

# Report on Default of Credit Card Clients Dataset

Masayoshi Sato

2021/6/26

## Contents

Introduction . . . . .	2
Packages and Dataset . . . . .	2
Data Exploration . . . . .	3
1 Outcome . . . . .	6
2 “LIMIT_BAL” . . . . .	6
3 “SEX” . . . . .	7
4 “EDUCATION” . . . . .	9
5 “MARRIAGE” . . . . .	11
6 “AGE” . . . . .	13
7 “PAY” . . . . .	14
8 “BILL_AMT” . . . . .	17
9 “PAY_AMT” . . . . .	19
Data Preparation . . . . .	21
Model analysis . . . . .	24
1 Baseline prediction . . . . .	24
2 Logistic regression . . . . .	25
3 Decision tree default model . . . . .	36
4 Decision tree further tuning . . . . .	38
5 Random forest default . . . . .	41
6 Random forest cross validation . . . . .	43
Evaluation . . . . .	45
Conclusion . . . . .	47

## Introduction

Finance thought to be a field where machine learning can be effective. It is surrounded by uncertainties, such as economic downturn, a collapse of markets. Individuals are no exception, we are not sure that a person is credible enough to lend money. They might be a deadbeat, or struggle in a significant debt, even though they look credible. This is the reason why machine learning plays a role in predicting uncertain economic future.

In this paper, we will deal with a familiar problem. Based on objective data, can we predict effectively whether a credit user will pay their debt or result in default? Traditionally, finding a credible borrower have been a skill and experience nurtured by financial institutions, like banks, credit company. Instead, we will look into open data, and using machine learning models, such as logistic regression, decision tree, and random forest. We will fit the data with these models and find the model which will show the most accurate prediction.

The dataset we use, “Default of Credit Card Clients Dataset” is stored in Kaggle website. It was collected in Taiwan in 2005. It has 24 variables, such as age, education, and payment condition. Outcome has two results, “0” non-default, “1” default.

Our goal is to find a classification model which predicts the most accurate outcome. However, its distribution of outcome is imbalanced. Namely, the number of default clients are small compared to non- default clients. To address the issue, we will use other criteria, balanced accuracy.

We will use three machine learning models, logistic regression, decision tree, and random forest. If necessary, we will tune their parameters to find the best solution. Overall procedures are as follows :

1. Data exploration and data cleansing
2. Splitting the dataset into train\_set, validation\_set, and test\_set
3. Applying models, logistic regression, decision tree, and random forest
4. Considering models performance, and evaluating

This paper is written as a final assignment in “HarvardX PH125.9x Data Science: Capstone.”

## Packages and Dataset

In this paper, we use R packages, “tidyverse<sup>1</sup>”, “DataExplorer<sup>2</sup>”, “gridExtra<sup>3</sup>”, “rpart<sup>4</sup>”, “caret<sup>5</sup>”, and “ranger<sup>6</sup>”.

We use a dataset stored in Kaggle<sup>7</sup> website. In the description, it says, “This dataset contains information on default payments, demographic factors, credit data, history of payment, and bill statements of credit card clients in Taiwan from April 2005 to September 2005.” It is CSV file.

Kaggle requires registration to download the data. For the sake of convenience, the data file is stored in my GitHub repository<sup>8</sup>.

---

<sup>1</sup><https://cran.r-project.org/web/packages/tidyverse/index.html>

<sup>2</sup><https://cran.r-project.org/web/packages/DataExplorer/vignettes/dataexplorer-intro.html>

<sup>3</sup><https://cran.r-project.org/web/packages/gridExtra/index.html>

<sup>4</sup><https://cran.r-project.org/web/packages/rpart/rpart.pdf>

<sup>5</sup><https://topepo.github.io/caret/>

<sup>6</sup><https://cran.r-project.org/web/packages/ranger/ranger.pdf>

Compared to “randomForest”, “ranger” is very quick and easy to operate.

<sup>7</sup><https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>

<sup>8</sup>[https://github.com/masa951125/Final\\_project/raw/main/UCI\\_Credit\\_Card.csv](https://github.com/masa951125/Final_project/raw/main/UCI_Credit_Card.csv)

## Data Exploration

First, we need to check the downloaded dataset. It has 30000 rows and 25 columns.

We look into the data using “str” and “summary” function.

```
## tibble [30,000 x 25] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ ID : num [1:30000] 1 2 3 4 5 6 7 8 9 10 ...
## $ LIMIT_BAL : num [1:30000] 20000 120000 90000 50000 50000 50000 50000 100000 140000 ...
## $ SEX : num [1:30000] 2 2 2 2 1 1 1 2 2 1 ...
## $ EDUCATION : num [1:30000] 2 2 2 2 2 1 1 2 3 3 ...
## $ MARRIAGE : num [1:30000] 1 2 2 1 1 2 2 2 1 2 ...
## $ AGE : num [1:30000] 24 26 34 37 57 37 29 23 28 35 ...
## $ PAY_0 : num [1:30000] 2 -1 0 0 -1 0 0 0 0 -2 ...
## $ PAY_2 : num [1:30000] 2 2 0 0 0 0 0 -1 0 -2 ...
## $ PAY_3 : num [1:30000] -1 0 0 0 -1 0 0 -1 2 -2 ...
## $ PAY_4 : num [1:30000] -1 0 0 0 0 0 0 0 0 -2 ...
## $ PAY_5 : num [1:30000] -2 0 0 0 0 0 0 0 0 -1 ...
## $ PAY_6 : num [1:30000] -2 2 0 0 0 0 0 -1 0 -1 ...
## $ BILL_AMT1 : num [1:30000] 3913 2682 29239 46990 8617 ...
## $ BILL_AMT2 : num [1:30000] 3102 1725 14027 48233 5670 ...
## $ BILL_AMT3 : num [1:30000] 689 2682 13559 49291 35835 ...
## $ BILL_AMT4 : num [1:30000] 0 3272 14331 28314 20940 ...
## $ BILL_AMT5 : num [1:30000] 0 3455 14948 28959 19146 ...
## $ BILL_AMT6 : num [1:30000] 0 3261 15549 29547 19131 ...
## $ PAY_AMT1 : num [1:30000] 0 0 1518 2000 2000 ...
## $ PAY_AMT2 : num [1:30000] 689 1000 1500 2019 36681 ...
## $ PAY_AMT3 : num [1:30000] 0 1000 1000 1200 10000 657 38000 0 432 0 ...
## $ PAY_AMT4 : num [1:30000] 0 1000 1000 1100 9000 ...
## $ PAY_AMT5 : num [1:30000] 0 0 1000 1069 689 ...
## $ PAY_AMT6 : num [1:30000] 0 2000 5000 1000 679 ...
## $ default.payment.next.month: num [1:30000] 1 1 0 0 0 0 0 0 0 0 ...
## - attr(*, "spec")=
## .. cols(
## .. ID = col_double(),
## .. LIMIT_BAL = col_double(),
## .. SEX = col_double(),
## .. EDUCATION = col_double(),
## .. MARRIAGE = col_double(),
## .. AGE = col_double(),
## .. PAY_0 = col_double(),
## .. PAY_2 = col_double(),
## .. PAY_3 = col_double(),
## .. PAY_4 = col_double(),
## .. PAY_5 = col_double(),
## .. PAY_6 = col_double(),
## .. BILL_AMT1 = col_double(),
## .. BILL_AMT2 = col_double(),
## .. BILL_AMT3 = col_double(),
## .. BILL_AMT4 = col_double(),
## .. BILL_AMT5 = col_double(),
## .. BILL_AMT6 = col_double(),
## .. PAY_AMT1 = col_double(),
## .. PAY_AMT2 = col_double(),
```

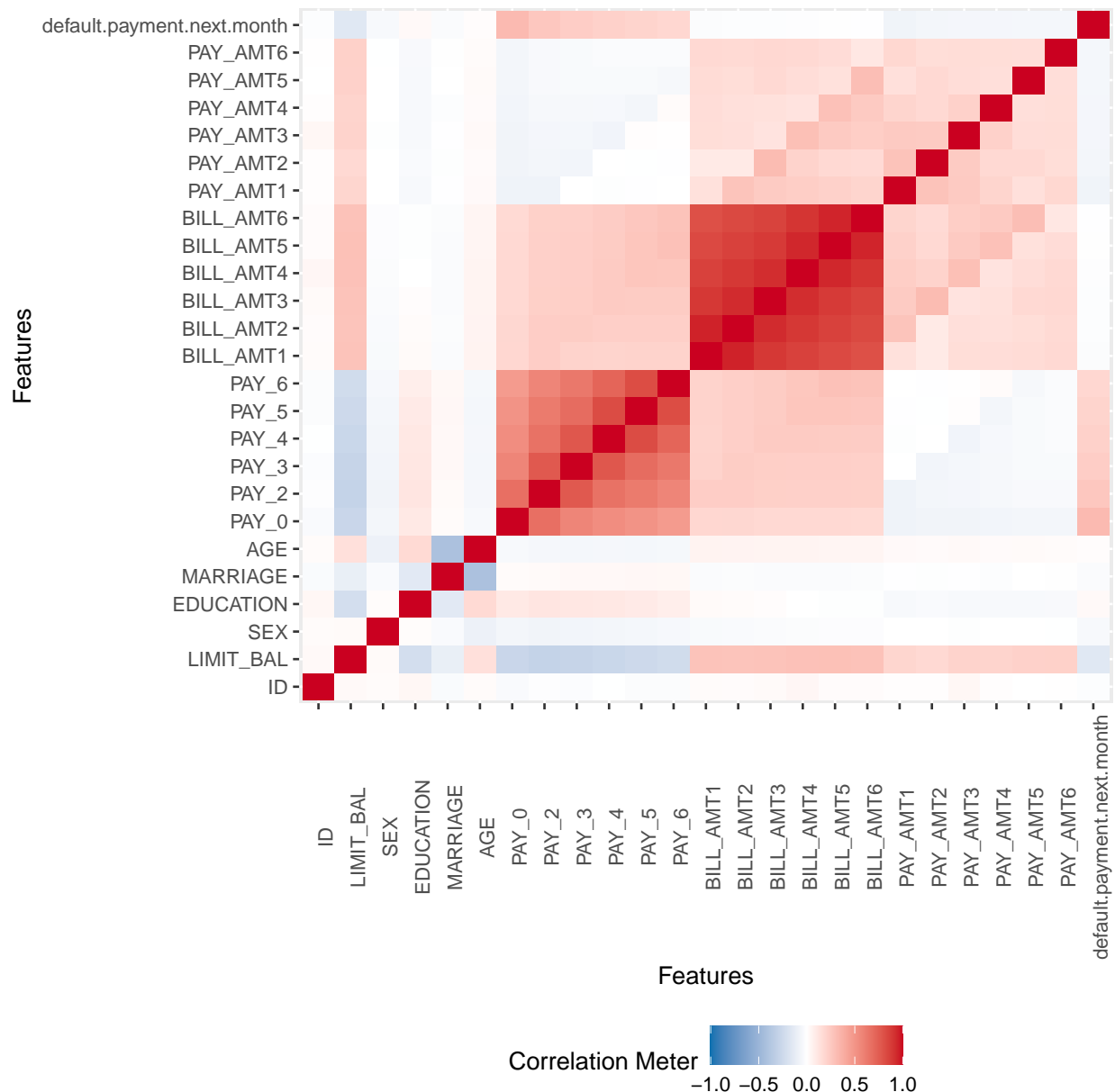
```
## .. PAY_AMT3 = col_double(),
## .. PAY_AMT4 = col_double(),
## .. PAY_AMT5 = col_double(),
## .. PAY_AMT6 = col_double(),
## .. default.payment.next.month = col_double()
## .. )
```

```
##          ID          LIMIT_BAL          SEX          EDUCATION
## Min.   :    1   Min.   : 10000   Min.   :1.000   Min.   :0.000
## 1st Qu.: 7501   1st Qu.: 50000   1st Qu.:1.000   1st Qu.:1.000
## Median :15000   Median : 140000   Median :2.000   Median :2.000
## Mean   :15000   Mean   : 167484   Mean   :1.604   Mean   :1.853
## 3rd Qu.:22500   3rd Qu.: 240000   3rd Qu.:2.000   3rd Qu.:2.000
## Max.   :30000   Max.   :1000000   Max.   :2.000   Max.   :6.000
## MARRIAGE          AGE          PAY_0          PAY_2
## Min.   :0.000   Min.   :21.00   Min.   : -2.0000   Min.   : -2.0000
## 1st Qu.:1.000   1st Qu.:28.00   1st Qu.: -1.0000   1st Qu.: -1.0000
## Median :2.000   Median :34.00   Median : 0.0000   Median : 0.0000
## Mean   :1.552   Mean   :35.49   Mean   : -0.0167   Mean   : -0.1338
## 3rd Qu.:2.000   3rd Qu.:41.00   3rd Qu.: 0.0000   3rd Qu.: 0.0000
## Max.   :3.000   Max.   :79.00   Max.   : 8.0000   Max.   : 8.0000
## PAY_3          PAY_4          PAY_5          PAY_6
## Min.   : -2.0000   Min.   : -2.0000   Min.   : -2.0000   Min.   : -2.0000
## 1st Qu.: -1.0000   1st Qu.: -1.0000   1st Qu.: -1.0000   1st Qu.: -1.0000
## Median : 0.0000   Median : 0.0000   Median : 0.0000   Median : 0.0000
## Mean   : -0.1662   Mean   : -0.2207   Mean   : -0.2662   Mean   : -0.2911
## 3rd Qu.: 0.0000   3rd Qu.: 0.0000   3rd Qu.: 0.0000   3rd Qu.: 0.0000
## Max.   : 8.0000   Max.   : 8.0000   Max.   : 8.0000   Max.   : 8.0000
## BILL_AMT1       BILL_AMT2       BILL_AMT3       BILL_AMT4
## Min.   : -165580   Min.   : -69777   Min.   : -157264   Min.   : -170000
## 1st Qu.: 3559     1st Qu.: 2985     1st Qu.: 2666     1st Qu.: 2327
## Median : 22382     Median : 21200     Median : 20089     Median : 19052
## Mean   : 51223     Mean   : 49179     Mean   : 47013     Mean   : 43263
## 3rd Qu.: 67091     3rd Qu.: 64006     3rd Qu.: 60165     3rd Qu.: 54506
## Max.   : 964511     Max.   : 983931     Max.   : 1664089     Max.   : 891586
## BILL_AMT5       BILL_AMT6       PAY_AMT1       PAY_AMT2
## Min.   : -81334     Min.   : -339603   Min.   : 0         Min.   : 0
## 1st Qu.: 1763       1st Qu.: 1256     1st Qu.: 1000     1st Qu.: 833
## Median : 18105     Median : 17071     Median : 2100     Median : 2009
## Mean   : 40311     Mean   : 38872     Mean   : 5664     Mean   : 5921
## 3rd Qu.: 50191     3rd Qu.: 49198     3rd Qu.: 5006     3rd Qu.: 5000
## Max.   : 927171     Max.   : 961664     Max.   : 873552     Max.   : 1684259
## PAY_AMT3       PAY_AMT4       PAY_AMT5       PAY_AMT6
## Min.   : 0         Min.   : 0         Min.   : 0.0       Min.   : 0.0
## 1st Qu.: 390       1st Qu.: 296     1st Qu.: 252.5     1st Qu.: 117.8
## Median : 1800     Median : 1500     Median : 1500.0     Median : 1500.0
## Mean   : 5226     Mean   : 4826     Mean   : 4799.4     Mean   : 5215.5
## 3rd Qu.: 4505     3rd Qu.: 4013     3rd Qu.: 4031.5     3rd Qu.: 4000.0
## Max.   : 896040     Max.   : 621000     Max.   : 426529.0     Max.   : 528666.0
## default.payment.next.month
## Min.   :0.0000
## 1st Qu.:0.0000
## Median :0.0000
## Mean   :0.2212
```

```
## 3rd Qu.:0.0000
## Max.    :1.0000
```

“Default.payment.next.month” is an outcome of the dataset which has values of 0 and 1. Other features seem to be either numerical or categorical data. “SEX”, “EDUCATION”, “MARRIAGE” and “PAY\_0” -“PAY\_6” look like categorical data, as their values are limited number of integers. Other features seem to be numerical. From checking its summary, we understand there are no NAs in the dataset.

