# Predicting the Status of Credit

2023-04-02

Lu Zheng: Model Building, Model Selection

Yuxin Yao: Model Building, Model Selection

Zhiquan Cui 1005835857: Exploratory Data Analysis, Model Validation and Diagnostics

# Contents

# 1 Introduction

# 2 Exploratory Data Analysis

## 2.1 Load Data & Inspect Variables

```
# Read the data
data <- read.csv("Credit.csv")
# Check the number of observations and number of variables
n <- nrow(data)
m <- ncol(data)
n
```

```
## [1] 1000
```

```
m
```

```
## [1] 21
```

```
# Check the data
kable(head(data[, 1:8]), format = "latex", align=rep("c", 8), booktabs=TRUE)
```

| status | duration | credit_history | purpose | amount | savings | employment_duration | installment_rate |
|--------|----------|----------------|---------|--------|---------|---------------------|------------------|
| 1 | 18 | 4 | 2 | 1049 | 1 | 2 | 4 |
| 1 | 9 | 4 | 0 | 2799 | 1 | 3 | 2 |
| 2 | 12 | 2 | 9 | 841 | 2 | 4 | 2 |
| 1 | 12 | 4 | 0 | 2122 | 1 | 3 | 3 |
| 1 | 12 | 4 | 0 | 2171 | 1 | 3 | 4 |
| 1 | 10 | 4 | 0 | 2241 | 1 | 2 | 1 |

```
kable(head(data[, 9:14]), format = "latex", align=rep("c", 6), booktabs=TRUE)
```

| personal_status_sex | other_debtors | present_residence | property | age | other_installment_plans |
|---------------------|---------------|-------------------|----------|-----|-------------------------|
| 2 | 1 | 4 | 2 | 21 | 3 |
| 3 | 1 | 2 | 1 | 36 | 3 |
| 2 | 1 | 4 | 1 | 23 | 3 |
| 3 | 1 | 2 | 1 | 39 | 3 |
| 3 | 1 | 4 | 2 | 38 | 1 |
| 3 | 1 | 3 | 1 | 48 | 3 |

```
kable(head(data[, 15:21]), format = "latex", align=rep("c", 7), booktabs=TRUE)
```

| housing | number_credits | job | people_liable | telephone | foreign_worker | credit_risk |
|---------|----------------|-----|---------------|-----------|----------------|-------------|
| 1 | 1 | 3 | 2 | 1 | 2 | 1 |
| 1 | 2 | 3 | 1 | 1 | 2 | 1 |
| 1 | 1 | 2 | 2 | 1 | 2 | 1 |
| 1 | 2 | 2 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 | 1 | 1 | 1 |
| 1 | 2 | 2 | 1 | 1 | 1 | 1 |

```r
# Check invalid or missing values
anyNA(data)
```

```
## [1] FALSE
```

```r
# Check the data type of each column
sapply(data, class)
```

```
##              status            duration        credit_history
##           "integer"           "integer"           "integer"
##             purpose              amount             savings
##           "integer"           "integer"           "integer"
##  employment_duration     installment_rate    personal_status_sex
##           "integer"           "integer"           "integer"
##        other_debtors    present_residence            property
##           "integer"           "integer"           "integer"
##                 age other_installment_plans            housing
##           "integer"           "integer"           "integer"
##      number_credits                 job        people_liable
##           "integer"           "integer"           "integer"
##           telephone       foreign_worker         credit_risk
##           "integer"           "integer"           "integer"
```

As we can see from the above outputs, there is no NaN values so the data is clean. And all of the columns
are of type integer. Some of them are quantitative variable while some of them are qualitative variables.
Here is a summary of the variables:

- status: status of the debtor's checking account with the bank (categorical)
- duration: credit duration in months (quantitative)
- credit_history: history of compliance with previous or concurrent credit contracts (categorical)
- purpose: purpose for which the credit is needed (categorical)
- amount: credit amount in DM (quantitative; result of monotonic transformation; actual data and type
  of transformation unknown)
- savings: debtor's savings (categorical)
- employment_duration: duration of debtor's employment with current employer (ordinal; discretized
  quantitative)
- installment_rate: credit installments as a percentage of debtor's disposable income (ordinal; discretized
  quantitative)
- personal_status_sex: combined information on sex and marital status (categorical)
- other_debtors: is there another debtor or a guarantor for the credit? (categorial)
- present_residence: length of time (in years) the debtor lives in the present residence (ordinal; dis-
  cretized quantitative)
- property: the debtor's most valuable property (ordinal)
- age: age in years (quantitative)
- other_installment_plans: installment plans from providers other than the credit-giving bank (cate-
  gorical)
- housing: type of housing the debtor lives in (categorical)
- number_credits: number of credits including the current one the debtor has (or had) at the bank
  (ordinal; discretized quantitative)
- job: quality of debtor's job (ordinal)
- people_liable: number of persons who financially depend on the debtor (binary; discretized quantita-
  tive)

- telephone: is there a telephone landline registered on the debtor's name? (binary)
- foreign_ worker: is the debtor a foreign worker? (binary)
- credit_risk: has the credit contract been complied with (good) or not (bad)? (binary)

We can see that the **quantitative variables** include duration, amount and age, while **qualitative variables** include status, credit_history, purpose, savings, employment_duration, installment_rate, personal_status_sex, other_debtors, present_residence, property, other_installment_plans, housing, number_credits, job, people_liable, telephone, foreign_worker and credit_risk.

## 2.2  Univariate Data Analysis & Visualization

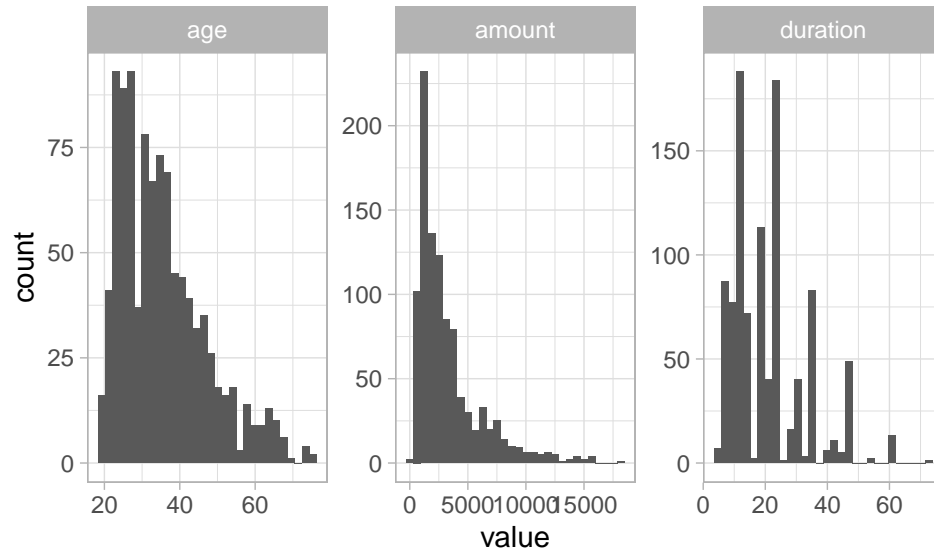### 2.2.1  Histogram of Quantitative Variables

