

Report on Default of Credit Card Clients Dataset

Masayoshi Sato

2021/6/26

Contents

Introduction	2
Data set and library	2
Data Exploration	2
1 Outcome	5
2 “LIMIT_BAL”	6
3 “SEX”	7
4 “EDUCATION”	8
5 “MARRIAGE”	10
6 “AGE”	12
7 “PAY”	13
8 “BILL_AMT”	16
9 “PAY_AMT”	18
Data Preparation	20
Model analysis	23
1 Baseline prediction	23
2 Logistic regression	24
3 Decision tree default model	35
4 Decision tree further tuning	37
5 Random forest default	40
6 Random forest cross validation	42
Evaluation	44
Conclusion	46

Introduction

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

Data set and library

In this paper, we use libraries as follows: tidyverse, gridExtra, caret, rpart, pROC, DataExplorer, ranger.

Data is stored in my GitHub repository. We will use the direct link from my GitHub repository.

Data Exploration

First, we need to check the downloaded dataset.

```
## 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()
## .. )
```

```
summary(original_default)
```

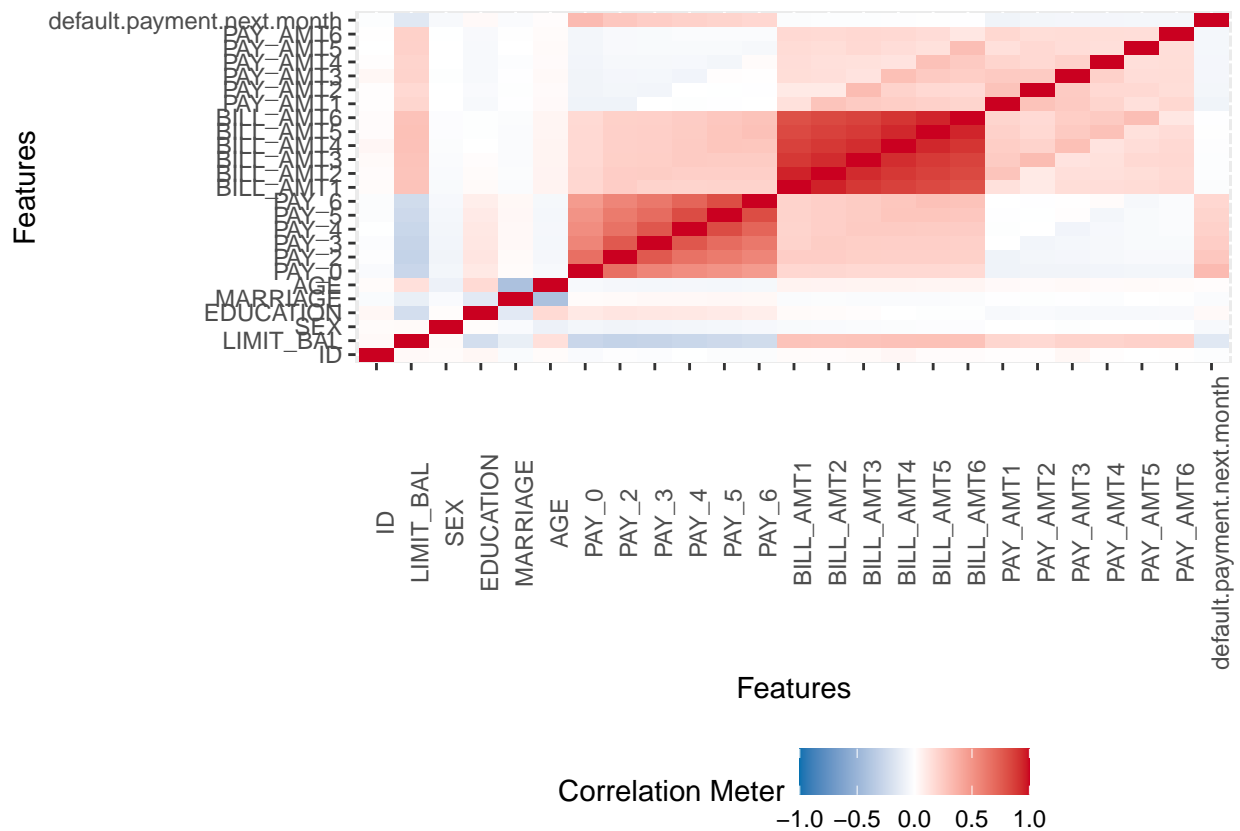
```
##          ID          LIMIT_BAL          SEX          EDUCATION
## Min.      :    1    Min.      : 10000    Min.      :1.000    Min.      :0.000
## 1st Qu.: 7501    1st Qu.:  50000    1st Qu.:1.000    1st Qu.:1.000
## Median :15000    Median : 140000    Median :2.000    Median :2.000
## Mean     :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
```

No NAs.