How these predictors are correlated? We use “plot\_correlation” function to investigate this.



Takeaways from this are;

1. Outcome (default.payment.next.month) has a strong positive correlation with PAY.
2. Overall, LIMIT\_BAL has a relatively strong correlation with other factors (except SEX).
3. EDUCATION, MARRIAGE, AGE have relatively strong correlation with one another.

4. EDUCATION and AGE have a relatively weak correlation with PAY and BILL\_AMT respectively.
5. PAY and BILL\_AMT, BILL\_AMT and PAY\_AMT have strong correlation.

Then, we will look into these features further.

## 1 Outcome

First, we look into the outcome, “default.payment.next.month”. The data description says, “Default payment, 1=yes, 0=no.”<sup>9</sup> We show the proportion of “0”, “1”.

```
##
##      0      1
## 0.7788 0.2212
```

This means that if we predict all the outcome as “0”, we will get 77.9% accuracy. We need to take into account this fact. We change the name, “default.payment.next.month”, to “DEFAULT” for the sake of convenience. Also, we change this numeric variable into factor.

## 2 “LIMIT\_BAL”

This is an “amount of given credit in NT dollars (includes individual and family/supplementary credit)”<sup>10</sup> It is clearly numeric data.

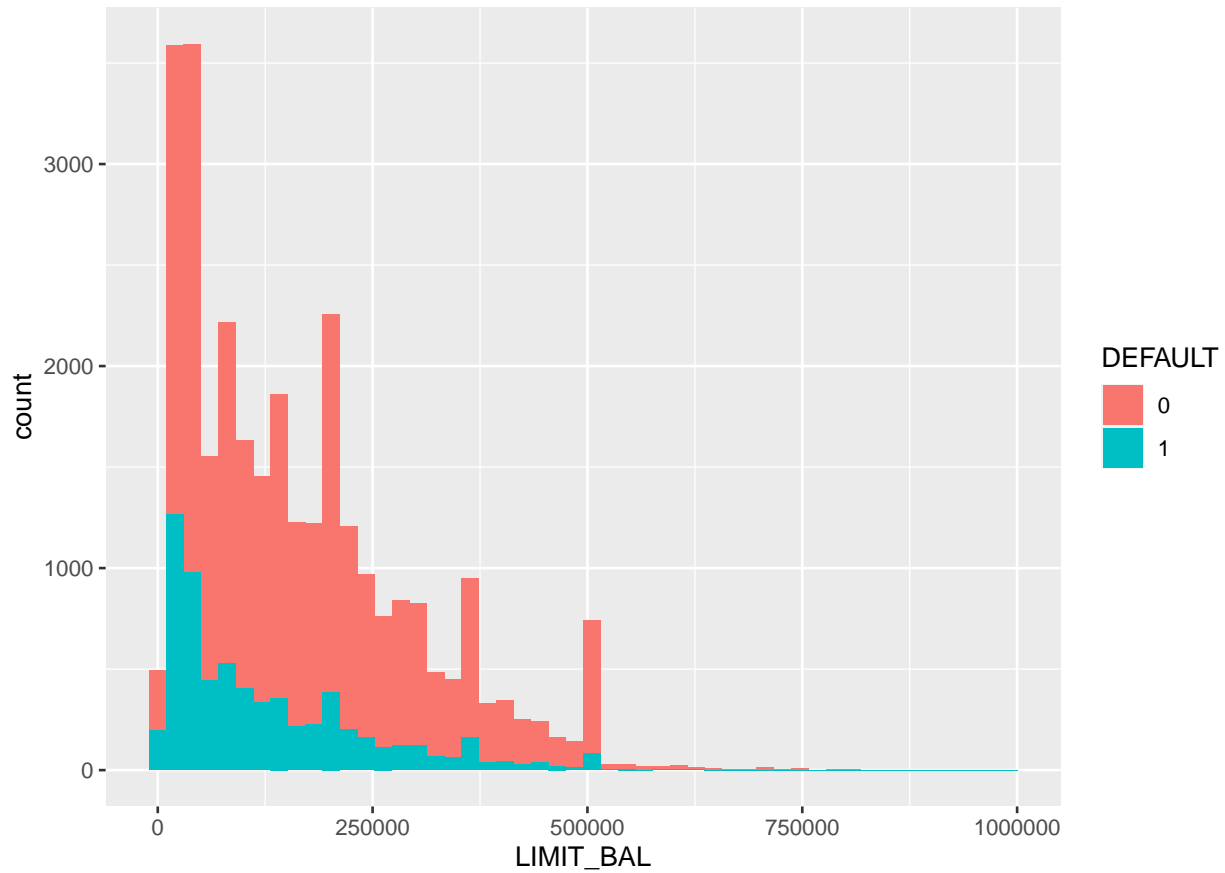
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  10000   50000  140000  167484  240000 1000000
```

We draw its distribution filling the proportion of default.

---

<sup>9</sup><https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>

<sup>10</sup><https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset> NT stands for “New Taiwan”.



Distribution is skewed right. Default clients seem to be gathered around lower range of LIMIT\_BAL values.

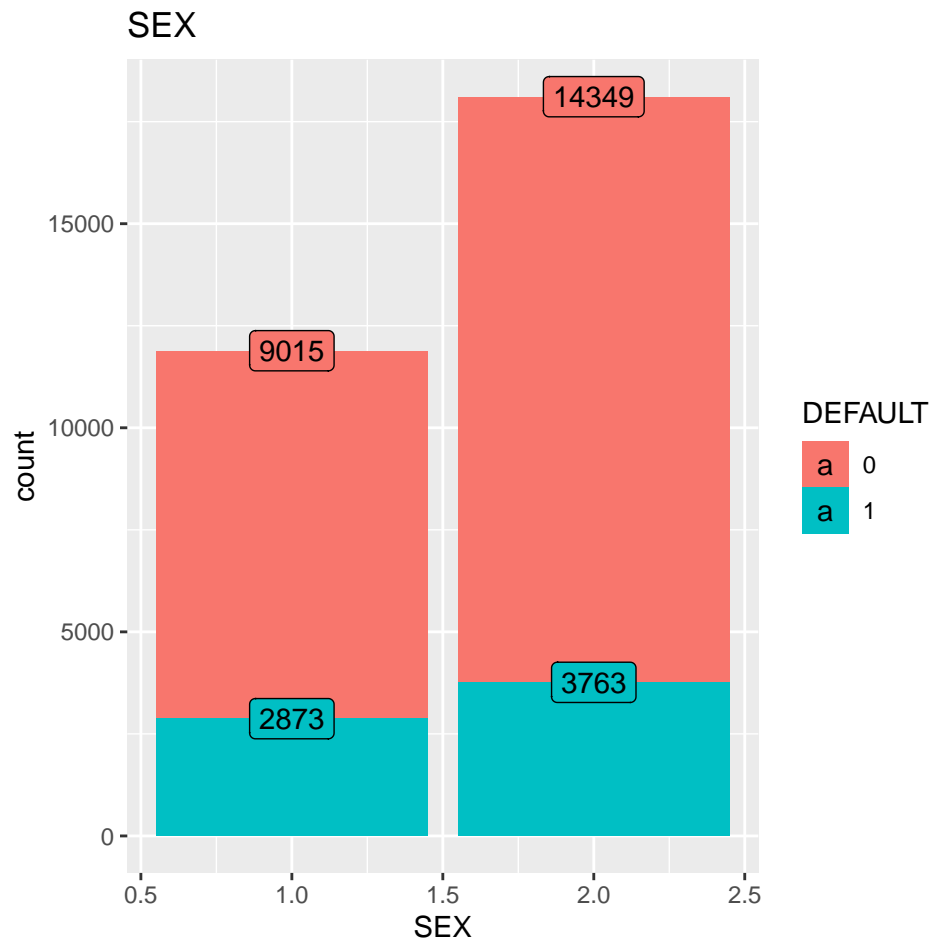
### 3 “SEX”

The values “1”, “2” correspond to male and female respectively<sup>11</sup>. Male is 40% and female is 60%.

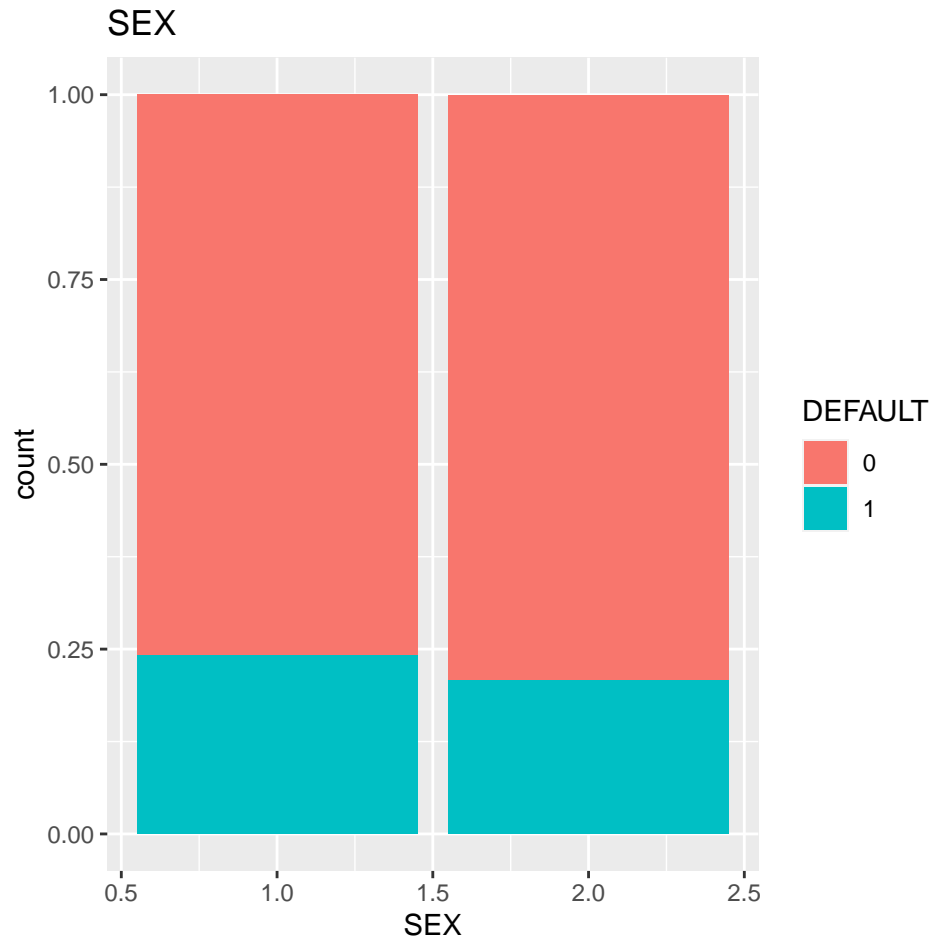
```
##
##           1           2
## 0.3962667 0.6037333
```

We draw its distribution and proportion in terms of default rates.

<sup>11</sup><https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>





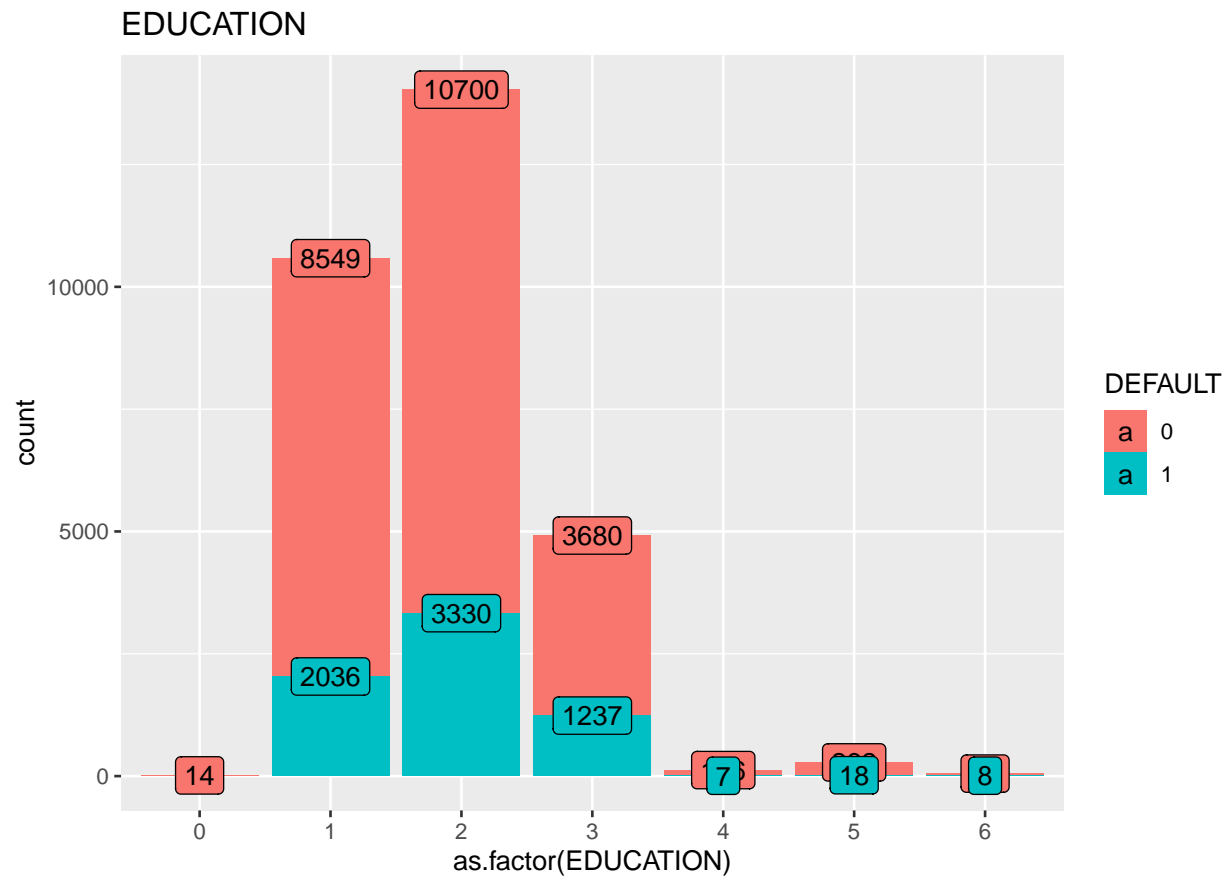


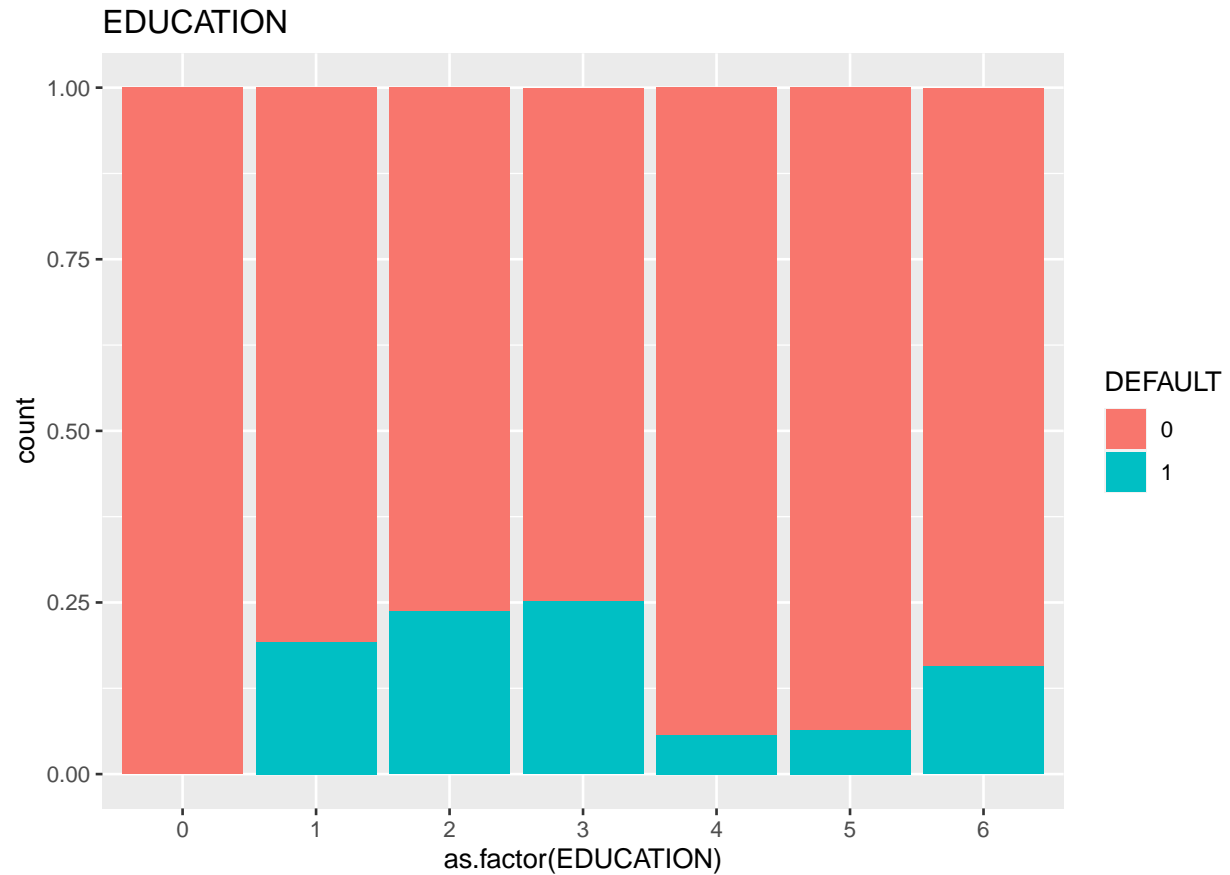
There seemed to be little difference between genders. This categorical variable is somewhat irrelevant to the outcome.

#### 4 “EDUCATION”

In this variable, values are “1”, “2”, “3”, “4”, “5”, “6”. They are categorical values. The numbers have meanings as follows ; 1=graduate school, 2=university, 3=high school, 4=others, 5=unknown, 6=unknown<sup>12</sup>. We plot its distribution and stacked bar graph.

<sup>12</sup><https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>



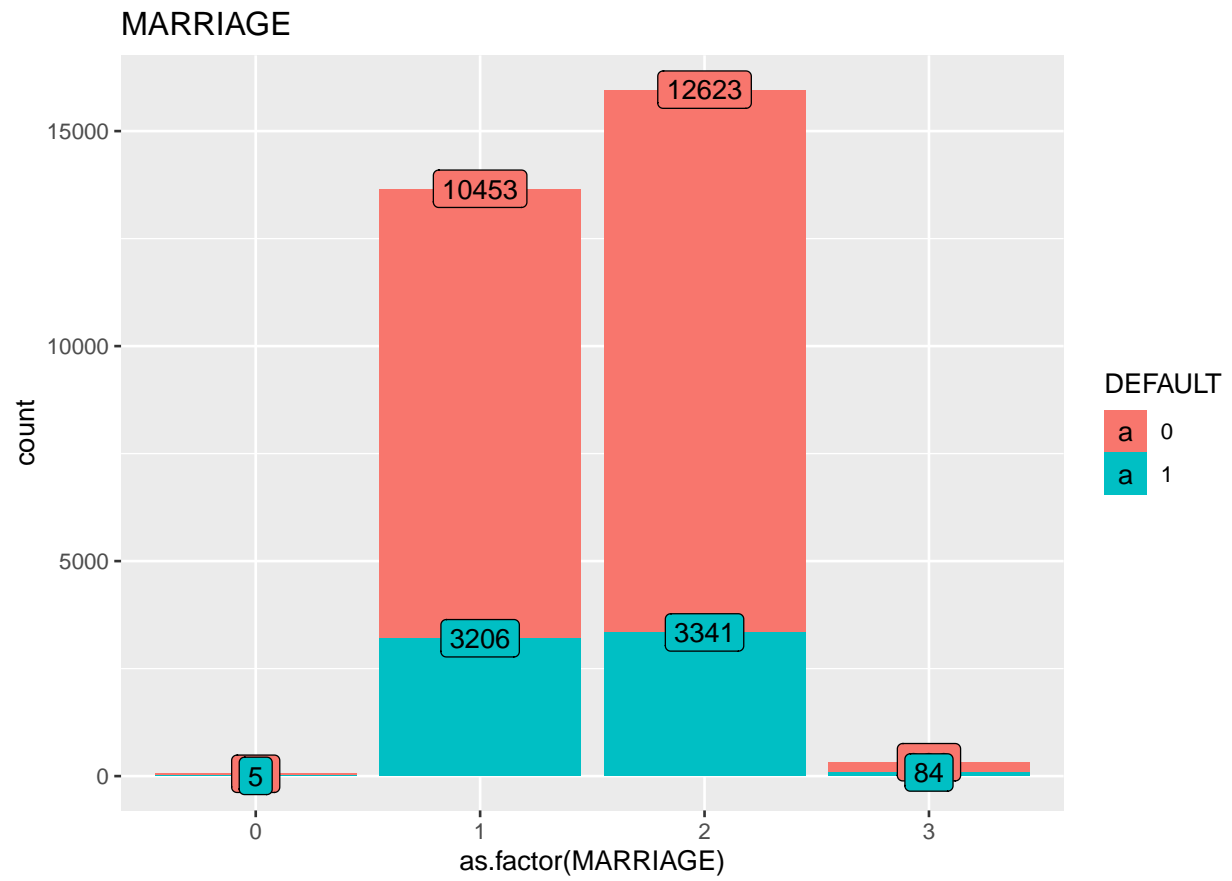


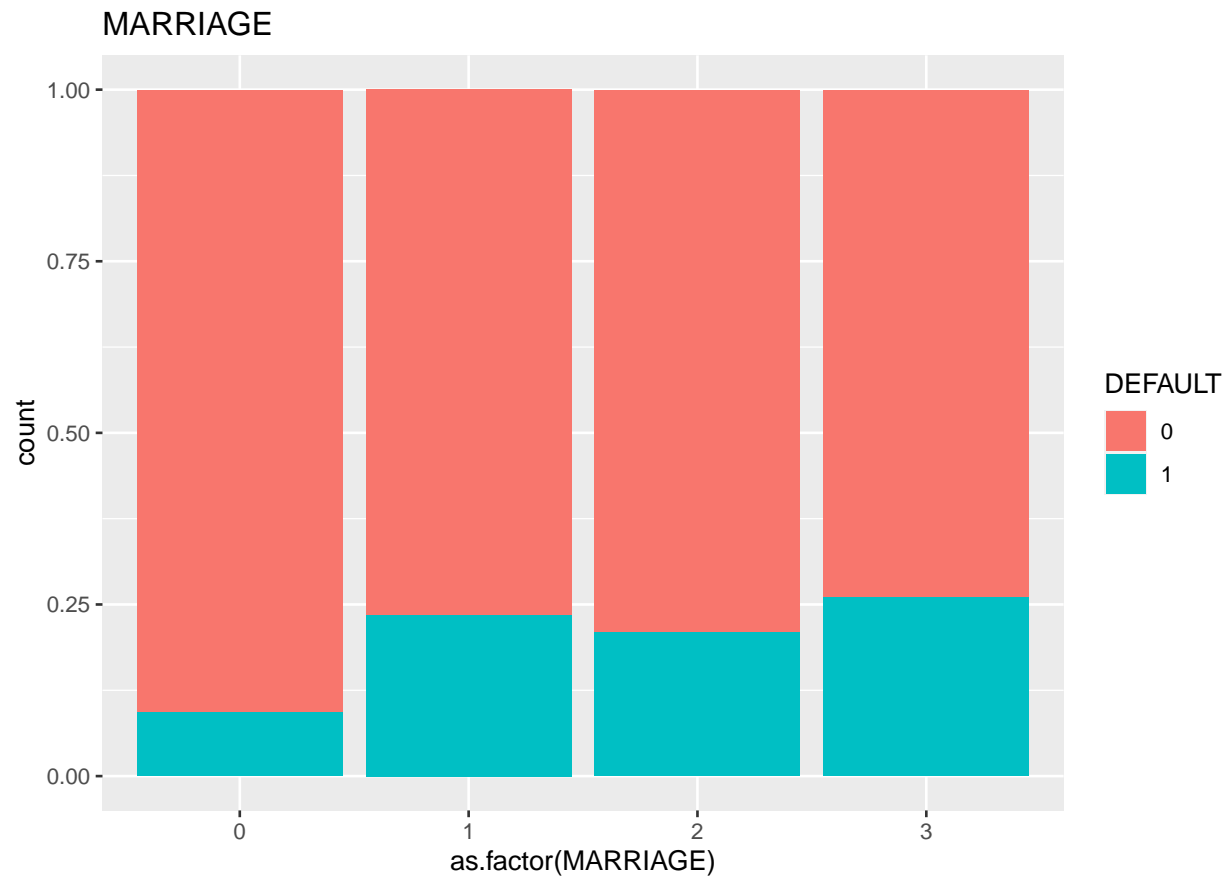
People whose final education is high school have relatively high default rate. On the other hand, people whose final education is graduate school have low default rate.

## 5 “MARRIAGE”

Kaggle’s data explanation says;

marital status. 1=married, 2=single, 3=others. Categorical data. Plot. Stack bar graph





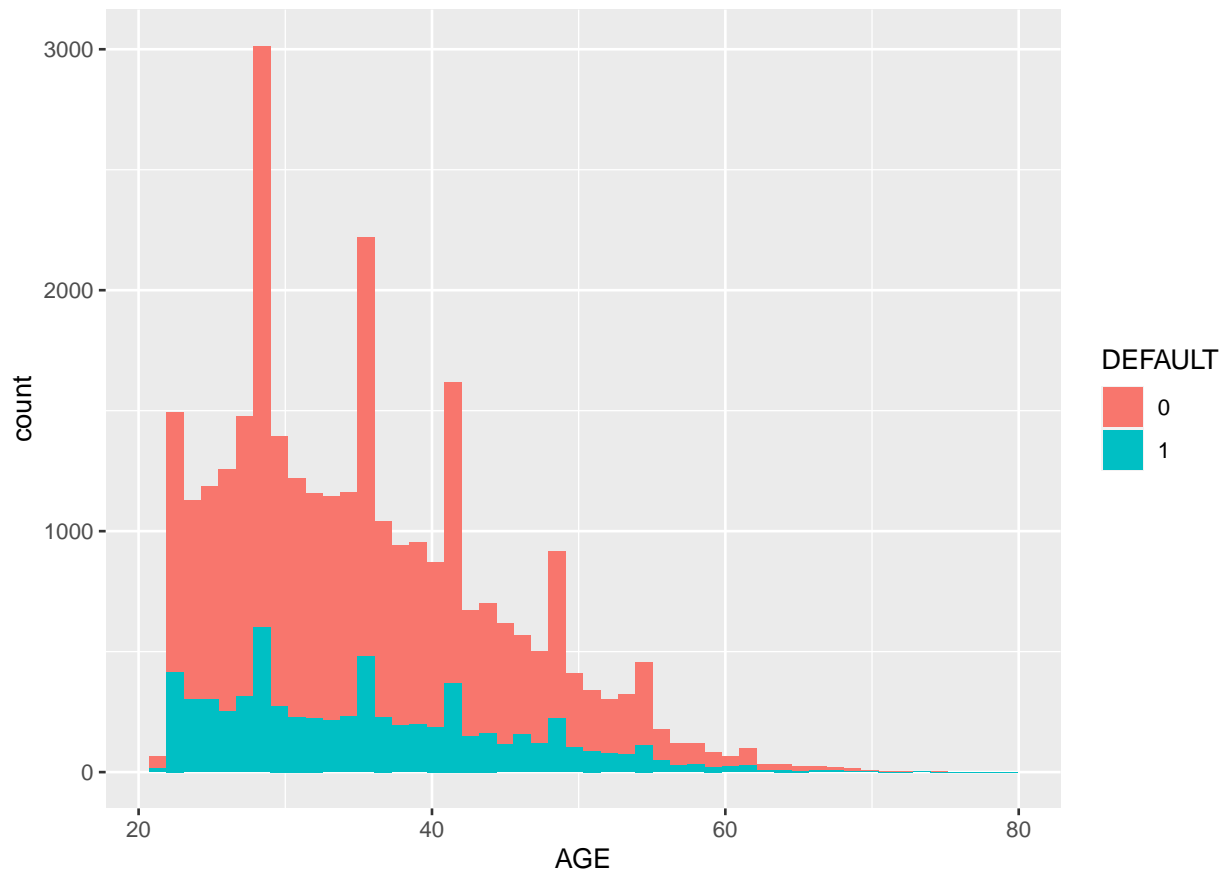
There seems to be little difference among the groups.

## 6 "AGE"

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	21.00	28.00	34.00	35.49	41.00	79.00

numeric data

Plot.



## 7 “PAY”

Kaggle’s data explanation says;

PAY\_0 means repayment status in September, 2005.

-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. Regarding values from PAY\_2 to PAY\_6, the scales are the same as PAY\_0. As the number increases, the date of repayment status goes back in time by a month until April, 2005 which is PAY\_6.

. PAY\_0.