First we will perform univariate analysis on each of the variables and look at their distribution. Here is the summary statistics:

```
quant_vars <- c("duration", "amount", "age")
qual_vars <- c("status", "credit_history", "purpose", "savings", "employment_duration",
               "installment_rate", "personal_status_sex", "other_debtors", "present_residence",
               "property", "other_installment_plans", "housing", "number_credits", "job",
               "people_liable", "telephone", "foreign_worker", "credit_risk")
summary(data[, quant_vars])
```

```
##     duration         amount            age
## Min.   : 4.0   Min.   :  250   Min.   :19.00
## 1st Qu.:12.0   1st Qu.: 1366   1st Qu.:27.00
## Median :18.0   Median : 2320   Median :33.00
## Mean   :20.9   Mean   : 3271   Mean   :35.54
## 3rd Qu.:24.0   3rd Qu.: 3972   3rd Qu.:42.00
## Max.   :72.0   Max.   :18424   Max.   :75.00
```

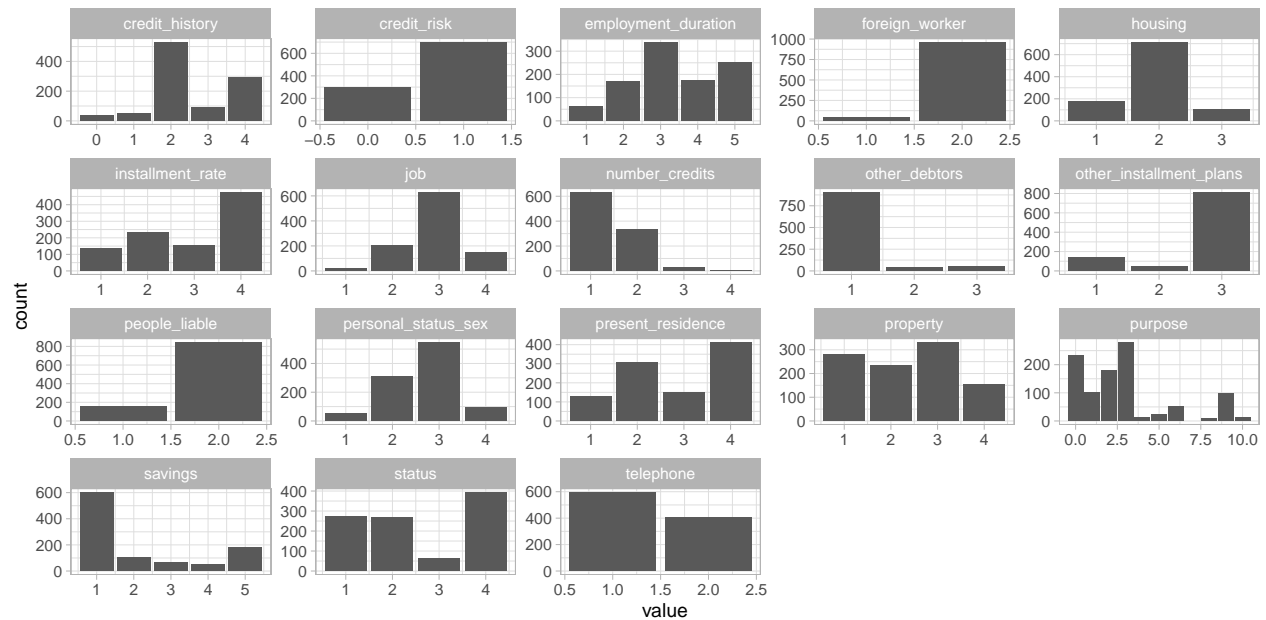Next, let us check the histograms of the quantitative variables:

```
data[, quant_vars] %>%
  gather() %>%
  ggplot(aes(value)) +
    facet_wrap(~ key, scales = "free") +
    geom_histogram() +
    theme_light()
```

### 2.2.2 Barplot of Qualitative Variables

Then, let us check the barplots of qualitative variables:

```
data[, qual_vars] %>%
  gather() %>%
  ggplot(aes(value)) +
    facet_wrap(~ key, scales = "free") +
    geom_bar() +
    theme_light()
```
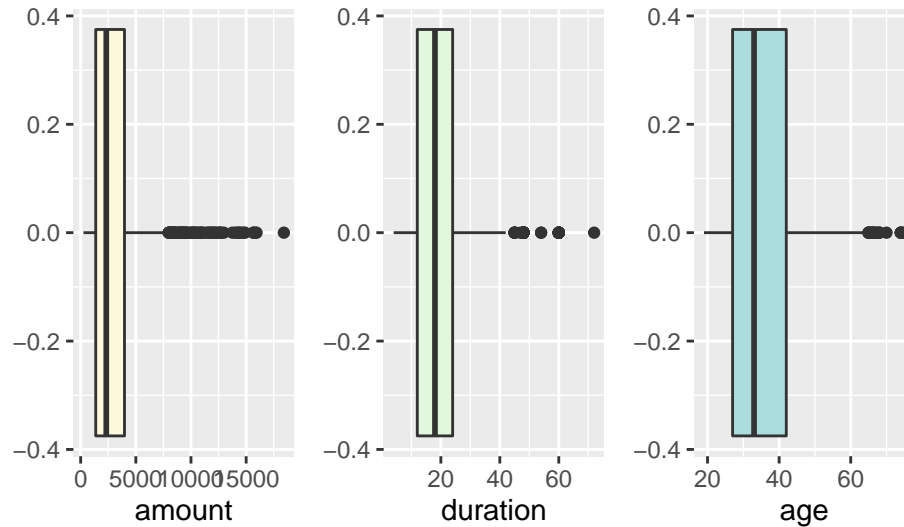


As we can see, the response variable credit risk is a binary variable while we have more than 2 predictors. This indicates that it is a good idea to use Multiple Logistic Regression as our model.

### 2.2.3 Boxplot of Quantitative Variables

After checking the histograms and barplots, we will check the boxplots of the quantitative variables. Here we will not check barplots for qualitative variables because it only makes sense to examine the median, first and third quartiles and maximum value for quantitative variables.

```
g1 <- ggplot(data, aes(x = amount)) + geom_boxplot(fill="#FEF8DD")
g2 <- ggplot(data, aes(x = duration)) + geom_boxplot(fill="#E1F8DC")
g3 <- ggplot(data, aes(x = age)) + geom_boxplot(fill="#ACDDDE")
g1 + g2 + g3
```



From the above box plots, we can see that there are a few outliers for the variable amount. If we look at the histogram of variable amount, we can see that it is a right skewed distribution with a long right tail, which results in these outliers.