Correlation.

```
plot_correlation(original_default)
```



1 Outcome

Kaggle's data explanation says

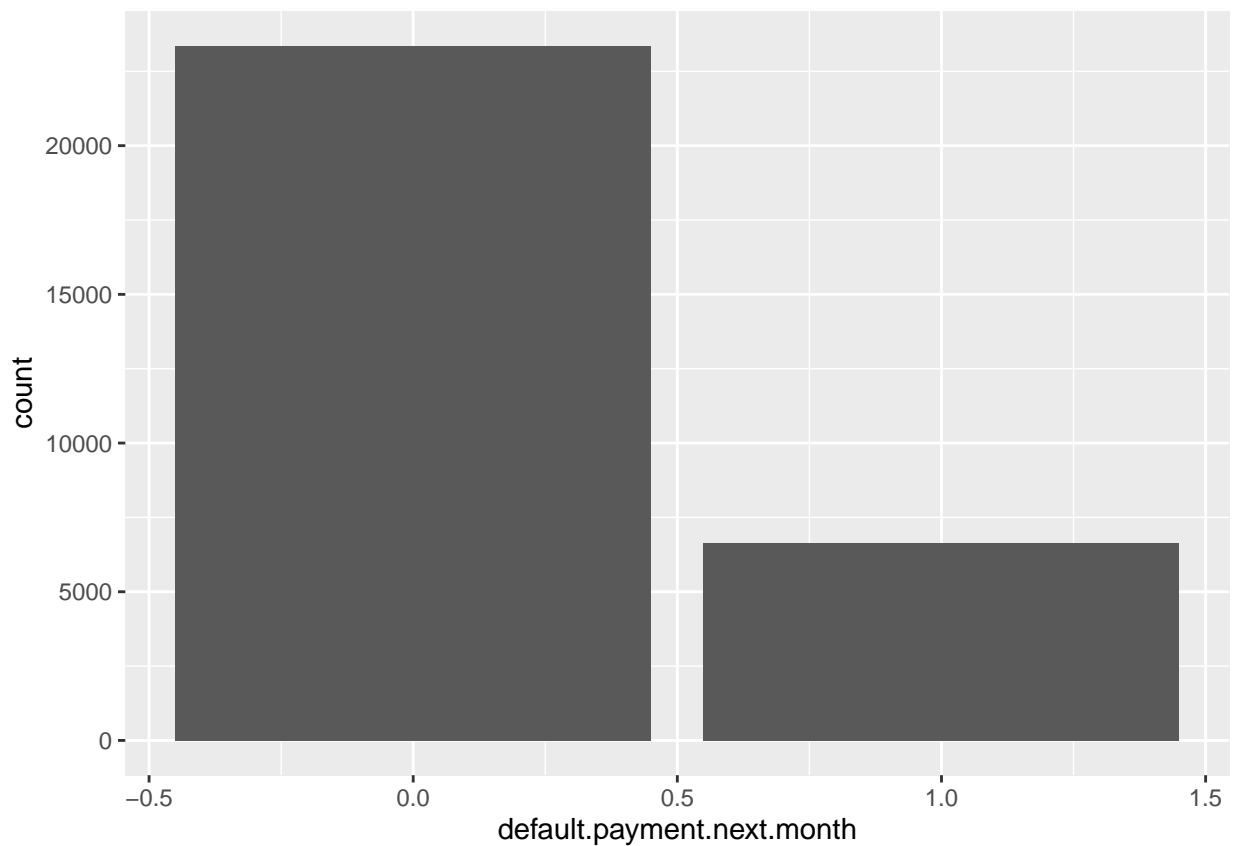
Default payment, 1=yes, 0=no

```
summary(original_default$default.payment.next.month)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.0000  0.0000  0.0000  0.2212  0.0000  1.0000
```

Show distribution graph.

```
ggplot(data=original_default, aes(default.payment.next.month)) +geom_bar()
```



Show proportion of 0,1

```
prop.table(table(original_default$default.payment.next.month))
```

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

Change name of “default.payment.next.month”

```
n <-which(names(original_default)=="default.payment.next.month")
names(original_default)[n] <- "DEFAULT"
```

Change outcome into factor

```
original_default$DEFAULT <- as.factor(original_default$DEFAULT)
```

2 “LIMIT_BAL”

Kaggle’s data explanation says

Amount of given credit in NT dollars (includes individual and family/supplementary credit