```
summary(original_default$PAY_0)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -2.0000 -1.0000   0.0000 -0.0167  0.0000   8.0000
```

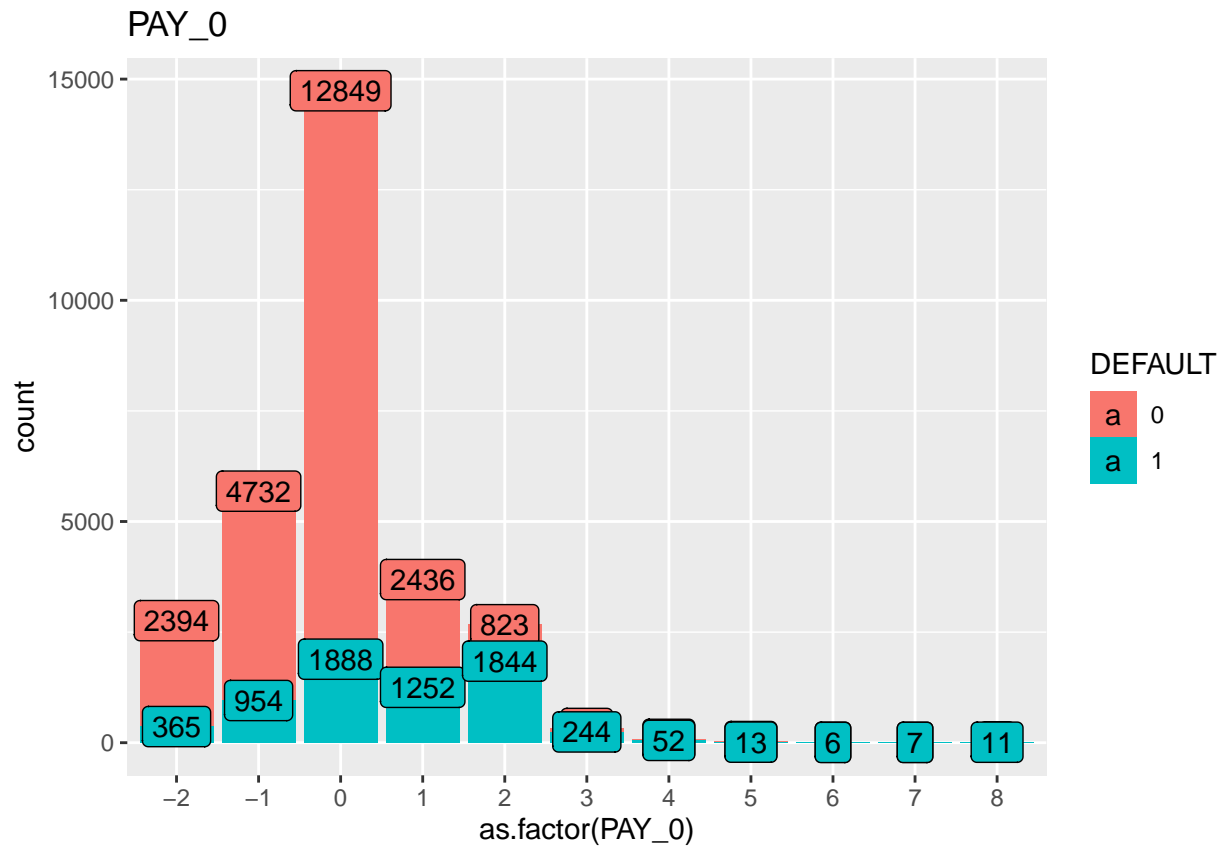
```
unique(original_default$PAY_0)
```

