

BT2103 Project 3

Chen Haoli, Lo Zhi Hao, Luah Jun Yang, Toh Zhan Ting

Table of Contents

- 1) Exploratory data analysis on dataset ("card")
- 2) Data pre-processing
- 3) Feature Selection
- 4) Model Selection
- 5) Model Evaluation
- 6) Areas for improvement on dataset

Introduction to the Dataset and Problem Statement

The cash and credit card debt problem that Taiwan's credit card issuers experienced in recent years is predicted to peak in the third quarter of 2006 (Chou, 2006). Taiwan's card-issuing banks over-issued cash and credit cards to unqualified applicants in an effort to gain market dominance. In addition, most cardholders, regardless of their capacity to pay back, abused their cards for consumption and racked up large credit and cash-card debt. The crisis damaged consumer confidence in finance, and therefore presents a significant problem for both banks and cardholders.

Crisis management and risk prediction take place upstream and downstream in a mature financial system, respectively. The main goal of risk prediction is to lessen the harm and uncertainty caused by corporate performance or individual customer credit risk by using financial information, such as business financial statements, customer transaction and repayment histories, etc.

Therefore, with extensive data collected from the period of April to September in 2005, this report aims to build a predictive model to accurately forecast clients who have tendencies to default on the bank and thus ensures the profitability and stability of banks. In Finance, a default occurs when a borrower doesn't fulfill the terms of the loan. In this situation, default would occur if the cardholder failed to pay the credit card account within a given month.

In the dataset, there are 23 explanatory variables, together with a dependent variable of client's default status. The detailed description is as follows:

X1: Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her family (supplementary) credit.

X2: Gender (1 = male; 2 = female).

X3: Education (1 = graduate school; 2 = university; 3 = high school; 4 = others).

X4: Marital status (1 = married; 2 = single; 3 = others).

X5: Age (year).

X6–X11: History of past payment. We tracked the past monthly payment records (from April to September, 2005) as follows: X6 = the repayment status in September, 2005; X7 = the repayment status in August, 2005; . . . ; X11 = the repayment status in April, 2005. The measurement scale for the repayment status is: -1 = pay duly; 1 = payment delay for one month; 2 = payment delay for two months; . . . ; 8 = payment delay for eight months; 9 = payment delay for nine months and above.

X12–X17: Amount of bill statement (NT dollar). X12 = amount of bill statement in September, 2005; X13 = amount of bill statement in August, 2005; . . . ;X17 = amount of bill statement in April, 2005.

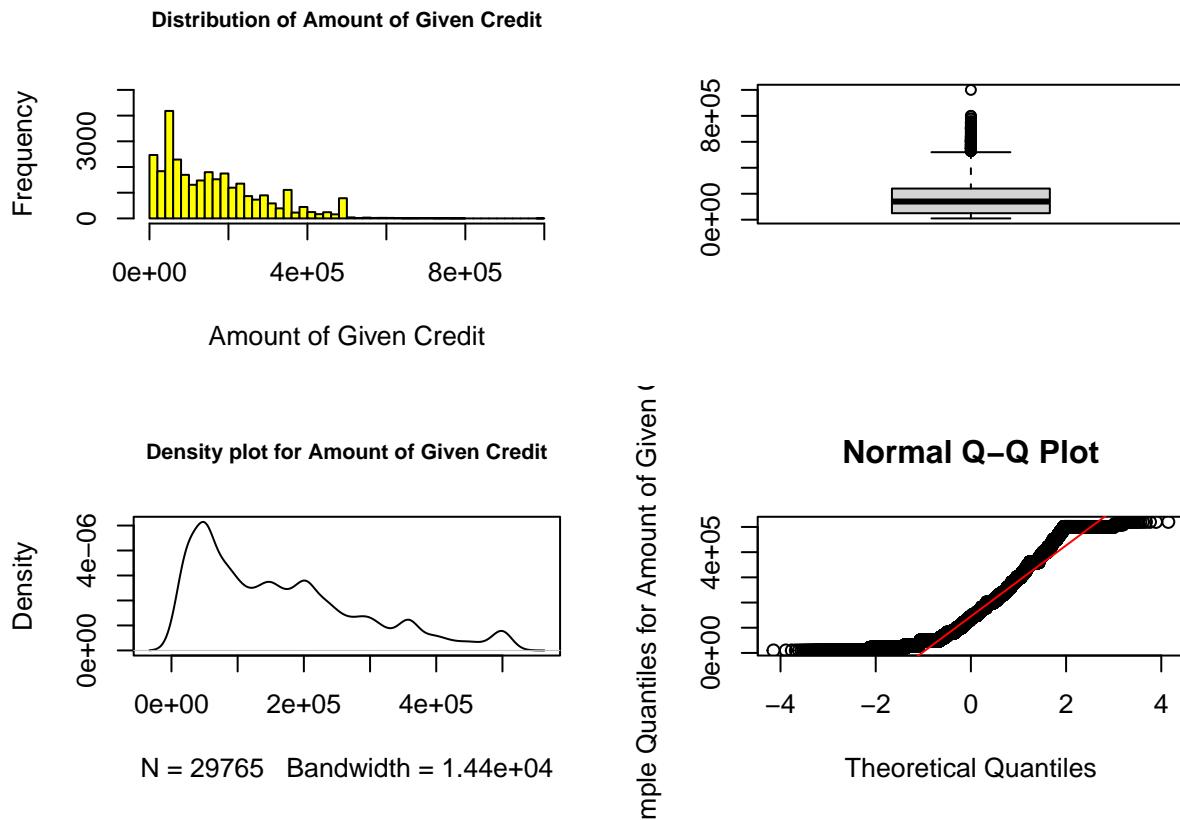
X18–X23: Amount of previous payment (NT dollar). X18 = amount paid in September, 2005; X19 = amount paid in August, 2005; . . . ;X23 = amount paid in April, 2005.

Exploratory Data Analysis

As all of the columns in the data are originally categorized as characters, in order to proceed with the data pre-processing and data visualization, we decided to transform the columns that are considered as continuous and numeric back to a numeric data type. After further research from reading through the description of our data, we identified x1, x5, x6, x7, x8, x9, x10, x11, x12, x13, x14, x15, x16, x17, x18, x19, x20, x21, x22, x23 as the data columns that we may need to transform. x2, x3, x4 are factors. For variables x3 and x4, we dropped the factor levels 0, which are indicative of N.A values for those variables Education and Marital Status respectively and they only amount to a small number of outliers. The other non-classified factor levels for these 2 variables are parked under the ‘others’ factor so that we will not exclude too many data points.

By using summary statistics and graphical representations, we are conducting preliminary analyses on the data in order to find trends, identify anomalies, inconsistencies, and missing values to test hypotheses and verify assumptions.

Amount of Given Credit (X1)

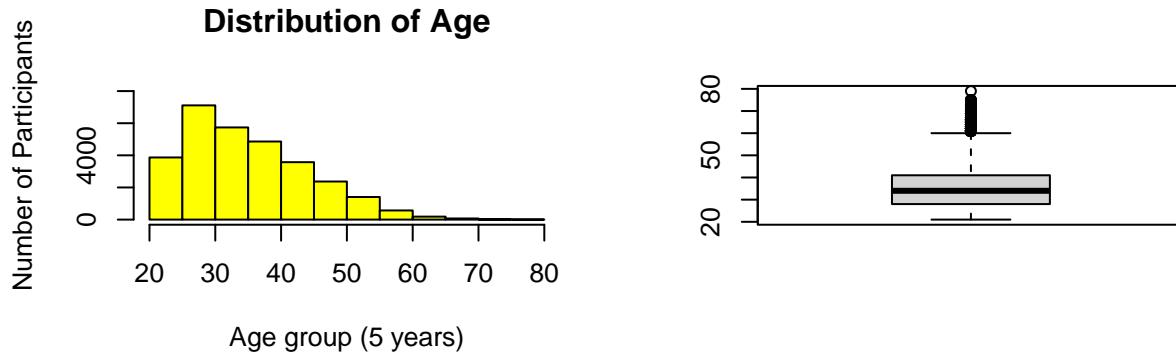


From the boxplot, we observed points lie outside $Q3 + 1.5 \text{IQR}$ from the mean. We consider these points as outliers and we will remove these points from dataset. From the density plot, we can observe that the data is positively skewed and it is not normally distributed.

From the qqplot, we can observe that the data is not normally distributed since the line and plots are not truly aligned.

Since the amount of given credit is positively skewed, we will use log transformation to improve the distribution of the data to normality.

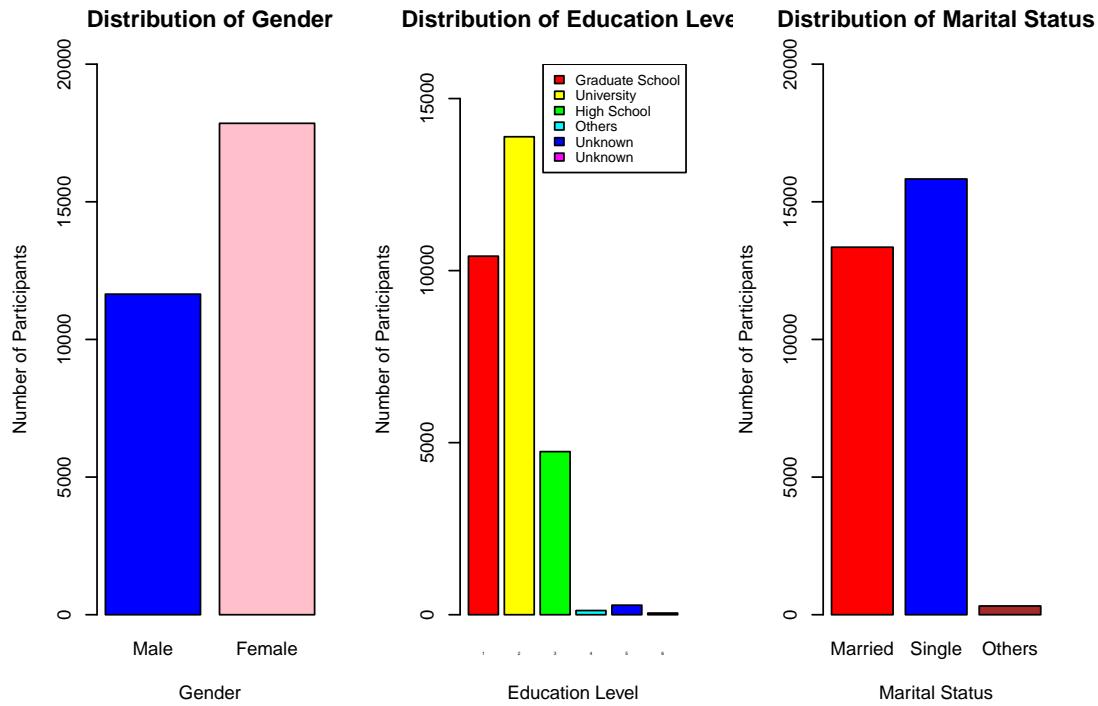
Age (X5)



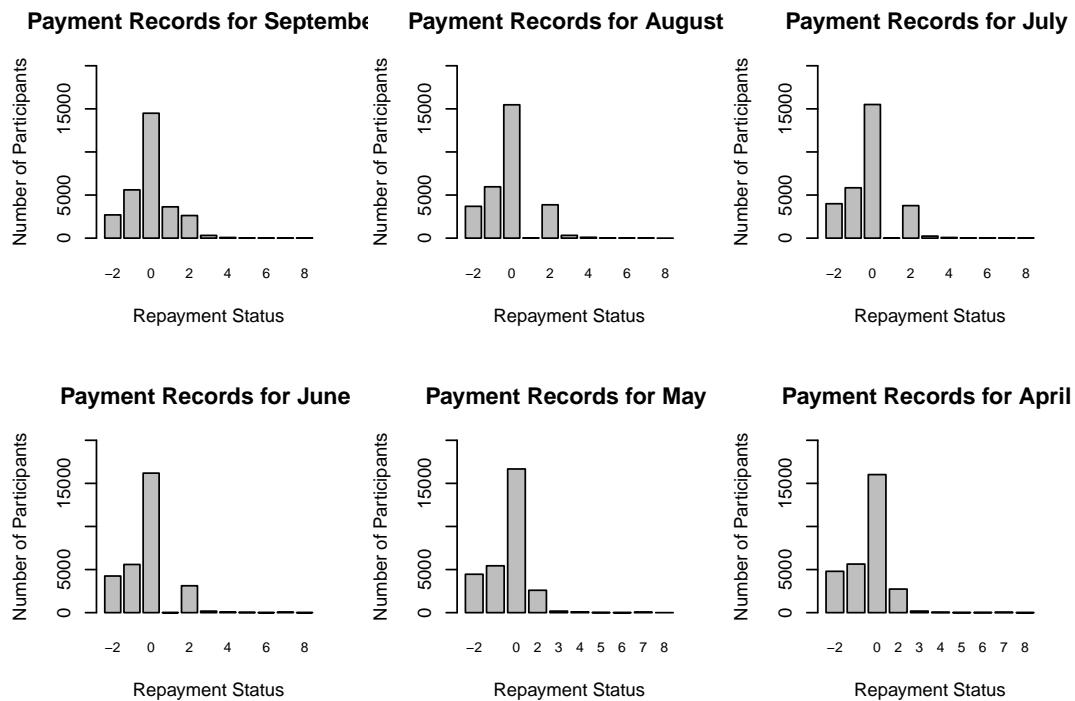
From the boxplot, we observed points lie outside $Q3 + 1.5IQR$ from the mean. We consider these points as outliers and we will remove these points from dataset

Gender (X2), Education (X3), and Marital Status (X5)

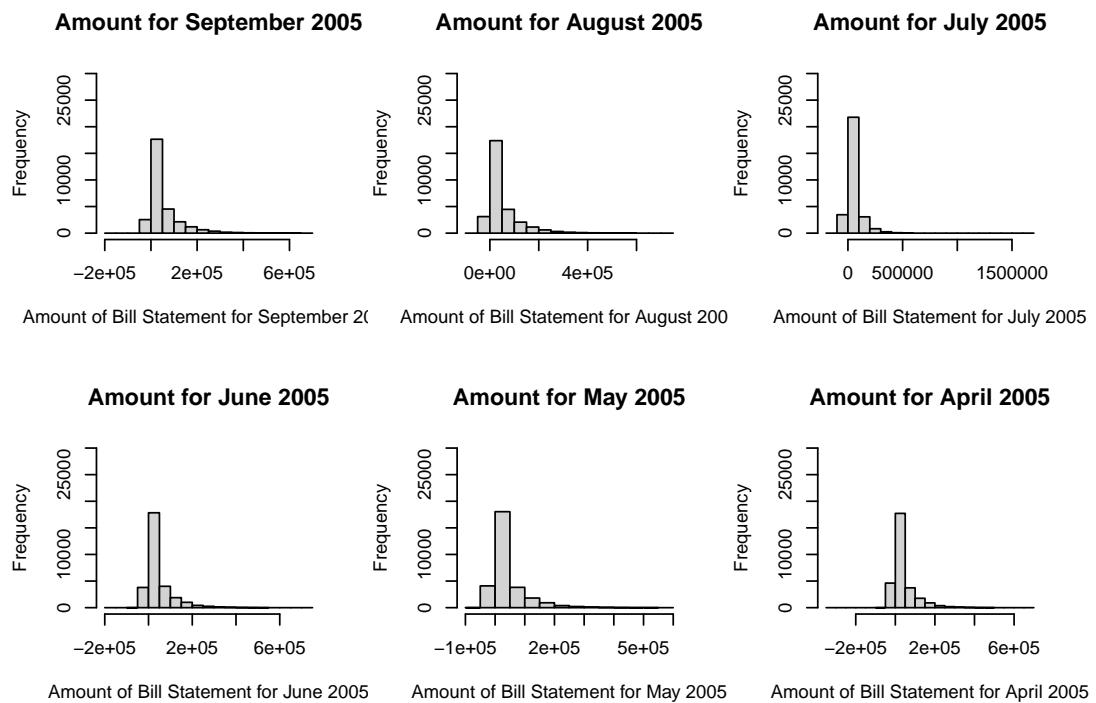
```
## Male Female
## 11645 17851
## Graduate School University High School Others Unknown
## 10419 13890 4740 122 277
## Unknown
## 48
## Married Single Others
## 13350 15828 318
```



Monthly Payment Records for 2005 (X6 - X11)

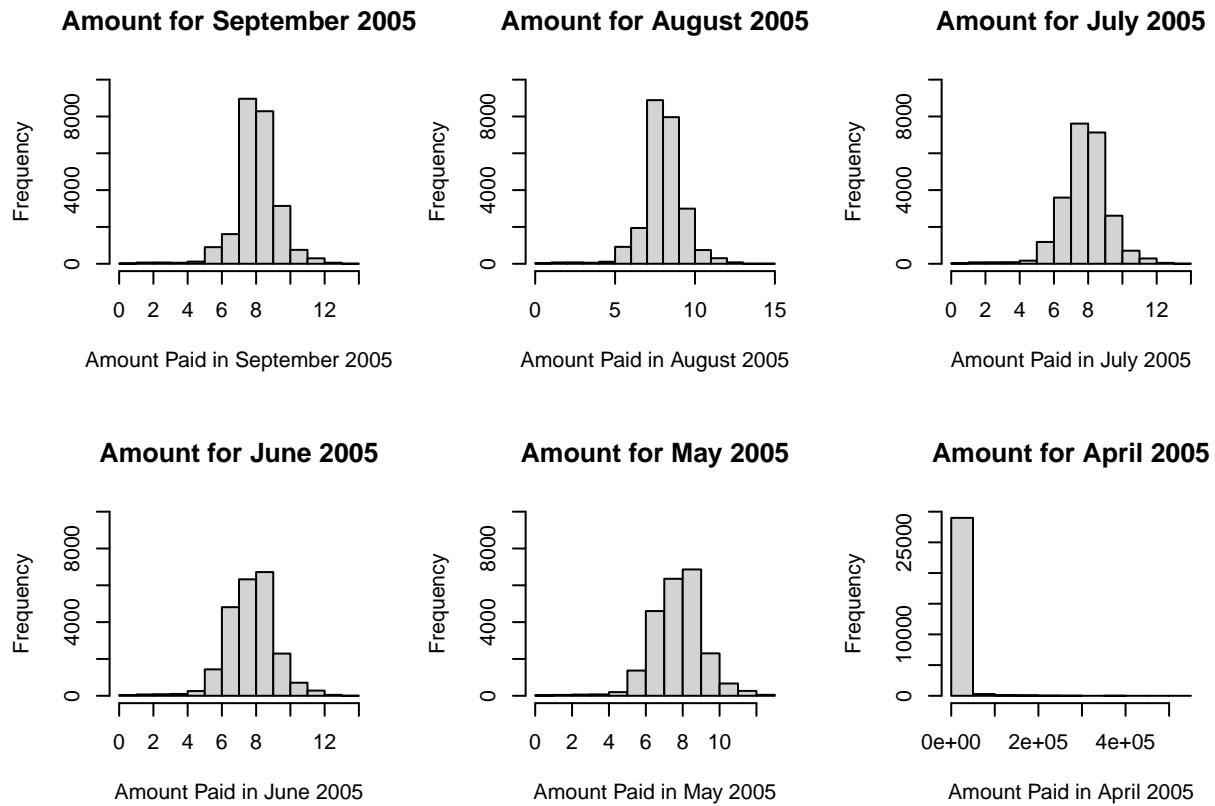


Monthly Bill Distributions for 2005 (X12 - X17)



Amount of previous payment (X18 - X23)

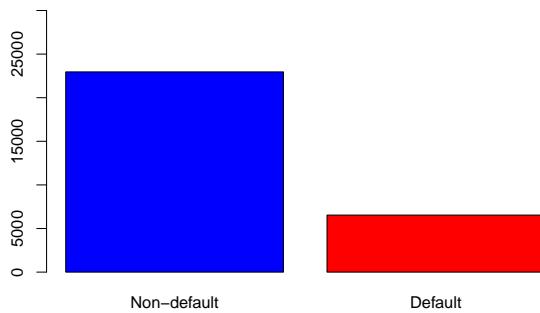
We will perform log transformation for Amount Paid in order to make it resemble a normal distribution.



Distribution of Default and Non-default data

```
## Non-default      Default
##        22957       6539
```

Distribution of Default and Non-default Customers



Data Pre-Processing

As we noticed from the data visualisation shown in our part above, there is a significant issue with the data we gathered from our dataset, as some of the data are not described in our original cited data source.

For example, for distribution of education level and marital status, we noticed some entries that does not belong to any factor levels that were described in the original data source. Similarly, for the history of payment records, we also realized the occurrence of erroneous factor levels such as -2 and 0, which are not described for in the original data source.

First, we are going to change the names of each columns for better understanding and easier interpretation in our following chapters.