### 2.2.4 Sample Odds of Binary Variables

For binary variables people_liable, telephone, foreign_worker and credit_risk, we can calculate and interpret the sample odds:

```
binary_var <- c("Statistics", "people_liable", "telephone", "foreign_worker", "credit_risk")
odds <- c("Sample Odds")
for (var in binary_var[2:5]) {
  if (var == "credit_risk") {y <- sum(data[, var] == 1)}
  else {y <- sum(data[, var] == 2)}
  n <- length(data[, var])
  odds <- append(odds, round(y / (n - y), 2))
}
kable(data.frame(t(odds)), col.names = binary_var, format = "latex") %>%
  kable_styling(position = "center", latex_options = "hold_position") %>% row_spec(0, bold = TRUE)
```

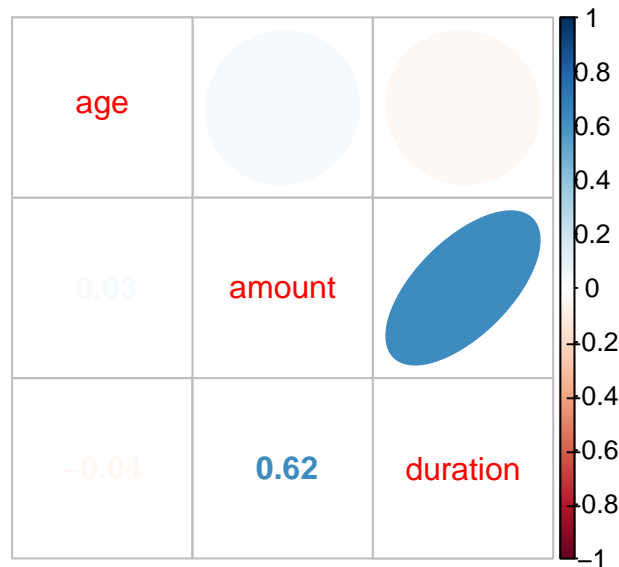| Statistics | people_liable | telephone | foreign_worker | credit_risk |
|---|---|---|---|---|
| Sample Odds | 5.45 | 0.68 | 26.03 | 2.33 |

7

Based on our sample, the estimated probability of a person to have good credit is 2.33 times as likely as having a bad credit. Similarly, the estimated probability of a person to have a telephone landline registered on his/her name is 0.68 times as likely as not having such a telephone landline.

## 2.3 Multivariate Data Analysis & Visualization

### 2.3.1 Quantitative Variable

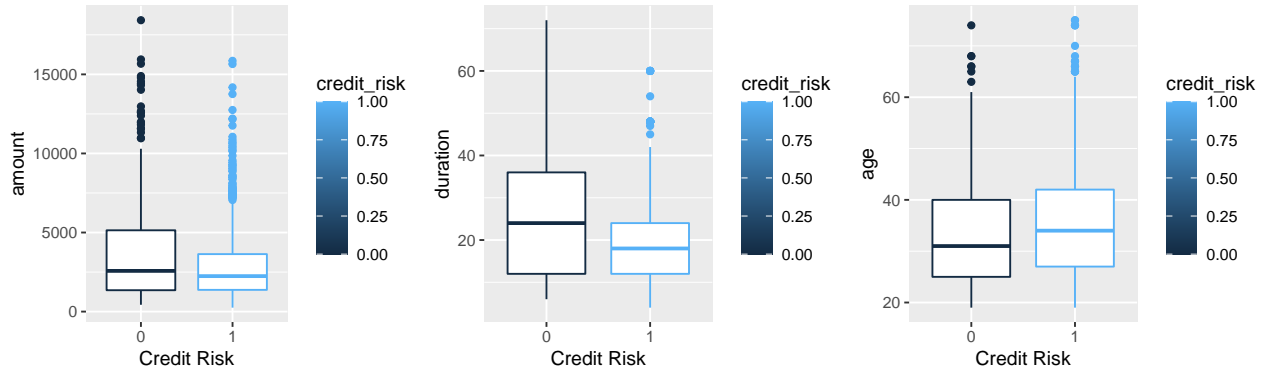First, let us look at the correlation plots of the quantitative variables.

```
corrplot.mixed(cor(data[quant_vars]), lower='number', upper='ellipse', order='AOE')
```

From the above correlation plot, we can see that the correlation coefficient between amount and duration is as high as 0.62, which indicates a strong positive correlation between the two variables. This also makes sense intuitively because the longer credit duration one has in months, he/she will have a higher chance to build up his/her credit and obtain a higher credit amount. Similarly, if one has a high credit amount, then he/she is more likely to have a long credit duration. In order to avoid multicollinearity, we will consider droping one of amount and duration in our model. However, before making a decision, we shall examine the side by side box plots.

```
g1 <- ggplot(data, aes(x=as.factor(credit_risk), y=amount, color=credit_risk)) +
    geom_boxplot() + xlab("Credit Risk")
g2 <- ggplot(data, aes(x=as.factor(credit_risk), y=duration, color=credit_risk)) +
    geom_boxplot() + xlab("Credit Risk")
g3 <- ggplot(data, aes(x=as.factor(credit_risk), y=age, color=credit_risk)) +
    geom_boxplot() + xlab("Credit Risk")

grid.arrange(g1, g2, g3, nrow=1)
```

From the above side by side box plots, we can see that for variables duration and age, there are significant differences on the box plots between two levels of credit risks. This indicates a significant association between credit risk and these two variables. However, we don't see a significant difference between two credit risk levels for variable amount.

Therefore, we will drop the variable amount.

### 2.3.2 Qualitative Variables

After examining the quantitative variables, we will now look at the qualitative variables. Since they are not continuous and numeric data, we should not use the same methodology as above. Instead, we will use Pearson's Chi-sq Test of Indepence and Cramer's V designed for qualitative variables to examine the data.

```
Pearson_chisq_test <- data.frame(matrix(0, ncol = length(qual_vars),
                                 nrow = length(qual_vars)), row.names = qual_vars)
colnames(Pearson_chisq_test) <- qual_vars
for (var in qual_vars) {
  for (var_2 in qual_vars) {
    test <- chisq.test(table(data[, var], data[, var_2]), simulate.p.value = TRUE)
    Pearson_chisq_test[var, var_2] <- test$p.value
  }
}
kable(Pearson_chisq_test[, 1:6], format = "latex", booktabs=TRUE) %>%
  kable_styling(font_size = 6, latex_options = "hold_position")
```