```
##  [1]  2 -1  0 -2  1  3  4  8  7  5  6
```

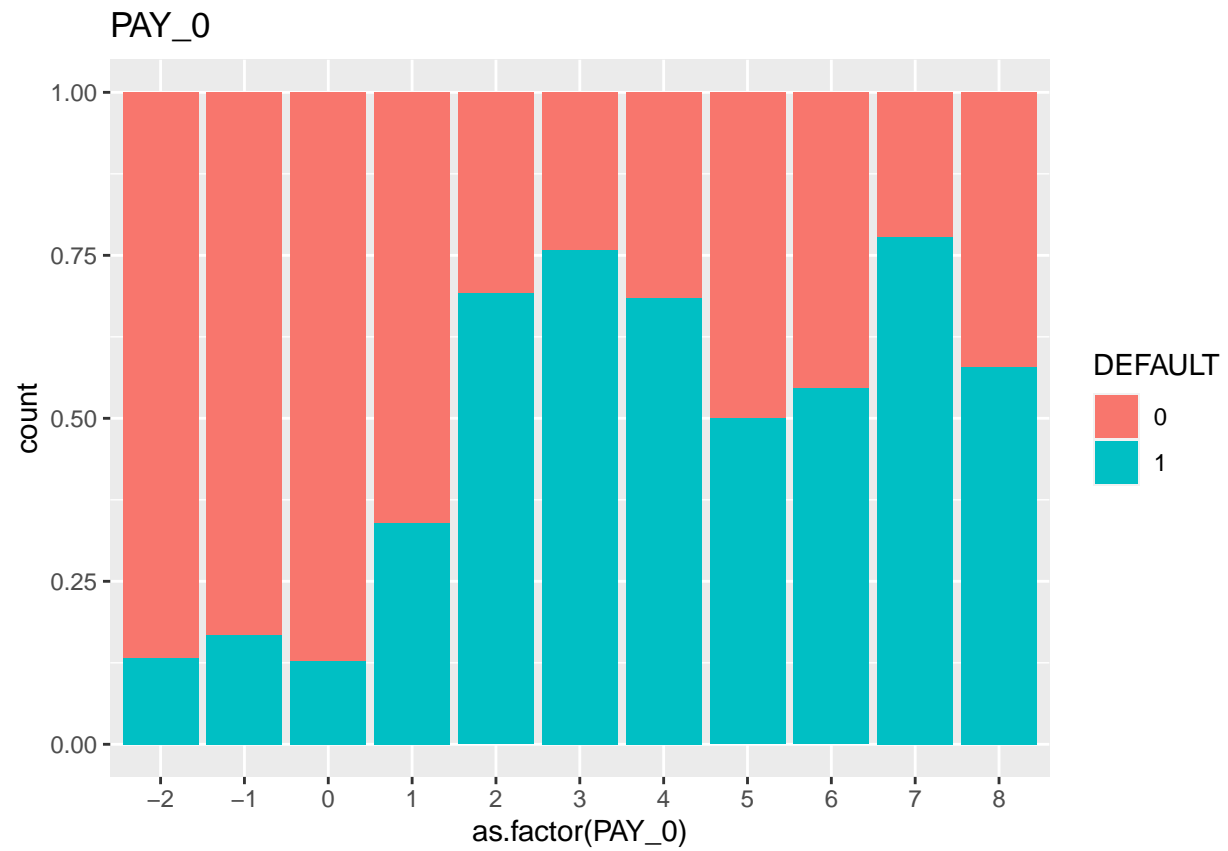
They are categorical data.

Plot.

```
original_default %>% ggplot(aes(x=as.factor(PAY_0), fill= DEFAULT)) +
  geom_bar() +
  ggtitle("PAY_0")+
  stat_count(aes(label = ..count..), geom = "label")# illustrate numbers
```

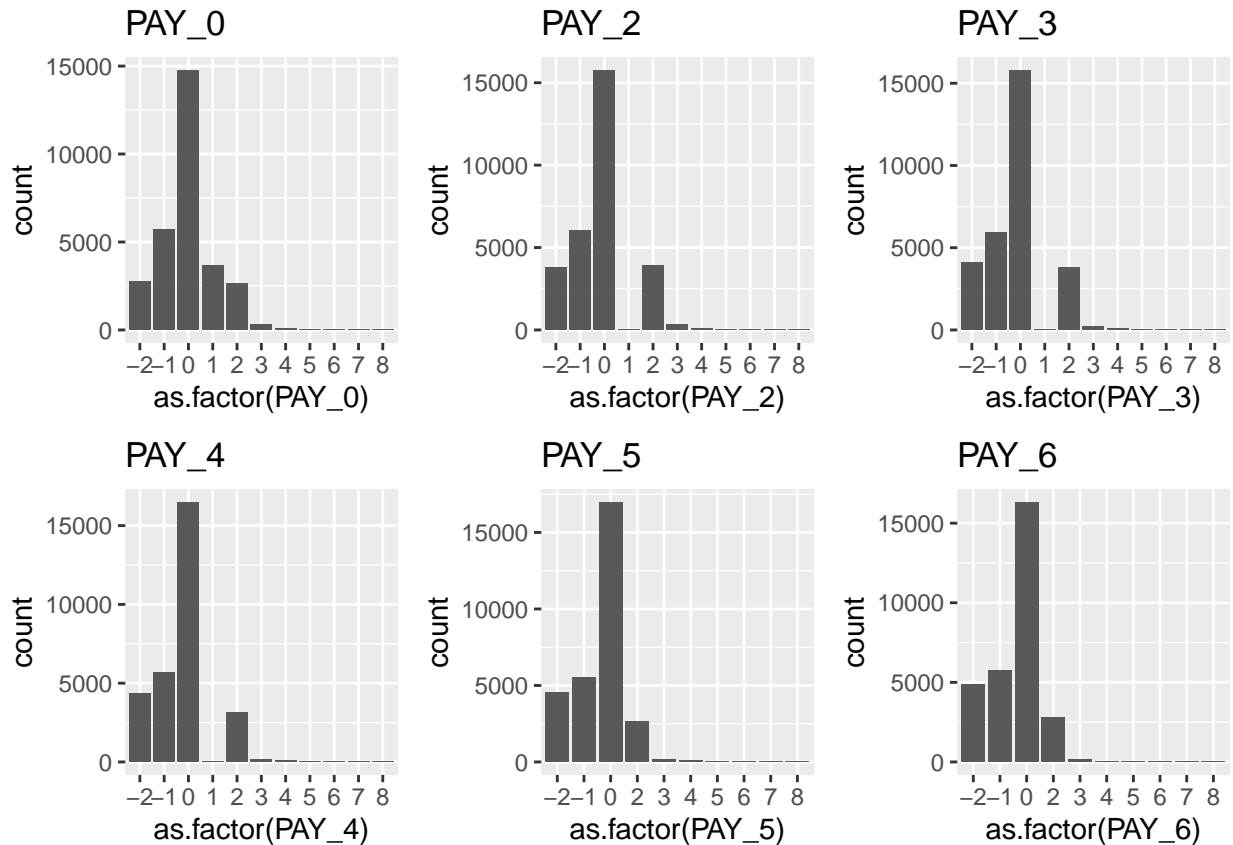


Stack bar graph PAY\_0.



PAY\_2 ~ PAY\_6 's structures are almost as same as PAY\_0. Show distribution.





## 8 “BILL\_AMT”

Kaggle’s data explanation says;

BILL\_AMT1 is an amount of bill statement in September, 2005 (NT dollar). Likewise PAY, BILL\_AMT goes back in time by a month from August to April, 2005 which is BILL\_AMT6.

```
summary(original_default$BILL_AMT1)
```

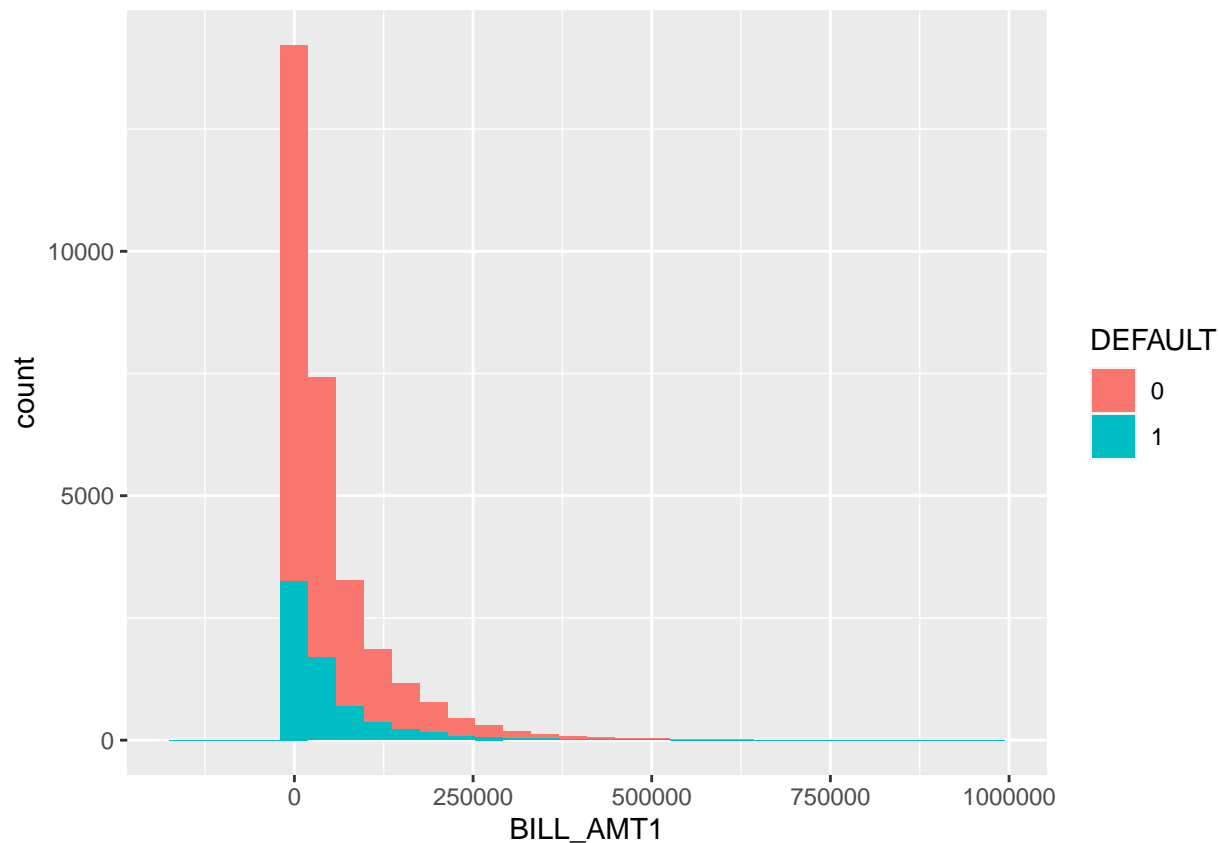
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -165580   3559    22382   51223   67091   964511
```

These are numerical data.

Here is BILL\_AMT1’s plot.

```
ggplot(data=original_default, aes(BILL_AMT1, fill= DEFAULT)) +geom_histogram()
```

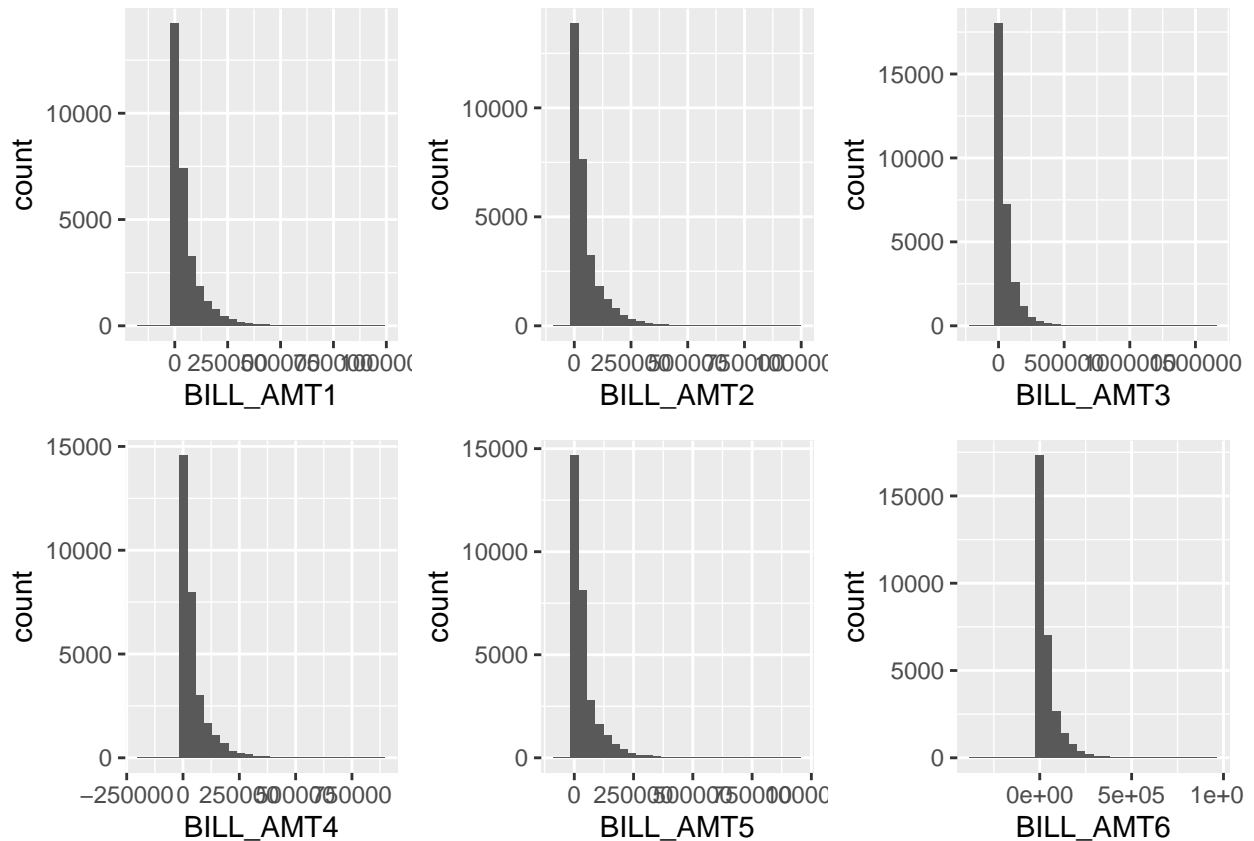
```
## ‘stat_bin()’ using ‘bins = 30’. Pick better value with ‘binwidth’.
```



From BILL\_AMT1 to BILL\_AMT6, their structures are almost the same as are shown in following plots.

```
b1 <- ggplot(data=original_default, aes(BILL_AMT1)) +geom_histogram()
b2 <- ggplot(data=original_default, aes(BILL_AMT2)) +geom_histogram()
b3 <- ggplot(data=original_default, aes(BILL_AMT3)) +geom_histogram()
b4 <- ggplot(data=original_default, aes(BILL_AMT4)) +geom_histogram()
b5 <- ggplot(data=original_default, aes(BILL_AMT5)) +geom_histogram()
b6 <- ggplot(data=original_default, aes(BILL_AMT6)) +geom_histogram()
grid.arrange(b1,b2,b3,b4,b5,b6, nrow=2, ncol=3)
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



## 9 “PAY\_AMT”

Kaggle’s data explanation says;

PAY\_AMT1 is an amount of previous payment in September, 2005 (NT dollar). Likewise BILL\_AMT, PAY\_AMT goes back in time by a month from August to April, 2005 which is PAY\_AMT6.