Moving on we are going to perform some data manipulation on some of the data anomalies that we spotted earlier. In that sense, we are going to manipulate and change the data points that are originally not described in our cited data source.

Data Manipulation for Marital Status (X2)

As we can see from the visualisation for distribution of marital status for our clients, there are categories that are not described for in our original data cited source, namely 0. For our team, as we consider 0 values as possible null values where our clients did not provide their marital status information properly, we decided to drop rows with 0 values recorded for marital status.

```
##  
##      1      2      3  
## 13350 15828   318
```

Data Manipulation for Education Status (X3)

As we can see from the visualisation for distribution of education level for our clients, there are categories that are not described for in our original data cited source, namely 0, 5 and 6. For our team, as we consider 0 as a possible null value where our clients did not provide their academic credentials, we decided to drop those observations with 0 recorded as their education level. Meanwhile, as 5 and 6 may be education levels that are higher than or not included in the provided options (such as PhD), we decided to manipulate those data to add them into the 'others' category, 4 that is provided.

```
##  
##      1      2      3      4  
## 10419 13890  4740    447
```

Feature Selection

Make use of feature selection methodologies to select the most relevant independent variables to create a prediction model.

For our project, as it is a large dataset with around 30,000 observations, for each type of feature selection method, we are going to perform a test in order to determine the best features in that category, namely:

Forward and Backward variable selection

Filter Method - ANOVA for continuous data and Chi Squared Test for categorical data

Wrapper Method - Boruta Method

Embedded Method - Information Gain

Forward Variable Selection:

- Start with model containing no possible explanatory variable and for each variable in turn, we will investigate effect of adding variable from the current model.
- 5 most significant variables will be used for our subsequent models

```
## (Intercept) Amt_credit Gender Edu_Lvl Marital Age PayRec_Sep PayRec_Aug
## 1          1          0      0      0      0      0      1      0
## 2          1          1      0      0      0      0      1      0
## 3          1          1      0      0      0      0      1      0
## 4          1          1      0      0      0      0      1      1
## 5          1          1      0      0      0      1      0      1
## PayRec_Jul PayRec_Jun PayRec_May PayRec_Apr BillAmt_Sep BillAmt_Aug
## 1          0          0      0      0      0      0
## 2          0          0      0      0      0      0
## 3          0          0      0      0      1      0
## 4          0          0      0      0      1      0
## 5          0          0      0      0      1      0
## BillAmt_Jul BillAmt_Jun BillAmt_May BillAmt_Apr PaidAmt_Sep PaidAmt_Aug
## 1          0          0      0      0      0      0
## 2          0          0      0      0      0      0
## 3          0          0      0      0      0      0
## 4          0          0      0      0      0      0
## 5          0          0      0      0      0      0
## PaidAmt_Jul PaidAmt_Jun PaidAmt_May PaidAmt_Apr   rsq adjr2      cp      bic
## 1          0          0      0      0  0.106 0.106 659.541 -3290.498
## 2          0          0      0      0  0.113 0.113 422.683 -3514.801
## 3          0          0      0      0  0.117 0.116 315.037 -3612.826
## 4          0          0      0      0  0.121 0.120 179.525 -3738.945
## 5          0          0      0      0  0.122 0.122 129.422 -3780.503
##           rss
## 1 4548.953
## 2 4512.917
## 3 4496.375
## 4 4475.628
## 5 4467.768
```

Backward Variable Selection:

- Start with model containing all possible explanatory variables and for each variable in turn, we will investigate effect of removing that variable from the current model.
- 5 most significant variables will be used for our subsequent models

```
## (Intercept) Amt_credit Gender Edu_Lvl Marital Age PayRec_Sep PayRec_Aug
## 1          1          0      0      0      0      0      1      0
## 2          1          0      0      0      0      0      1      0
## 3          1          0      0      0      0      0      1      1
## 4          1          1      0      0      0      0      1      1
## 5          1          1      0      0      0      1      1      1
## PayRec_Jul PayRec_Jun PayRec_May PayRec_Apr BillAmt_Sep BillAmt_Aug
## 1          0          0      0      0      0      0
## 2          0          0      0      0      1      0
## 3          0          0      0      0      1      0
## 4          0          0      0      0      1      0
## 5          0          0      0      0      1      0
## BillAmt_Jul BillAmt_Jun BillAmt_May BillAmt_Apr PaidAmt_Sep PaidAmt_Aug
## 1          0          0      0      0      0      0
```

```

## 1      0      0      0      0      0      0
## 2      0      0      0      0      0      0
## 3      0      0      0      0      0      0
## 4      0      0      0      0      0      0
## 5      0      0      0      0      0      0
##   PaidAmt_Jul PaidAmt_Jun PaidAmt_May PaidAmt_Apr   rsq adjr2     cp      bic
## 1      0          0          0          0  0.106 0.106 659.541 -3290.498
## 2      0          0          0          0  0.113 0.113 426.521 -3511.017
## 3      0          0          0          0  0.119 0.119 230.592 -3696.520
## 4      0          0          0          0  0.121 0.120 179.525 -3738.945
## 5      0          0          0          0  0.122 0.122 130.981 -3778.950
##      rss
## 1 4548.953
## 2 4513.496
## 3 4483.635
## 4 4475.628
## 5 4468.003

```

From both the Forward and Backward Variable Selection, we observed similarities in the significant variables selected in both models (i.e.,Amt_credit, PayRec_Sep, PayRec_Aug, PayRec_Jul). However, we also observed that Marital Status & Age are ranked differently in both models. Thus, in order to increase accuracy of our subsequent models, we decided to include both Marital & Age.

Splitting into train and test set

First, we will split our dataset into train and test set before actually performing the feature selection models on our data. This is to prevent any leakage of information from our test set into our training set leading to biased and overfitted models.

Anova Test

```

##      Variable      p-value
## 2      Age 0.000000e+00
## 14 PaidAmt_Apr 0.000000e+00
## 8 BillAmt_Apr 5.897689e-20
## 10 PaidAmt_Aug 4.891178e-15
## 11 PaidAmt_Jul 8.157781e-13
## 12 PaidAmt_Jun 1.371457e-12
## 9 PaidAmt_Sep 1.665843e-11
## 1 Amt_Credit 2.898041e-11
## 13 PaidAmt_May 2.794714e-10
## 3 BillAmt_Sep 6.482189e-03
## 4 BillAmt_Aug 1.310158e-02
## 5 BillAmt_Jul 7.495121e-02
## 6 BillAmt_Jun 1.744674e-01
## 7 BillAmt_May 1.864872e-01

```

Chi Square test

Next, for categorical variables, Chi Squared Method is used to identify the most important features.

```

##      Statistics      p-value      V3
## 1      Gender 6.713353 3.485089e-02
## 2      Edu_Lvl 85.508305 2.492442e-05
## 7 PayRec_Jun 983.964863 4.943317e-206

```

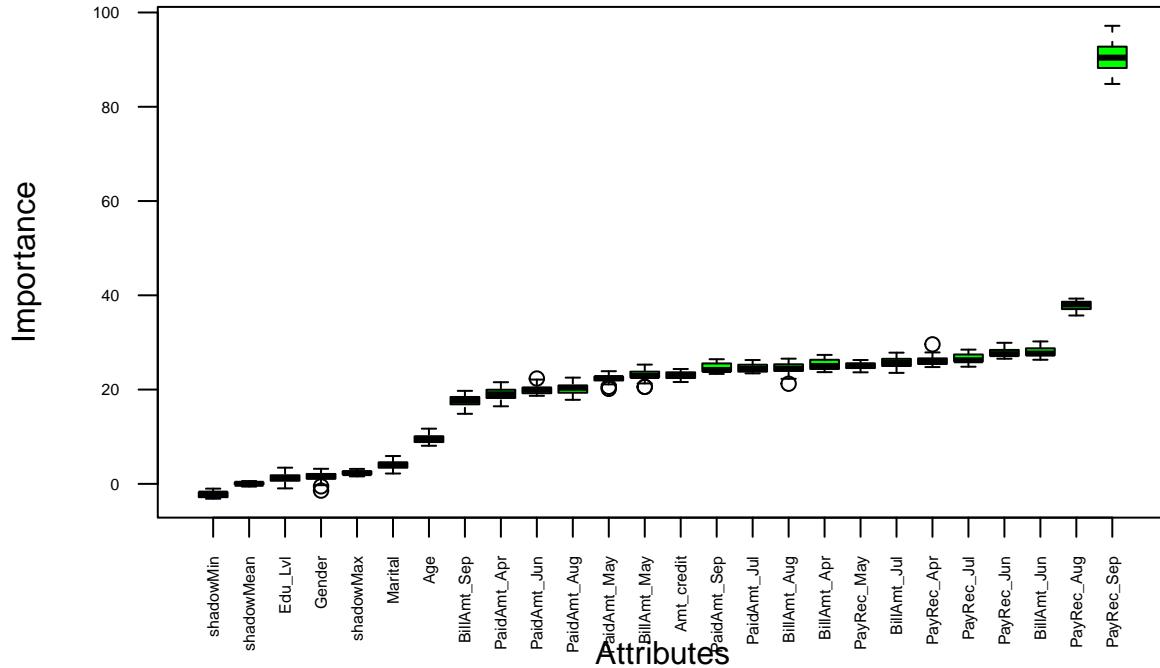
```

## 6 PayRec_Jul 1054.926738 2.457447e-221
## 5 PayRec_Aug 1131.704234 7.713010e-237
## 4 PayRec_Sep 1303.027160 8.499636e-274
## 3 Marital 1812.283792 0.000000e+00
## 8 PayRec_May 11924.409892 8.638387e-01

```

Wrapper method for feature selection

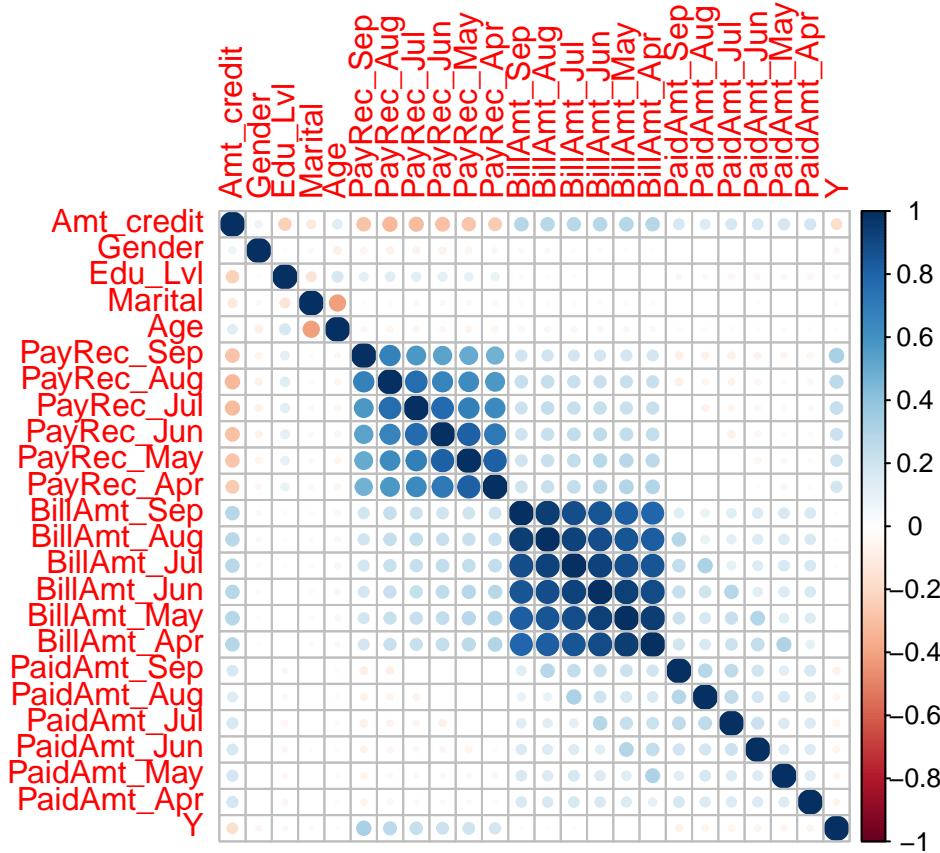
For wrapper method, the package and method that we are going to use is the Boruta Method that utilises random forest decision tree model in computing for the importance of each feature.



Embedded method for feature selection

We utilise the information gain function from the FSelectorRcpp package to inspect and identify the important features in our model.

Correlation



The highly correlated attributes are BillAmt_Jun, BillAmt_May, BillAmt_Jul, BillAmt_Apr, BillAmt_Aug, PayRec_Jun, PayRec_Jul, PayRec_Aug, and PayRec_May.

Model Selection

For our models, a way to identify the effectiveness and accuracy of each of our models is to evaluate each of the models based on accuracy, null accuracy, ROC and AUC values and harmonic mean. Each of these methods have their own merits and advantages, and by computing the confusion matrix and compute each of these values, we can get a deeper understanding on how each of our models are performing.

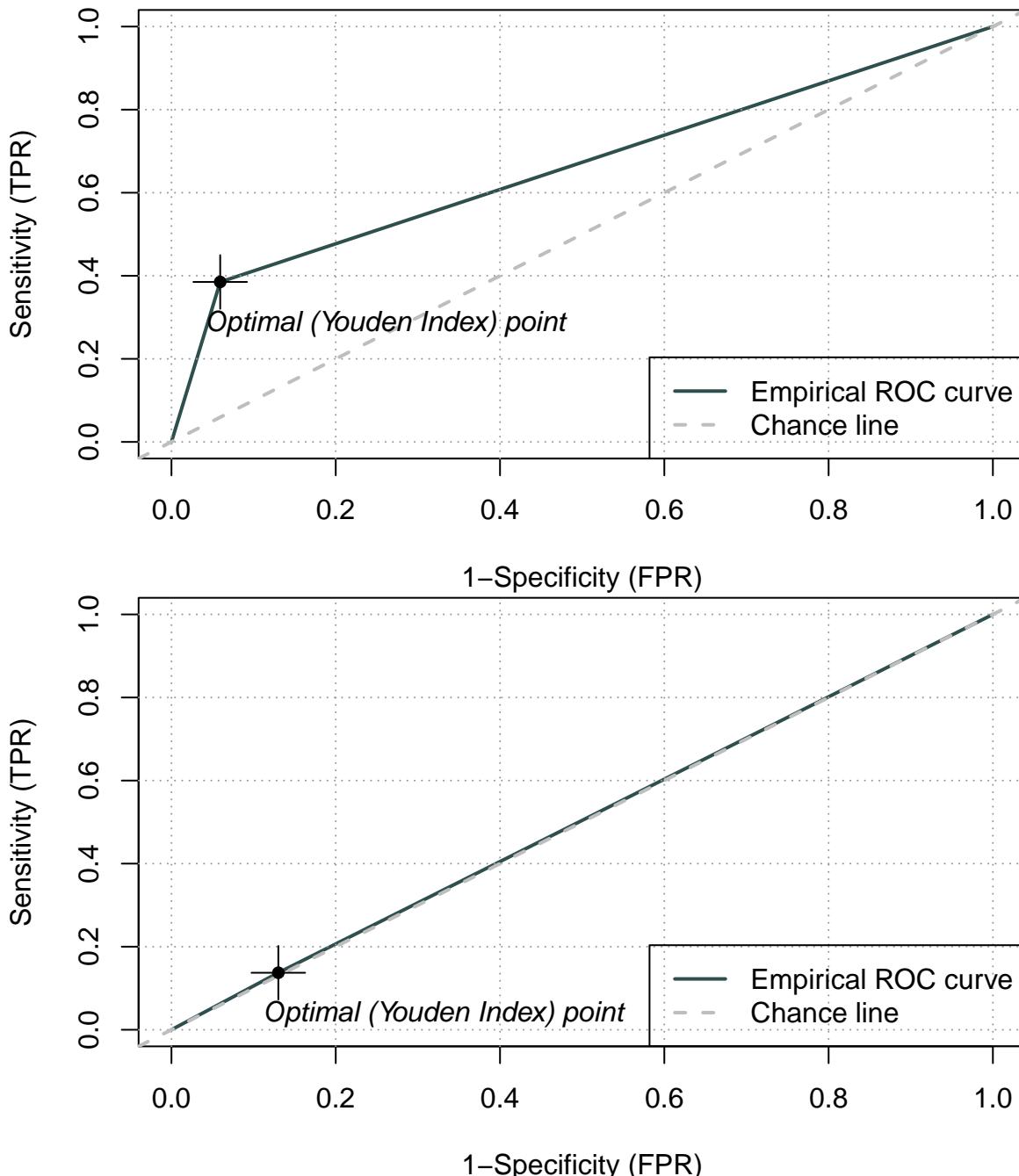
For null accuracy, which is the accuracy that could be achieved by always predicting the most frequent class, is used as a metrics for reference for the effectiveness of the model on overall. If a model performs better than the null accuracy significantly, it suggests that the model itself is effective and should be utilised in our decision making.

For harmonic mean, as it is a function of both recall and precision, it therefore strives and performs better in imbalanced datasets, where the accuracy score may be affected by the large number of data samples on one side. Therefore, we decided to include harmonic mean as a way to evaluate our model performance.

For ROC/AUC curves, the curve measures the sensitivity and specificity of the model, and provides us with a clearer view and understanding on our final model. It also serves as a way to prevent overfitting from happening in our model. Essentially, this visualisation method provides us with a clearer understanding on the performance of our models, and thus is included as a way to evaluate our model.

From our calculations, our null accuracy for the dataset is at 0.778.

Logistic Regression



Evaluation for Logistic Model

This is a model which we build using the logistics regression classification method. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

Logistics regression Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Test Accuracy: 0.708

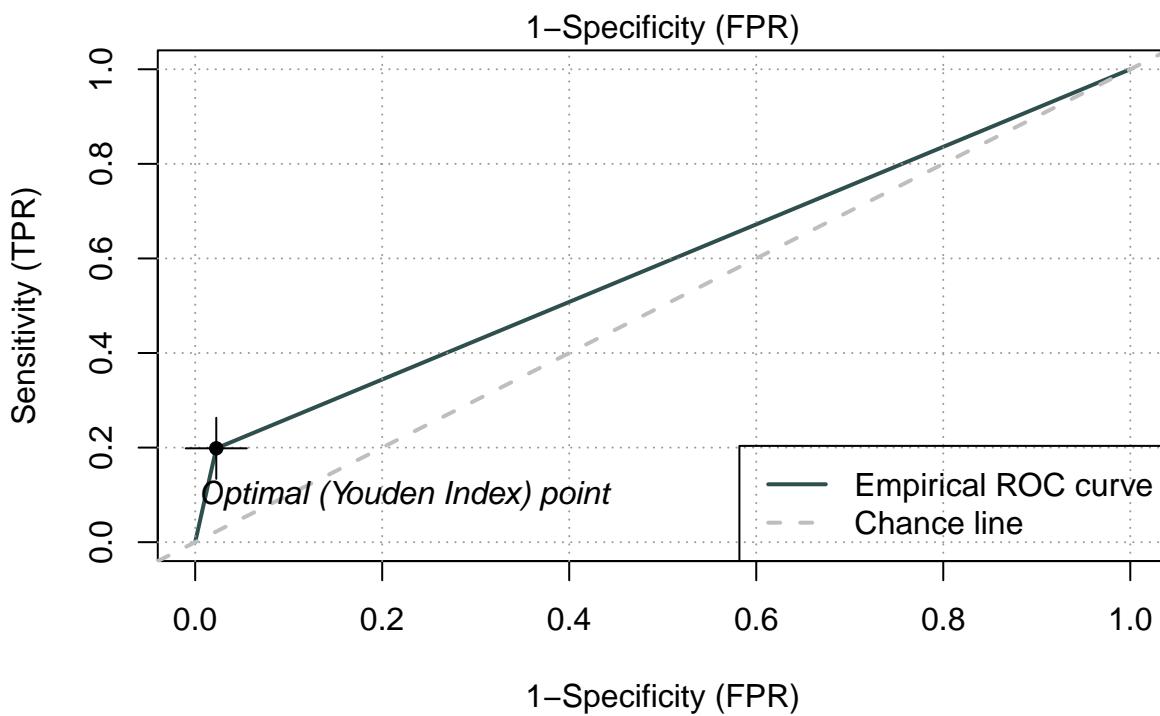
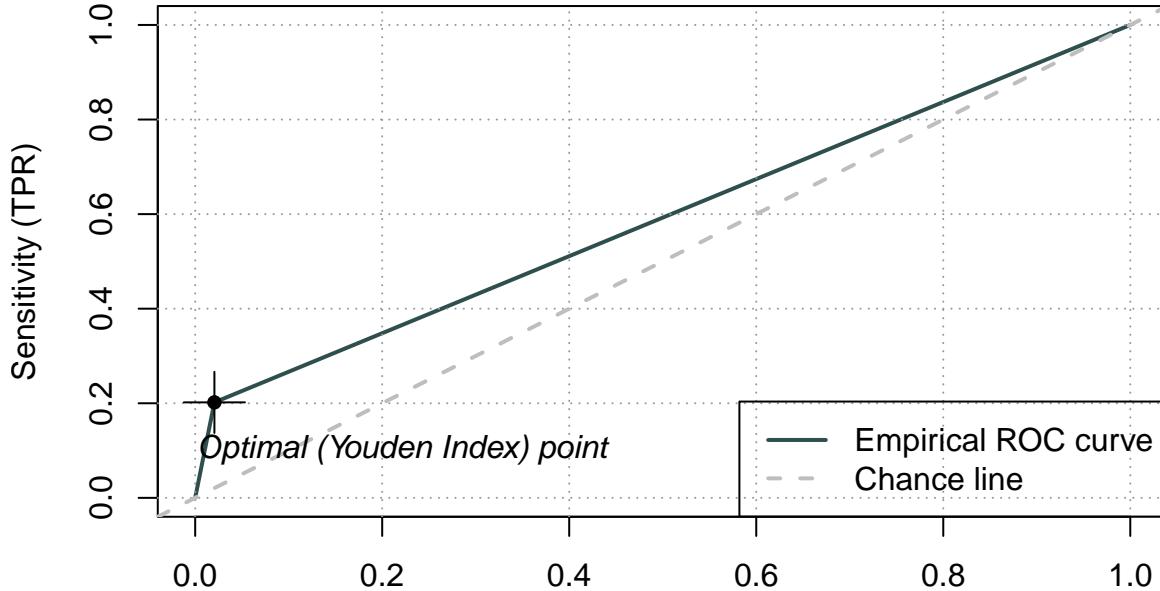
Train AUC: 0.6629

Test AUC: 0.5036

Test Harmonic Mean: 0.237

From the results, although the observed accuracy is quite high at 0.708, the harmonic mean is quite low at 0.237, which suggests a possible issue with the actual performance of our model. Furthermore, the AUC in test dataset is very low, which indicates that the model does not perform better by a significant margin relative to random selection.