| | status | credit_history | purpose | savings | employment_duration | installment_rate |
|---|---|---|---|---|---|---|
| status | 0.0004998 | 0.0004998 | 0.0004998 | 0.0004998 | 0.0069965 | 0.4517741 |
| credit_history | 0.0004998 | 0.0004998 | 0.0004998 | 0.1924038 | 0.0009995 | 0.6996502 |
| purpose | 0.0004998 | 0.0004998 | 0.0004998 | 0.0569715 | 0.0124938 | 0.0009995 |
| savings | 0.0004998 | 0.1939030 | 0.0459770 | 0.0004998 | 0.0264868 | 0.3318341 |
| employment_duration | 0.0099950 | 0.0004998 | 0.0084958 | 0.0234883 | 0.0004998 | 0.0014993 |
| installment_rate | 0.4482759 | 0.6846577 | 0.0014993 | 0.3388306 | 0.0004998 | 0.0004998 |
| personal_status_sex | 0.1454273 | 0.0159920 | 0.0004998 | 0.4822589 | 0.0004998 | 0.0004998 |
| other_debtors | 0.0024988 | 0.0499750 | 0.0004998 | 0.0179910 | 0.0849575 | 0.9800100 |
| present_residence | 0.0019990 | 0.1034483 | 0.0044978 | 0.1219390 | 0.0004998 | 0.4277861 |
| property | 0.0434783 | 0.0884558 | 0.0004998 | 0.0979510 | 0.0004998 | 0.4442779 |
| other_installment_plans | 0.2878561 | 0.0004998 | 0.0029985 | 0.9990005 | 0.2898551 | 0.6331834 |
| housing | 0.0019990 | 0.0174913 | 0.0004998 | 0.8125937 | 0.0004998 | 0.1039480 |
| number_credits | 0.0414793 | 0.0004998 | 0.2953523 | 0.1014493 | 0.0004998 | 0.4387806 |
| job | 0.0509745 | 0.3928036 | 0.0004998 | 0.3553223 | 0.0004998 | 0.0674663 |
| people_liable | 0.1159420 | 0.0589705 | 0.0029985 | 0.8875562 | 0.0389805 | 0.1139430 |
| telephone | 0.0799600 | 0.2628686 | 0.0004998 | 0.0604698 | 0.0004998 | 0.4532734 |
| foreign_worker | 0.1104448 | 0.3958021 | 0.0054973 | 0.9220390 | 0.4762619 | 0.0034983 |
| credit_risk | 0.0004998 | 0.0004998 | 0.0004998 | 0.0004998 | 0.0014993 | 0.1289355 |

```
kable(Pearson_chisq_test[, 7:12], format = "latex", booktabs=TRUE) %>%
  kable_styling(font_size = 6, latex_options = "hold_position")
```

|  | personal_status_sex | other_debtors | present_residence | property | other_installment_plans | housing |
|---|---|---|---|---|---|---|
| status | 0.1369315 | 0.0019990 | 0.0009995 | 0.0399800 | 0.2733633 | 0.0054973 |
| credit_history | 0.0164918 | 0.0569715 | 0.1114443 | 0.0844578 | 0.0004998 | 0.0154923 |
| purpose | 0.0004998 | 0.0009995 | 0.0024988 | 0.0004998 | 0.0039980 | 0.0004998 |
| savings | 0.4822589 | 0.0204898 | 0.1244378 | 0.0914543 | 0.9990005 | 0.8165917 |
| employment_duration | 0.0004998 | 0.0874563 | 0.0004998 | 0.0004998 | 0.2793603 | 0.0004998 |
| installment_rate | 0.0004998 | 0.9795102 | 0.4372814 | 0.4527736 | 0.6236882 | 0.1214393 |
| personal_status_sex | 0.0004998 | 0.6076962 | 0.0004998 | 0.0004998 | 0.4597701 | 0.0004998 |
| other_debtors | 0.6156922 | 0.0004998 | 0.6491754 | 0.0004998 | 0.2293853 | 0.1419290 |
| present_residence | 0.0004998 | 0.6501749 | 0.0004998 | 0.0004998 | 0.7566217 | 0.0004998 |
| property | 0.0004998 | 0.0004998 | 0.0004998 | 0.0004998 | 0.0039980 | 0.0004998 |
| other_installment_plans | 0.4412794 | 0.2233883 | 0.7316342 | 0.0039980 | 0.0004998 | 0.0024988 |
| housing | 0.0004998 | 0.1384308 | 0.0004998 | 0.0004998 | 0.0049975 | 0.0004998 |
| number_credits | 0.0539730 | 0.8370815 | 0.0069965 | 0.1159420 | 0.0219890 | 0.0044978 |
| job | 0.0374813 | 0.0474763 | 0.8985507 | 0.0004998 | 0.0344828 | 0.0004998 |
| people_liable | 0.0004998 | 0.3398301 | 0.2773613 | 0.0344828 | 0.0429785 | 0.0019990 |
| telephone | 0.0504748 | 0.0599700 | 0.0239880 | 0.0004998 | 0.2738631 | 0.0009995 |
| foreign_worker | 0.0789605 | 0.0029985 | 0.3823088 | 0.0004998 | 0.8570715 | 0.0254873 |
| credit_risk | 0.0224888 | 0.0394803 | 0.8740630 | 0.0004998 | 0.0004998 | 0.0004998 |

```
kable(Pearson_chisq_test[, 13:17], format = "latex", booktabs=TRUE) %>%
  kable_styling(font_size = 6, latex_options = "hold_position")
```

|  | number_credits | job | people_liable | telephone | foreign_worker |
|---|---|---|---|---|---|
| status | 0.0404798 | 0.0484758 | 0.1104448 | 0.0809595 | 0.1149425 |
| credit_history | 0.0004998 | 0.3958021 | 0.0419790 | 0.2918541 | 0.4017991 |
| purpose | 0.2943528 | 0.0004998 | 0.0014993 | 0.0004998 | 0.0054973 |
| savings | 0.1069465 | 0.3503248 | 0.8910545 | 0.0694653 | 0.9130435 |
| employment_duration | 0.0009995 | 0.0004998 | 0.0519740 | 0.0004998 | 0.4797601 |
| installment_rate | 0.4357821 | 0.0629685 | 0.0974513 | 0.4737631 | 0.0034983 |
| personal_status_sex | 0.0449775 | 0.0399800 | 0.0004998 | 0.0379810 | 0.0759620 |
| other_debtors | 0.8530735 | 0.0534733 | 0.3278361 | 0.0514743 | 0.0024988 |
| present_residence | 0.0049975 | 0.9025487 | 0.2828586 | 0.0229885 | 0.3563218 |
| property | 0.1309345 | 0.0004998 | 0.0294853 | 0.0004998 | 0.0004998 |
| other_installment_plans | 0.0159920 | 0.0264868 | 0.0454773 | 0.2643678 | 0.8595702 |
| housing | 0.0074963 | 0.0004998 | 0.0009995 | 0.0014993 | 0.0309845 |
| number_credits | 0.0004998 | 0.0019990 | 0.0049975 | 0.0599700 | 0.9150425 |
| job | 0.0024988 | 0.0004998 | 0.0009995 | 0.0004998 | 0.0204898 |
| people_liable | 0.0039980 | 0.0004998 | 0.0004998 | 0.6681659 | 0.0179910 |
| telephone | 0.0669665 | 0.0004998 | 0.6566717 | 0.0004998 | 0.0289855 |
| foreign_worker | 0.9205397 | 0.0244878 | 0.0189905 | 0.0269865 | 0.0004998 |
| credit_risk | 0.4392804 | 0.6041979 | 1.0000000 | 0.2713643 | 0.0074963 |