```
summary(original_default$PAY_AMT1)
```

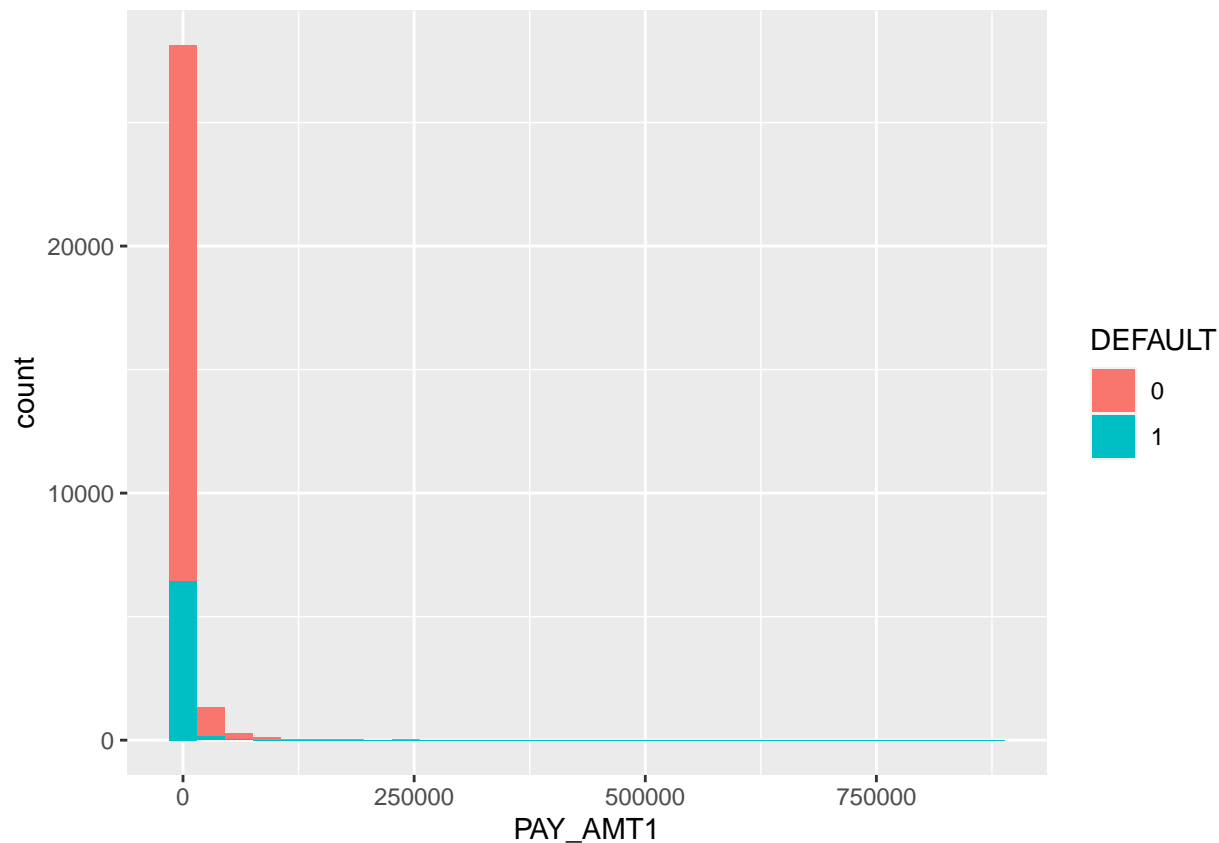
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         0    1000    2100   5664   5006 873552
```

They are numerical data.

Here is PAY\_AMT1’s plot.

```
ggplot(data=original_default, aes(PAY_AMT1, fill= DEFAULT)) +geom_histogram()
```

```
## ‘stat_bin()’ using ‘bins = 30’. Pick better value with ‘binwidth’.
```

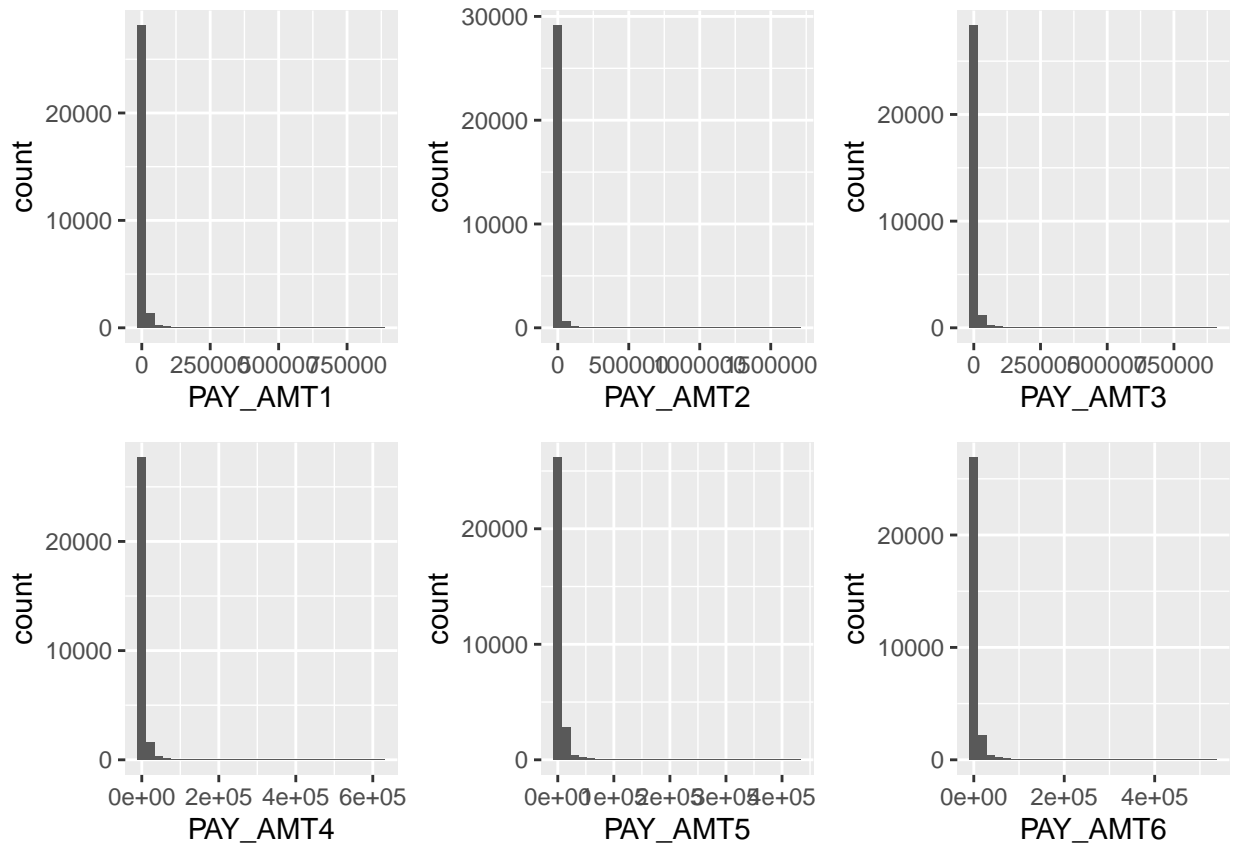


From PAY\_AMT1 to PAY\_AMT6, their structures are almost the same as are shown in following plots.

```
p1 <- ggplot(data=original_default, aes(PAY_AMT1)) +geom_histogram()
p2 <- ggplot(data=original_default, aes(PAY_AMT2)) +geom_histogram()
p3 <- ggplot(data=original_default, aes(PAY_AMT3)) +geom_histogram()
p4 <- ggplot(data=original_default, aes(PAY_AMT4)) +geom_histogram()
p5 <- ggplot(data=original_default, aes(PAY_AMT5)) +geom_histogram()
p6 <- ggplot(data=original_default, aes(PAY_AMT6)) +geom_histogram()

grid.arrange(p1,p2,p3,p4,p5,p6, nrow=2, ncol=3)
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



## Data Preparation

Remove ID

```
original_default <- original_default %>% select(-ID)
```

Categorical data, change numeric to factor. SEX, EDUCATION, MARRIAGE, PAY\_0~PAY\_6 are categorical data

```
original_default <- original_default %>%
  mutate(SEX = as.factor(SEX),
         EDUCATION = as.factor(EDUCATION),
         MARRIAGE = as.factor(MARRIAGE),
         PAY_0 = as.factor(PAY_0),
         PAY_2 = as.factor(PAY_2),
         PAY_3 = as.factor(PAY_3),
         PAY_4 = as.factor(PAY_4),
         PAY_5 = as.factor(PAY_5),
         PAY_6 = as.factor(PAY_6) )
```

Scaling. We use “scale” function to standardize predictors. Categorical data columns. we assume these can be defined as factors.

```

cat_col <- c("SEX", "EDUCATION", "MARRIAGE",
            "PAY_0", "PAY_2", "PAY_3", "PAY_4", "PAY_5", "PAY_6", "DEFAULT")

#all columns
all_col <- names(original_default)

#numerical data columns
num_col <- all_col[~which(all_col %in% cat_col)]

#scaling numerical data
original_default[num_col] <-original_default %>% select(-all_of(cat_col)) %>% scale()

```

Check the dataset.

```
str(original_default)
```

```

## tibble [30,000 x 24] (S3: tbl_df/tbl/data.frame)
## $ LIMIT_BAL: num [1:30000] -1.137 -0.366 -0.597 -0.905 -0.905 ...
## $ SEX      : Factor w/ 2 levels "1","2": 2 2 2 2 1 1 1 2 2 1 ...
## $ EDUCATION: Factor w/ 7 levels "0","1","2","3",...: 3 3 3 3 3 2 2 3 4 4 ...
## $ MARRIAGE  : Factor w/ 4 levels "0","1","2","3": 2 3 3 2 2 3 3 3 2 3 ...
## $ AGE       : num [1:30000] -1.246 -1.029 -0.161 0.164 2.334 ...
## $ PAY_0     : Factor w/ 11 levels "-2","-1","0",...: 5 2 3 3 2 3 3 3 3 1 ...
## $ PAY_2     : Factor w/ 11 levels "-2","-1","0",...: 5 5 3 3 3 3 3 2 3 1 ...
## $ PAY_3     : Factor w/ 11 levels "-2","-1","0",...: 2 3 3 3 2 3 3 2 5 1 ...
## $ PAY_4     : Factor w/ 11 levels "-2","-1","0",...: 2 3 3 3 3 3 3 3 3 1 ...
## $ PAY_5     : Factor w/ 10 levels "-2","-1","0",...: 1 3 3 3 3 3 3 3 3 2 ...
## $ PAY_6     : Factor w/ 10 levels "-2","-1","0",...: 1 4 3 3 3 3 3 2 3 2 ...
## $ BILL_AMT1: num [1:30000] -0.6425 -0.6592 -0.2986 -0.0575 -0.5786 ...
## $ BILL_AMT2: num [1:30000] -0.6474 -0.6667 -0.4939 -0.0133 -0.6113 ...
## $ BILL_AMT3: num [1:30000] -0.668 -0.6392 -0.4824 0.0328 -0.1612 ...
## $ BILL_AMT4: num [1:30000] -0.672 -0.622 -0.45 -0.232 -0.347 ...
## $ BILL_AMT5: num [1:30000] -0.663 -0.606 -0.417 -0.187 -0.348 ...
## $ BILL_AMT6: num [1:30000] -0.653 -0.598 -0.392 -0.157 -0.331 ...
## $ PAY_AMT1  : num [1:30000] -0.342 -0.342 -0.25 -0.221 -0.221 ...
## $ PAY_AMT2  : num [1:30000] -0.227 -0.214 -0.192 -0.169 1.335 ...
## $ PAY_AMT3  : num [1:30000] -0.297 -0.24 -0.24 -0.229 0.271 ...
## $ PAY_AMT4  : num [1:30000] -0.308 -0.244 -0.244 -0.238 0.266 ...
## $ PAY_AMT5  : num [1:30000] -0.314 -0.314 -0.249 -0.244 -0.269 ...
## $ PAY_AMT6  : num [1:30000] -0.2934 -0.1809 -0.0121 -0.2371 -0.2552 ...
## $ DEFAULT   : Factor w/ 2 levels "0","1": 2 2 1 1 1 1 1 1 1 1 ...

```

```
summary(original_default)
```

```

##   LIMIT_BAL   SEX   EDUCATION MARRIAGE   AGE
## Min.   : -1.2138 1:11888 0:   14  0:   54 Min.   : -1.5715
## 1st Qu.: -0.9055 2:18112 1:10585 1:13659 1st Qu.: -0.8121
## Median : -0.2118          2:14030 2:15964 Median : -0.1612
## Mean    :  0.0000          3: 4917  3:  323 Mean    :  0.0000
## 3rd Qu.:  0.5589          4:  123          3rd Qu.:  0.5982
## Max.    :  6.4164          5:  280          Max.    :  4.7207
##                                6:   51

```

```
##      PAY_0      PAY_2      PAY_3      PAY_4
## 0      :14737  0      :15730  0      :15764  0      :16455
## -1     : 5686 -1      : 6050 -1      : 5938 -1      : 5687
## 1      : 3688 2       : 3927 -2      : 4085 -2      : 4348
## -2     : 2759 -2      : 3782 2       : 3819 2       : 3159
## 2      : 2667 3       :  326 3       :  240 3       :  180
## 3      :  322 4       :   99 4       :   76 4       :   69
## (Other): 141 (Other):  86 (Other):  78 (Other): 102
##      PAY_5      PAY_6      BILL_AMT1      BILL_AMT2
## 0      :16947  0      :16286  Min.    :-2.9443  Min.    :-1.6713
## -1     : 5539 -1      : 5740  1st Qu.: -0.6473  1st Qu.: -0.6490
## -2     : 4546 -2      : 4895  Median  :-0.3917  Median  :-0.3931
## 2      : 2626 2       : 2766  Mean     : 0.0000  Mean     : 0.0000
## 3      :  178 3       :  184  3rd Qu.:  0.2155  3rd Qu.:  0.2083
## 4      :   84 4       :   49  Max.    :12.4028  Max.    :13.1334
## (Other):  80 (Other):  80
##      BILL_AMT3      BILL_AMT4      BILL_AMT5      BILL_AMT6
## Min.    :-2.9456  Min.    :-3.3150  Min.    :-2.0008  Min.    :-6.3551
## 1st Qu.: -0.6395  1st Qu.: -0.6363  1st Qu.: -0.6340  1st Qu.: -0.6316
## Median  :-0.3882  Median  :-0.3763  Median  :-0.3653  Median  :-0.3661
## Mean     : 0.0000  Mean     : 0.0000  Mean     : 0.0000  Mean     : 0.0000
## 3rd Qu.:  0.1896  3rd Qu.:  0.1748  3rd Qu.:  0.1625  3rd Qu.:  0.1734
## Max.    :23.3178  Max.    :13.1865  Max.    :14.5872  Max.    :15.4950
##
##      PAY_AMT1      PAY_AMT2      PAY_AMT3      PAY_AMT4
## Min.    :-0.3419  Min.    :-0.25699  Min.    :-0.29680  Min.    :-0.30806
## 1st Qu.: -0.2816  1st Qu.: -0.22083  1st Qu.: -0.27465  1st Qu.: -0.28916
## Median  :-0.2152  Median  :-0.16979  Median  :-0.19456  Median  :-0.21231
## Mean     : 0.0000  Mean     : 0.00000  Mean     : 0.00000  Mean     : 0.00000
## 3rd Qu.: -0.0397  3rd Qu.: -0.03998  3rd Qu.: -0.04093  3rd Qu.: -0.05188
## Max.    :52.3983  Max.    :72.84177  Max.    :50.59444  Max.    :39.33152
##
##      PAY_AMT5      PAY_AMT6      DEFAULT
## Min.    :-0.31413  Min.    :-0.29338  0:23364
## 1st Qu.: -0.29760  1st Qu.: -0.28675  1: 6636
## Median  :-0.21595  Median  :-0.20900
## Mean     : 0.00000  Mean     : 0.00000
## 3rd Qu.: -0.05026  3rd Qu.: -0.06837
## Max.    :27.60317  Max.    :29.44461
##
```

Splitting into train\_set, validation\_set, test\_set.

First we split data into test\_set, and default. Test\_set will be only used as evaluation. We use “createDataPartition” function in “caret” package. Set seed 2021.