SVM



Evaluation for SVM model

This is a model which we build using the support vector machine classification method. C-classification, as well as radial kernel is used as our final settings for the model, based on our repeated testing on the test dataset. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

SVM regression Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

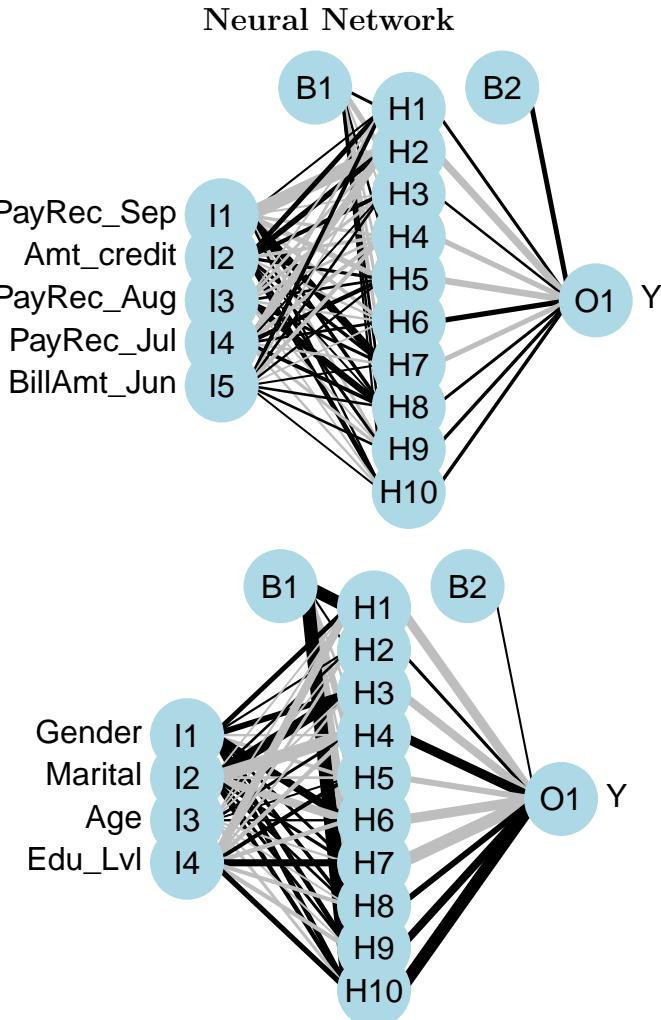
Test Accuracy: 0.805

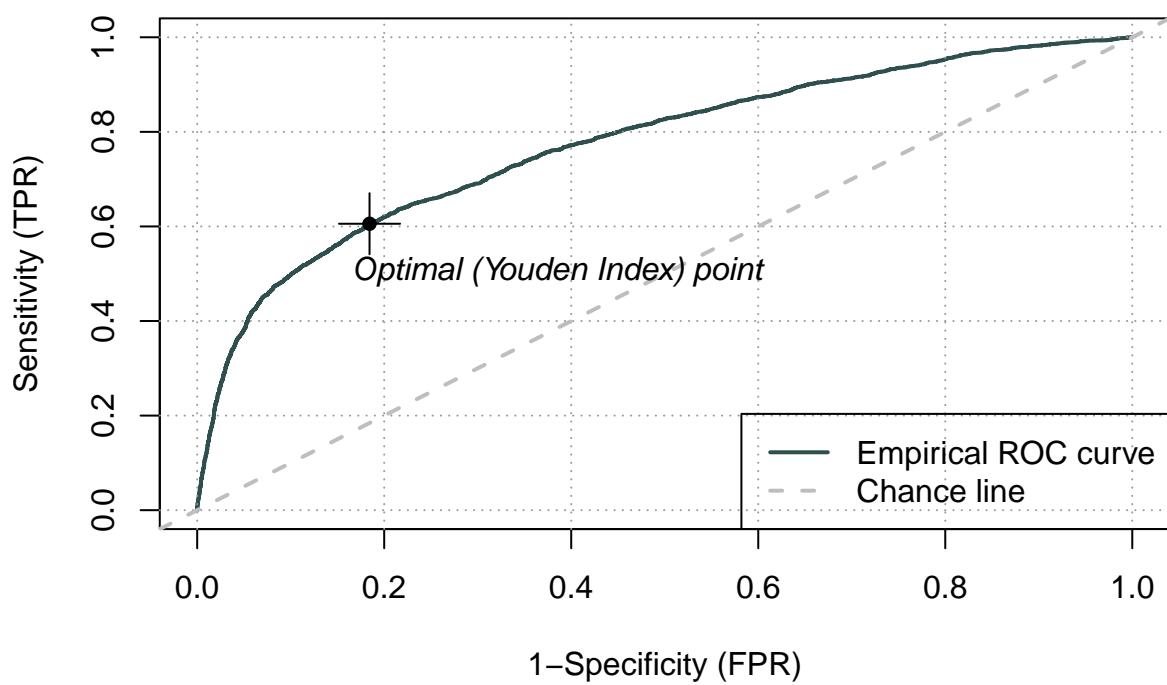
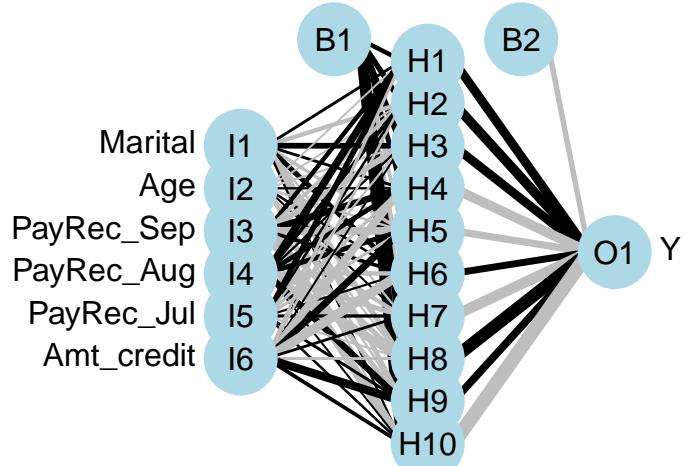
Train AUC: 0.591

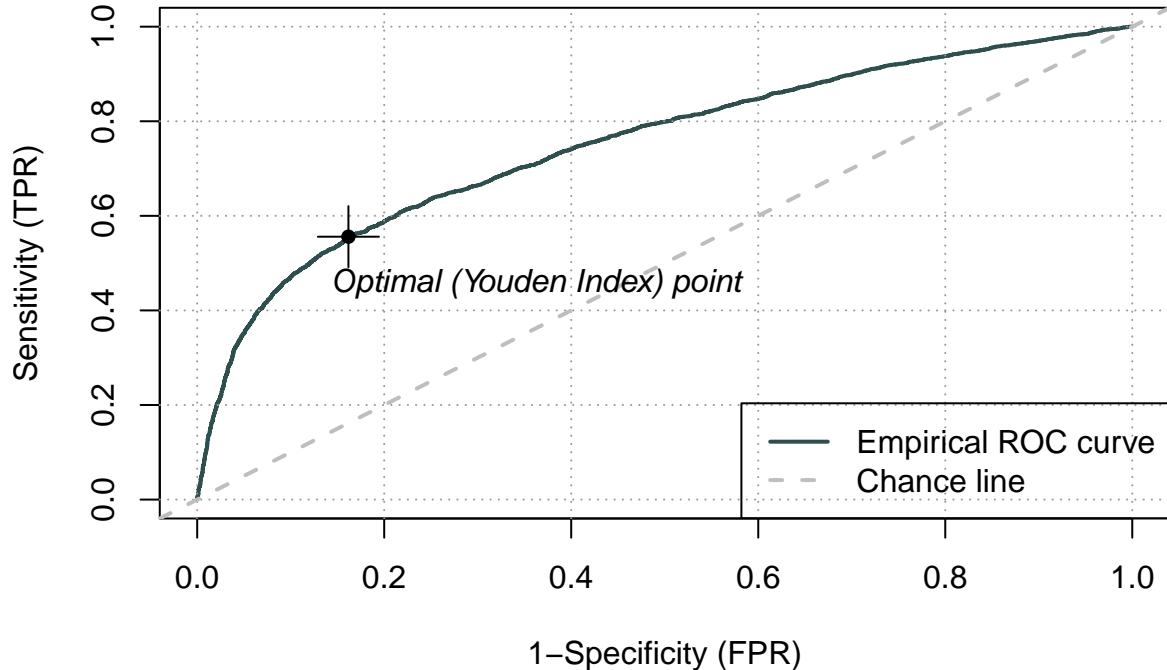
Test AUC: 0.588

Test Harmonic Mean: 0.330

From the results above, although the observed accuracy is quite high as well at 0.805, the harmonic mean for our test dataset is again quite low at 0.330, which suggests a possible issue with the actual performance of our model. However, from a AUC review, the train and test datasets achieved better results than that of the logistics regression model.







Evaluation for Neural Network model

This is a model which we build using the neural network classification method. A value of 1000 as max iterations to run, 10 hidden nodes and a decay factor of 0.08 is set for the model to prevent overfitting and achieve the best results. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

Neural Network Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Test Accuracy: 0.818

Train AUC: 0.653

Test AUC: 0.647

Test Harmonic Mean: 0.501

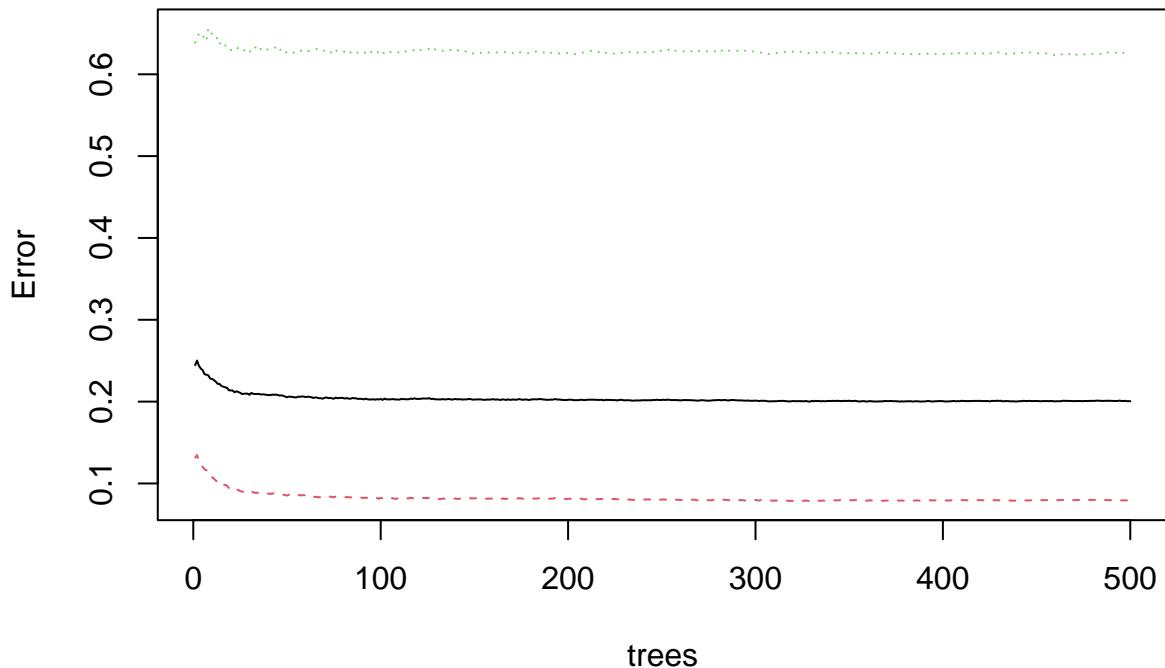
From the results above, as we can see, the test accuracy is again quite high albeit still remaining quite close to the null accuracy. We can see that the harmonic mean still remains low, and it suggests that the model may not perform as good as we think it is. For AUC, the AUC for test and train set are close, which suggests overfitting is not an issue. However, we still believe that we can create a better model make the predictions.

Random Forest Classifier

In this case, $\text{mtry} = 4$ is the best mtry as it has least OOB error. Coincidentally, $\text{mtry} = 4$ was also used as default mtry .

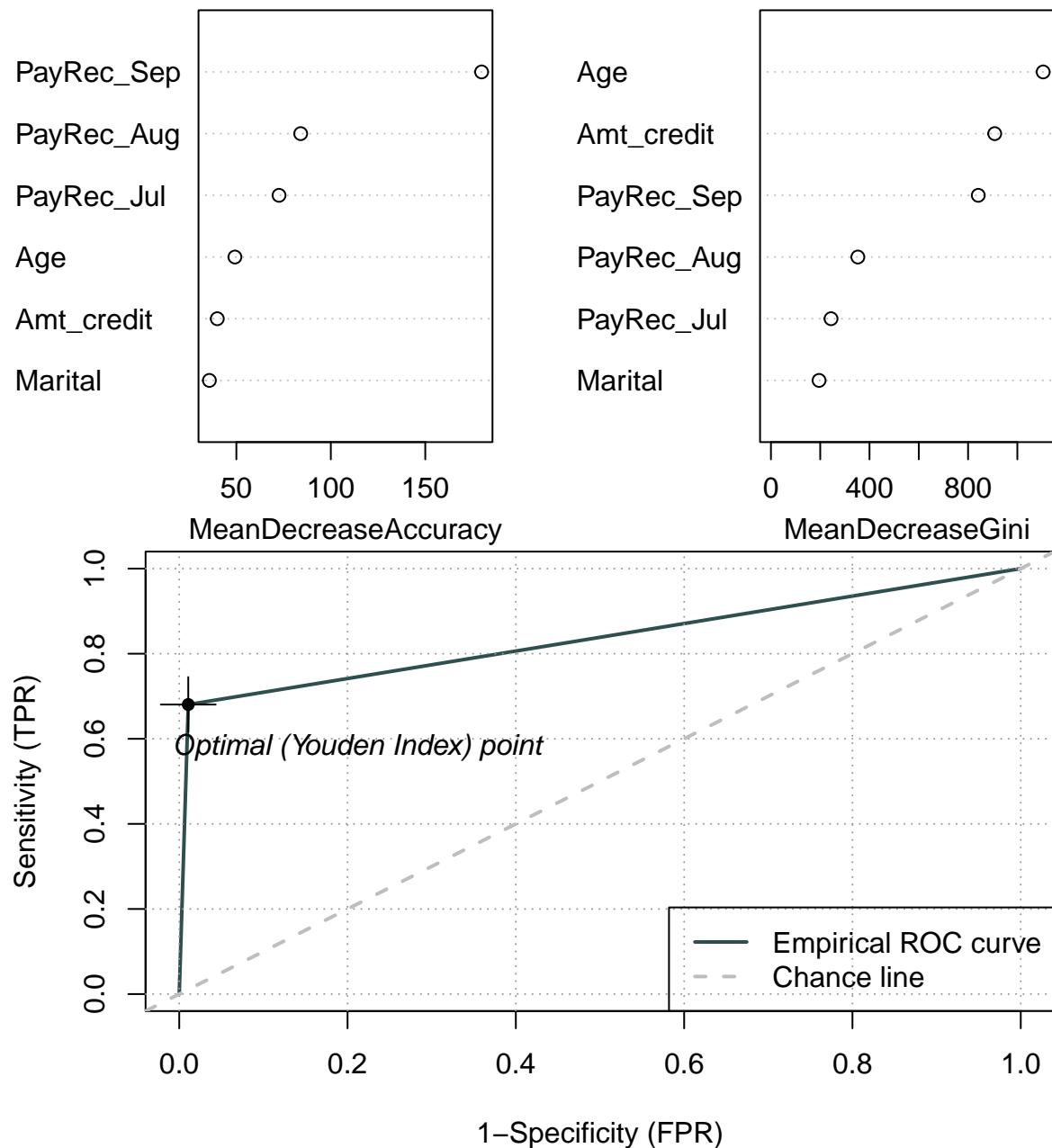
```
##      y_pred
##      0     1
##  0 10507  976
##  1 2065   1200
```

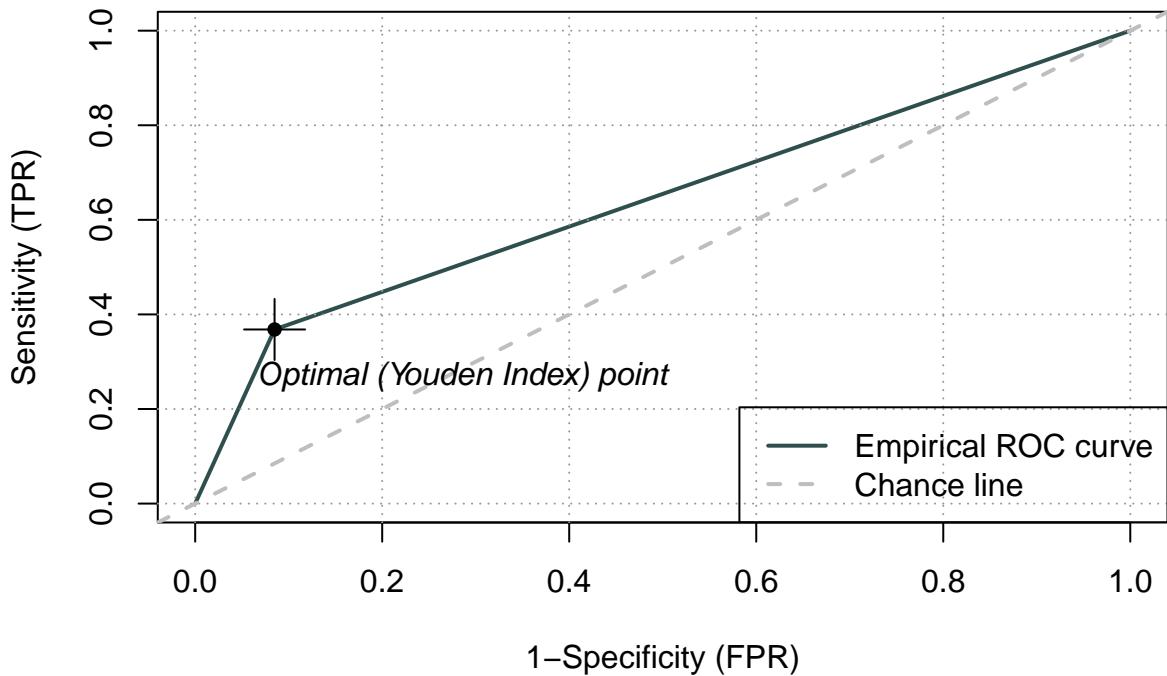
tuned_classifier_RF



```
##          0      1 MeanDecreaseAccuracy MeanDecreaseGini
## Marital    40.89120 -3.973139      35.72078      195.3995
## Age        49.26334  4.949118      49.23292     1104.6684
## Amt_credit 33.85993 16.085824      39.90170      908.1640
## PayRec_Sep 131.88687 69.726807     179.82329     841.4174
## PayRec_Aug  81.05604 -7.127203      84.04086     352.7161
## PayRec_Jul  54.28743 37.671283      72.59057     243.8465
```

tuned_classifier_RF





Evaluation for Random Forest model

This is a model which we build using the random forest decision tree classification method. From our repeated training, the number of nodes is set at 4, with TRUE set for gini importance and number of trees at 500. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

Seed 120

Random Forest with all variables:

Accuracy: 0.861

Train AUC: 1

Test AUC: 0.6241

Random Forest Chi Square with $Y \sim \text{Edu_Lvl} + \text{PayRec_May} + \text{PayRec_Jun} + \text{PayRec_Jul} + \text{PayRec_Aug}$:

Accuracy: 0.846

Train AUC: 0.5633

Test AUC: 0.5413

Random Forest first 5 of Boruta with $Y \sim \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul} + \text{BillAmt_Jun} + \text{PayRec_Jun}$:

Accuracy: 0.845

Train AUC: 0.753

Test AUC: 0.618

Seed 2

Random Forest Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Accuracy: 0.795

Train AUC: 0.832

Test AUC: 0.642

From the evaluation above, we can notice that our random forest model is consistent throughout different seed values, and thus makes it a reliable model that we can trust.

