      1436
##      1440
##      1444
##      1448
##      1452
##      1456
##      1460
##      1464
##      1468
##      1472
##      1476
##      1480
##      1484
##      1488
##      1492
##      1496
##      1500
##      1504
##      1508
##      1512
##      1516
##      1520
##      1524
##      1528
##      1532
##      1536
##      1540
##      1544
##      1548
##      1552
##      1556
##      1560
##      1564
##      1568
##      1572
##      1576
##      1580
##      1584
##      1588
##      1592
##      1596
##      1600
##      1604
##      1608
##      1612
##      1616
##      1620
##      1624
##      1628
##      1632
##      1636
##      1640
##      1644
##      1648
##      1652
##      1656
##      1660
##      1664
##      1668
##      1672
##      1676
##      1680
##      1684
##      1688
##      1692
##      1696
##      1700
##      1704
##      1708
##      1712
##      1716
##      1720
##      1724
##      1728
##      1732
##      1736
##      1740
##      1744
##      1748
##      1752
##      1756
##      1760
##      1764
##      1768
##      1772
##      1776
##      1780
##      1784
##      1788
##      1792
##      1796
##      1800
##      1804
##      1808
##      1812
##      1816
##      1820
##      1824
##      1828
##      1832
##      1836
##      1840
##      1844
##      1848
##      1852
##      1856
##      1860
##      1864
##      1868
##      1872
##      1876
##      1880
##      1884
##      1888
##      1892
##      1896
##      1900
##      1904
##      1908
##      1912
##      1916
##      1920
##      1924
##      1928
##      1932
##      1936
##      1940
##      1944
##      1948
##      1952
##      1956
##      1960
##      1964
##      1968
##      1972
##      1976
##      1980
##      1984
##      1988
##      1992
##      1996
##      2000
##      2004
##      2008
##      2012
##      2016
##      2020
##      2024
##      2028
##      2032
##      2036
##      2040
##      2044
##      2048
##      2052
##      2056
##      2060
##      2064
##      2068
##      2072
##      2076
##      2080
##      2084
##      2088
##      2092
##      2096
##      2100
##      2104
##      2108
##      2112
##      2116
##      2120
##      2124
##      2128
##      2132
##      2136
##      2140
##      2144
##      2148
##      2152
##      2156
##      2160
##      2164
##      2168
##      2172
##      2176
##      2180
##      2184
##      2188
##      2192
##      2196
##      2200
##      2204
##      2208
##      2212
##      2216
##      2220
##      2224
##      2228
##      2232
##      2236
##      2240
##      2244
##      2248
##      2252
##      2256
##      2260
##      2264
##      2268
##      2272
##      2276
##      2280
##      2284
##      2288
##      2292
##      2296
##      2300
##      2304
##      2308
##      2312
##      2316
##      2320
##      2324
##      2328
##      2332
##      2336
##      2340
##      2344
##      2348
##      2352
##      2356
##      2360
##      2364
##      2368
##      2372
##      2376
##      2380
##      2384
##      2388
##      2392
##      2396
##      2400
##      2404
##      2408
##      2412
##      2416
##      2420
##      2424
##      2428
##      2432
##      2436
##      2440
##      2444
##      2448
##      2452
##      2456
##      2460
##      2464
##      2468
##      2472
##      2476
##      2480
##      2484
##      2488
##      2492
##      2496
##      2500
##      2504
##      2508
##      2512
##      2516
##      2520
##      2524
##      2528
##      2532
##      2536
##      2540
##      2544
##      2548
##      2552
##      2556
##      2560
##      2564
##      2568
##      2572
##      2576
##      2580
##      2584
##      2588
##      2592
##      2596
##      2600
##      2604
##      2608
##      2612
##      2616
##      2620
##      2624
##      2628
##      2632
##      2636
##      2640
##      2644
##      2648
##      2652
##      2656
##      2660
##      2664
##      2668
##      2672
##      2676
##      2680
##      2684
##      2688
##      2692
##      2696
##      2700
##      2704
##      2708
##      2712
##      2716
##      2720
##      2724