```
set.seed(2021, sample.kind = "Rounding")
```

```
## Warning in set.seed(2021, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used
```

```
index_1 <- createDataPartition(original_default$DEFAULT, p=0.2, list=F, times=1)
test_set <- original_default[index_1,]
default <- original_default[-index_1,]
```

As we tune hyperparameters, we split default into train\_set and validation\_set. Validation set will be used when tuning models.

```
set.seed(2021, sample.kind = "Rounding")
```

```
## Warning in set.seed(2021, sample.kind = "Rounding"): non-uniform 'Rounding'  
## sampler used
```

```
index_2 <- createDataPartition(default$DEFAULT, p=0.2, list=F, times=1)  
validation_set <- default[index_2,]  
train_set <- default[-index_2,]
```

Check default ratio.

```
#train_set  
prop.table(table(train_set$DEFAULT))
```

```
##  
##      0      1  
## 0.7788311 0.2211689
```

```
#validation_set  
prop.table(table(validation_set$DEFAULT))
```

```
##  
##      0      1  
## 0.7787961 0.2212039
```

```
#test_set  
prop.table(table(test_set$DEFAULT))
```

```
##  
##      0      1  
## 0.7787035 0.2212965
```

Almost similar ratio.

## Model analysis

### 1 Baseline prediction

All predicted as non\_default make factor vectors.

```
base_pred <- factor(numeric(length(test_set$DEFAULT)), levels=c("0", "1"))
```

Confusion matrix.



```
confusionMatrix(base_pred, test_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 4673 1328
##           1    0    0
##
##           Accuracy : 0.7787
##           95% CI : (0.768, 0.7892)
##       No Information Rate : 0.7787
##       P-Value [Acc > NIR] : 0.5074
##
##           Kappa : 0
##
##  Mcnemar's Test P-Value : <2e-16
##
##           Sensitivity : 1.0000
##           Specificity : 0.0000
##       Pos Pred Value : 0.7787
##       Neg Pred Value :      NaN
##           Prevalence : 0.7787
##       Detection Rate : 0.7787
##   Detection Prevalence : 1.0000
##       Balanced Accuracy : 0.5000
##
##       'Positive' Class : 0
##
```

We need to find models which exceed these values(except sensitivity). In this model, sensitivity is 1, but specificity is 0. This means the credit company falsely give credit to a person who fail to repay a debt. The loss for the company would be huge.

evaluation method

as this is a classification problem, we calculate accuracy using confusion matrix. However, as is shown in this baseline prediction, default rate is imbalanced. As well as accuracy, we will pay attention to specificity and balanced accuracy.

## 2 Logistic regression

As this is a classification, we use logistic regression. we use “glm” function. There are 24 predictors in the train\_set. We use “step regression” to find the best logistic regression model.

Stepwise regression explanation. First we make null-model and full-model.

```
#a null model with no predictors
null_model <- glm(DEFAULT~1, data = train_set, family = binomial(link = "logit"))

#a full model using all of the potential predictors
full_model <- glm(DEFAULT~., data = train_set, family = binomial(link = "logit"))
```

Forward and backward stepwise algorithm.

```
step_md1 <- step(null_model,
  scope = list(lower = null_model, upper = full_model),
  direction = "both")
```

```
## Start: AIC=20289.81
```

```
## DEFAULT ~ 1
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

##		Df	Deviance	AIC
##	+ PAY_0	10	17383	17405
##	+ PAY_2	10	18439	18461
##	+ PAY_3	10	18834	18856
##	+ PAY_4	10	18980	19002
##	+ PAY_5	9	19077	19097
##	+ PAY_6	9	19232	19252
##	+ LIMIT_BAL	1	19768	19772
##	+ PAY_AMT2	1	20063	20067
##	+ PAY_AMT1	1	20085	20089
##	+ PAY_AMT3	1	20112	20116
##	+ PAY_AMT5	1	20164	20168
##	+ PAY_AMT4	1	20177	20181
##	+ EDUCATION	6	20169	20183
##	+ PAY_AMT6	1	20220	20224
##	+ SEX	1	20257	20261
##	+ MARRIAGE	3	20269	20277
##	+ BILL_AMT1	1	20279	20283
##	+ BILL_AMT3	1	20284	20288
##	+ BILL_AMT2	1	20284	20288
##	+ BILL_AMT4	1	20285	20289
##	<none>		20288	20290
##	+ BILL_AMT5	1	20286	20290
##	+ BILL_AMT6	1	20286	20290
##	+ AGE	1	20288	20292
##				

```
## Step: AIC=17404.59
```

```
## DEFAULT ~ PAY_0
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

##		Df	Deviance	AIC
##	+ PAY_4	10	17107	17149
##	+ PAY_5	9	17114	17154
##	+ PAY_3	10	17128	17170
##	+ PAY_6	9	17137	17177
##	+ LIMIT_BAL	1	17178	17202
##	+ PAY_2	9	17243	17283
##	+ PAY_AMT2	1	17294	17318
##	+ PAY_AMT3	1	17312	17336
##	+ PAY_AMT1	1	17321	17345
##	+ PAY_AMT5	1	17329	17353
##	+ EDUCATION	6	17320	17354

```
## + PAY_AMT4      1      17338 17362
## + PAY_AMT6      1      17352 17376
## + SEX            1      17364 17388
## + MARRIAGE       3      17366 17394
## + BILL_AMT5      1      17377 17401
## + BILL_AMT6      1      17377 17401
## + BILL_AMT4      1      17379 17403
## + BILL_AMT3      1      17379 17403
## + BILL_AMT1      1      17380 17404
## <none>           17383 17405
## + BILL_AMT2      1      17381 17405
## + AGE            1      17381 17405
## - PAY_0          10      20288 20290
```

```
##
```

```
## Step: AIC=17148.59
```

```
## DEFAULT ~ PAY_0 + PAY_4
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
##           Df Deviance   AIC
## + LIMIT_BAL  1      16957 17001
## + PAY_6      9      16990 17050
## + PAY_AMT2   1      17019 17063
## + PAY_5      9      17027 17087
## + PAY_3     10      17027 17089
## + PAY_AMT1   1      17051 17095
## + PAY_AMT5   1      17062 17106
## + PAY_AMT3   1      17067 17111
## + PAY_2      9      17052 17112
## + EDUCATION  6      17058 17112
## + PAY_AMT4   1      17075 17119
## + PAY_AMT6   1      17082 17126
## + SEX        1      17091 17135
## + MARRIAGE    3      17089 17137
## + BILL_AMT6   1      17095 17139
## + BILL_AMT5   1      17096 17140
## + BILL_AMT4   1      17099 17143
## + BILL_AMT3   1      17101 17145
## + BILL_AMT1   1      17104 17148
## + BILL_AMT2   1      17105 17149
## <none>        17107 17149
## + AGE         1      17105 17149
## - PAY_4       10      17383 17405
## - PAY_0       10      18980 19002
```

```
##
```

```
## Step: AIC=17001.21
```

```
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL
```

```
##
```

```
##           Df Deviance   AIC
## + PAY_6      9      16855 16917
## + PAY_5      9      16884 16946
## + PAY_3     10      16895 16959
## + PAY_AMT2   1      16916 16962
## + PAY_AMT1   1      16934 16980
```

```

## + PAY_2      9      16919 16981
## + EDUCATION  6      16925 16981
## + BILL_AMT2  1      16936 16982
## + MARRIAGE   3      16932 16982
## + BILL_AMT1  1      16937 16983
## + PAY_AMT5   1      16941 16987
## + SEX        1      16943 16989
## + BILL_AMT3  1      16943 16989
## + PAY_AMT3   1      16944 16990
## + BILL_AMT4  1      16946 16992
## + PAY_AMT4   1      16947 16993
## + AGE        1      16950 16996
## + BILL_AMT5  1      16950 16996
## + BILL_AMT6  1      16952 16998
## + PAY_AMT6   1      16952 16998
## <none>              16957 17001
## - LIMIT_BAL  1      17107 17149
## - PAY_4      10     17178 17202
## - PAY_0      10     18714 18738
##
## Step:  AIC=16916.88
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6
##
##           Df Deviance   AIC
## + PAY_AMT2  1      16818 16882
## + PAY_3     10     16803 16885
## + BILL_AMT2  1      16828 16892
## + BILL_AMT1  1      16830 16894
## + MARRIAGE   3      16829 16897
## + EDUCATION  6      16823 16897
## + PAY_AMT1   1      16834 16898
## + BILL_AMT3  1      16837 16901
## + PAY_2      9      16823 16903
## + BILL_AMT4  1      16840 16904
## + SEX        1      16842 16906
## + PAY_AMT3   1      16843 16907
## + PAY_AMT5   1      16844 16908
## + BILL_AMT5  1      16845 16909
## + PAY_AMT4   1      16847 16911
## + AGE        1      16847 16911
## + PAY_5      9      16831 16911
## + BILL_AMT6  1      16848 16912
## + PAY_AMT6   1      16851 16915
## <none>              16855 16917
## - PAY_6      9      16957 17001
## - PAY_4     10     16975 17017
## - LIMIT_BAL  1      16990 17050
## - PAY_0     10     18464 18506
##
## Step:  AIC=16881.94
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##		Df	Deviance	AIC
##	+ BILL_AMT3	1	16780	16846
##	+ BILL_AMT2	1	16785	16851
##	+ BILL_AMT1	1	16787	16853
##	+ BILL_AMT4	1	16792	16858
##	+ PAY_3	10	16777	16861
##	+ MARRIAGE	3	16792	16862
##	+ EDUCATION	6	16788	16864
##	+ BILL_AMT5	1	16800	16866
##	+ BILL_AMT6	1	16805	16871
##	+ SEX	1	16805	16871
##	+ PAY_AMT1	1	16806	16872
##	+ PAY_2	9	16790	16872
##	+ AGE	1	16810	16876
##	+ PAY_5	9	16794	16876
##	+ PAY_AMT5	1	16811	16877
##	+ PAY_AMT3	1	16812	16878
##	+ PAY_AMT4	1	16813	16879
##	<none>		16818	16882
##	+ PAY_AMT6	1	16816	16882
##	- PAY_AMT2	1	16855	16917
##	- PAY_6	9	16916	16962
##	- LIMIT_BAL	1	16912	16974
##	- PAY_4	10	16943	16987
##	- PAY_0	10	18402	18446

##

## Step: AIC=16845.65

## DEFAULT ~ PAY\_0 + PAY\_4 + LIMIT\_BAL + PAY\_6 + PAY\_AMT2 + BILL\_AMT3

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##		Df	Deviance	AIC
##	+ PAY_AMT1	1	16758	16826
##	+ MARRIAGE	3	16755	16827
##	+ EDUCATION	6	16749	16827
##	+ PAY_3	10	16744	16830
##	+ SEX	1	16767	16835
##	+ PAY_AMT5	1	16768	16836
##	+ PAY_AMT3	1	16770	16838
##	+ PAY_AMT4	1	16773	16841
##	+ AGE	1	16773	16841
##	+ PAY_5	9	16757	16841
##	+ PAY_2	9	16758	16842
##	+ BILL_AMT6	1	16774	16842
##	+ BILL_AMT5	1	16774	16842
##	+ PAY_AMT6	1	16776	16844
##	<none>		16780	16846
##	+ BILL_AMT4	1	16778	16846
##	+ BILL_AMT2	1	16780	16848

```

## + BILL_AMT1  1    16780 16848
## - BILL_AMT3  1    16818 16882
## - PAY_AMT2   1    16837 16901
## - PAY_6      9    16882 16930
## - PAY_4     10    16895 16941
## - LIMIT_BAL  1    16909 16973
## - PAY_0     10    18342 18388
##
## Step:  AIC=16826.41
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##      PAY_AMT1

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##           Df Deviance   AIC
## + PAY_3     10    16720 16808
## + MARRIAGE   3    16734 16808
## + EDUCATION  6    16729 16809
## + SEX        1    16746 16816
## + PAY_AMT5   1    16749 16819
## + AGE        1    16752 16822
## + PAY_AMT3   1    16752 16822
## + PAY_5      9    16736 16822
## + BILL_AMT6  1    16753 16823
## + BILL_AMT5  1    16753 16823
## + PAY_AMT4   1    16753 16823
## + PAY_AMT6   1    16756 16826
## <none>          16758 16826
## + BILL_AMT4  1    16757 16827
## + PAY_2      9    16741 16827
## + BILL_AMT1  1    16758 16828
## + BILL_AMT2  1    16758 16828
## - PAY_AMT1   1    16780 16846
## - BILL_AMT3  1    16806 16872
## - PAY_AMT2   1    16807 16873
## - PAY_6      9    16860 16910
## - PAY_4     10    16873 16921
## - LIMIT_BAL  1    16879 16945
## - PAY_0     10    18295 18343
##
## Step:  AIC=16807.47
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##      PAY_AMT1 + PAY_3
##
##           Df Deviance   AIC
## + MARRIAGE   3    16695 16789
## + EDUCATION  6    16691 16791
## + SEX        1    16707 16797
## + PAY_AMT5   1    16710 16800
## + PAY_AMT3   1    16713 16803
## + AGE        1    16713 16803
## + PAY_5      9    16697 16803

```

```

## + BILL_AMT6 1 16714 16804
## + BILL_AMT5 1 16714 16804
## + PAY_AMT4 1 16714 16804
## + PAY_AMT6 1 16717 16807
## <none> 16720 16808
## + BILL_AMT4 1 16718 16808
## + BILL_AMT1 1 16719 16809
## + BILL_AMT2 1 16719 16809
## + PAY_2 9 16712 16818
## - PAY_3 10 16758 16826
## - PAY_AMT1 1 16744 16830
## - PAY_AMT2 1 16756 16842
## - PAY_4 10 16778 16846
## - BILL_AMT3 1 16761 16847
## - PAY_6 9 16813 16883
## - LIMIT_BAL 1 16829 16915
## - PAY_0 10 18044 18112
##
## Step: AIC=16788.84
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## PAY_AMT1 + PAY_3 + MARRIAGE
##
## Df Deviance AIC
## + EDUCATION 6 16667 16773
## + SEX 1 16682 16778
## + PAY_AMT5 1 16685 16781
## + PAY_5 9 16672 16784
## + PAY_AMT3 1 16688 16784
## + BILL_AMT6 1 16689 16785
## + BILL_AMT5 1 16689 16785
## + PAY_AMT4 1 16690 16786
## + PAY_AMT6 1 16692 16788
## <none> 16695 16789
## + BILL_AMT4 1 16693 16789
## + AGE 1 16694 16790
## + BILL_AMT1 1 16694 16790
## + BILL_AMT2 1 16695 16791
## + PAY_2 9 16686 16798
## - MARRIAGE 3 16720 16808
## - PAY_3 10 16734 16808
## - PAY_AMT1 1 16719 16811
## - PAY_AMT2 1 16731 16823
## - PAY_4 10 16753 16827
## - BILL_AMT3 1 16735 16827
## - PAY_6 9 16789 16865
## - LIMIT_BAL 1 16809 16901
## - PAY_0 10 18015 18089
##
## Step: AIC=16772.76
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## PAY_AMT1 + PAY_3 + MARRIAGE + EDUCATION
##
## Df Deviance AIC
## + SEX 1 16654 16762