Overall, our final model that we would recommend would be a model with Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables in our model, and the algorithm for the model that we are using is the random forest decision tree algorithm.

Final Review and Discussion

As we run the model using Marital + Age + Amt_credit + PayRec_Sep + PayRec_Aug + PayRec_Jul, we observed the train AUC & test AUC values closer to 1. A model which has AUC near to the 1 which means it has a good measure of separability. By analogy, the higher the AUC, the better the model is at forecasting clients who have tendencies to default on the bank.

Overall, from all the models we run as shown above, despite the increasing harmonic means and AUCs we achieved across models, we are still unable to achieve a good prediction score for all the models.

From our repeated model trainings, we realise that taking a huge number of variables in as predictors does not only makes the model more complex, but also may overfit to the test dataset. On the other hand, too few predictors may lead to a model that would not classify the data well. Therefore, based on the importance of each variable, we decided to come up with the variables as stated above as our best features to predict the model.

This is due to a few factors, namely the small number of predictors we selected for our models to achieve more efficient and quick predictions without wastage of resources, potential overfitting issues of using too much predictors as shown in our random forest run for a model with all the variables included, as well as some potential data integrity issues that arises from the unexplained values that are found in our data visualisation. Unbalanced nature of our dataset also suggests that there may be harder to predict due to the lower number of default clients as well.

However, there is no denial that our model is not perfect and there are ways to improve on it. For example, the lower score for both AUC, as well as harmonic means suggests that there are still ways to improve on the model. We believe that with better knowledge and understanding of the data, we will be able to create a better model, as compared to the satisfactory model we have now.

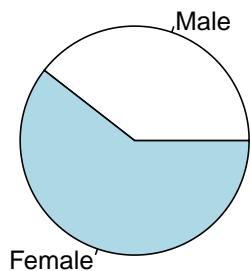
While reviewing our project, we discovered that the dataset is unbalanced which could result in bias in the model. This can be shown through the figures below.

Breakdown of Customers based on default status

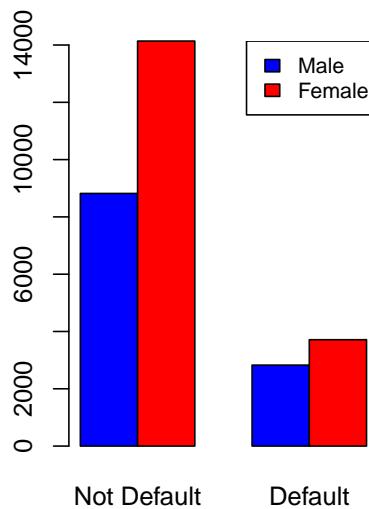
```
## Not Default      Default  
##      22957       6539
```

Breakdown of Customers based on Gender and Limit Balance

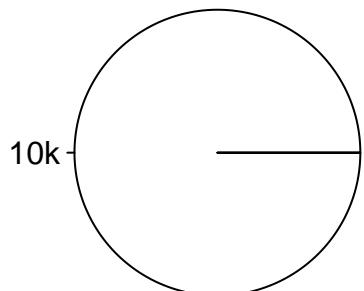
Credit Card Clients By Gender



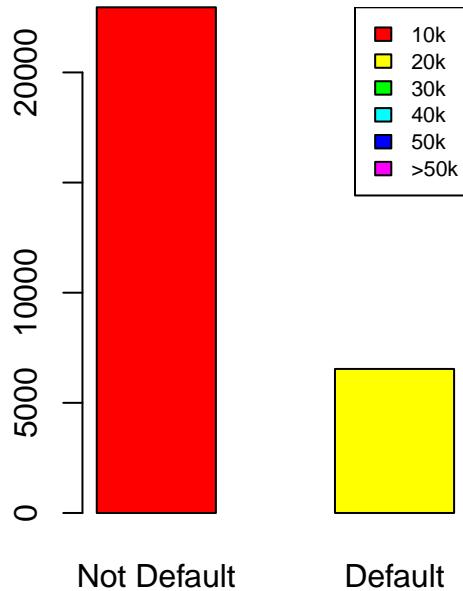
Defaulters vs Non-defaulters based on Gender



Credit Card Clients By Limit Balance



Defaulters vs Non-defaulters based on Limit Balance



BT2103 Project 3

Chen Haoli, Lo Zhi Hao, Luah Jun Yang, Toh Zhan Ting

Table of Contents

- 1) Exploratory data analysis on dataset ("card")
- 2) Data pre-processing
- 3) Feature Selection
- 4) Model Selection
- 5) Model Evaluation
- 6) Areas for improvement on dataset

Introduction to the Dataset and Problem Statement

The cash and credit card debt problem that Taiwan's credit card issuers experienced in recent years is predicted to peak in the third quarter of 2006 (Chou, 2006). Taiwan's card-issuing banks over-issued cash and credit cards to unqualified applicants in an effort to gain market dominance. In addition, most cardholders, regardless of their capacity to pay back, abused their cards for consumption and racked up large credit and cash-card debt. The crisis damaged consumer confidence in finance, and therefore presents a significant problem for both banks and cardholders.

Crisis management and risk prediction take place upstream and downstream in a mature financial system, respectively. The main goal of risk prediction is to lessen the harm and uncertainty caused by corporate performance or individual customer credit risk by using financial information, such as business financial statements, customer transaction and repayment histories, etc.

Therefore, with extensive data collected from the period of April to September in 2005, this report aims to build a predictive model to accurately forecast clients who have tendencies to default on the bank and thus ensures the profitability and stability of banks. In Finance, a default occurs when a borrower doesn't fulfill the terms of the loan. In this situation, default would occur if the cardholder failed to pay the credit card account within a given month.

In the dataset, there are 23 explanatory variables, together with a dependent variable of client's default status. The detailed description is as follows:

X1: Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her family (supplementary) credit.

X2: Gender (1 = male; 2 = female).

X3: Education (1 = graduate school; 2 = university; 3 = high school; 4 = others).

X4: Marital status (1 = married; 2 = single; 3 = others).

X5: Age (year).

X6–X11: History of past payment. We tracked the past monthly payment records (from April to September, 2005) as follows: X6 = the repayment status in September, 2005; X7 = the repayment status in August, 2005; . . . ; X11 = the repayment status in April, 2005. The measurement scale for the repayment status is: -1 = pay duly; 1 = payment delay for one month; 2 = payment delay for two months; . . . ; 8 = payment delay for eight months; 9 = payment delay for nine months and above.

X12–X17: Amount of bill statement (NT dollar). X12 = amount of bill statement in September, 2005; X13 = amount of bill statement in August, 2005; . . . ;X17 = amount of bill statement in April, 2005.

X18–X23: Amount of previous payment (NT dollar). X18 = amount paid in September, 2005; X19 = amount paid in August, 2005; . . . ;X23 = amount paid in April, 2005.

```
data_processed <- data[-1, ]  
#head(data_processed)  
#str(data_processed)
```

Exploratory Data Analysis

As all of the columns in the data are originally categorized as characters, in order to proceed with the data pre-processing and data visualization, we decided to transform the columns that are considered as continuous and numeric back to a numeric data type. After further research from reading through the description of our data, we identified x1, x5, x6, x7, x8, x9, x10, x11, x12, x13, x14, x15, x16, x17, x18, x19, x20, x21, x22, x23 as the data columns that we may need to transform. x2, x3, x4 are factors. For variables x3 and x4, we dropped the factor levels 0, which are indicative of N.A values for those variables Education and Marital Status respectively and they only amount to a small number of outliers. The other non-classified factor levels for these 2 variables are parked under the ‘others’ factor so that we will not exclude too many data points.

```
card <- data_processed %>% mutate(X1 = as.numeric(as.character(X1))) %>%  
  mutate(X2 = as.numeric(as.character(X2))) %>%  
  mutate(X3 = as.numeric(as.character(X3))) %>%  
  mutate(X4 = as.numeric(as.character(X4))) %>%  
  mutate(X5 = as.numeric(as.character(X5))) %>%  
  mutate(X5 = as.numeric(as.character(X5))) %>%  
  mutate(X6 = as.numeric(as.character(X6))) %>%  
  mutate(X7 = as.numeric(as.character(X7))) %>%  
  mutate(X8 = as.numeric(as.character(X8))) %>%  
  mutate(X9 = as.numeric(as.character(X9))) %>%  
  mutate(X10 = as.numeric(as.character(X10))) %>%  
  mutate(X11 = as.numeric(as.character(X11))) %>%  
  mutate(X12 = as.numeric(as.character(X12))) %>%  
  mutate(X13 = as.numeric(as.character(X13))) %>%  
  mutate(X14 = as.numeric(as.character(X14))) %>%  
  mutate(X15 = as.numeric(as.character(X15))) %>%  
  mutate(X16 = as.numeric(as.character(X16))) %>%  
  mutate(X17 = as.numeric(as.character(X17))) %>%  
  mutate(X18 = as.numeric(as.character(X18))) %>%  
  mutate(X19 = as.numeric(as.character(X19))) %>%  
  mutate(X20 = as.numeric(as.character(X20))) %>%  
  mutate(X21 = as.numeric(as.character(X21))) %>%  
  mutate(X22 = as.numeric(as.character(X22))) %>%  
  mutate(X23 = as.numeric(as.character(X23))) %>%  
  mutate(Y = as.numeric(Y))  
  
## Dropping factor levels that are unused  
card <- filter(card, X3 != 0)  
card <- filter(card, X4 != 0)  
  
## tests to make sure the values are changed accordingly  
#str(card)  
#class(card$X1)
```

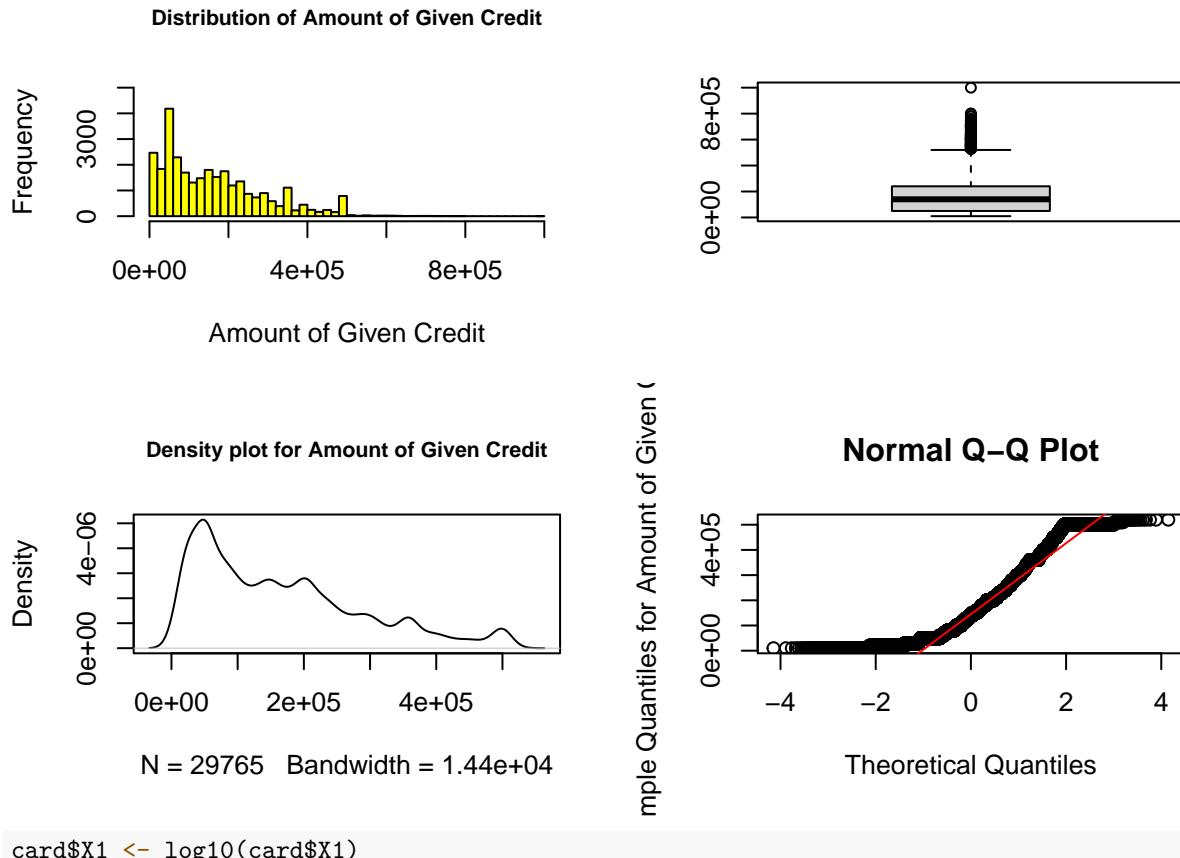
By using summary statistics and graphical representations, we are conducting preliminary analyses on the data in order to find trends, identify anomalies, inconsistencies, and missing values to test hypotheses and verify assumptions.

Amount of Given Credit (X1)

```
## X1 -> Amount of given credit
par(mfrow = c(2, 2))
hist(card$X1, main = "Distribution of Amount of Given Credit", xlim = c(0, 1000000),
     ylim = c(0, 5000), breaks = 50, col = "yellow", xlab = "Amount of Given Credit",
     ylab = "Frequency", cex.main = 0.8)
boxplot(card$X1)

card<-filter(card, X1 < (quantile(card$X1, 0.75) + 1.5*IQR(card$X1)))
plot(density(card$X1),main="Density plot for Amount of Given Credit", cex.main = 0.8)

qqnorm(card$X1,ylab="Sample Quantiles for Amount of Given Credit")
qqline(card$X1,col="red")
```



```
card$X1 <- log10(card$X1)
```

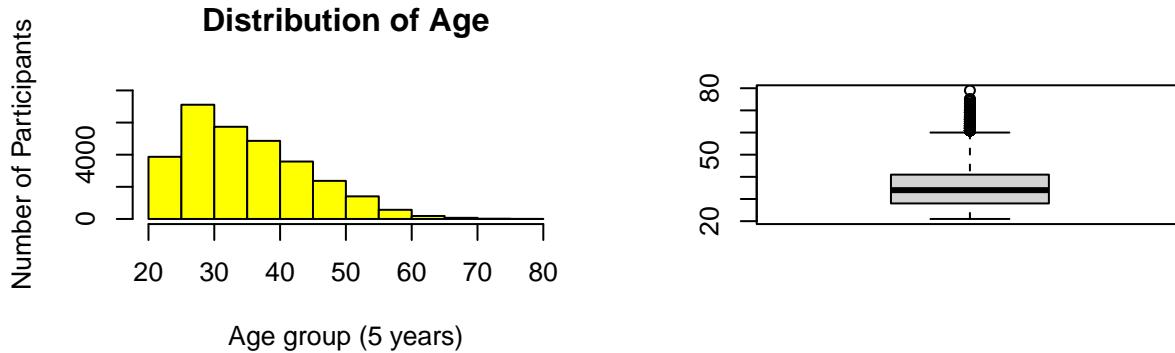
From the boxplot, we observed points lie outside $Q3 + 1.5\text{IQR}$ from the mean. We consider these points as outliers and we will remove these points from dataset From the density plot, we can observe that the data is positively skewed and it is not normally distributed.

From the qqplot, we can observe that the data is not normally distributed since the line and plots are not truly aligned.

Since the amount of given credit is positively skewed, we will use log transformation to improve the distribution of the data to normality.

Age (X5)

```
par(mfrow = c(2, 2))
## X5 -> Age
hist(card$X5, main = "Distribution of Age", ylim = c(0, 8000), breaks = 16,
     xlab = "Age group (5 years)", ylab = "Number of Participants", col = "yellow")
boxplot(card$X5)
card<-filter(card, X5 < quantile(card$X5,0.75) + 1.5 *IQR(card$X5))
```



From the boxplot, we observed points lie outside $Q3 + 1.5\text{IQR}$ from the mean. We consider these points as outliers and we will remove these points from dataset

Gender (X2), Education (X3), and Marital Status (X5)

```
## X2 -> gender
par(mfrow = c(1, 3))
gender_table <- table(card$X2)
names(gender_table) <- c("Male", "Female")
gender_table

## Male Female
## 11645 17851

barplot(table(card$X2), names.arg = c("Male", "Female"), col = c("blue", "pink"),
        ylim = c(0, 20000), ylab = "Number of Participants",
        xlab = "Gender", main = "Distribution of Gender")
# 11813 records are males, while 18020 are females

education_table <- table(card$X3)
names(education_table) <- c("Graduate School", "University",
                            "High School", "Others", "Unknown", "Unknown")
education_table

## Graduate School      University      High School      Others      Unknown
##           10419          13890           4740            122           277
##             Unknown
##                 48

barplot(table(card$X3),
       ylim = c(0, 16000),
       col = rainbow(length(c("Graduate School", "University", "High School",
                             "Others", "Unknown", "Unknown"))),
       main = "Distribution of Education Level", ylab = "Number of Participants",
       xlab = "Education Level", cex.names = 0.3)
```

```

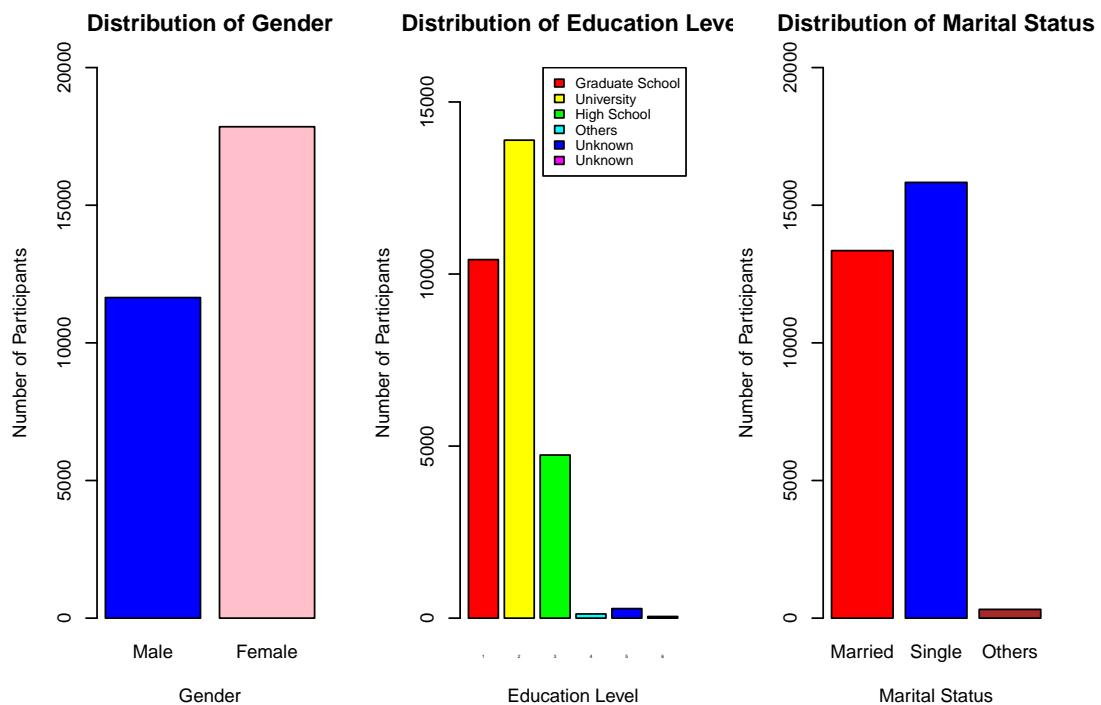
legend("topright", c("Graduate School", "University", "High School", "Others", "Unknown",
                     "Unknown"),
      fill = rainbow(length(c("Graduate School", "University", "High School",
                             "Others", "Unknown", "Unknown"))),
      cex = 0.7)

marital_table <- table(card$X4)
names(marital_table) <- c("Married", "Single", "Others")
marital_table

## Married Single Others
## 13350 15828 318

barplot(table(card$X4), main = "Distribution of Marital Status", ylim = c(0, 20000),
       names.arg = c("Married", "Single", "Others"), col = c("red", "blue", "brown"),
       ylab = "Number of Participants", xlab = "Marital Status")