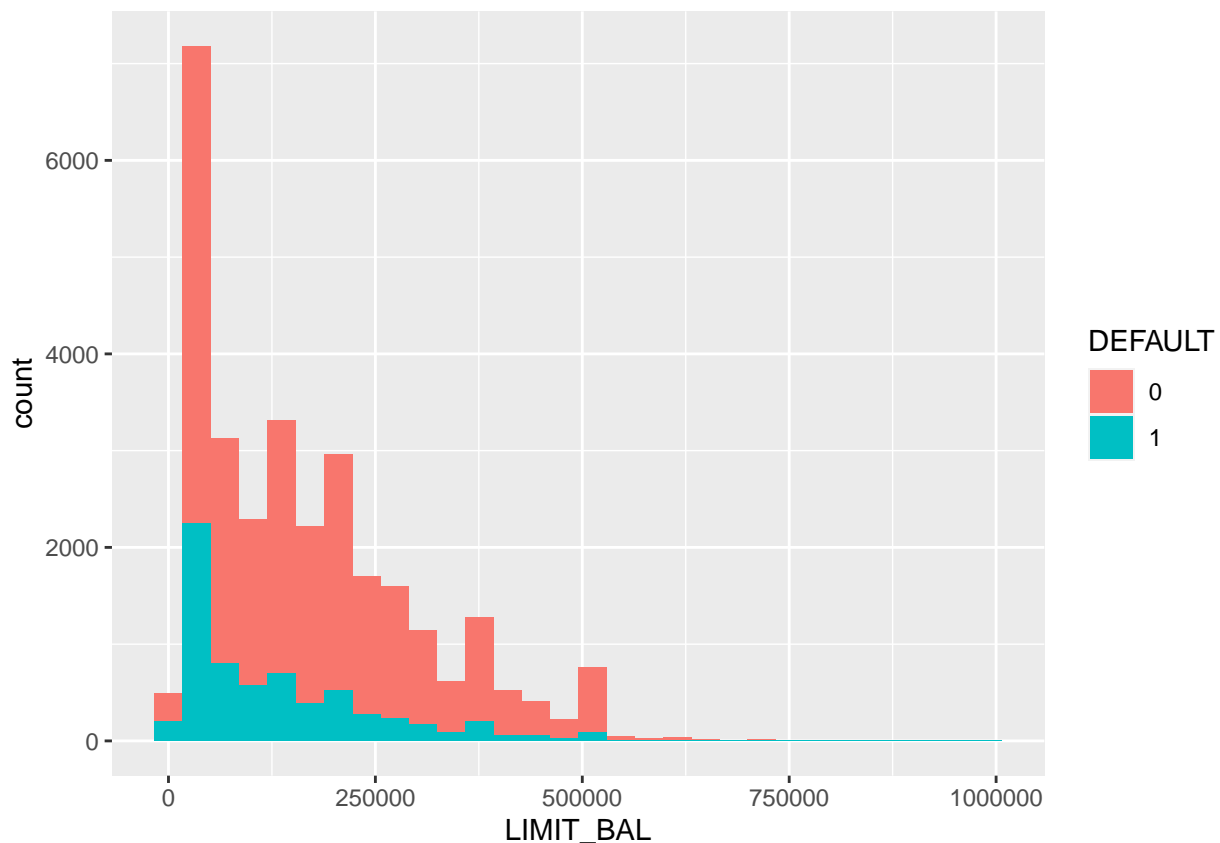
```
summary(original_default$LIMIT_BAL)
```

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

numeric data

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

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



Distribution is skewed right

3 “SEX”

Kaggle’s data explanation says

1=male, 2=female

```
summary(original_default$SEX)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1.000   1.000   2.000   1.604   2.000   2.000
```

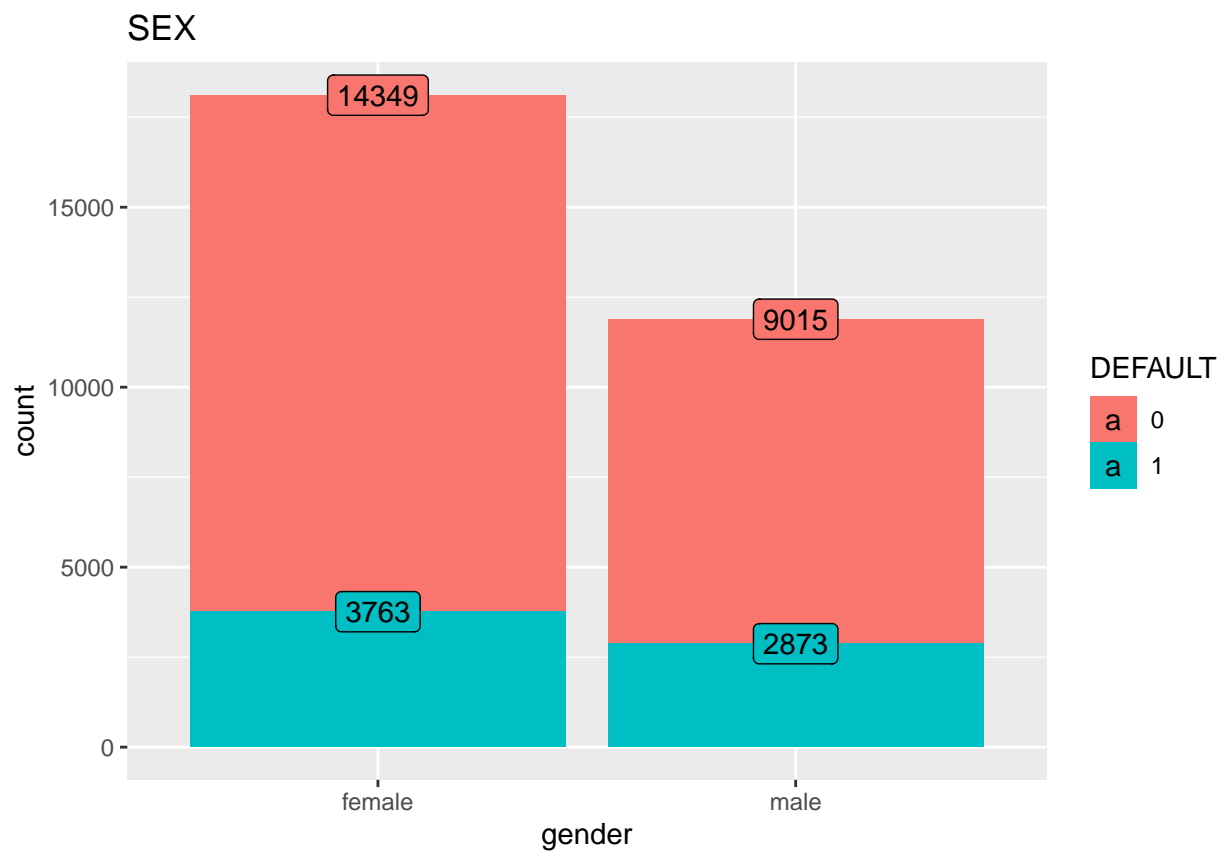
```
unique(original_default$SEX)
```

```
## [1] 2 1
```