Based on the above table, we conclude that the following predictors are dependent to most of the predictors with $\alpha = 0.05$ according to Pearson's Chi-sq Test of Independence, and we consider dropping these predictors:

- job
- credit_history
- purpose
- employment_duration
- housing
- people_liable

Also, we can see that the following predictors have very weak association with the response variable:

- installment_rate

- personal_status_sex
- other_debtors
- present_residence
- number_credits
- job
- people_liable
- telephone
- foreign_worker

To summarize, the variables we will use in model building are:

- status
- duration
- savings
- property
- age
- other_installment_plans

# 3  Model Building and Model Selection

## 3.1  Load the data

```r
data.credit = read.csv("Credit.csv")
# Transform categorical variables
data.credit$credit_risk = as.factor(data.credit$credit_risk)
data.credit$status  = as.factor(data.credit$status)
data.credit$savings = as.factor(data.credit$savings)
data.credit$property = as.ordered(data.credit$property)
data.credit$other_installment_plans = as.factor(data.credit$other_installment_plans)
```

## 3.2  Split the data into training set and testing set

```r
set.seed(1006742107)

n = nrow(data.credit)
index = sample(n, round(0.75 * n), replace = FALSE)
traindata = data.credit[index, ]
testdata = data.credit[-index, ]
```

## 3.3  Main effect model

### 3.3.1  Training model

```r
step(glm(credit_risk ~ 1, family = binomial, data = traindata), scope =
      ~status + duration + savings + property + age +
      other_installment_plans, direction = "forward", test = "Chisq")
```

### 3.3.1.1 Forward method

```
## Start:  AIC=921.66
## credit_risk ~ 1
##
##                            Df Deviance    AIC     LRT  Pr(>Chi)
## + status                    3   828.43 836.43 91.225 < 2.2e-16 ***
## + duration                  1   887.29 891.29 32.368 1.276e-08 ***
## + savings                   4   892.66 902.66 27.003 1.985e-05 ***
## + property                  3   903.76 911.76 15.899  0.001189 **
## + age                       1   909.12 913.12 10.537  0.001170 **
## + other_installment_plans   2   909.77 915.77  9.889  0.007123 **
## <none>                          919.66 921.66
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=836.43
## credit_risk ~ status
##
##                            Df Deviance    AIC     LRT  Pr(>Chi)
## + duration                  1   801.74 811.74 26.6915 2.387e-07 ***
## + property                  3   811.90 825.90 16.5377 0.0008796 ***
## + savings                   4   814.60 830.60 13.8362 0.0078365 **
## + other_installment_plans   2   820.32 832.32  8.1125 0.0173134 *
## + age                       1   823.40 833.40  5.0337 0.0248586 *
## <none>                          828.43 836.43
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=811.74
## credit_risk ~ status + duration
##
##                            Df Deviance    AIC     LRT Pr(>Chi)
## + savings                   4   787.09 805.09 14.6529 0.005478 **
## + other_installment_plans   2   794.74 808.74  7.0039 0.030138 *
## + age                       1   797.00 809.00  4.7411 0.029450 *
## + property                  3   794.81 810.81  6.9303 0.074154 .
## <none>                          801.74 811.74
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=805.09
## credit_risk ~ status + duration + savings
##
##                            Df Deviance    AIC    LRT Pr(>Chi)
## + other_installment_plans   2   779.60 801.60 7.4878  0.02366 *
## + age                       1   783.28 803.28 3.8072  0.05103 .
## + property                  3   780.56 804.56 6.5277  0.08858 .
## <none>                          787.09 805.09
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=801.6
## credit_risk ~ status + duration + savings + other_installment_plans
```

```
##
##            Df Deviance    AIC     LRT Pr(>Chi)
## + age       1    775.22 799.22 4.3864  0.03623 *
## <none>           779.60 801.60
## + property  3    774.53 802.53 5.0741  0.16645
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=799.22
## credit_risk ~ status + duration + savings + other_installment_plans +
##     age
##
##            Df Deviance    AIC     LRT Pr(>Chi)
## + property  3    767.65 797.65 7.5621  0.05598 .
## <none>           775.22 799.22
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=797.65
## credit_risk ~ status + duration + savings + other_installment_plans +
##     age + property


##
## Call:  glm(formula = credit_risk ~ status + duration + savings + other_installment_plans +
##     age + property, family = binomial, data = traindata)
##
## Coefficients:
##             (Intercept)                    status2                    status3
##                -0.72146                    0.45381                    0.85533
##                 status4                   duration                   savings2
##                 1.75141                   -0.03046                    0.25614
##                 savings3                   savings4                   savings5
##                 0.13951                    1.49757                    0.72922
## other_installment_plans2  other_installment_plans3                        age
##                 0.19640                    0.59292                    0.02236
##              property.L                 property.Q                 property.C
##                -0.58463                   -0.16008                   -0.06641
##
## Degrees of Freedom: 749 Total (i.e. Null);   735 Residual
## Null Deviance:        919.7
## Residual Deviance: 767.7     AIC: 797.7
```

```
step(glm(credit_risk ~status + duration + savings + property + age +
         other_installment_plans, family = binomial, data = traindata), test = "Chisq")
```

### 3.3.1.2   Backward method

```
## Start:  AIC=797.65
## credit_risk ~ status + duration + savings + property + age +
##     other_installment_plans
##
```

```
##                              Df Deviance    AIC    LRT  Pr(>Chi)
## <none>                          767.65 797.65
## - property                   3  775.22 799.22  7.562  0.055984 .
## - other_installment_plans    2  774.23 800.23  6.577  0.037312 *
## - age                        1  774.53 802.53  6.874  0.008744 **
## - savings                    4  781.43 803.43 13.779  0.008036 **
## - duration                   1  784.38 812.38 16.727 4.316e-05 ***
## - status                     3  832.72 856.72 65.065 4.857e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


##
## Call:  glm(formula = credit_risk ~ status + duration + savings + property +
##      age + other_installment_plans, family = binomial, data = traindata)
##
## Coefficients:
##              (Intercept)                    status2                    status3
##                 -0.72146                    0.45381                    0.85533
##                  status4                   duration                   savings2
##                  1.75141                   -0.03046                    0.25614
##                  savings3                   savings4                   savings5
##                  0.13951                    1.49757                    0.72922
##                property.L                 property.Q                 property.C
##                 -0.58463                   -0.16008                   -0.06641
##                      age  other_installment_plans2  other_installment_plans3
##                  0.02236                    0.19640                    0.59292
##
## Degrees of Freedom: 749 Total (i.e. Null);  735 Residual
## Null Deviance:       919.7
## Residual Deviance: 767.7     AIC: 797.7
```