```



Monthly Payment Records for 2005 (X6 - X11)

```

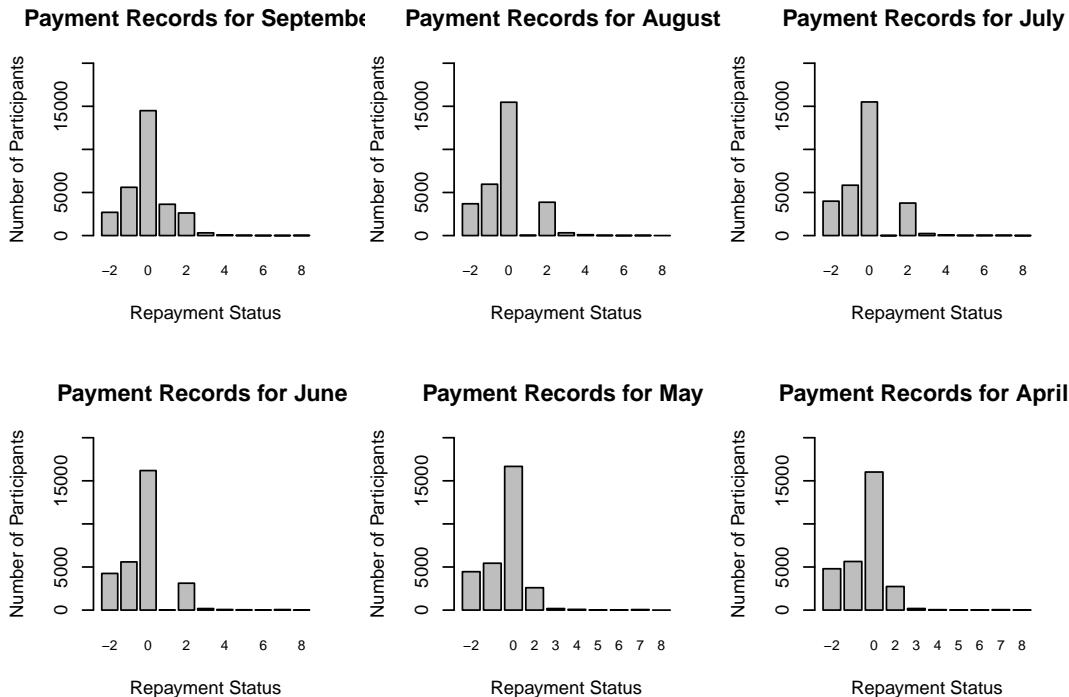
## X6 - 11 -> History of past payment. past monthly payment records (from April to September, 2005)
par(mfrow = c(2, 3))
## Tabulate the data
df <- data.frame(table(card$X6), table(card$X7), table(card$X8), table(card$X9))
## Visualisation
barplot(table(card$X6), main = "Payment Records for September ",
        ylab = "Number of Participants", xlab = "Repayment Status",
        ylim = c(0,20000), cex.names = 0.7)
barplot(table(card$X7), main = "Payment Records for August ",
        ylab = "Number of Participants", xlab = "Repayment Status",
        ylim = c(0,20000), cex.names = 0.7)
barplot(table(card$X8), main = "Payment Records for July ",
        ylab = "Number of Participants", xlab = "Repayment Status",
        ylim = c(0,20000), cex.names = 0.7)

```

```

ylab = "Number of Participants", xlab = "Repayment Status",
ylim = c(0,20000), cex.names = 0.7)
barplot(table(card$X9), main = "Payment Records for June ",
       ylab = "Number of Participants", xlab = "Repayment Status",
       ylim = c(0,20000), cex.names = 0.7)
barplot(table(card$X10), main = "Payment Records for May ",
       ylab = "Number of Participants", xlab = "Repayment Status",
       ylim = c(0,20000), cex.names = 0.7)
barplot(table(card$X11), main = "Payment Records for April ",
       ylab = "Number of Participants", xlab = "Repayment Status",
       ylim = c(0,20000), cex.names = 0.7)

```



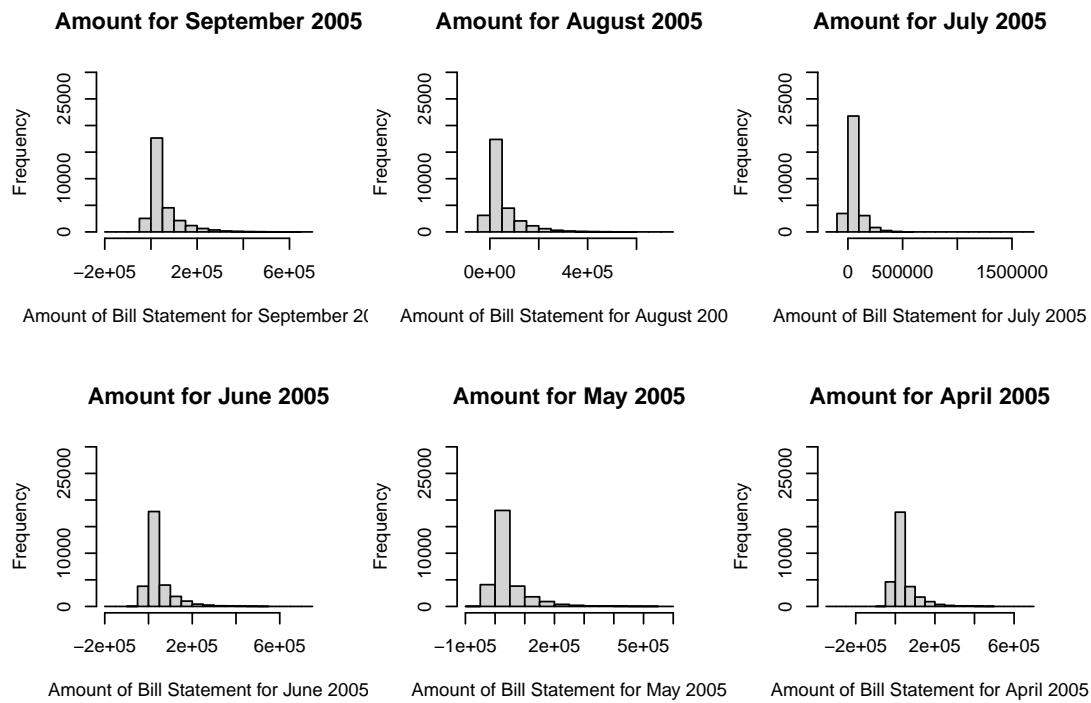
Monthly Bill Distributions for 2005 (X12 - X17)

```

## X12 - X17 -> Amount of bill statement (NT dollar).
## X12 = amount of bill statement in September 2005;
## X13 = amount of bill statement in August 2005; . . .
## X17 = amount of bill statement in April, 2005
par(mfrow = c(2, 3))
## Doing Visualisation for the data
hist(card$X12, xlab = "Amount of Bill Statement for September 2005",
      main = "Amount for September 2005", ylim = c(0, 30000))
hist(card$X13, xlab = "Amount of Bill Statement for August 2005",
      main = "Amount for August 2005", ylim = c(0, 30000))
hist(card$X14, xlab = "Amount of Bill Statement for July 2005",
      main = "Amount for July 2005", ylim = c(0, 30000))
hist(card$X15, xlab = "Amount of Bill Statement for June 2005",
      main = "Amount for June 2005", ylim = c(0, 30000))
hist(card$X16, xlab = "Amount of Bill Statement for May 2005",
      main = "Amount for May 2005", ylim = c(0, 30000))
hist(card$X17, xlab = "Amount of Bill Statement for April 2005",
      main = "Amount for April 2005", ylim = c(0, 30000))

```

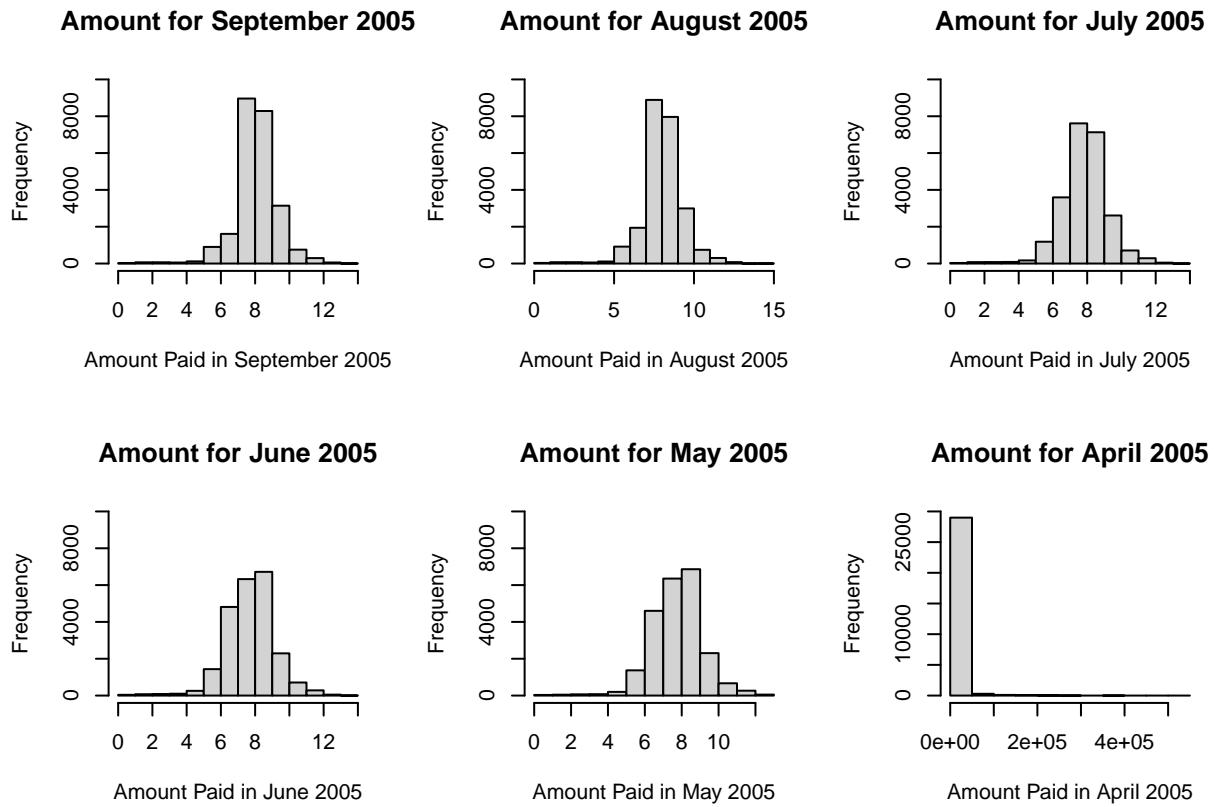
```
main = "Amount for April 2005", ylim = c(0, 30000))
```



Amount of previous payment (X18 - X23)

We will perform log transformation for Amount Paid in order to make it resemble a normal distribution.

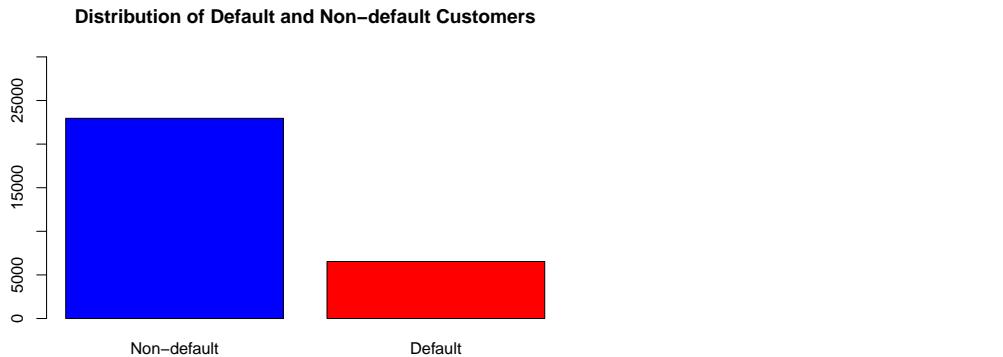
```
## X18 - X23 -> Amount of previous payment (NT dollar).
## X18 = amount paid in September, 2005;
## X19 = amount paid in August, 2005; . . .
## X23 = amount paid in April, 2005
par(mfrow = c(2, 3))
## Log Transformed Data Visualisation In order to make it resemble a normal distribution
hist(log(card$X18), xlab = "Amount Paid in September 2005",
     main = "Amount for September 2005", ylim = c(0, 10000))
hist(log(card$X19), xlab = "Amount Paid in August 2005",
     main = "Amount for August 2005", ylim = c(0, 10000))
hist(log(card$X20), xlab = "Amount Paid in July 2005",
     main = "Amount for July 2005", ylim = c(0, 10000))
hist(log(card$X21), xlab = "Amount Paid in June 2005",
     main = "Amount for June 2005", ylim = c(0, 10000))
hist(log(card$X22), xlab = "Amount Paid in May 2005",
     main = "Amount for May 2005", ylim = c(0, 10000))
hist(card$X23, xlab = "Amount Paid in April 2005",
     main = "Amount for April 2005", ylim = c(0, 30000))
```



Distribution of Default and Non-default data

```
default_table <- table(card$Y)
names(default_table) <- c("Non-default", "Default")
default_table
```

```
## Non-default      Default
##          22957       6539
barplot(table(card$Y), ylim = c(0, 30000),
        main = "Distribution of Default and Non-default Customers",
        names.arg = c("Non-default", "Default"), col = c("blue","red"))
```



Data Pre-Processing

As we noticed from the data visualisation shown in our part above, there is a significant issue with the data we gathered from our dataset, as some of the data are not described in our original cited data source.

For example, for distribution of education level and marital status, we noticed some entries that does not belong to any factor levels that were described in the original data source. Similarly, for the history of payment records, we also realized the occurrence of erroneous factor levels such as -2 and 0, which are not described for in the original data source.

First, we are going to change the names of each columns for better understanding and easier interpretation in our following chapters.

Moving on we are going to perform some data manipulation on some of the data anomalies that we spotted earlier. In that sense, we are going to manipulate and change the data points that are originally not described in our cited data source.

```
## Renaming the column names
card_renamed <- card %>% rename(Amt_credit = X1) %>%
  rename(Gender = X2) %>%
  rename(Edu_Lvl = X3) %>%
  rename(Marital = X4) %>%
  rename(Age = X5) %>%
  rename(PayRec_Sep = X6) %>%
  rename(PayRec_Aug = X7) %>%
  rename(PayRec_Jul = X8) %>%
  rename(PayRec_Jun = X9) %>%
  rename(PayRec_May = X10) %>%
  rename(PayRec_Apr = X11) %>%
  rename(BillAmt_Sep = X12) %>%
  rename(BillAmt_Aug = X13) %>%
  rename(BillAmt_Jul = X14) %>%
  rename(BillAmt_Jun = X15) %>%
  rename(BillAmt_May = X16) %>%
  rename(BillAmt_Apr = X17) %>%
  rename(PaidAmt_Sep = X18) %>%
  rename(PaidAmt_Aug = X19) %>%
  rename(PaidAmt_Jul = X20) %>%
  rename(PaidAmt_Jun = X21) %>%
  rename(PaidAmt_May = X22) %>%
  rename(PaidAmt_Apr = X23)
#(card_renamed)

## Manipulating the erroneous data that are observed for EDUCATION LEVEL
card_renamed <- card_renamed %>% filter(Edu_Lvl != 0)
## Checking whether 0 values are removed
#nrow(card_renamed)
#nrow(card)

## Adding 5, 6 into 4
card_renamed$Edu_Lvl[card_renamed$Edu_Lvl == 5] <- 4
card_renamed$Edu_Lvl[card_renamed$Edu_Lvl == 6] <- 4
## Checking whether 5, 6 values are manipulated
card_renamed <- card_renamed[, -1]
#table(card_renamed$Edu_Lvl)
```

Data Manipulation for Marital Status (X2)

As we can see from the visualisation for distribution of marital status for our clients, there are categories that are not described for in our original data cited source, namely 0. For our team, as we consider 0 values as possible null values where our clients did not provide their marital status information properly, we decided to drop rows with 0 values recorded for marital status.

```
## Manipulating the erroneous data that are observed MARITAL STATUS
card_renamed <- card_renamed %>% filter(Marital != 0)
## Checking whether 0 values are removed
table(card_renamed$Marital)

##
##      1      2      3
## 13350 15828   318
```

Data Manipulation for Education Status (X3)

As we can see from the visualisation for distribution of education level for our clients, there are categories that are not described for in our original data cited source, namely 0, 5 and 6. For our team, as we consider 0 as a possible null value where our clients did not provide their academic credentials, we decided to drop those observations with 0 recorded as their education level. Meanwhile, as 5 and 6 may be education levels that are higher than or not included in the provided options (such as PhD), we decided to manipulate those data to add them into the 'others' category, 4 that is provided.

```
## Manipulating the erroneous data that are observed
card_renamed <- card_renamed %>% filter(Edu_Lvl != 0)

## Checking whether 0 values are removed
nrow(card_renamed)
nrow(card)

## Adding 5, 6 into 4
card_renamed$Edu_Lvl[card_renamed$Edu_Lvl == 5] <- 4
card_renamed$Edu_Lvl[card_renamed$Edu_Lvl == 6] <- 4

## Checking whether 5, 6 values are manipulated
table(card_renamed$Edu_Lvl)

##
##      1      2      3      4
## 10419 13890  4740   447
```

Feature Selection

Make use of feature selection methodologies to select the most relevant independent variables to create a prediction model.

For our project, as it is a large dataset with around 30,000 observations, for each type of feature selection method, we are going to perform a test in order to determine the best features in that category, namely:

Forward and Backward variable selection

Filter Method - ANOVA for continuous data and Chi Squared Test for categorical data

Wrapper Method - Boruta Method

Embedded Method - Information Gain

Forward Variable Selection:

- Start with model containing no possible explanatory variable and for each variable in turn, we will investigate effect of adding variable from the current model.
- 5 most significant variables will be used for our subsequent models

```
forwardSearch<-regsubsets(Y~., data=card_renamed,nbest=1, nvmax=5,method = "forward")

info_Forward <- summary(forwardSearch)

cbind(info_Forward$which, round(cbind(rsq=info_Forward$rsq,adjr2=info_Forward$adjr2,
                                     cp=info_Forward$cp, bic=info_Forward$bic,
                                     rss=info_Forward$rss), 3))

##   (Intercept) Amt_credit Gender Edu_Lvl Marital Age PayRec_Sep PayRec_Aug
## 1           1          0     0      0     0    0       1         0
## 2           1          1     0      0     0    0       1         0
## 3           1          1     0      0     0    0       1         0
## 4           1          1     0      0     0    0       1         1
## 5           1          1     0      0     1    0       1         1
##   PayRec_Jul PayRec_Jun PayRec_May PayRec_Apr BillAmt_Sep BillAmt_Aug
## 1           0          0        0        0       0         0
## 2           0          0        0        0       0         0
## 3           0          0        0        0       1         0
## 4           0          0        0        0       1         0
## 5           0          0        0        0       1         0
##   BillAmt_Jul BillAmt_Jun BillAmt_May BillAmt_Apr PaidAmt_Sep PaidAmt_Aug
## 1           0          0        0        0       0         0
## 2           0          0        0        0       0         0
## 3           0          0        0        0       0         0
## 4           0          0        0        0       0         0
## 5           0          0        0        0       0         0
##   PaidAmt_Jul PaidAmt_Jun PaidAmt_May PaidAmt_Apr   rsq adjr2      cp      bic
## 1           0          0        0        0   0 0.106 0.106 659.541 -3290.498
## 2           0          0        0        0   0 0.113 0.113 422.683 -3514.801
## 3           0          0        0        0   0 0.117 0.116 315.037 -3612.826
## 4           0          0        0        0   0 0.121 0.120 179.525 -3738.945
## 5           0          0        0        0   0 0.122 0.122 129.422 -3780.503
##   rss
## 1 4548.953
## 2 4512.917
## 3 4496.375
## 4 4475.628
## 5 4467.768
```

Backward Variable Selection:

- Start with model containing all possible explanatory variables and for each variable in turn, we will investigate effect of removing that variable from the current model.
- 5 most significant variables will be used for our subsequent models

```
backwardSearch<-regsubsets(Y~., data=card_renamed,nbest=1, nvmax=5,method = "backward")

info_Backward <- summary(backwardSearch)

cbind(info_Backward$which, round(cbind(rsq=info_Backward$rsq,adjr2=info_Backward$adjr2,
```

```

          cp=info_Backward$cp, bic=info_Backward$bic,
          rss=info_Backward$rss), 3))

##   (Intercept) Amt_credit Gender Edu_Lvl Marital Age PayRec_Sep PayRec_Aug
## 1           1         0     0     0     0     0       1       0
## 2           1         0     0     0     0     0       1       0
## 3           1         0     0     0     0     0       1       1
## 4           1         1     0     0     0     0       1       1
## 5           1         1     0     0     0     1       1       1
##   PayRec_Jul PayRec_Jun PayRec_May PayRec_Apr BillAmt_Sep BillAmt_Aug
## 1           0         0       0       0       0       0
## 2           0         0       0       0       1       0
## 3           0         0       0       0       1       0
## 4           0         0       0       0       1       0
## 5           0         0       0       0       1       0
##   BillAmt_Jul BillAmt_Jun BillAmt_May BillAmt_Apr PaidAmt_Sep PaidAmt_Aug
## 1           0         0       0       0       0       0
## 2           0         0       0       0       0       0
## 3           0         0       0       0       0       0
## 4           0         0       0       0       0       0
## 5           0         0       0       0       0       0
##   PaidAmt_Jul PaidAmt_Jun PaidAmt_May PaidAmt_Apr   rsq adjr2      cp      bic
## 1           0         0       0       0   0.106 0.106 659.541 -3290.498
## 2           0         0       0       0   0.113 0.113 426.521 -3511.017
## 3           0         0       0       0   0.119 0.119 230.592 -3696.520
## 4           0         0       0       0   0.121 0.120 179.525 -3738.945
## 5           0         0       0       0   0.122 0.122 130.981 -3778.950
##   rss
## 1 4548.953
## 2 4513.496
## 3 4483.635
## 4 4475.628
## 5 4468.003

```

From both the Forward and Backward Variable Selection, we observed similarities in the significant variables selected in both models (i.e., Amt_credit, PayRec_Sep, PayRec_Aug, PayRec_Jul). However, we also observed that Marital Status & Age are ranked differently in both models. Thus, in order to increase accuracy of our subsequent models, we decided to include both Marital & Age.