categorical data. to make a plot, introducing new character vector.

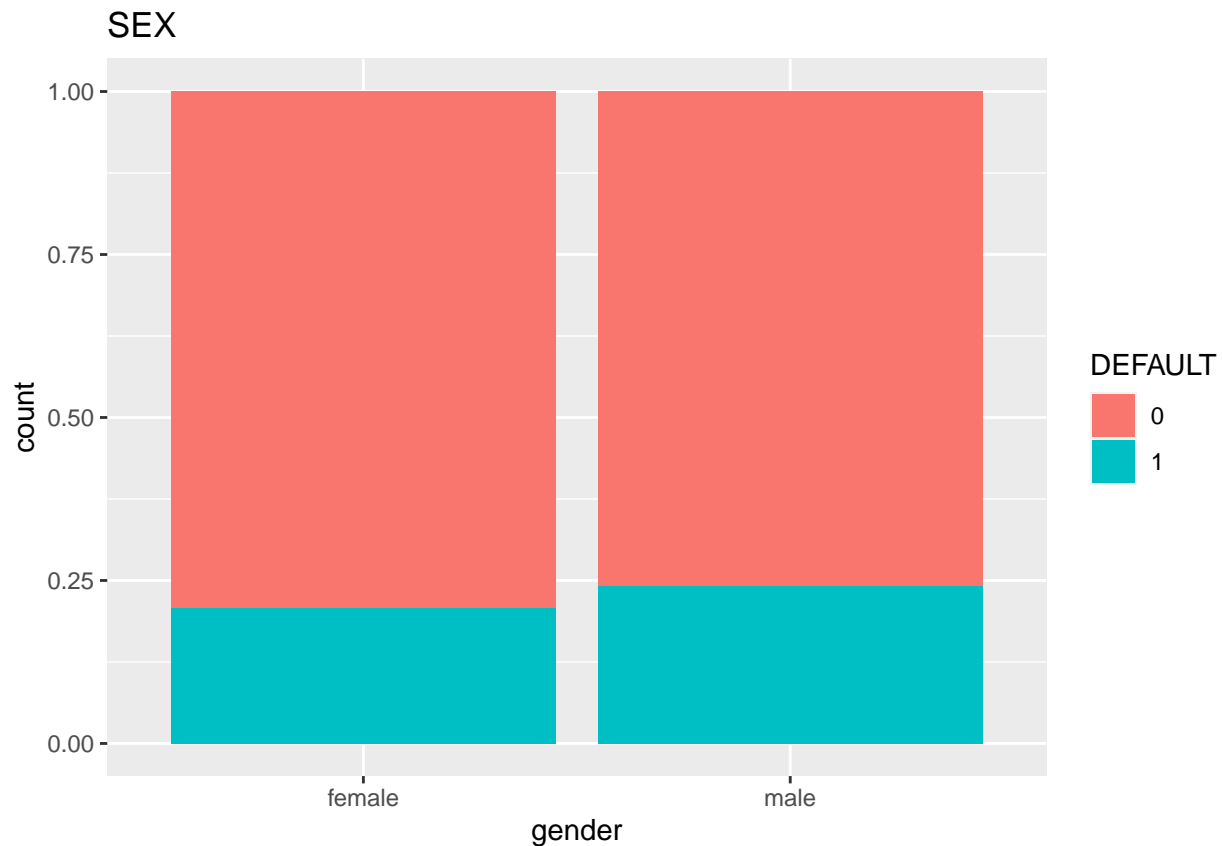
```
gender <- ifelse(original_default$SEX == 1, "male", "female")
```

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



To make stacked bar graph.

```
original_default %>% ggplot(aes(x=gender, fill= DEFAULT)) +
  geom_bar(position="fill") +
  ggtitle("SEX")
```



There seemed to be little difference between genders.