##      2728
##      2732
##      2736
##      2740
##      2744
##      2748
##      2752
##      2756
##      2760
##      2764
##      2768
##      2772
##      2776
##      2780
##      2784
##      2788
##      2792
##      2796
##      2800
##      2804
##      2808
##      2812
##      2816
##      2820
##      2824
##      2828
##      2832
##      2836
##      2840
##      2844
##      2848
##      2852
##      2856
##      2860
##      2864
##      2868
##      2872
##      2876
##      2880
##      2884
##      2888
##      2892
##      2896
##      2900
##      2904
##      2908
##      2912
##      2916
##      2920
##      2924
##      2928
##      2932
##      2936
##      2940
##      2944
##      2948
##      2952
##      2956
##      2960
##      2964
##      2968
##      2972
##      2976
##      2980
##      2984
##      2988
##      2992
##      2996
##      3000
##      3004
##      3008
##      3012
##      3016
##      3020
##      3024
##      3028
##      3032
##      3036
##      3040
##      3044
##      3048
##      3052
##      3056
##      3060
##      3064
##      3068
##      3072
##      3076
##      3080
##      3084
##      3088
##      3092
##      3096
##      3100
##      3104
##      3108
##      3112
##      3116
##      3120
##      3124
##      3128
##      3132
##      3136
##      3140
##      3144
##      3148
##      3152
##      3156
##      3160
##      3164
##      3168
##      3172
##      3176
##      3180
##      3184
##      3188
##      3192
##      3196
##      3200
##      3204
##      3208
##      3212
##      3216
##      3220
##      3224
##      3228
##      3232
##      3236
##      3240
##      3244
##      3248
##      3252
##      3256
##      3260
##      3264
##      3268
##      3272
##      3276
##      3280
##      3284
##      3288
##      3292
##      3296
##      3300
##      3304
##      3308
##      3312
##      3316
##      3320
##      3324
##      3328
##      3332
##      3336
##      3340
##      3344
##      3348
##      3352
##      3356
##      3360
##      3364
##      3368
##      3372
##      3376
##      3380
##      3384
##      3388
##      3392
##      3396
##      3400
##      3404
##      3408
##      3412
##      3416
##      3420
##      3424
##      3428
##      3432
##      3436
##      3440
##      3444
##      3448
##      3452
##      3456
##      3460
##      3464
##      3468
##      3472
##      3476
##      3480
##      3484
##      3488
##      3492
##      3496
##      3500
##      3504
##      3508
##      3512
##      3516
##      3520
##      3524
##      3528
##      3532
##      3536
##      3540
##      3544
##      3548
##      3552
##      3556
##      3560
##      3564
##      3568
##      3572
##      3576
##      3580
##      3584
##      3588
##      3592
##      3596
##      3600
##      3604
##      3608
##      3612
##      3616
##      3620
##      3624
##      3628
##      3632
##      3636
##      3640
##      3644
##      3648
##      3652
##      3656
##      3660
##      3664
##      3668
##      3672
##      3676
##      3680
##      3684
##      3688
##      3692
##      3696
##      3700
##      3704
##      3708
##      3712
##      3716
##      3720
##      3724
##      3728
##      3732
##      3736
##      3740
##      3744
##      3748
##      3752
##      3756
##      3760
##      3764
##      3768
##      3772
##      3776
##      3780
##      3784
##      3788
##      3792
##      3796
##      3800
##      3804
##      3808
##      3812
##      3816
##      3820
##      3824
##      3828
##      3832
##      3836
##      3840
##      3844
##      3848
##      3852
##      3856
##      3860
##      3864
##      3868
##      3872
##      3876
##      3880
##      3884
##      3888
##      3892
##      3896
##      3900
##      3904
##      3908
##      3912
##      3916
##      3920
##      3924
##      3928
##      3932
##      3936
##      3940
##      3944
##      3948
##      3952
##      3956
##      3960
##      3964
##      3968
##      3972
##      3976
##      3980
##      3984
##      3988
##      3992
##      3996
##      4000
##      4004
##      4008
##      4012
##      4016