From above coding, we could find that both forward selection and backward elimination choose the model:
glm(credit_risk ~status + duration + savings + property + age + other_installment_plans, family = binomial, data = traindata)

$$logit(\hat{\pi}) = -0.72+0.45{\cdot}S_1+0.86{\cdot}S_2+1.75{\cdot}S_3-0.03{\cdot}D+0.26{\cdot}SV_1+0.14{\cdot}SV_2+1.50SV_3+0.73SV_4-0.58{\cdot}P_L-0.16{\cdot}P_Q$$

$$-0.07 \cdot P_C + 0.02 \cdot A + 0.20 \cdot O_1 + 0.59 \cdot O_2$$

where * $S_i$'s are dummy variables for status * D is duration * $SV$'s are dummy variables for savings * $P_i$'s are dummy variables for property * A is age * $O_i$'s are dummy variables for other_installment_plans
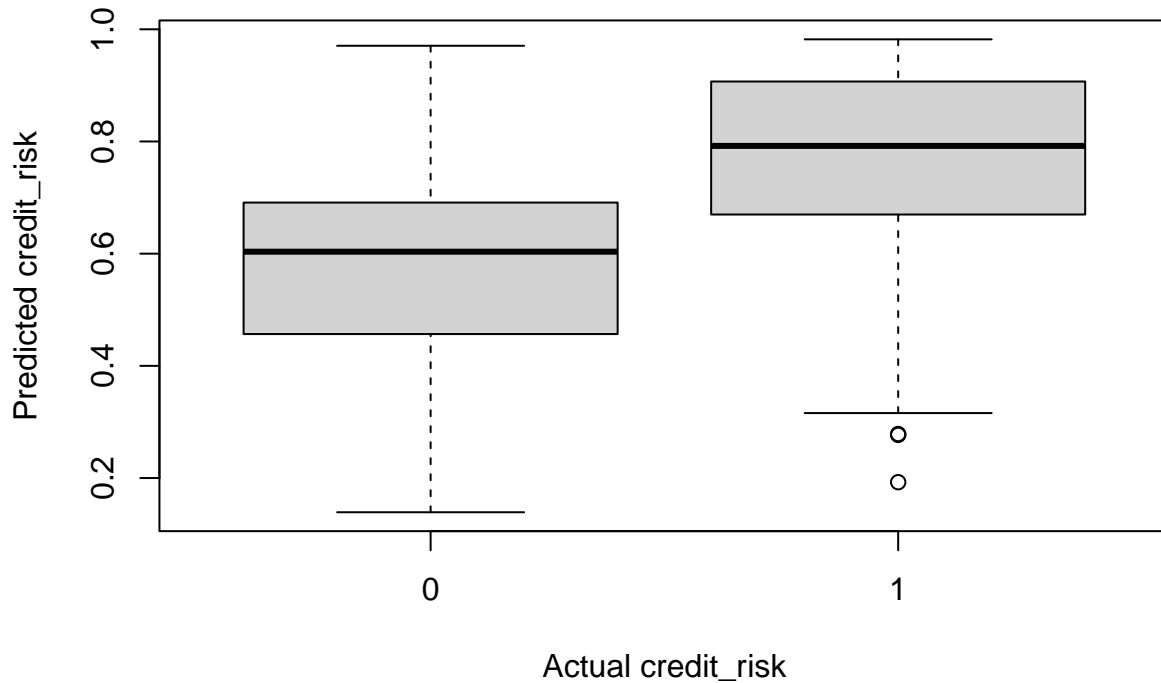
```
bestmodel.1 = glm(credit_risk ~status + duration + savings + property + age + other_installment_plans,
summary(bestmodel.1)
```

```
##
## Call:
## glm(formula = credit_risk ~ status + duration + savings + property +
##      age + other_installment_plans, family = binomial, data = traindata)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -2.7517   -0.9595    0.4902    0.8265    1.8215
##
## Coefficients:
##                           Estimate Std. Error z value Pr(>|z|)
## (Intercept)              -0.721459   0.436812  -1.652  0.09861 .
## status2                   0.453812   0.222048   2.044  0.04098 *
## status3                   0.855325   0.366951   2.331  0.01976 *
## status4                   1.751411   0.235658   7.432 1.07e-13 ***
## duration                 -0.030462   0.007502  -4.061 4.89e-05 ***
## savings2                  0.256136   0.302805   0.846  0.39762
## savings3                  0.139505   0.392170   0.356  0.72205
## savings4                  1.497575   0.657083   2.279  0.02266 *
## savings5                  0.729215   0.264897   2.753  0.00591 **
## property.L               -0.584630   0.213990  -2.732  0.00629 **
## property.Q               -0.160077   0.189526  -0.845  0.39832
## property.C               -0.066411   0.173897  -0.382  0.70254
## age                       0.022359   0.008708   2.568  0.01024 *
## other_installment_plans2  0.196405   0.440276   0.446  0.65553
## other_installment_plans3  0.592919   0.240369   2.467  0.01364 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 919.66  on 749  degrees of freedom
## Residual deviance: 767.65  on 735  degrees of freedom
## AIC: 797.65
##
## Number of Fisher Scoring iterations: 5
```

### 3.3.2 Testing model

```
pred.1 = predict(bestmodel.1, newdata = testdata)
plot(testdata$credit_risk, inv.logit(pred.1), xlab = "Actual credit_risk", ylab = "Predicted credit_risk
```

From the plot, we find that the main effect model can describe the actual data fairly well.