```

```

## + PAY_AMT5      1      16657 16765
## + PAY_5         9      16644 16768
## + BILL_AMT6     1      16661 16769
## + PAY_AMT3      1      16661 16769
## + BILL_AMT5     1      16661 16769
## + PAY_AMT4      1      16662 16770
## + PAY_AMT6      1      16664 16772
## <none>          16667 16773
## + BILL_AMT4     1      16665 16773
## + AGE           1      16666 16774
## + BILL_AMT1     1      16666 16774
## + BILL_AMT2     1      16666 16774
## + PAY_2         9      16658 16782
## - EDUCATION     6      16695 16789
## - MARRIAGE       3      16691 16791
## - PAY_3        10      16705 16791
## - PAY_AMT1      1      16690 16794
## - PAY_AMT2      1      16702 16806
## - PAY_4        10      16724 16810
## - BILL_AMT3     1      16708 16812
## - PAY_6         9      16760 16848
## - LIMIT_BAL     1      16774 16878
## - PAY_0        10      17986 18072
##
## Step:  AIC=16762.09
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          PAY_AMT1 + PAY_3 + MARRIAGE + EDUCATION + SEX
##
##           Df Deviance   AIC
## + PAY_AMT5      1      16644 16754
## + PAY_5         9      16631 16757
## + PAY_AMT3      1      16648 16758
## + BILL_AMT6     1      16648 16758
## + PAY_AMT4      1      16649 16759
## + BILL_AMT5     1      16649 16759
## + PAY_AMT6      1      16651 16761
## <none>          16654 16762
## + BILL_AMT4     1      16653 16763
## + BILL_AMT1     1      16653 16763
## + AGE           1      16654 16764
## + BILL_AMT2     1      16654 16764
## + PAY_2         9      16645 16771
## - SEX           1      16667 16773
## - EDUCATION     6      16682 16778
## - PAY_3        10      16691 16779
## - MARRIAGE       3      16678 16780
## - PAY_AMT1      1      16678 16784
## - PAY_AMT2      1      16689 16795
## - PAY_4        10      16711 16799
## - BILL_AMT3     1      16695 16801
## - PAY_6         9      16746 16836
## - LIMIT_BAL     1      16761 16867
## - PAY_0        10      17974 18062
##

```



```

## Step: AIC=16754.29
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## PAY_AMT1 + PAY_3 + MARRIAGE + EDUCATION + SEX + PAY_AMT5
##
##          Df Deviance   AIC
## + PAY_5      9    16621 16749
## + BILL_AMT5   1    16639 16751
## + PAY_AMT3    1    16640 16752
## + PAY_AMT4    1    16640 16752
## + PAY_AMT6    1    16642 16754
## + BILL_AMT6   1    16642 16754
## <none>          16644 16754
## + BILL_AMT4   1    16643 16755
## + BILL_AMT1   1    16644 16756
## + BILL_AMT2   1    16644 16756
## + AGE         1    16644 16756
## - PAY_AMT5    1    16654 16762
## + PAY_2       9    16635 16763
## - SEX         1    16657 16765
## - EDUCATION   6    16673 16771
## - PAY_3      10    16682 16772
## - MARRIAGE     3    16668 16772
## - PAY_AMT1    1    16665 16773
## - PAY_AMT2    1    16675 16783
## - PAY_4      10    16702 16792
## - BILL_AMT3   1    16689 16797
## - PAY_6       9    16733 16825
## - LIMIT_BAL   1    16743 16851
## - PAY_0      10    17963 18053
##
## Step: AIC=16748.96
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## PAY_AMT1 + PAY_3 + MARRIAGE + EDUCATION + SEX + PAY_AMT5 +
## PAY_5

```

```

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```

```

##          Df Deviance   AIC
## + PAY_AMT3    1    16615 16745
## + BILL_AMT5   1    16616 16746
## + PAY_AMT4    1    16618 16748
## + BILL_AMT6   1    16619 16749
## + PAY_AMT6    1    16619 16749
## <none>          16621 16749
## + BILL_AMT4   1    16619 16749
## + BILL_AMT2   1    16621 16751
## + AGE         1    16621 16751
## + BILL_AMT1   1    16621 16751
## - PAY_5       9    16644 16754
## - PAY_AMT5    1    16631 16757
## + PAY_2       9    16612 16758
## - SEX         1    16634 16760
## - EDUCATION   6    16649 16765
## - PAY_4      10    16658 16766

```

```
## - PAY_3      10      16659 16767
## - PAY_AMT1   1      16641 16767
## - MARRIAGE   3      16646 16768
## - PAY_AMT2   1      16652 16778
## - PAY_6      9      16668 16778
## - BILL_AMT3  1      16665 16791
## - LIMIT_BAL  1      16719 16845
## - PAY_0      10      17914 18022
##
## Step:  AIC=16744.78
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          PAY_AMT1 + PAY_3 + MARRIAGE + EDUCATION + SEX + PAY_AMT5 +
##          PAY_5 + PAY_AMT3
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
##          Df Deviance  AIC
## + BILL_AMT5  1      16612 16744
## + PAY_AMT4   1      16613 16745
## <none>                16615 16745
## + PAY_AMT6   1      16613 16745
## + BILL_AMT6  1      16613 16745
## + BILL_AMT2  1      16614 16746
## + AGE        1      16614 16746
## + BILL_AMT1  1      16615 16747
## + BILL_AMT4  1      16615 16747
## - PAY_AMT3   1      16621 16749
## - PAY_5      9      16640 16752
## - PAY_AMT5   1      16624 16752
## + PAY_2      9      16606 16754
## - SEX        1      16628 16756
## - PAY_4      10      16649 16759
## - PAY_AMT1   1      16633 16761
## - EDUCATION  6      16643 16761
## - PAY_3      10      16653 16763
## - MARRIAGE   3      16639 16763
## - PAY_AMT2   1      16642 16770
## - PAY_6      9      16661 16773
## - BILL_AMT3  1      16660 16788
## - LIMIT_BAL  1      16706 16834
## - PAY_0      10      17905 18015
##
```

```
## Step:  AIC=16743.74
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          PAY_AMT1 + PAY_3 + MARRIAGE + EDUCATION + SEX + PAY_AMT5 +
##          PAY_5 + PAY_AMT3 + BILL_AMT5
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
##          Df Deviance  AIC
## <none>                16612 16744
## + PAY_AMT6   1      16610 16744
## - BILL_AMT5  1      16615 16745
```

```
## + BILL_AMT4 1 16611 16745
## + BILL_AMT2 1 16611 16745
## + PAY_AMT4 1 16611 16745
## + AGE 1 16611 16745
## + BILL_AMT1 1 16612 16746
## + BILL_AMT6 1 16612 16746
## - PAY_AMT3 1 16616 16746
## - PAY_5 9 16636 16750
## - PAY_AMT5 1 16621 16751
## + PAY_2 9 16603 16753
## - SEX 1 16624 16754
## - PAY_4 10 16646 16758
## - EDUCATION 6 16640 16760
## - PAY_AMT1 1 16630 16760
## - MARRIAGE 3 16636 16762
## - PAY_3 10 16650 16762
## - BILL_AMT3 1 16635 16765
## - PAY_6 9 16657 16771
## - PAY_AMT2 1 16642 16772
## - LIMIT_BAL 1 16699 16829
## - PAY_0 10 17904 18016
```

Predict by using `validation_set`. First we predict probabilities and then classify them using cut-off 0.5.

```
step_prob <- predict(step_mdl, validation_set, type="response")
step_pred <- ifelse(step_prob > 0.5, 1, 0)
```

To show accuracy we use `confusionMatrix` function in `caret` library.

```
confusionMatrix(as.factor(step_pred), validation_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 3567  714
##           1  172  348
##
##           Accuracy : 0.8155
##           95% CI : (0.8042, 0.8263)
##           No Information Rate : 0.7788
##           P-Value [Acc > NIR] : 2.329e-10
##
##           Kappa : 0.3446
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9540
##           Specificity : 0.3277
##           Pos Pred Value : 0.8332
##           Neg Pred Value : 0.6692
##           Prevalence : 0.7788
##           Detection Rate : 0.7430
```

```
## Detection Prevalence : 0.8917
## Balanced Accuracy : 0.6408
##
## 'Positive' Class : 0
##
```

Make a table.

```
results <- tibble(method = "logistic regression",
  Accuracy = confusionMatrix(as.factor(step_pred), validation_set$DEFAULT)$overall[1],
  Sensitivity = confusionMatrix(as.factor(step_pred), validation_set$DEFAULT)$byClass[1],
  Specificity = confusionMatrix(as.factor(step_pred), validation_set$DEFAULT)$byClass[2],
  Balanced_Accuracy = confusionMatrix(as.factor(step_pred), validation_set$DEFAULT)$byClass[3])

results %>% knitr::kable()
```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
logistic regression	0.8154551	0.9539984	0.3276836	0.640841

### 3 Decision tree default model

Use CART classification and regression tree. Rpart ~ using default minsplit=20, cp=0.01.

```
set.seed(2021, sample.kind = "Rounding")
```

```
## Warning in set.seed(2021, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used
```

```
rpart_md1 <- rpart(DEFAULT ~ ., data = train_set)
```

Predict.

```
rpart_pred <- predict(rpart_md1, validation_set, type="class")
```

Confusion Matrix.

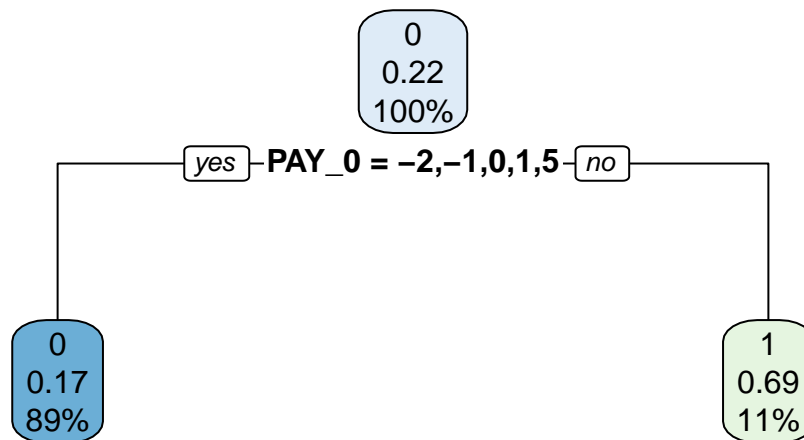
```
confusionMatrix(rpart_pred, validation_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 3597  736
##           1  142  326
##
##           Accuracy : 0.8171
##           95% CI : (0.8059, 0.828)
##           No Information Rate : 0.7788
##           P-Value [Acc > NIR] : 3.487e-11
```

```
##
##           Kappa : 0.3363
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9620
##           Specificity : 0.3070
##           Pos Pred Value : 0.8301
##           Neg Pred Value : 0.6966
##           Prevalence : 0.7788
##           Detection Rate : 0.7492
##           Detection Prevalence : 0.9025
##           Balanced Accuracy : 0.6345
##
##           'Positive' Class : 0
##
```

Draw decision tree `rpart.plot` is good function to show decision tree clearly.

```
rpart.plot(rpart_md1)
```



Find used features.

```
rpart_md1$variable.importance
```

```
##      PAY_0      PAY_4      PAY_5      PAY_6      PAY_3      PAY_2
```

```
## 1000.94794 38.19276 36.20872 26.78453 25.29650 21.82443
```

This model illustrates that PAY\_0 is overwhelmingly important.

Make a table

```
results <- bind_rows(  
  results,  
  tibble(method="CART default",  
    Accuracy = confusionMatrix(rpart_pred,  
                               validation_set$DEFAULT)$overall[1],  
    Sensitivity = confusionMatrix(rpart_pred,  
                                 validation_set$DEFAULT)$byClass[1],  
    Specificity = confusionMatrix(rpart_pred,  
                                 validation_set$DEFAULT)$byClass[2],  
    Balanced_Accuracy = confusionMatrix(rpart_pred,  
                                       validation_set$DEFAULT)$byClass[11]))  
  
results %>% knitr::kable()
```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
logistic regresion	0.8154551	0.9539984	0.3276836	0.640841
CART default	0.8171214	0.9620219	0.3069680	0.634495

#### 4 Decision tree further tuning

We use “train” function in “caret” package. and tune cp. Cross validation

rpart ~tuning using smaller cp, less than 0.01

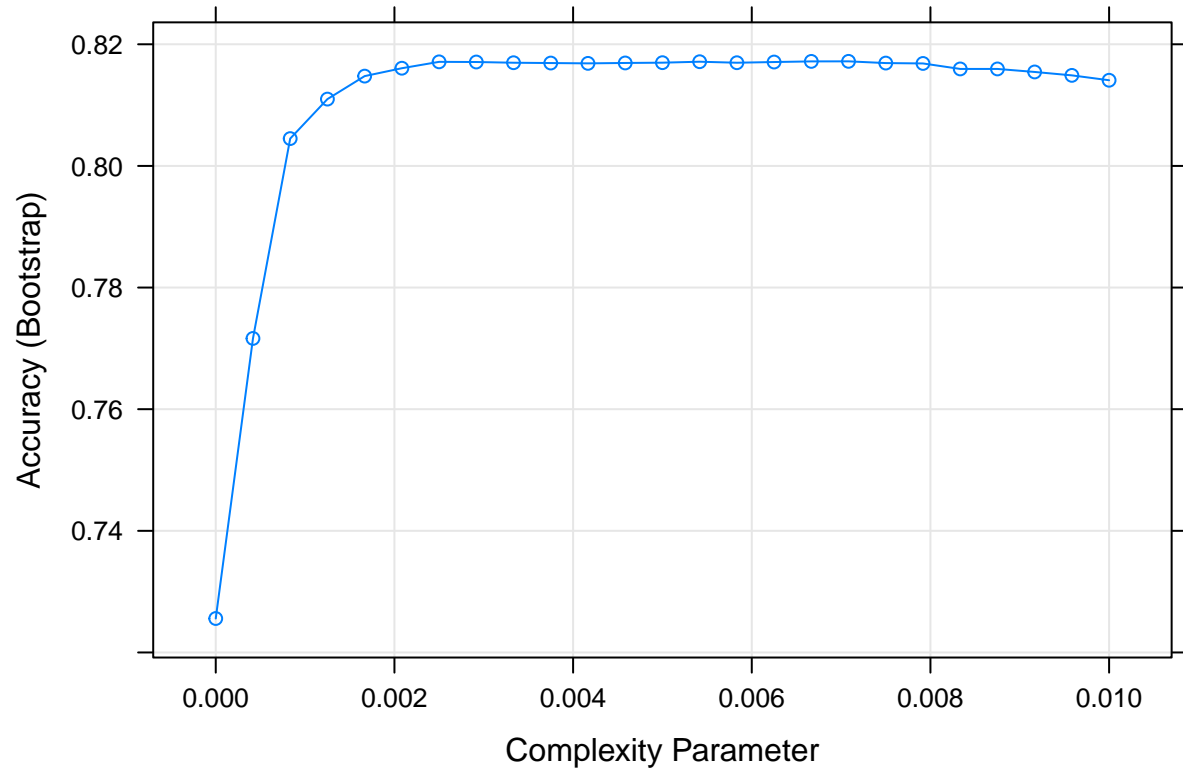
```
set.seed(2021, sample.kind = "Rounding")
```

```
## Warning in set.seed(2021, sample.kind = "Rounding"): non-uniform 'Rounding'  
## sampler used
```