Splitting into train and test set

First, we will split our dataset into train and test set before actually performing the feature selection models on our data. This is to prevent any leakage of information from our test set into our training set leading to biased and overfitted models.

```

## Splitting into train and test set

n <- length(card_renamed$Y)
set.seed(123)
index <- 1:nrow(card_renamed)
trainindex <- sample(index, trunc(n)/2)
train.data <- card_renamed[trainindex,]
test.data <- card_renamed[-trainindex,]
ntrain <- length(train.data)
ntest <- length(test.data)

```

```

#nrow(train.data)
#nrow(test.data)

#head(train.data)
#head(test.data)

```

Anova Test

```

## Anova test for each continuous variable
## Select continuous variables
cont.vars <- train.data[, c(1:2, 6, 13:24)]
#cont.vars
aov.stats <- data.frame(nrow = 14, ncol = 2)
colnames(aov.stats) <- c("Variable", "p-value")
colnames <- c("Amt_Credit", "Age", "BillAmt_Sep", "BillAmt_Aug", "BillAmt_Jul",
             "BillAmt_Jun", "BillAmt_May", "BillAmt_Apr", "PaidAmt_Sep",
             "PaidAmt_Aug", "PaidAmt_Jul", "PaidAmt_Jun", "PaidAmt_May", "PaidAmt_Apr")

col <- ncol(cont.vars) - 1
for (I in 1:col) {
  x <- cont.vars[, I + 1]
  #plot(x ~ train.data$Y)
  x.aov <- aov(x ~ train.data$Y)
  tests <- summary(x.aov)[[1]][1,5]
  #str(tests)
  #p_value <- tests[[1]]$'Pr(>F)'
  aov.stats[I, 1] <- colnames[I]
  aov.stats[I, 2] <- tests
}
aov.stats[order(aov.stats[,2]), ]

```

##	Variable	p-value
## 2	Age	0.000000e+00
## 14	PaidAmt_Apr	0.000000e+00
## 8	BillAmt_Apr	5.897689e-20
## 10	PaidAmt_Aug	4.891178e-15
## 11	PaidAmt_Jul	8.157781e-13
## 12	PaidAmt_Jun	1.371457e-12
## 9	PaidAmt_Sep	1.665843e-11
## 1	Amt_Credit	2.898041e-11
## 13	PaidAmt_May	2.794714e-10
## 3	BillAmt_Sep	6.482189e-03
## 4	BillAmt_Aug	1.310158e-02
## 5	BillAmt_Jul	7.495121e-02
## 6	BillAmt_Jun	1.744674e-01
## 7	BillAmt_May	1.864872e-01
## Variables Amt_Credit, PaidAmt_Sep, PaidAmt_Jun, PaidAmt_Jul, PaidAmt_Apr, PaidAmt_May,		
## PaidAmt_Aug should be choosen according to ANOVA tests		

Chi Square test

Next, for categorical variables, Chi Squared Method is used to identify the most important features.

```

## Chi square test is used for each categorical variables

## Select categorical variables
cat.vars <- train.data[, c(3:5, 7:12)]
chi.stats <- data.frame(nrow = 9, ncol = 3)
colnames(chi.stats) <- c("Statistics", "p-value")
colnames <- c("Gender", "Edu_Lvl", "Marital", "PayRec_Sep", "PayRec_Aug",
             "PayRec_Jul", "PayRec_Jun", "PayRec_May", "PayRec_Apr")

col <- ncol(cat.vars) - 1
for (I in 1:col) {
  x <- cat.vars[, I + 1]
  tbl <- table(x, train.data$Y)
  #print(tbl)
  chi2res <- chisq.test(tbl)
  #print(chi2res)
  chi.stats[I, 1] <- colnames[I]
  chi.stats[I, 2] <- chi2res$statistic
  chi.stats[I, 3] <- chi2res$p.value
}

chi.stats[order(chi.stats[,2]), ]

##   Statistics      p-value      V3
## 1    Gender  6.713353 3.485089e-02
## 2    Edu_Lvl  85.508305 2.492442e-05
## 7 PayRec_Jun  983.964863 4.943317e-206
## 6 PayRec_Jul  1054.926738 2.457447e-221
## 5 PayRec_Aug  1131.704234 7.713010e-237
## 4 PayRec_Sep  1303.027160 8.499636e-274
## 3   Marital  1812.283792 0.000000e+00
## 8 PayRec_May 11924.409892 8.638387e-01

## Variables PayRec_Aug, PayRec_Jul, PayRec_Jun, PayRec_May
## should be chosen according to Chi Square test

```

Wrapper method for feature selection

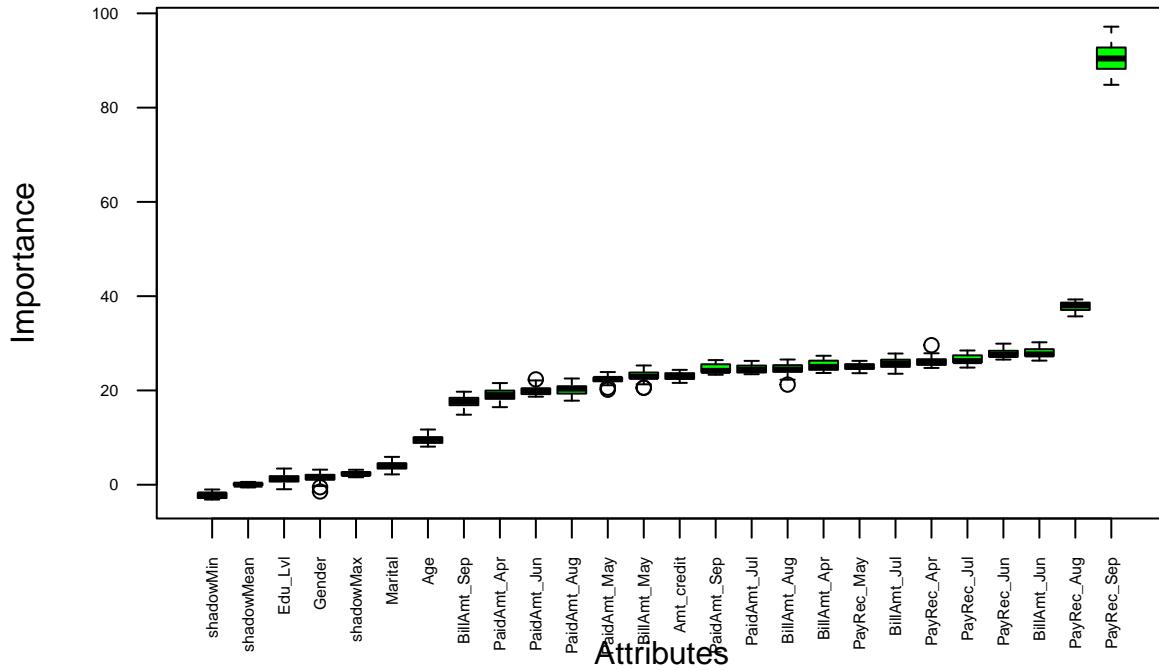
For wrapper method, the package and method that we are going to use is the Boruta Method that utilises random forest decision tree model in computing for the importance of each feature.

```

## Running the Boruta method for feature importance
set.seed(123)
##idpt_vars <- train.data[, 2:23]
##head(idpt_vars)
## Running the feature selection Boruta Method
boruta <- Boruta(Y ~ ., data = train.data, doTrace = 2, maxRuns = 20)

## print results and plot
#print(boruta)
plot(boruta, las = 2, cex.axis = 0.5)

```



Embedded method for feature selection

We utilise the information gain function from the FSelectorRcpp package to inspect and identify the important features in our model.

```
## running the information gain algorithm to identify importance of each features
IG.Fselector <- information_gain(Y ~ ., data = train.data)
#print(IG.Fselector[order(IG.Fselector[, 2]), ])
```

Correlation

```
# Correlation
card_renamed_corr <- card_renamed %>%
  mutate(Gender = as.numeric(as.character(Gender))) %>%
  mutate(Edu_Lvl = as.numeric(as.character(Edu_Lvl))) %>%
  mutate(Marital = as.numeric(as.character(Marital))) %>%
  mutate(Y = as.numeric(as.character(Y)))

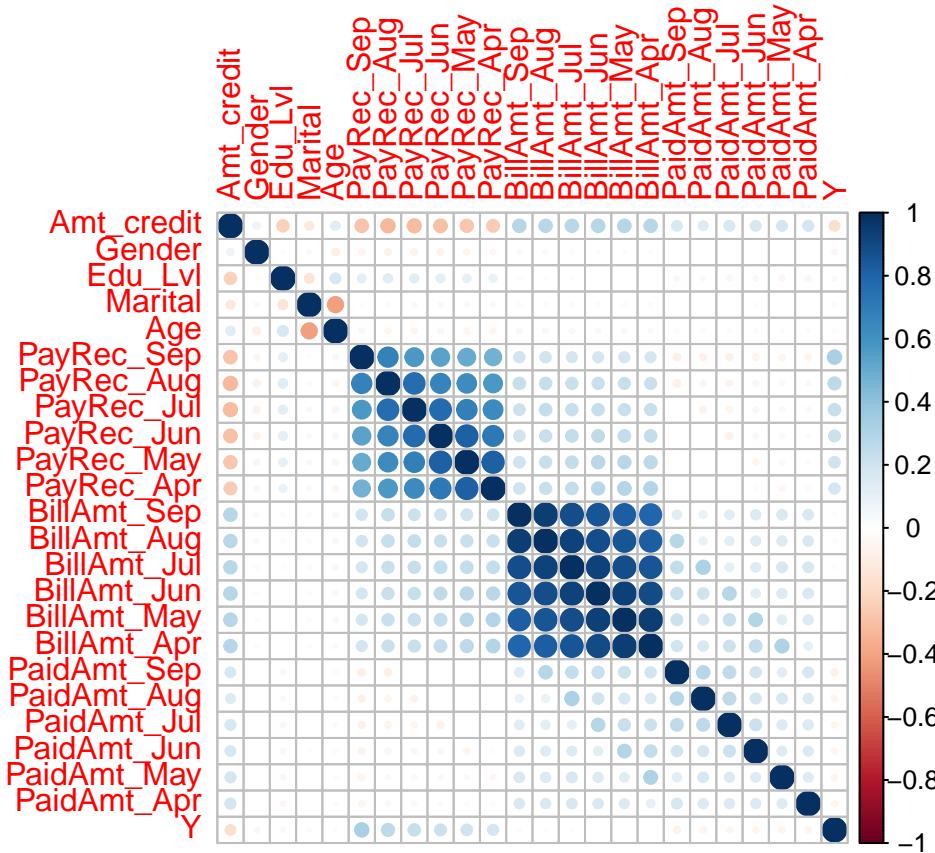
correlationMatrix <- cor(card_renamed_corr)

# summarize the correlation matrix
#print(correlationMatrix)

# find attributes that are highly correlated (ideally >0.75, but used 0.5 here)
highlyCorrelated <- findCorrelation(correlationMatrix, cutoff = 0.5)

# print indexes of highly correlated attributes
#print(highlyCorrelated)

corrplot(correlationMatrix)
```



The highly correlated attributes are BillAmt_Jun, BillAmt_May, BillAmt_Jul, BillAmt_Apr, BillAmt_Aug, PayRec_Jun, PayRec_Jul, PayRec_Aug, PayRec_May.

Model Selection

For our models, a way to identify the effectiveness and accuracy of each of our models is to evaluate each of the models based on accuracy, null accuracy, ROC and AUC values and harmonic mean. Each of these methods have their own merits and advantages, and by computing the confusion matrix and compute each of these values, we can get a deeper understanding on how each of our models are performing.

For null accuracy, which is the accuracy that could be achieved by always predicting the most frequent class, is used as a metrics for reference for the effectiveness of the model on overall. If a model performs better than the null accuracy significantly, it suggests that the model itself is effective and should be utilised in our decision making.

For harmonic mean, as it is a function of both recall and precision, it therefore strives and performs better in imbalanced datasets, where the accuracy score may be affected by the large number of data samples on one side. Therefore, we decided to include harmonic mean as a way to evaluate our model performance.

For ROC/AUC curves, the curve measures the sensitivity and specificity of the model, and provides us with a clearer view and understanding on our final model. It also serves as a way to prevent overfitting from happening in our model. Essentially, this visualisation method provides us with a clearer understanding on the performance of our models, and thus is included as a way to evaluate our model.

From our calculations, our null accuracy for the dataset is at 0.778.

```
harmonic_mean <- function(table) {
  tab <- data.frame(table)
  return ((1/(0.5 * ((tab[1,3]+tab[2,3])/tab[1,3] + (tab[3,3]+tab[4,3])/tab[4,3])))
```

```

## Computing null accuracy
## training dataset





```

Logistic Regression

```

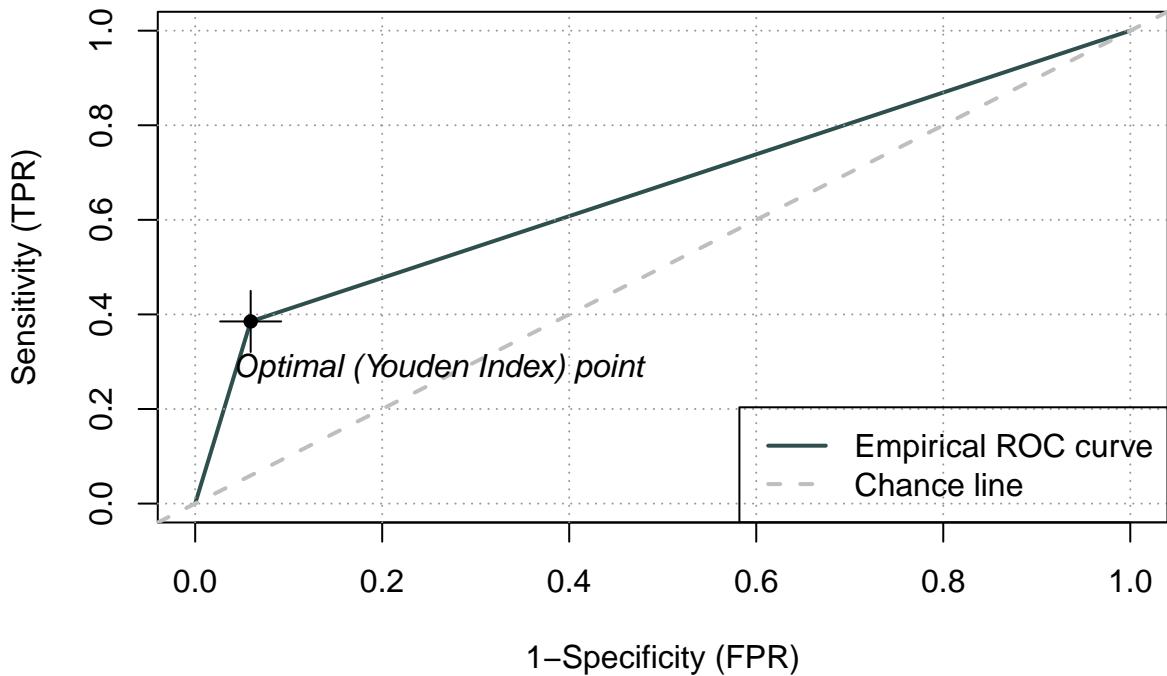
log_model5 <- glm(Y ~ Marital + Age + PayRec_Sep + PayRec_Aug + PayRec_Jul + PayRec_Jun +
                     BillAmt_Sep + PaidAmt_Sep, data = train.data, family = binomial)
#summary(log_model5)

## prediction on the test dataset (Logistics regression model 5)
train.log_model5_pred <- predict(log_model5, data = train.data, type = "response")
test.log_model5_pred <- predict(log_model5, data = test.data, type = "response")

## Finding the optimal cutoff
optcut <- optimalCutoff(train.data$Y, train.log_model5_pred,
                         optimiseFor = "misclasserror")
#print(optcut)

## Classifying based on the prediction model (train dataset)
train.binpred_glm5 <- ifelse(train.log_model5_pred < optcut ,0,1)


```

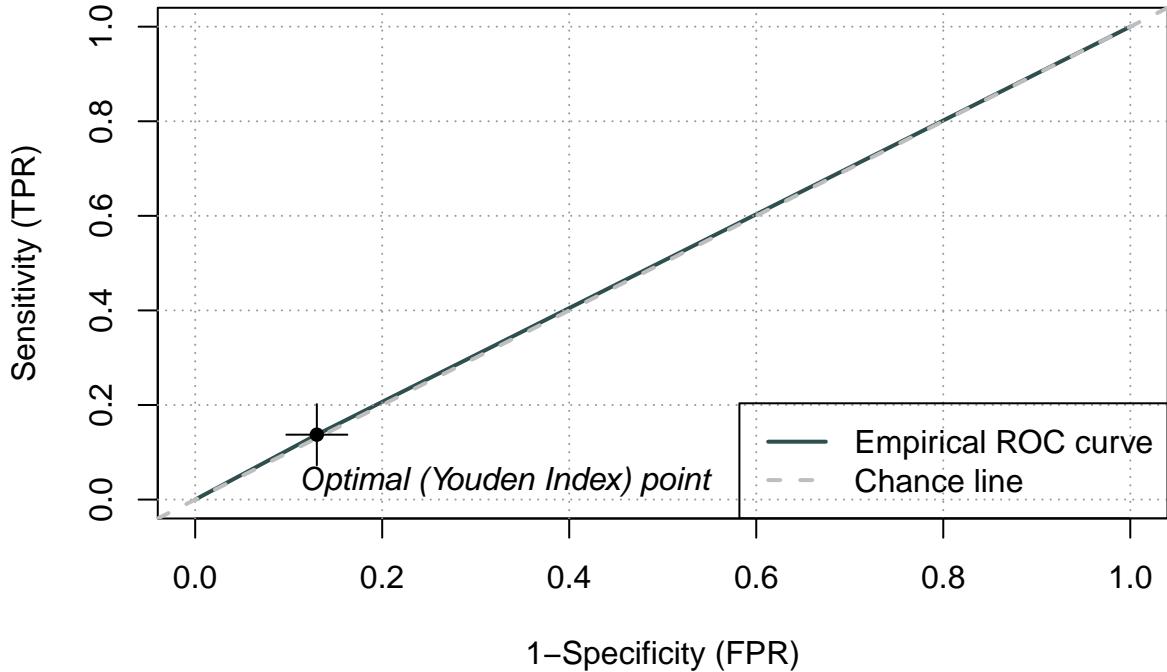


```

## Classifying based on the prediction model (train dataset)
train_table_log <- table(pred=train.binpred_glm5,actual=train.data$Y)
#sprintf("Train Accuracy: %f", mean(train.binpred_glm5 == train.data$Y))
#sprintf("Train Harmonic Mean: %f", harmonic_mean(train_table_log))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_train0 <- F1_Score(train.data$Y, train.binpred_glm5, positive = "0")
#sprintf("train f1score : %f", f1_train0)
f1_train1 <- F1_Score(train.data$Y, train.binpred_glm5, positive = "1")
#sprintf("train f1score : %f", f1_train1)

## for test set
test.roc5 <- rocit(test.binpred_glm5, test.data$Y)
#test.roc5
plot(test.roc5)

```



```
## Classifying based on the prediction model (test dataset)
test_table_log <- table(pred=test.binpred_glm5,actual=test.data$Y)
sprintf("Test Accuracy: %f", mean(test.binpred_glm5 == test.data$Y))
sprintf("Test Harmonic Mean: %f", harmonic_mean(test_table_log))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_test0 <- F1_Score(test.data$Y, test.binpred_glm5, positive = "0")
sprintf("test f1score : %f", f1_test0)
f1_test1 <- F1_Score(test.data$Y, test.binpred_glm5, positive = "1")
sprintf("test f1score : %f", f1_test1)
```

Evaluation for Logistic Model

This is a model which we build using the logistics regression classification method. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

Logistics regression Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Test Accuracy: 0.708

Train AUC: 0.6629

Test AUC: 0.5036

Test Harmonic Mean: 0.237

From the results, although the observed accuracy is quite high at 0.708, the harmonic mean is quite low at 0.237, which suggests a possible issue with the actual performance of our model. Furthermore, the AUC in test dataset is very low, which indicates that the model does not perform better by a significant margin relative to random selection.