4 “EDUCATION”

Kaggle’s data explanation says;

1=graduate school, 2=university, 3=high school, 4=others, 5=unknown, 6=unknown

```
summary(original_default$EDUCATION)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.000   1.000   2.000   1.853   2.000   6.000
```

```
unique(original_default$EDUCATION)
```

```
## [1] 2 1 3 5 4 6 0
```

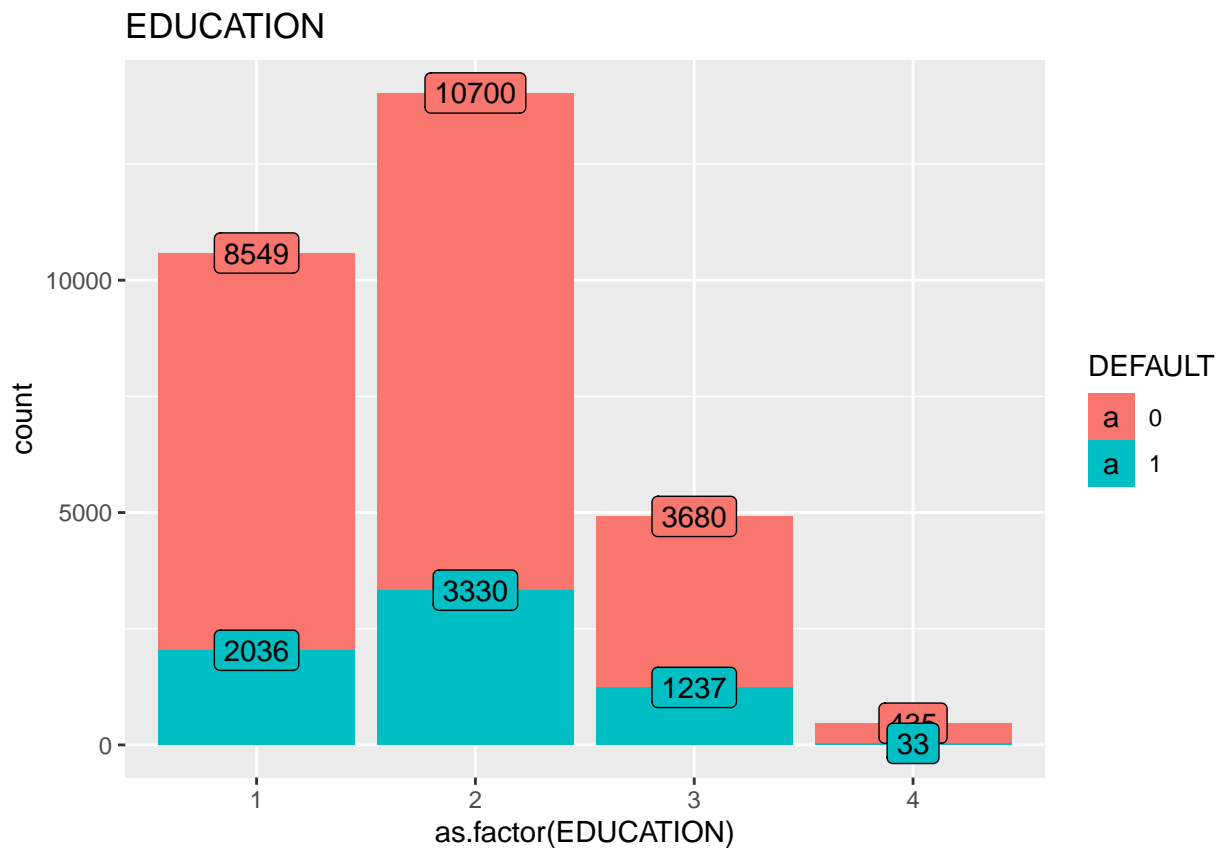
categorical data.