```
rpart_tuned_md1 <- train(DEFAULT ~ .,  
  method = "rpart",  
  tuneGrid = data.frame(cp = seq(0, 0.01, len = 25)),  
  control = rpart.control(minsplit = 0),  
  data = train_set)
```

Plot cp.

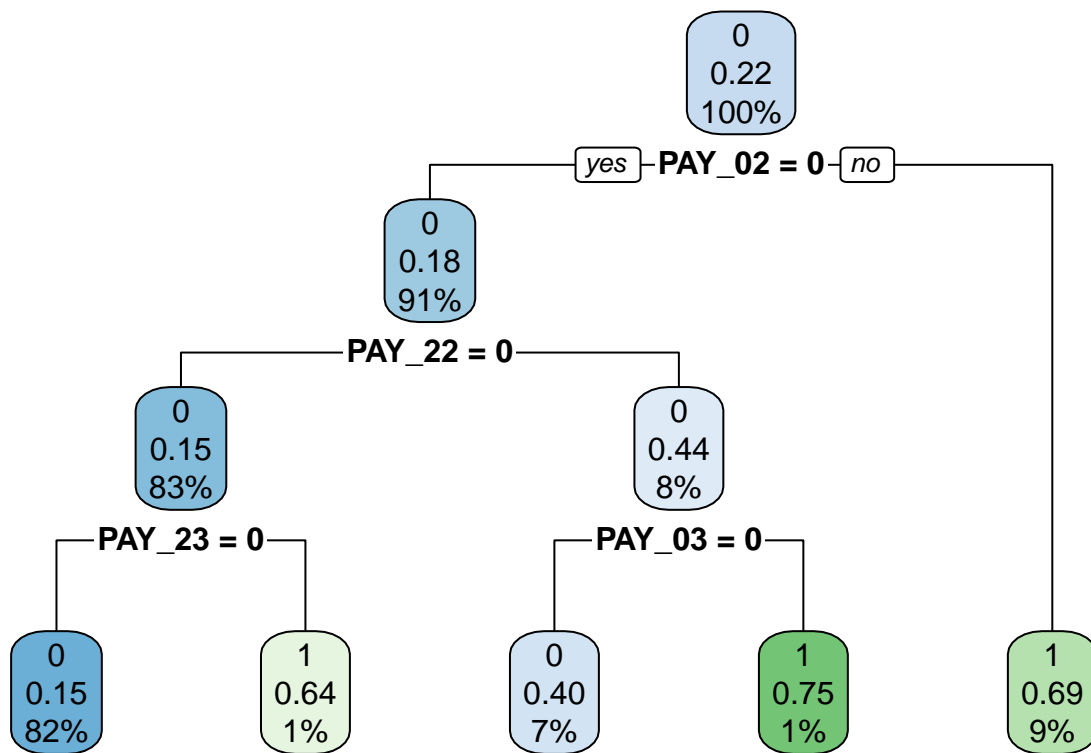
```
plot(rpart_tuned_md1)
```



```
opt_cp <- rpart_tuned_md1$bestTune
```

Draw decision tree. using rpart.plot.

```
rpart.plot(rpart_tuned_md1$finalModel)
```



Note: numeric values are scaled

Prediction.

```
rpart_tuned_pred <- predict(rpart_tuned_mdl, validation_set)
```

Confusion matrix

```
confusionMatrix(rpart_tuned_pred, validation_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 3587  730
##           1  152  332
##
##           Accuracy : 0.8163
##           95% CI : (0.805, 0.8272)
##           No Information Rate : 0.7788
##           P-Value [Acc > NIR] : 9.111e-11
##
##           Kappa : 0.3378
##
##           McNemar's Test P-Value : < 2.2e-16
##
```



```
##           Sensitivity : 0.9593
##           Specificity : 0.3126
##           Pos Pred Value : 0.8309
##           Neg Pred Value : 0.6860
##           Prevalence : 0.7788
##           Detection Rate : 0.7471
##           Detection Prevalence : 0.8992
##           Balanced Accuracy : 0.6360
##
##           'Positive' Class : 0
##
```

Make a table.

```
results <- bind_rows(
  results,
  tibble(method="CART tuned cp",
    Accuracy = confusionMatrix(rpart_tuned_pred, validation_set$DEFAULT)$overall[1],
    Sensitivity = confusionMatrix(rpart_tuned_pred, validation_set$DEFAULT)$byClass[1],
    Specificity = confusionMatrix(rpart_tuned_pred, validation_set$DEFAULT)$byClass[2],
    Balanced_Accuracy = confusionMatrix(rpart_tuned_pred, validation_set$DEFAULT)$byClass[11]) )

results %>% knitr::kable()
```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
logistic regresion	0.8154551	0.9539984	0.3276836	0.6408410
CART default	0.8171214	0.9620219	0.3069680	0.6344950
CART tuned cp	0.8162883	0.9593474	0.3126177	0.6359826

## 5 Random forest default

Using “ranger”.

```
set.seed(2021, sample.kind = "Rounding")
```

```
## Warning in set.seed(2021, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used
```

```
rf_md1 <- ranger(
  formula = DEFAULT ~ .,
  data = train_set,
  probability = F)
```

Model details.

```
rf_md1
```

```
## Ranger result
##
```

```
## Call:
## ranger(formula = DEFAULT ~ ., data = train_set, probability = F)
##
## Type:                      Classification
## Number of trees:           500
## Sample size:               19198
## Number of independent variables: 23
## Mtry:                       4
## Target node size:          1
## Variable importance mode:   none
## Splitrule:                 gini
## OOB prediction error:      18.30 %
```

Prediction.

```
rf_pred <- predict(rf_mdl, validation_set)$predictions
```

Confusion matrix

```
confusionMatrix(rf_pred, validation_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 3564  698
##           1  175  364
##
##           Accuracy : 0.8182
##           95% CI : (0.807, 0.829)
##       No Information Rate : 0.7788
##       P-Value [Acc > NIR] : 1.018e-11
##
##           Kappa : 0.3593
##
##  McNemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9532
##           Specificity : 0.3427
##           Pos Pred Value : 0.8362
##           Neg Pred Value : 0.6753
##           Prevalence : 0.7788
##           Detection Rate : 0.7423
##       Detection Prevalence : 0.8877
##           Balanced Accuracy : 0.6480
##
##           'Positive' Class : 0
##
```

Make a table.

```

results <- bind_rows(
  results,
  tibble(method="random forest default",
    Accuracy = confusionMatrix(rf_pred, validation_set$DEFAULT)$overall[1],
    Sensitivity = confusionMatrix(rf_pred, validation_set$DEFAULT)$byClass[1],
    Specificity = confusionMatrix(rf_pred, validation_set$DEFAULT)$byClass[2],
    Balanced_Accuracy = confusionMatrix(rf_pred, validation_set$DEFAULT)$byClass[11]) )

results %>% knitr::kable()

```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
logistic regression	0.8154551	0.9539984	0.3276836	0.6408410
CART default	0.8171214	0.9620219	0.3069680	0.6344950
CART tuned cp	0.8162883	0.9593474	0.3126177	0.6359826
random forest default	0.8181629	0.9531960	0.3427495	0.6479728

## 6 Random forest cross validation

Grid search

```
modelLookup("ranger")
```

```

##      model      parameter      label forReg forClass probModel
## 1 ranger      mtry #Randomly Selected Predictors    TRUE    TRUE    TRUE
## 2 ranger      splitrule      Splitting Rule    TRUE    TRUE    TRUE
## 3 ranger min.node.size      Minimal Node Size    TRUE    TRUE    TRUE

```

Make a model.

```
set.seed(2021, sample.kind = "Rounding")
```

```

## Warning in set.seed(2021, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used

```

```

rf_cv_md1 <- train( DEFAULT~ .,
  data = train_set,
  method = 'ranger',
  metric = 'Accuracy',
  num.trees = 1000,
  tuneGrid = expand.grid(
    mtry = 3:10, splitrule = 'gini', min.node.size = 1),
  trControl = trainControl(method = 'cv', number = 5))

```

```

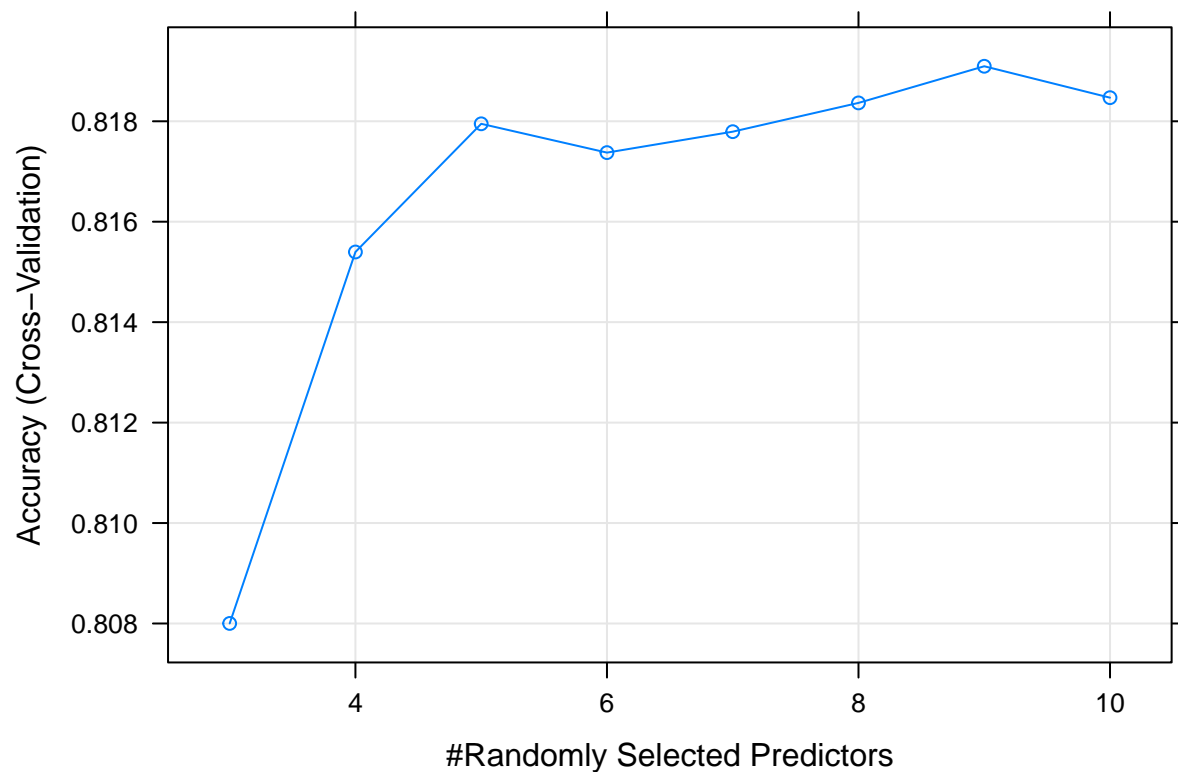
## Growing trees.. Progress: 92%. Estimated remaining time: 2 seconds.
## Growing trees.. Progress: 81%. Estimated remaining time: 7 seconds.
## Growing trees.. Progress: 79%. Estimated remaining time: 8 seconds.
## Growing trees.. Progress: 96%. Estimated remaining time: 1 seconds.
## Growing trees.. Progress: 89%. Estimated remaining time: 4 seconds.
## Growing trees.. Progress: 78%. Estimated remaining time: 8 seconds.

```

```
## Growing trees.. Progress: 100%. Estimated remaining time: 0 seconds.
## Growing trees.. Progress: 88%. Estimated remaining time: 4 seconds.
## Growing trees.. Progress: 77%. Estimated remaining time: 9 seconds.
## Growing trees.. Progress: 85%. Estimated remaining time: 5 seconds.
## Growing trees.. Progress: 81%. Estimated remaining time: 7 seconds.
## Growing trees.. Progress: 74%. Estimated remaining time: 10 seconds.
## Growing trees.. Progress: 97%. Estimated remaining time: 0 seconds.
## Growing trees.. Progress: 86%. Estimated remaining time: 5 seconds.
## Growing trees.. Progress: 76%. Estimated remaining time: 10 seconds.
## Growing trees.. Progress: 65%. Estimated remaining time: 16 seconds.
```

Plot.

```
plot(rf_cv_md1)
```



Prediction.

```
rf_cv_pred <- predict(rf_cv_md1, validation_set)
```

Confusion Matrix

```
confusionMatrix(rf_cv_pred, validation_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
```

```
##           Reference
## Prediction    0    1
##           0 3573  712
##           1  166  350
##
##           Accuracy : 0.8171
##           95% CI : (0.8059, 0.828)
##           No Information Rate : 0.7788
##           P-Value [Acc > NIR] : 3.487e-11
##
##           Kappa : 0.3495
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9556
##           Specificity : 0.3296
##           Pos Pred Value : 0.8338
##           Neg Pred Value : 0.6783
##           Prevalence : 0.7788
##           Detection Rate : 0.7442
##           Detection Prevalence : 0.8925
##           Balanced Accuracy : 0.6426
##
##           'Positive' Class : 0
##
```

Make a table.

```
results <- bind_rows( results,
                      tibble(
                        method="random forest tuned ",
                        Accuracy = confusionMatrix(rf_cv_pred, validation_set$DEFAULT)$overall[1],
                        Sensitivity = confusionMatrix(rf_cv_pred, validation_set$DEFAULT)$byClass[1],
                        Specificity = confusionMatrix(rf_cv_pred, validation_set$DEFAULT)$byClass[2],
                        Balanced_Accuracy= confusionMatrix(rf_cv_pred, validation_set$DEFAULT)$byClass[11]) )

results %>% knitr::kable()
```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
logistic regresion	0.8154551	0.9539984	0.3276836	0.6408410
CART default	0.8171214	0.9620219	0.3069680	0.6344950
CART tuned cp	0.8162883	0.9593474	0.3126177	0.6359826
random forest default	0.8181629	0.9531960	0.3427495	0.6479728
random forest tuned	0.8171214	0.9556031	0.3295669	0.6425850

## Evaluation

Best performance in terms of balanced accuracy is “random forest default model” Best performance in terms of accuracy is “CART default model” Then evaluate by using `test_set`.

```
final_pred_rpart <- predict(rpart_mdl, test_set, type="class")
confusionMatrix(final_pred_rpart, test_set$DEFAULT)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 4496  887
##           1  177  441
##
##           Accuracy : 0.8227
##           95% CI : (0.8128, 0.8323)
##       No Information Rate : 0.7787
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.3638
##
##  Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9621
##           Specificity : 0.3321
##       Pos Pred Value : 0.8352
##       Neg Pred Value : 0.7136
##           Prevalence : 0.7787
##       Detection Rate : 0.7492
##       Detection Prevalence : 0.8970
##       Balanced Accuracy : 0.6471
##
##       'Positive' Class : 0
##
```

```
final_pred_rf <- predict(rf_mdl, test_set)$predictions
confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass
```

```
##           Sensitivity           Specificity           Pos Pred Value
##           0.9492831           0.3554217           0.8382464
##       Neg Pred Value           Precision           Recall
##           0.6657264           0.8382464           0.9492831
##           F1           Prevalence           Detection Rate
##           0.8903161           0.7787035           0.7392101
##       Detection Prevalence       Balanced Accuracy
##           0.8818530           0.6523524
```

Make a table.

```
final_results <- tibble( method = "CART default",
                        Accuracy = confusionMatrix(final_pred_rpart, test_set$DEFAULT)$overall[1],
                        Sensitivity = confusionMatrix(final_pred_rpart, test_set$DEFAULT)$byClass[1],
                        Specificity = confusionMatrix(final_pred_rpart, test_set$DEFAULT)$byClass[2],
                        Balanced_Accuracy = confusionMatrix(final_pred_rpart, test_set$DEFAULT)$byClass[3],
                        Detection_Rate = confusionMatrix(final_pred_rf, test_set$DEFAULT)$overall[1],
                        Sensitivity = confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[1],
                        Specificity = confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[2],
                        Balanced_Accuracy = confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[3],
                        Detection_Rate = confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[4]
)

final_results <- bind_rows( final_results,
```

```

tibble( method ="Random forest default",
Accuracy =confusionMatrix(final_pred_rf, test_set$DEFAULT)$overall[1],
Sensitivity =confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[1],
Specificity =confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[2],
Balanced_Accuracy = confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass[3])

final_results %>% knitr::kable()

```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
CART default	0.8226962	0.9621228	0.3320783	0.6471006
Random forest default	0.8178637	0.9492831	0.3554217	0.6523524

## Conclusion

###