##      4020
##      4024
##      4028
##      4032
##      4036
##      4040
##      4044
##      4048
##      4052
##      4056
##      4060
##      4064
##      4068
##      4072
##      4076
##      4080
##      4084
##      4088
##      4092
##      4096
##      4100
##      4104
##      4108
##      4112
##      4116
##      4120
##      4124
##      4128
##      4132
##      4136
##      4140
##      4144
##      4148
##      4152
##      4156
##      4160
##      4164
##      4168
##      4172
##      4176
##      4180
##      4184
##      4188
##      4192
##      4196
##      4200
##      4204
##      4208
##      4212
##      4216
##      4220
##      4224
##      4228
##      4232
##      4236
##      4240
##      4244
##      4248
##      4252
##      4256
##      4260
##      4264
##      4268
##      4272
##      4276
##      4280
##      4284
##      4288
##      4292
##      4296
##      4300
##      4304
##      4308
##      4312
##      4316
##      4320
##      4324
##      4328
##      4332
##      4336
##      4340
##      4344
##      4348
##      4352
##      4356
##      4360
##      4364
##      4368
##      4372
##      4376
##      4380
##      4384
##      4388
##      4392
##      4396
##      4400
##      4404
##      4408
##      4412
##      4416
##      4420
##      4424
##      4428
##      4432
##      4436
##      4440
##      4444
##      4448
##      4452
##      4456
##      4460
##      4464
##      4468
##      4472
##      4476
##      4480
##      4484
##      4488
##      4492
##      4496
##      4500
##      4504
##      4508
##      4512
##      4516
##      4520
##      4524
##      4528
##      4532
##      4536
##      4540
##      4544
##      4548
##      4552
##      4556
##      4560
##      4564
##      4568
##      4572
##      4576
##      4580
##      4584
##      4588
##      4592
##      4596
##      4600
##      4604
##      4608
##      4612
##      4616
##      4620
##      4624
##      4628
##      4632
##      4636
##      4640
##      4644
##      4648
##      4652
##      4656
##      4660
##      4664
##      4668
##      4672
##      4676
##      4680
##      4684
##      4688
##      4692
##      4696
##      4700
##      4704
##      4708
##      4712
##      4716
##      4720
##      4724
##      4728
##      4732
##      4736
##      4740
##      4744
##      4748
##      4752
##      4756
##      4760
##      4764
##      4768
##      4772
##      4776
##      4780
##      4784
##      4788
##      4792
##      4796
##      4800
##      4804
##      4808
##      4812
##      4816
##      4820
##      4824
##      4828
##      4832
##      4836
##      4840
##      4844
##      4848
##      4852
##      4856
##      4860
##      4864
##      4868
##      4872
##      4876
##      4880
##      4884
##      4888
##      4892
##      4896
##      4900
##      4904
##      4908
##      4912
##      4916
##      4920
##      4924
##      4928
##      4932
##      4936
##      4940
##      4944
##      4948
##      4952
##      4956
##      4960
##      4964
##      4968
##      4972
##      4976
##      4980
##      4984
##      4988
##      4992
##      4996
##      5000
##      5004
##      5008
##      5012
##      5016
##      5020
##      5024
##      5028
##      5032
##      5036
##      5040
##      5044
##      5048
##      5052
##      5056
##      5060
##      5064
##      5068
##      5072
##      5076
##      5080
##      5084
##      5088
##      5092
##      5096
##      5100
##      5104
##      5108
##      5112
##      5116
##      5120
##      5124
##      5128
##      5132
##      5136
##      5140
##      5144
##      5148
##      5152
##      5156
##      5160
##      5164
##      5168
##      5172
##      5176
##      5180
##      5184
##      5188
##      5192
##      5196
##      5200
##      5204
##      5208
##      5212
##      5216
##      5220
##      5224
##      5228
##      5232
##      5236
##      5240
##      5244
##      5248
##      5252
##      5256
##      5260
##      5264
##      5268
##      5272
##      5276
##      5280
##      5284
##      5288
##      5292
##      5296
##      5300
##      5304
##      5
```