0 is not defined. 0,5 and 6 can be included into 4


```
original_default$EDUCATION <- ifelse( original_default$EDUCATION== 0|
                                       original_default$EDUCATION == 5|
                                       original_default$EDUCATION == 6, 4,
                                       original_default$EDUCATION)
```

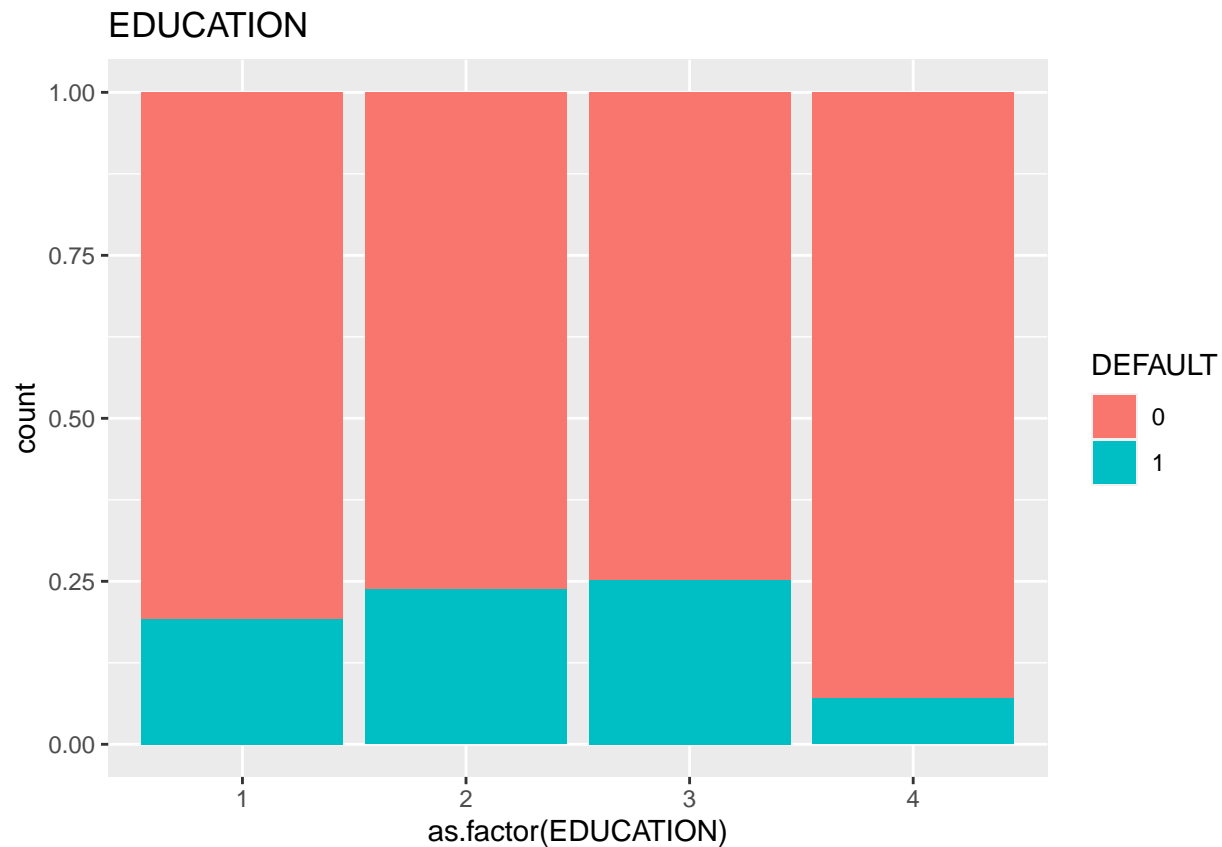
Plot.

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



Stacked bar graph.

```
original_default %>% ggplot(aes(x=as.factor(EDUCATION), fill= DEFAULT)) +
  geom_bar(position="fill") +
  ggtitle("EDUCATION")
```



4 is the smallest in terms of default rate. but its numbers are very small.

5 “MARRIAGE”

Kaggle’s data explanation says;

marital status. 1=married, 2=single, 3=others.

```
summary(original_default$ MARRIAGE)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.000   1.000   2.000   1.552   2.000   3.000
```

```
unique(original_default$ MARRIAGE)
```

```
## [1] 1 2 3 0
```