SVM

```
## Running the SVM model
svm_model <- svm(as.factor(Y) ~ PayRec_Sep + Amt_credit + PaidAmt_Sep + PaidAmt_Aug +
                  PaidAmt_Jun + PaidAmt_Apr + Edu_Lvl + PayRec_Jul, data = train.data,
```

```

        type = "C-classification", kernel = "linear")
#svm_model
#svm_model$SV
#svm_model$index

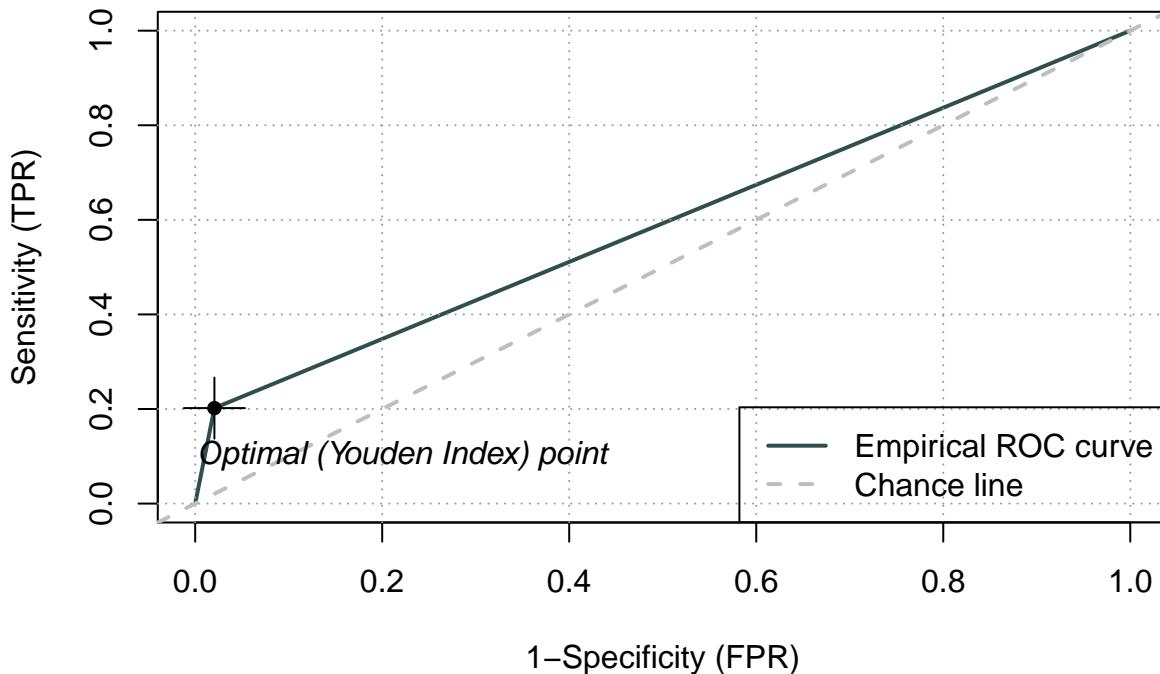
## testing and calculating the accuracy for the prediction model
results_train <- predict(svm_model, train.data)
results_test <- predict(svm_model, test.data)

## Classifying based on the prediction model (train dataset)
train_table_svm <- table(pred=results_train,actual=train.data$Y)
sprintf("Train Accuracy: %f", mean(results_train == train.data$Y))
sprintf("Train Harmonic Mean: %f", harmonic_mean(train_table_svm))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_train0 <- F1_Score(train.data$Y, results_train, positive = "0")
sprintf("Train f1score 0 : %f", f1_train0)
f1_train1 <- F1_Score(train.data$Y, results_train, positive = "1")
sprintf("Train f1score 1 : %f", f1_train1)

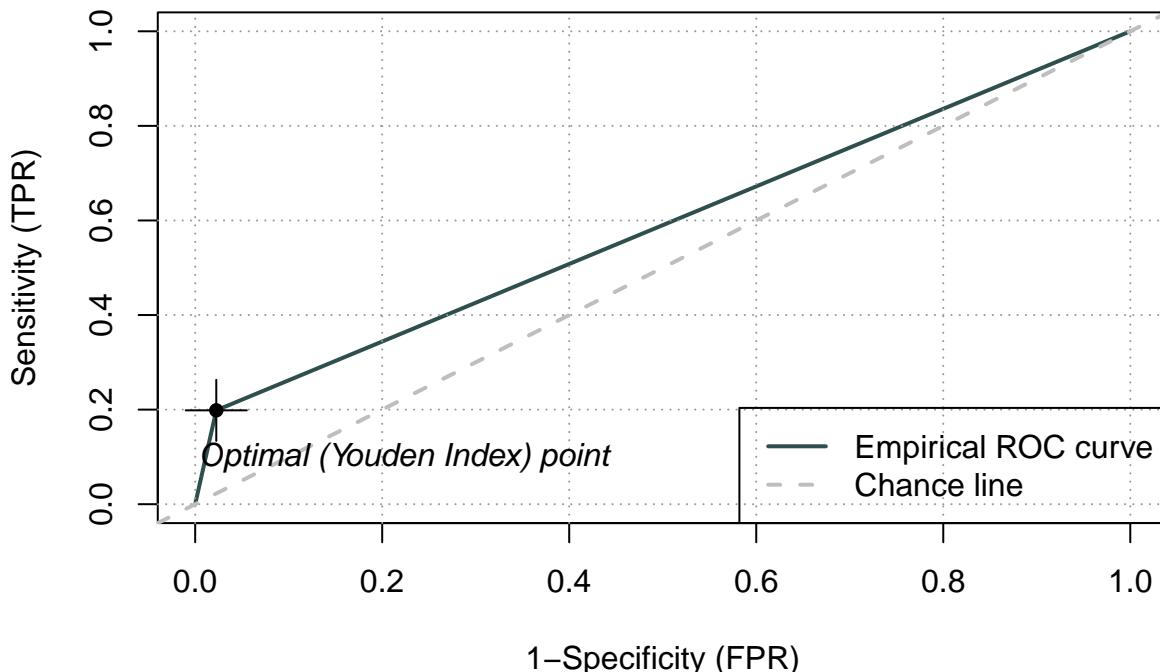
## Classifying based on the prediction model (test dataset)
test_table_svm <- table(pred=results_test,actual=test.data$Y)
sprintf("Test Accuracy: %f", mean(results_test == test.data$Y))
sprintf("Test Harmonic Mean: %f", harmonic_mean(test_table_svm))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_test0 <- F1_Score(train.data$Y, results_test, positive = "0")
sprintf("Test f1score 0 : %f", f1_train0)
f1_test1 <- F1_Score(train.data$Y, results_test, positive = "1")
sprintf("Test f1score 1 : %f", f1_train1)

## Using ROC curve for identifying accuracy
## for training set
train.roc <- rocit(as.numeric(results_train), train.data$Y)
#train.roc
plot(train.roc)

```



```
## for test set
test.roc <- rocit(as.numeric(results_test), test.data$Y)
#test.roc
plot(test.roc)
```



Evaluation for SVM model

This is a model which we build using the support vector machine classification method. C-classification, as well as radial kernel is used as our final settings for the model, based on our repeated testing on the test dataset. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

SVM regression Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Test Accuracy: 0.805

Train AUC: 0.591

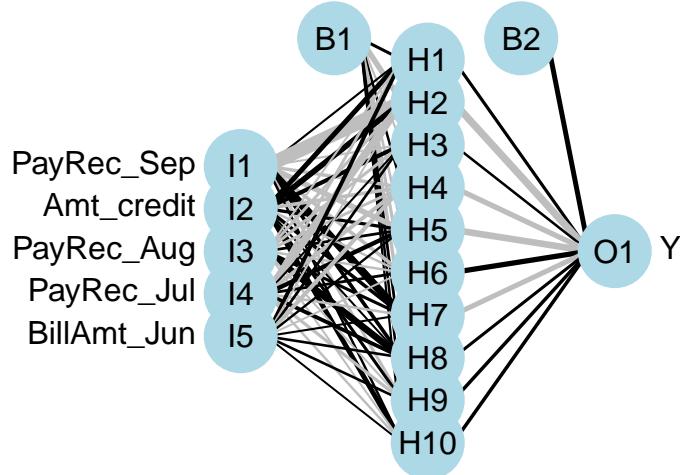
Test AUC: 0.588

Test Harmonic Mean: 0.330

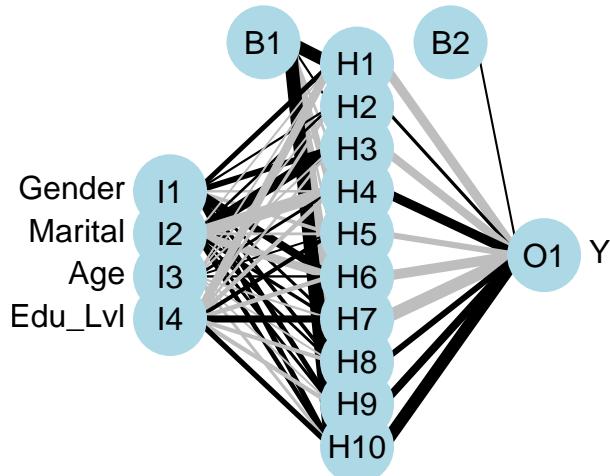
From the results above, although the observed accuracy is quite high as well at 0.805, the harmonic mean for our test dataset is again quite low at 0.330, which suggests a possible issue with the actual performance of our model. However, from a AUC review, the train and test datasets achieved better results than that of the logistics regression model.

Neural Network

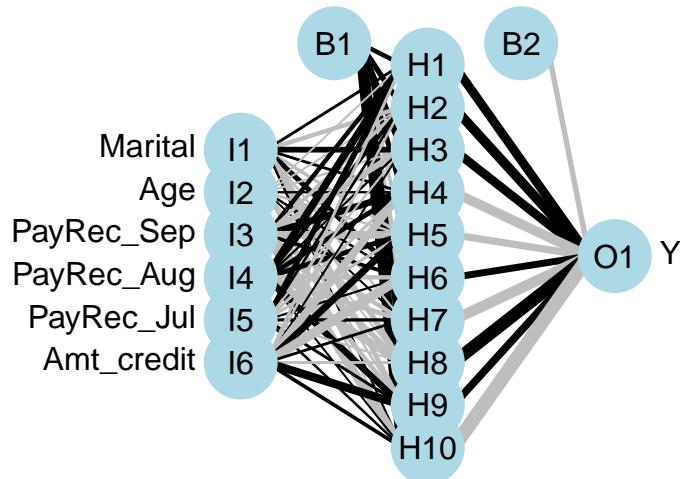
```
## Running the neural network model (best model)
set.seed(123)
nn_model <- nnet(Y ~ PayRec_Sep + Amt_credit + PayRec_Aug + PayRec_Jul + PayRec_Aug +
                  BillAmt_Jun, data = train.data, maxit = 1000, size = 10, decay = 0.08,
                  trace = FALSE)
plotnet(nn_model, pad_x = 0.5)
```



```
nn_model2 <- nnet(Y ~ Gender + Marital + Age + Edu_Lvl, data = train.data, maxit = 1000,
                     size = 10, decay = 0.08, trace = FALSE)
plotnet(nn_model2, pad_x = 0.5)
```



```
## running with forward selection
nn_model2 <- nnet(Y ~ Marital + Age + PayRec_Sep + PayRec_Aug + PayRec_Jul + Amt_credit,
                     data = train.data, maxit = 1000, size = 10, decay = 0.08, trace = FALSE)
plotnet(nn_model2, pad_x = 0.5)
```



```
## testing and calculating the accuracy for the prediction model (worst case)
results_train <- predict(nn_model2, train.data)
results_test <- predict(nn_model2, test.data)

## Classifying based on the prediction model (train dataset)
train_table_nn <- table(pred=results_train,actual=train.data$Y)
#train_table_nn
sprintf("Train Accuracy: %f", mean(results_train == train.data$Y))
sprintf("Train Harmonic Mean: %f", harmonic_mean(train_table_nn))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_train0 <- F1_Score(train.data$Y, results_train, positive = "0")
sprintf("train f1score 0 : %f", f1_train0)
f1_train1 <- F1_Score(train.data$Y, results_train, positive = "1")
sprintf("train f1score 1 : %f", f1_train1)

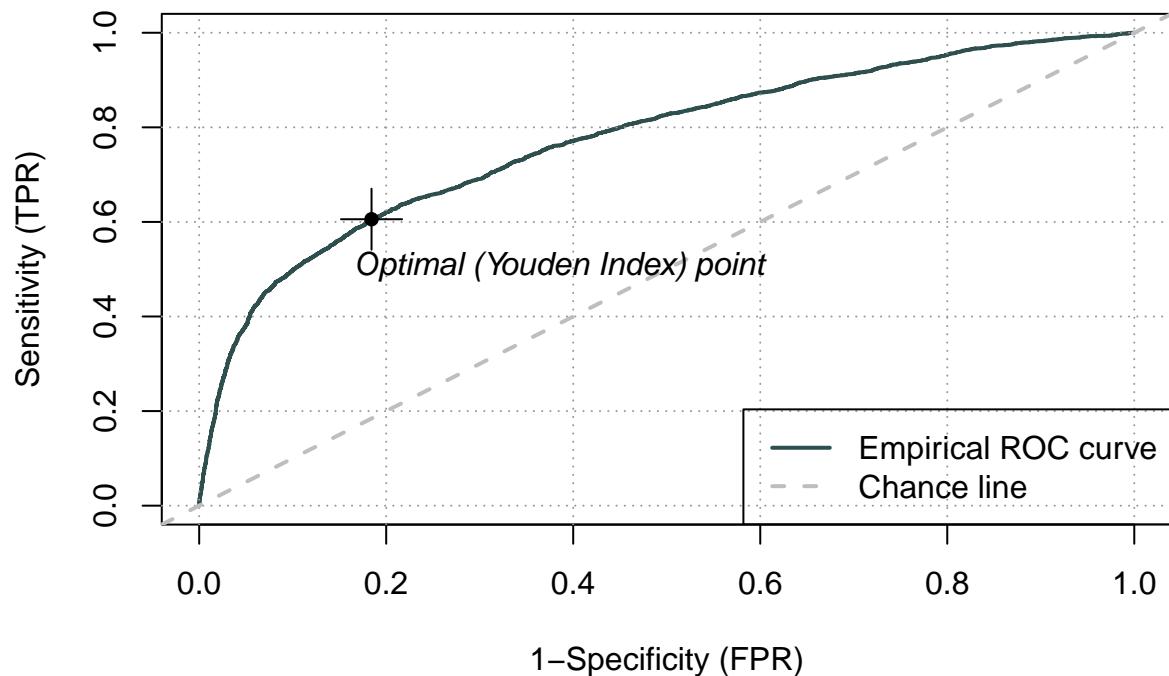
## Classifying based on the prediction model (test dataset)
test_table_nn <- table(pred=results_test,actual=test.data$Y)
```

```

 sprintf("Test Accuracy: %f", mean(results_test == test.data$Y))
 sprintf("Test Harmonic Mean: %f", harmonic_mean(test_table_nn))
 ## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
 ## as they bring losses to the bank)
 f1_test0 <- F1_Score(train.data$Y, results_test, positive = "0")
 sprintf("Test f1score 0 : %f", f1_train0)
 f1_test1 <- F1_Score(train.data$Y, results_test, positive = "1")
 sprintf("Test f1score 1 : %f", f1_train1)

## for training set
train.roc <- rocit(as.numeric(results_train), train.data$Y)
#train.roc
plot(train.roc)

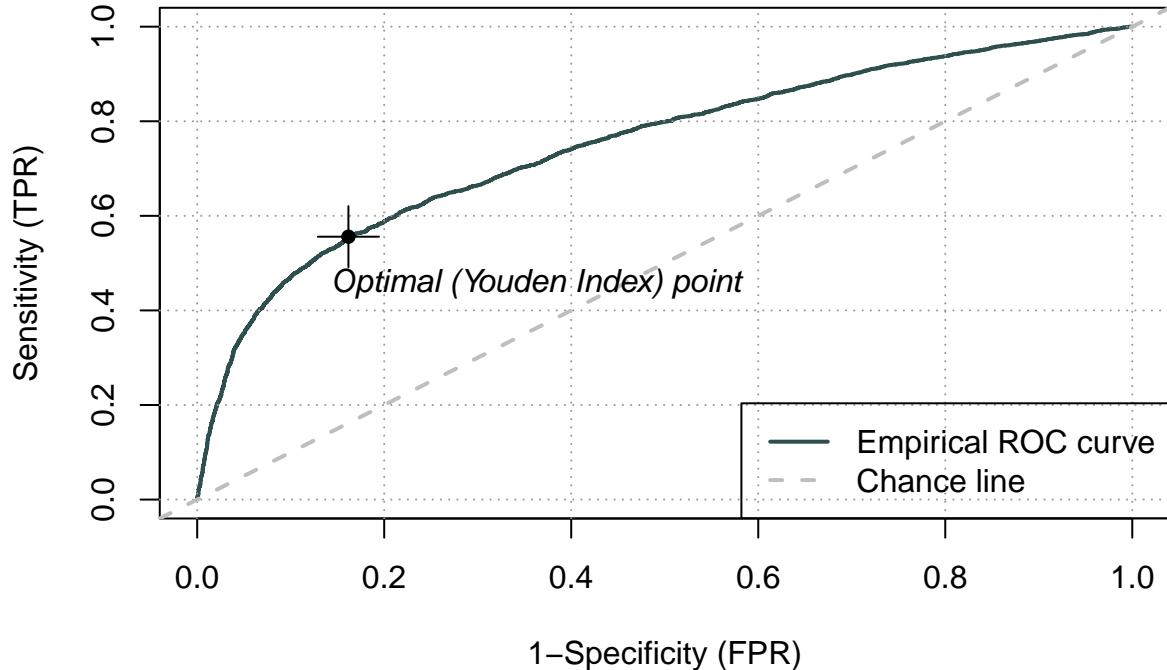
```



```

## for test set
test.roc <- rocit(as.numeric(results_test), test.data$Y)
#test.roc
plot(test.roc)

```



Evaluation for Neural Network model

This is a model which we build using the neural network classification method. A value of 1000 as max iterations to run, 10 hidden nodes and a decay factor of 0.08 is set for the model to prevent overfitting and achieve the best results. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

Neural Network Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Test Accuracy: 0.818

Train AUC: 0.653

Test AUC: 0.647

Test Harmonic Mean: 0.501

From the results above, as we can see, the test accuracy is again quite high albeit still remaining quite close to the null accuracy. We can see that the harmonic mean still remains low, and it suggests that the model may not perform as good as we think it is. For AUC, the AUC for test and train set are close, which suggests overfitting is not an issue. However, we still believe that we can create a better model make the predictions.