## 3.4 Interaction model

### 3.4.1 Training model

```
bestmodel.3 <- step(glm(credit_risk ~ 1, family = binomial, data = traindata), scope = ~status * durati
```

#### 3.4.1.1 Forward method

```
## Start:  AIC=921.66
## credit_risk ~ 1
##
##                          Df Deviance    AIC    LRT  Pr(>Chi)
## + status                  3   828.43 836.43 91.225 < 2.2e-16 ***
## + duration                1   887.29 891.29 32.368 1.276e-08 ***
## + savings                 4   892.66 902.66 27.003 1.985e-05 ***
## + property                3   903.76 911.76 15.899  0.001189 **
## + age                     1   909.12 913.12 10.537  0.001170 **
## + other_installment_plans 2   909.77 915.77  9.889  0.007123 **
## <none>                        919.66 921.66
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=836.43
## credit_risk ~ status
##
##                          Df Deviance    AIC    LRT  Pr(>Chi)
## + duration                1   801.74 811.74 26.6915 2.387e-07 ***
```

```
## + property                 3   811.90 825.90 16.5377 0.0008796 ***
## + savings                  4   814.60 830.60 13.8362 0.0078365 **
## + other_installment_plans  2   820.32 832.32  8.1125 0.0173134 *
## + age                      1   823.40 833.40  5.0337 0.0248586 *
## <none>                         828.43 836.43
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=811.74
## credit_risk ~ status + duration
##
##                           Df Deviance    AIC    LRT Pr(>Chi)
## + savings                  4   787.09 805.09 14.6529 0.005478 **
## + other_installment_plans  2   794.74 808.74  7.0039 0.030138 *
## + age                      1   797.00 809.00  4.7411 0.029450 *
## + property                 3   794.81 810.81  6.9303 0.074154 .
## <none>                         801.74 811.74
## + status:duration          3   797.97 813.97  3.7701 0.287383
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=805.09
## credit_risk ~ status + duration + savings
##
##                           Df Deviance    AIC    LRT Pr(>Chi)
## + duration:savings         4   773.61 799.61 13.4777 0.009163 **
## + other_installment_plans  2   779.60 801.60  7.4878 0.023662 *
## + age                      1   783.28 803.28  3.8072 0.051032 .
## + property                 3   780.56 804.56  6.5277 0.088578 .
## <none>                         787.09 805.09
## + status:duration          3   783.47 807.47  3.6232 0.305134
## + status:savings          12   769.40 811.40 17.6913 0.125392
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=799.61
## credit_risk ~ status + duration + savings + duration:savings
##
##                           Df Deviance    AIC    LRT Pr(>Chi)
## + other_installment_plans  2   765.10 795.10  8.5114  0.01418 *
## + age                      1   769.17 797.17  4.4427  0.03505 *
## + property                 3   766.99 798.99  6.6259  0.08483 .
## <none>                         773.61 799.61
## + status:duration          3   770.18 802.18  3.4347  0.32934
## + status:savings          12   760.49 810.49 13.1277  0.35983
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=795.1
## credit_risk ~ status + duration + savings + other_installment_plans +
##     duration:savings
##
##                                  Df Deviance    AIC    LRT Pr(>Chi)
## + age                             1   760.07 792.07  5.0305   0.0249 *
```