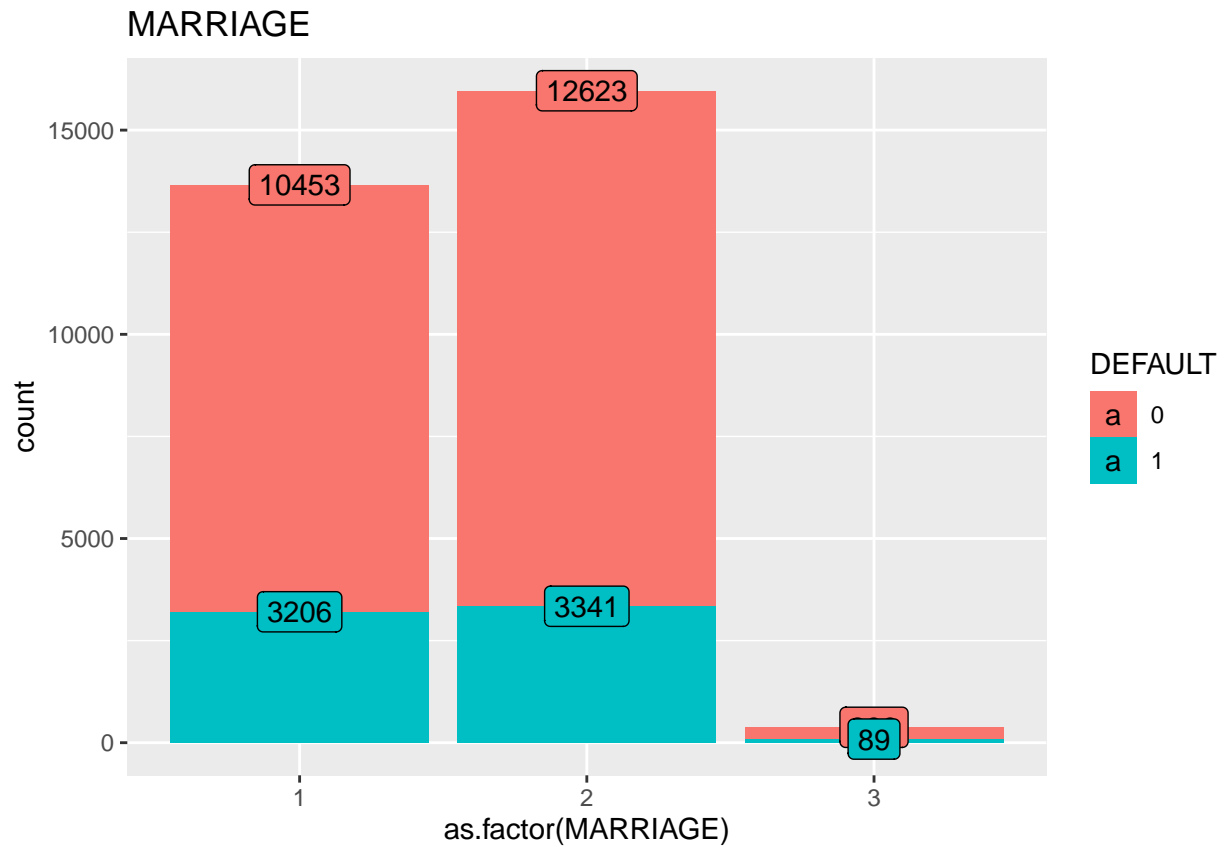
categorical data

0 is not defined. 0 can be included in 3.

```
original_default$MARRIAGE <- ifelse(original_default$MARRIAGE== 0, 3,
                                     original_default$MARRIAGE)
```

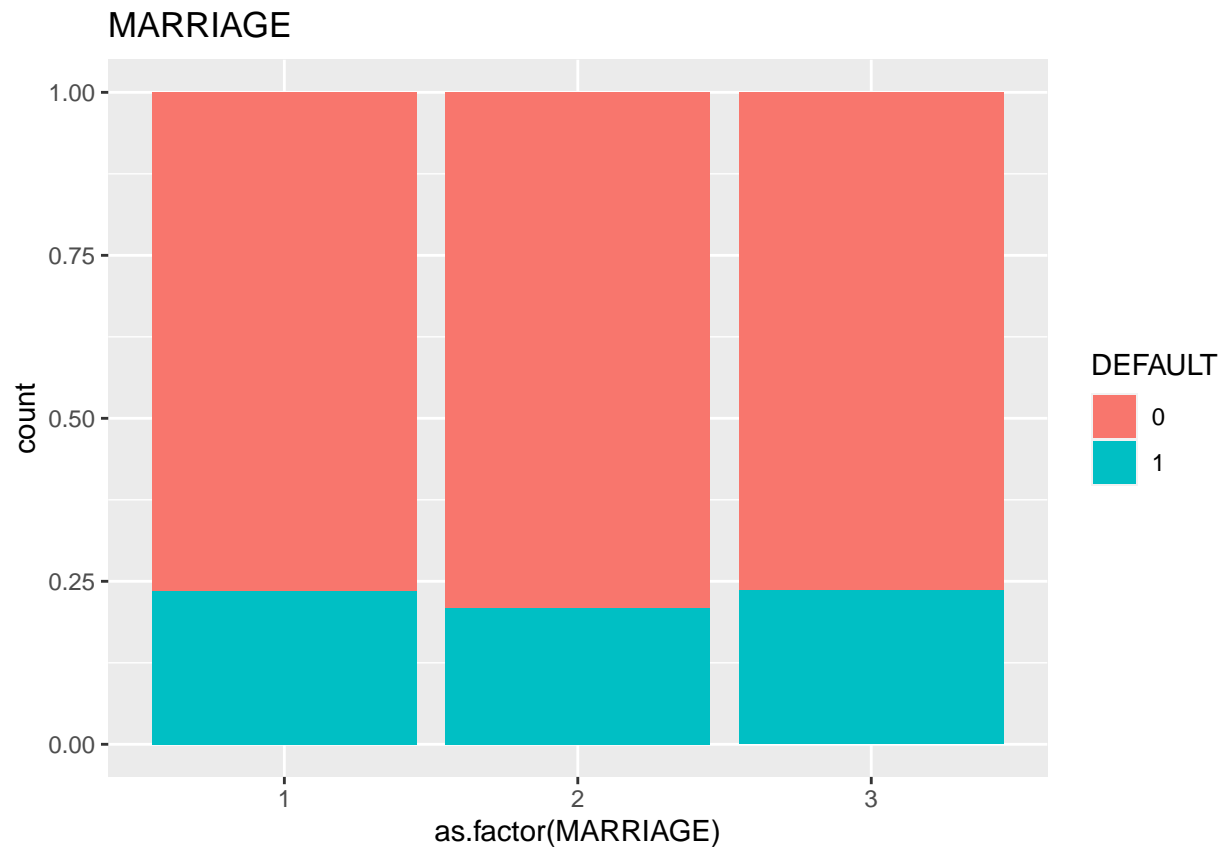
Plot.