Random Forest Classifier

```
set.seed(2) # Setting seed

# need to convert Y to a factor for random forest to work
train_data <- train.data %>%
  mutate(Y = as.factor(as.character(Y)))

# need to convert Y to a factor for random forest to work
test_data <- test.data %>%
  mutate(Y = as.factor(as.character(Y)))

classifier_RF <- randomForest(Y ~ Marital + Age + Amt_credit + PayRec_Sep + PayRec_Aug)
```

```

+ PayRec_Jul, data=train_data)

#classifier_RF
```

In this case, mtry = 4 is the best mtry as it has least OOB error. Coincidentally, mtry = 4 was also used as default mtry.

```

# Rebuild model using best mtry value

tuned_classifier_RF <- randomForest(Y ~ Marital + Age + Amt_credit + PayRec_Sep
+ PayRec_Aug + PayRec_Jul, data=train_data,
mtry=4, importance=TRUE, ntree=500)

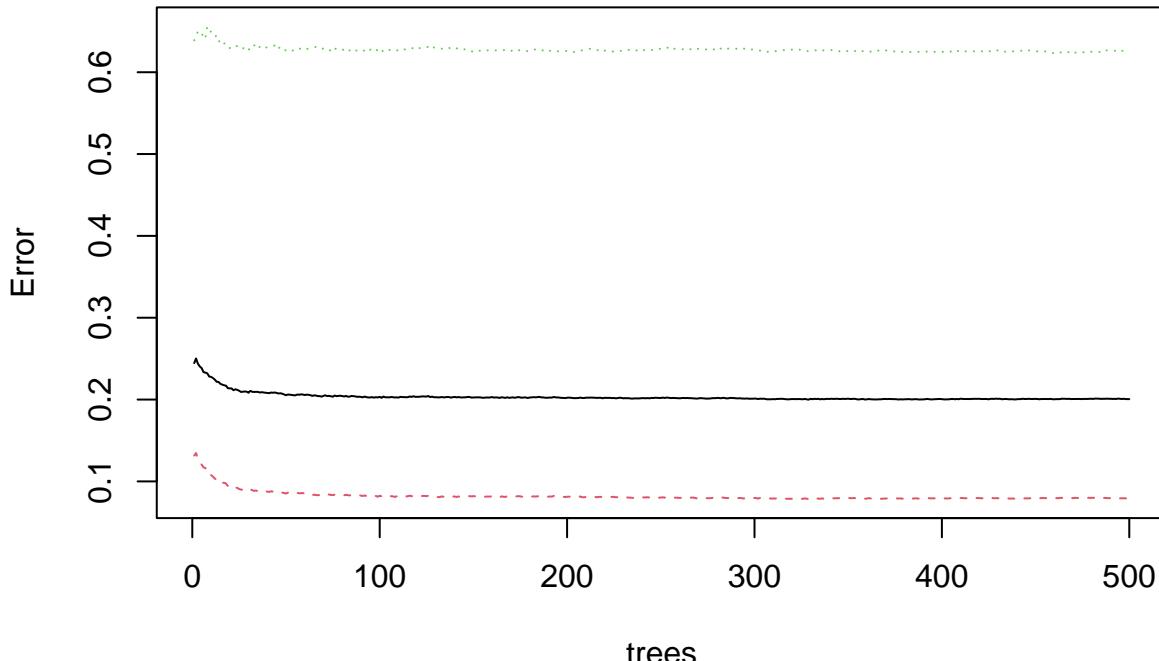
# Predicting the Test set results
y_pred = predict(tuned_classifier_RF, newdata = test_data[-24])

# Confusion Matrix
confusion_mtx = table(test_data[, 24], y_pred)
confusion_mtx

##      y_pred
##      0     1
##  0 10507   976
##  1 2065   1200

# Plotting model
plot(tuned_classifier_RF)
```

tuned_classifier_RF



```

# Importance plot
importance(tuned_classifier_RF)
```

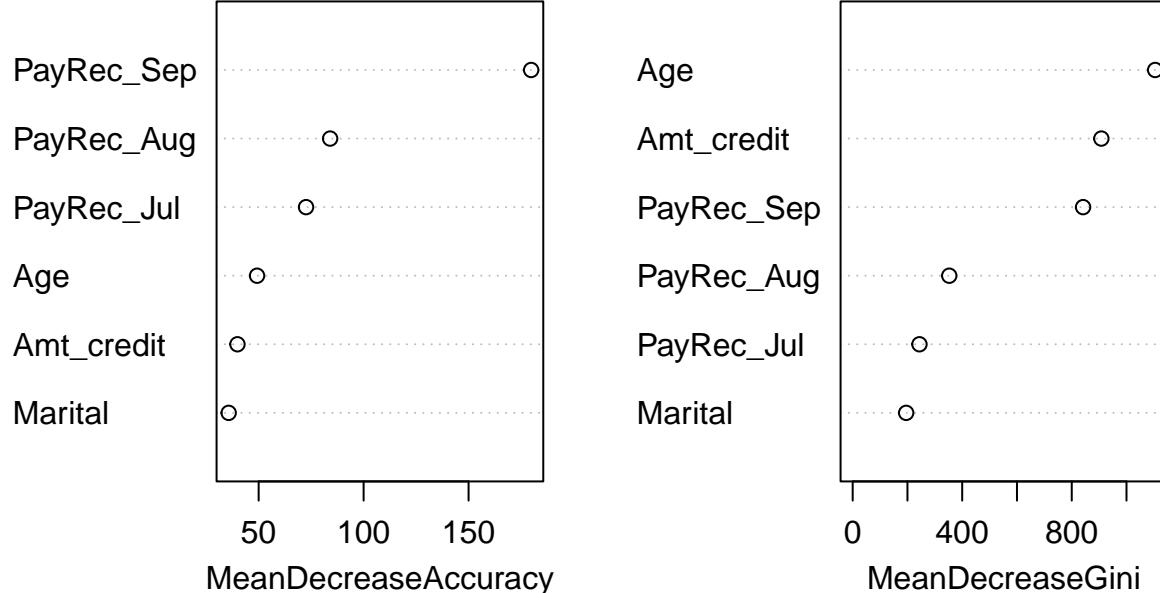
```

##          0      1 MeanDecreaseAccuracy MeanDecreaseGini
## Marital    40.89120 -3.973139      35.72078      195.3995
## Age        49.26334  4.949118      49.23292     1104.6684
## Amt_credit 33.85993 16.085824      39.90170      908.1640
## PayRec_Sep 131.88687 69.726807     179.82329     841.4174
## PayRec_Aug  81.05604 -7.127203      84.04086     352.7161
## PayRec_Jul  54.28743 37.671283      72.59057     243.8465

# Variable importance plot
varImpPlot(tuned_classifier_RF)

```

tuned_classifier_RF



```

# Random Forest Classifier AUC

## testing and calculating the accuracy for the prediction model
results_train_rf <- predict(tuned_classifier_RF, train_data)
results_test_rf <- predict(tuned_classifier_RF, test_data)

## Classifying based on the prediction model (train dataset)
train_table_rf <- table(pred=results_train_rf, actual=train_data$Y)
#train_table_rf
sprintf("Train Accuracy: %f", mean(results_train_rf == train_data$Y))
sprintf("Train Harmonic Mean: %f", harmonic_mean(train_table_rf))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_train0 <- F1_Score(train_data$Y, results_train_rf, positive = "0")
sprintf("train f1score 0 : %f", f1_train0)
f1_train1 <- F1_Score(train_data$Y, results_train_rf, positive = "1")
sprintf("train f1score 1 : %f", f1_train1)

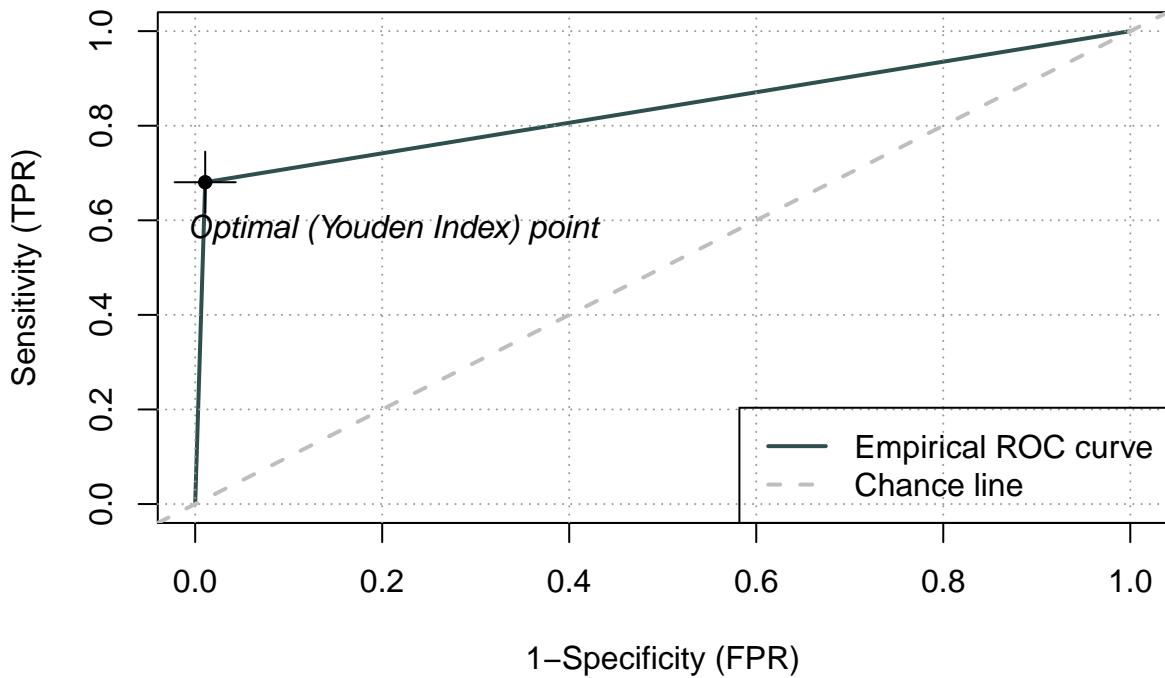
```

```

## Classifying based on the prediction model (test dataset)
test_table_rf <- table(pred=results_test_rf,actual=test_data$Y)
#test_table_rf
#sprintf("Test Accuracy: %f", mean(results_test_rf == test_data$Y))
#sprintf("Test Harmonic Mean: %f", harmonic_mean(test_table_rf))
## compute harmonic mean 0 - non default 1 default (focus on 1 (default),
## as they bring losses to the bank)
f1_test0 <- F1_Score(train_data$Y, results_test_rf, positive = "0")
#sprintf("test f1score 0 : %f", f1_train0)
f1_test1 <- F1_Score(train_data$Y, results_test_rf, positive = "1")
#sprintf("test f1score 1 : %f", f1_train1)

## for training set
train.roc_rf <- rocit(as.numeric(results_train_rf), train_data$Y)
#train.roc_rf
plot(train.roc_rf)

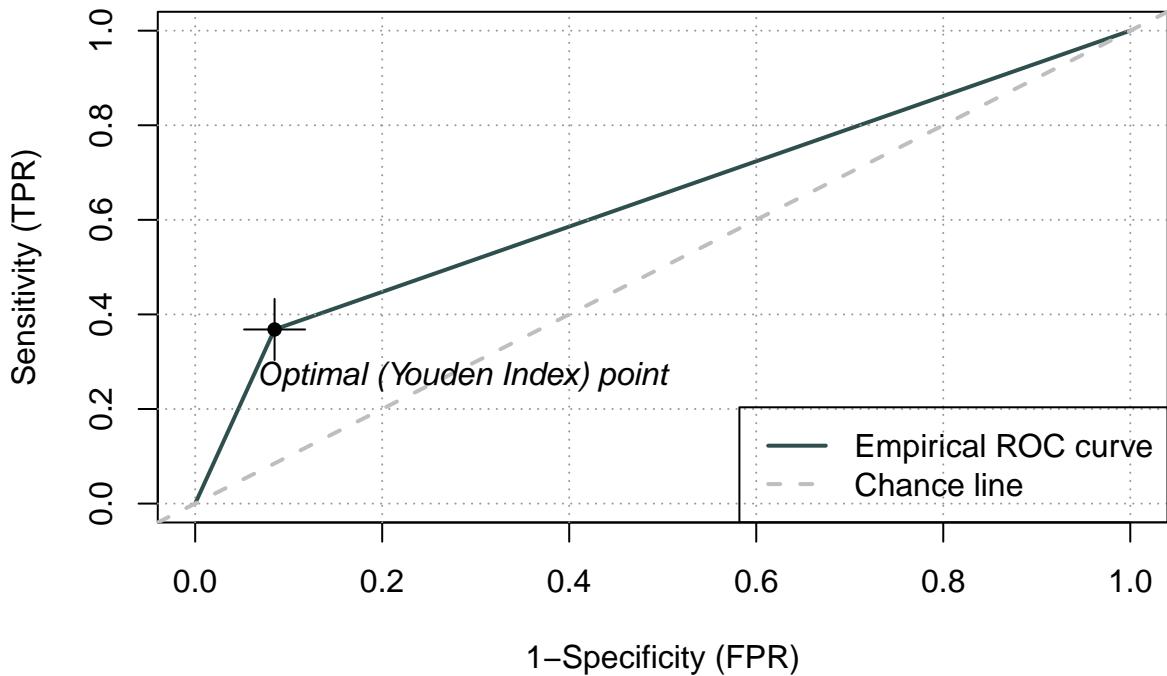
```



```

## for test set
test.roc_rf <- rocit(as.numeric(results_test_rf), test_data$Y)
#test.roc_rf
plot(test.roc_rf)

```



Evaluation for Random Forest model

This is a model which we build using the random forest decision tree classification method. From our repeated training, the number of nodes is set at 4, with TRUE set for gini importance and number of trees at 500. Using the variables that we selected on from earlier, which are Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables, our results for the model is as follows:

Seed 120

Random Forest with all variables:

Accuracy: 0.861

Train AUC: 1

Test AUC: 0.6241

Random Forest Chi Square with $Y \sim \text{Edu_Lvl} + \text{PayRec_May} + \text{PayRec_Jun} + \text{PayRec_Jul} + \text{PayRec_Aug}$:

Accuracy: 0.846

Train AUC: 0.5633

Test AUC: 0.5413

Random Forest first 5 of Boruta with $Y \sim \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul} + \text{BillAmt_Jun} + \text{PayRec_Jun}$:

Accuracy: 0.845

Train AUC: 0.753

Test AUC: 0.618

Seed 2

Random Forest Forward Backward with $Y \sim \text{Marital} + \text{Age} + \text{Amt_credit} + \text{PayRec_Sep} + \text{PayRec_Aug} + \text{PayRec_Jul}$:

Accuracy: 0.795

Train AUC: 0.832

Test AUC: 0.642

From the evaluation above, we can notice that our random forest model is consistent throughout different seed values, and thus makes it a reliable model that we can trust.

Overall, our final model that we would recommend would be a model with Marital, Age, Amt_credit, PayRec_Sep, PayRec_Aug and PayRec_Jul as the variables in our model, and the algorithm for the model that we are using is the random forest decision tree algorithm.

Final Review and Discussion

As we run the model using Marital + Age + Amt_credit + PayRec_Sep + PayRec_Aug + PayRec_Jul, we observed the train AUC & test AUC values closer to 1. A model which has AUC near to the 1 which means it has a good measure of separability. By analogy, the higher the AUC, the better the model is at forecasting clients who have tendencies to default on the bank.

Overall, from all the models we run as shown above, despite the increasing harmonic means and AUCs we achieved across models, we are still unable to achieve a good prediction score for all the models.

From our repeated model trainings, we realise that taking a huge number of variables in as predictors does not only makes the model more complex, but also may overfit to the test dataset. On the other hand, too few predictors may lead to a model that would not classify the data well. Therefore, based on the importance of each variable, we decided to come up with the variables as stated above as our best features to predict the model.

This is due to a few factors, namely the small number of predictors we selected for our models to achieve more efficient and quick predictions without wastage of resources, potential overfitting issues of using too much predictors as shown in our random forest run for a model with all the variables included, as well as some potential data integrity issues that arises from the unexplained values that are found in our data visualisation. Unbalanced nature of our dataset also suggests that there may be harder to predict due to the lower number of default clients as well.

However, there is no denial that our model is not perfect and there are ways to improve on it. For example, the lower score for both AUC, as well as harmonic means suggests that there are still ways to improve on the model. We believe that with better knowledge and understanding of the data, we will be able to create a better model, as compared to the satisfactory model we have now.

While reviewing our project, we discovered that the dataset is unbalanced which could result in bias in the model. This can be shown though the figures below.

```
# Convert LIMIT_BAL to factors
card_renamed$limit_bal_categorical <- ifelse(card_renamed$Amt_credit >= 0
                                              & card_renamed$Amt_credit < 10000, "1",
                                              ifelse(card_renamed$Amt_credit >= 10000
                                                    & card_renamed$Amt_credit < 20000, "2",
                                                    ifelse(card_renamed$Amt_credit >= 20000
                                                      & card_renamed$Amt_credit < 30000, "3",
                                                      ifelse(card_renamed$Amt_credit >= 30000
                                                        & card_renamed$Amt_credit < 40000, "4",
                                                        ifelse(card_renamed$Amt_credit >= 40000
                                                          & card_renamed$Amt_credit < 50000, "5", "6"))))

# Convert AGE to factors
card_renamed$age_categorical <- ifelse(card_renamed$Age >= 18
                                         & card_renamed$Age < 30, "1",
                                         ifelse(card_renamed$Age >= 30 & card_renamed$Age < 40, "2",
                                         ifelse(card_renamed$Age >= 40 & card_renamed$Age < 50, "3",
                                         ifelse(card_renamed$Age >= 50 & card_renamed$Age < 60, "4", "5"))))
```

Breakdown of Customers based on default status

```
default <- table(card_renamed$Y)
names(default) <- c("Not Default", "Default")
default

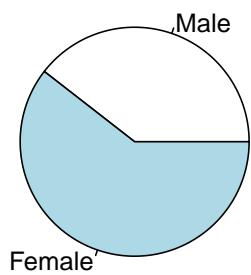
## Not Default      Default
##          22957       6539
```

Breakdown of Customers based on Gender and Limit Balance

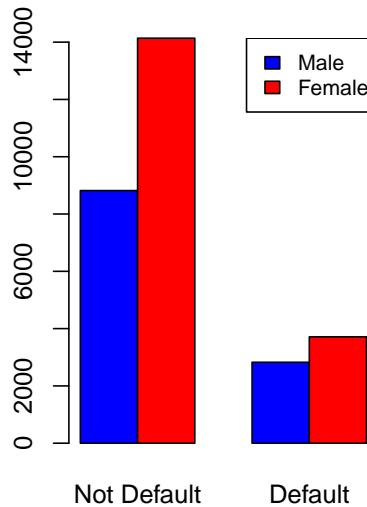
```
par(mfrow = c(1, 2))
# piechart of male to female clients
pie(table(card_renamed$Gender), c("Male", "Female"), main="Credit Card Clients By Gender",
    cex.main = 0.8)

# barplot showing defaulters vs non-defaulters based on gender
barplot(table(card_renamed$Gender, card_renamed$Y), beside = T, col = c("blue", "red"),
        names.arg = c("Not Default", "Default"),
        main = "Defaulters vs Non-defaulters based on Gender", cex.main = 0.8)
legend("topright", c("Male", "Female"), cex = 0.8, fill = c("blue", "red"))
```

Credit Card Clients By Gender



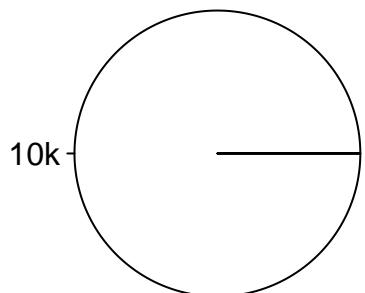
Defaulters vs Non-defaulters based on Gender



```
par(mfrow = c(1, 2))
# piechart of limit_balance
pie(table(card_renamed$limit_bal_categorical), c("10k", "20k", "30k", "40k", "50k", ">50k"),
    main="Credit Card Clients By Limit Balance", cex.main = 0.7)

# barplot showing defaulters vs non-defaulters based on age
barplot(table(card_renamed$limit_bal_categorical, card_renamed$Y), beside = T,
        col = rainbow(length(c("10k", "20k", "30k", "40k", "50k", ">50k"))),
        names.arg = c("Not Default", "Default"),
        main = "Defaulters vs Non-defaulters based on Limit Balance", cex.main = 0.7)
legend("topright", c("10k", "20k", "30k", "40k", "50k", ">50k"), cex = 0.7,
       fill = rainbow(length(c("10k", "20k", "30k", "40k", "50k", ">50k"))))
```

Credit Card Clients By Limit Balance



10k

Defaulters vs Non-defaulters based on Limit Balance