17

```
## <none>                                          765.10 795.10
## + property                             3     759.95 795.95   5.1531     0.1609
## + status:duration                      3     761.92 797.92   3.1767     0.3652
## + duration:other_installment_plans     2     765.05 799.05   0.0514     0.9746
## + savings:other_installment_plans      8     756.37 802.37   8.7308     0.3655
## + status:other_installment_plans       6     761.30 803.30   3.8027     0.7034
## + status:savings                      12     753.00 807.00  12.1014     0.4376
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=792.07
## credit_risk ~ status + duration + savings + other_installment_plans +
##     age + duration:savings
##
##                                        Df Deviance    AIC     LRT Pr(>Chi)
## + age:other_installment_plans           2     750.02 786.02 10.0493 0.006574 **
## + property                              3     752.24 790.24  7.8305 0.049647 *
## <none>                                        760.07 792.07
## + duration:age                          1     759.61 793.61  0.4620 0.496669
## + status:duration                       3     756.45 794.45  3.6226 0.305204
## + duration:other_installment_plans      2     760.04 796.04  0.0335 0.983389
## + status:age                            3     758.23 796.23  1.8404 0.606175
## + savings:age                           4     757.98 797.98  2.0917 0.718903
## + savings:other_installment_plans       8     750.31 798.31  9.7585 0.282388
## + status:other_installment_plans        6     756.14 800.14  3.9290 0.686284
## + status:savings                       12     747.44 803.44 12.6261 0.396794
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=786.02
## credit_risk ~ status + duration + savings + other_installment_plans +
##     age + duration:savings + other_installment_plans:age
##
##                                        Df Deviance    AIC     LRT Pr(>Chi)
## + property                              3     742.30 784.30  7.7188   0.0522 .
## <none>                                        750.02 786.02
## + duration:age                          1     749.76 787.76  0.2605   0.6098
## + status:duration                       3     745.89 787.89  4.1336   0.2474
## + status:age                            3     747.65 789.65  2.3749   0.4983
## + duration:other_installment_plans      2     749.97 789.97  0.0476   0.9765
## + savings:age                           4     747.70 791.70  2.3220   0.6768
## + status:other_installment_plans        6     745.83 793.83  4.1928   0.6506
## + savings:other_installment_plans       8     742.20 794.20  7.8179   0.4515
## + status:savings                       12     739.02 799.02 10.9992   0.5290
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=784.3
## credit_risk ~ status + duration + savings + other_installment_plans +
##     age + property + duration:savings + other_installment_plans:age
##
##                                        Df Deviance    AIC     LRT Pr(>Chi)
## + status:property                       9     722.44 782.44 19.8671  0.01875 *
## <none>                                        742.30 784.30
```

```
## + duration:age                      1   742.22 786.22  0.0787  0.77913
## + status:duration                   3   738.37 786.37  3.9359  0.26847
## + duration:property                 3   738.94 786.94  3.3634  0.33891
## + status:age                        3   740.05 788.05  2.2542  0.52135
## + duration:other_installment_plans  2   742.29 788.29  0.0141  0.99300
## + property:age                      3   741.47 789.47  0.8287  0.84258
## + savings:age                       4   739.79 789.79  2.5115  0.64258
## + property:other_installment_plans  6   737.35 791.35  4.9562  0.54944
## + status:other_installment_plans    6   738.23 792.23  4.0688  0.66737
## + savings:other_installment_plans   8   734.97 792.97  7.3372  0.50073
## + savings:property                 12   729.92 795.92 12.3800  0.41566
## + status:savings                   12   731.61 797.61 10.6915  0.55552
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=782.44
## credit_risk ~ status + duration + savings + other_installment_plans +
##     age + property + duration:savings + other_installment_plans:age +
##     status:property
##
##                                    Df Deviance    AIC     LRT Pr(>Chi)
## <none>                                722.44 782.44
## + duration:age                      1   722.35 784.35  0.0832  0.7730
## + status:age                        3   718.69 784.69  3.7428  0.2906
## + duration:property                 3   718.92 784.92  3.5143  0.3189
## + status:duration                   3   720.15 786.15  2.2823  0.5159
## + duration:other_installment_plans  2   722.22 786.22  0.2146  0.8983
## + savings:age                       4   719.83 787.83  2.6071  0.6256
## + property:age                      3   721.88 787.88  0.5560  0.9064
## + property:other_installment_plans  6   717.17 789.17  5.2635  0.5105
## + savings:other_installment_plans   8   715.61 791.61  6.8292  0.5552
## + status:other_installment_plans    6   720.08 792.08  2.3606  0.8837
## + savings:property                 12   709.66 793.66 12.7724  0.3858
## + status:savings                   12   710.09 794.09 12.3503  0.4180
```
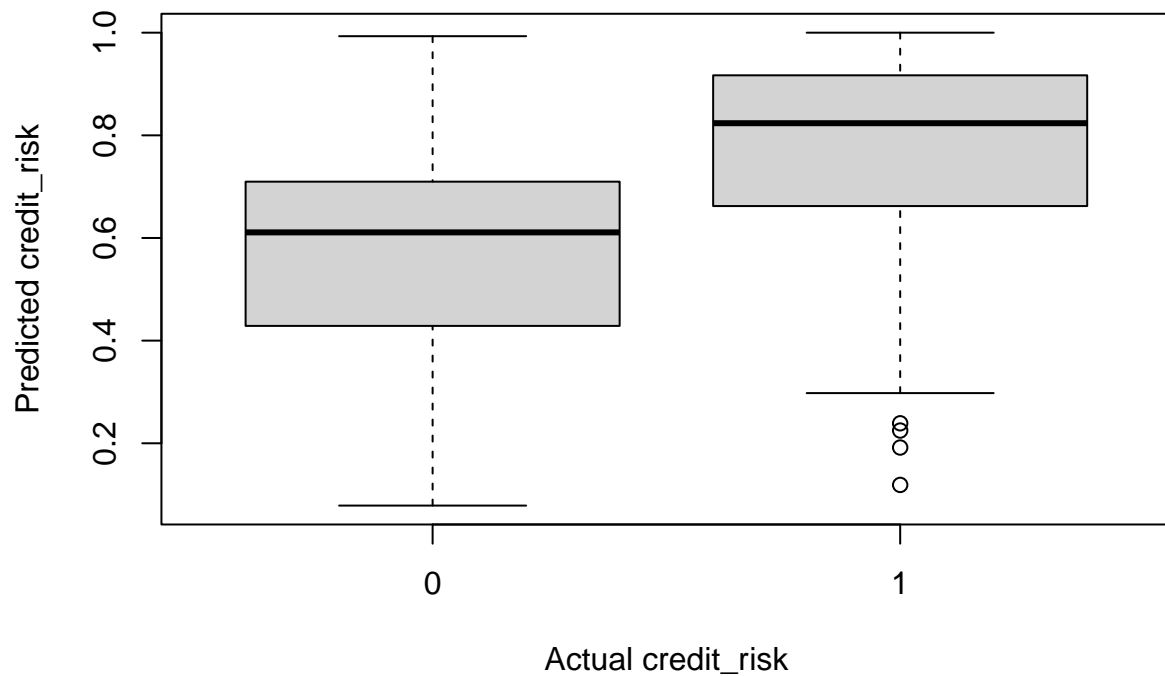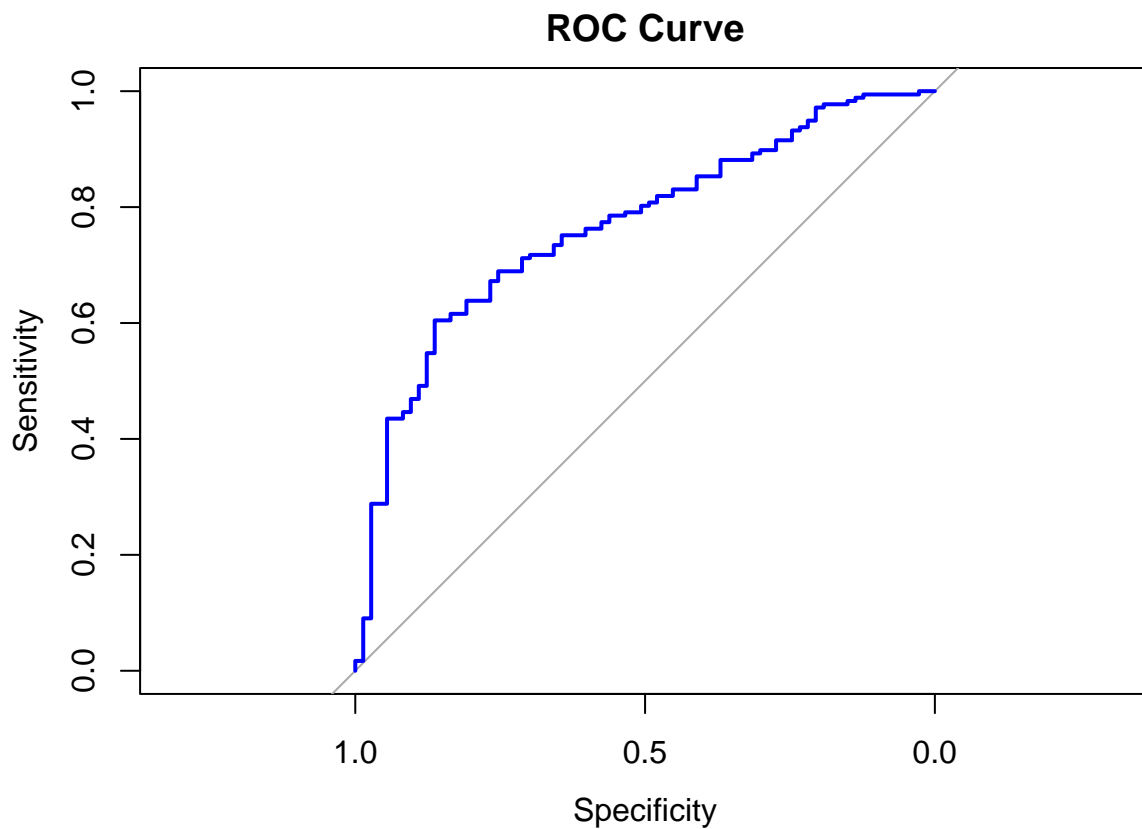
### 3.4.2  Testing model

```
pred.3 <- predict(bestmodel.3, newdata = testdata)
plot(testdata$credit_risk, inv.logit(pred.3), xlab = "Actual credit_risk", ylab = "Predicted credit_ris
```

```
roc(testdata$credit_risk~inv.logit(pred.3), plot=TRUE, main="ROC Curve", col="blue")
```

## ROC Curve



```
##
## Call:
## roc.formula(formula = testdata$credit_risk ~ inv.logit(pred.3),     plot = TRUE, main = "ROC Curve",
```

```
## 
## Data: inv.logit(pred.3) in 73 controls (testdata$credit_risk 0) < 177 cases (testdata$credit_risk 1)
## Area under the curve: 0.7659
```

```
auc(testdata$credit_risk~inv.logit(pred.3))
```

```
## Area under the curve: 0.7659
```