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



Stack bar graph

```
original_default %>% ggplot(aes(x=as.factor(MARRIAGE), fill= DEFAULT)) +
  geom_bar(position="fill") +
  ggtitle("MARRIAGE")
```



There seems to be little difference among the groups.

6 “AGE”

```
summary(original_default$AGE)
```

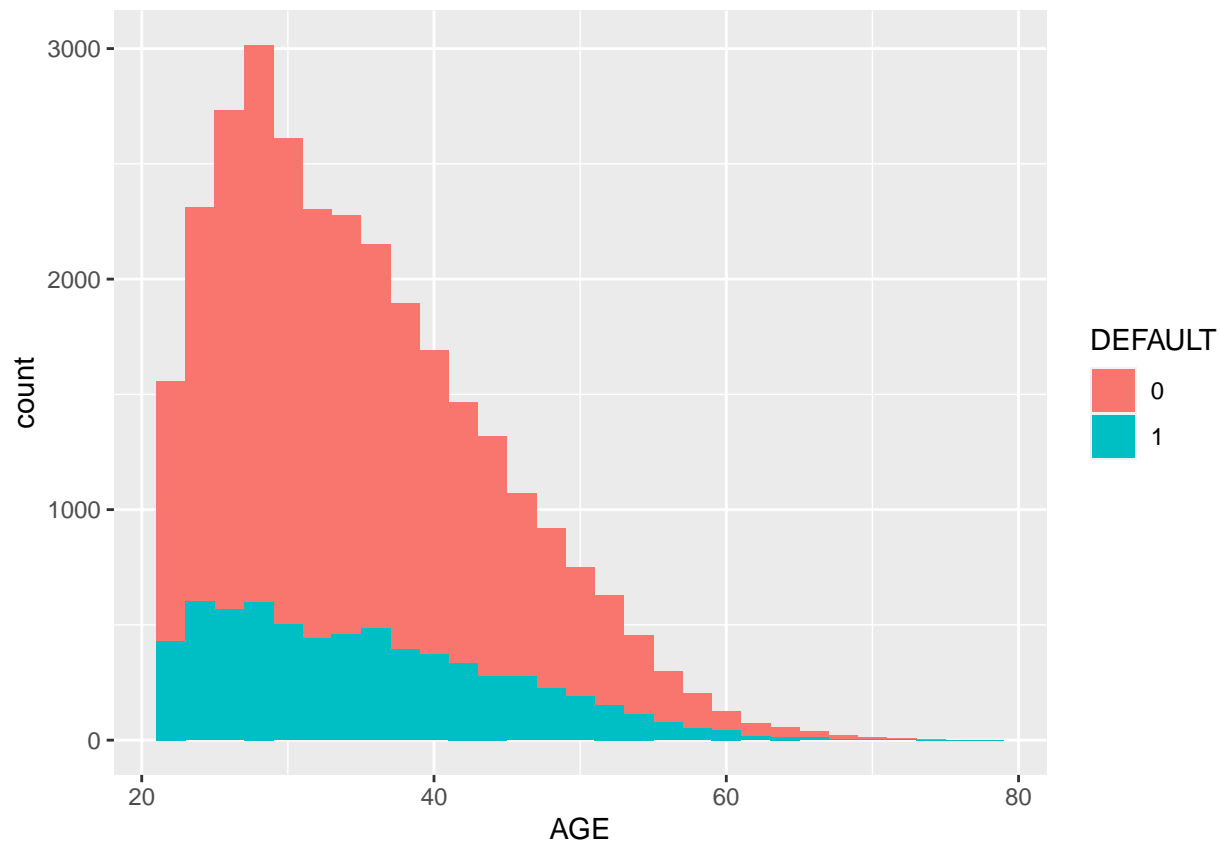
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  21.00   28.00   34.00   35.49   41.00   79.00
```

numeric data

Plot.

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

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



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
```

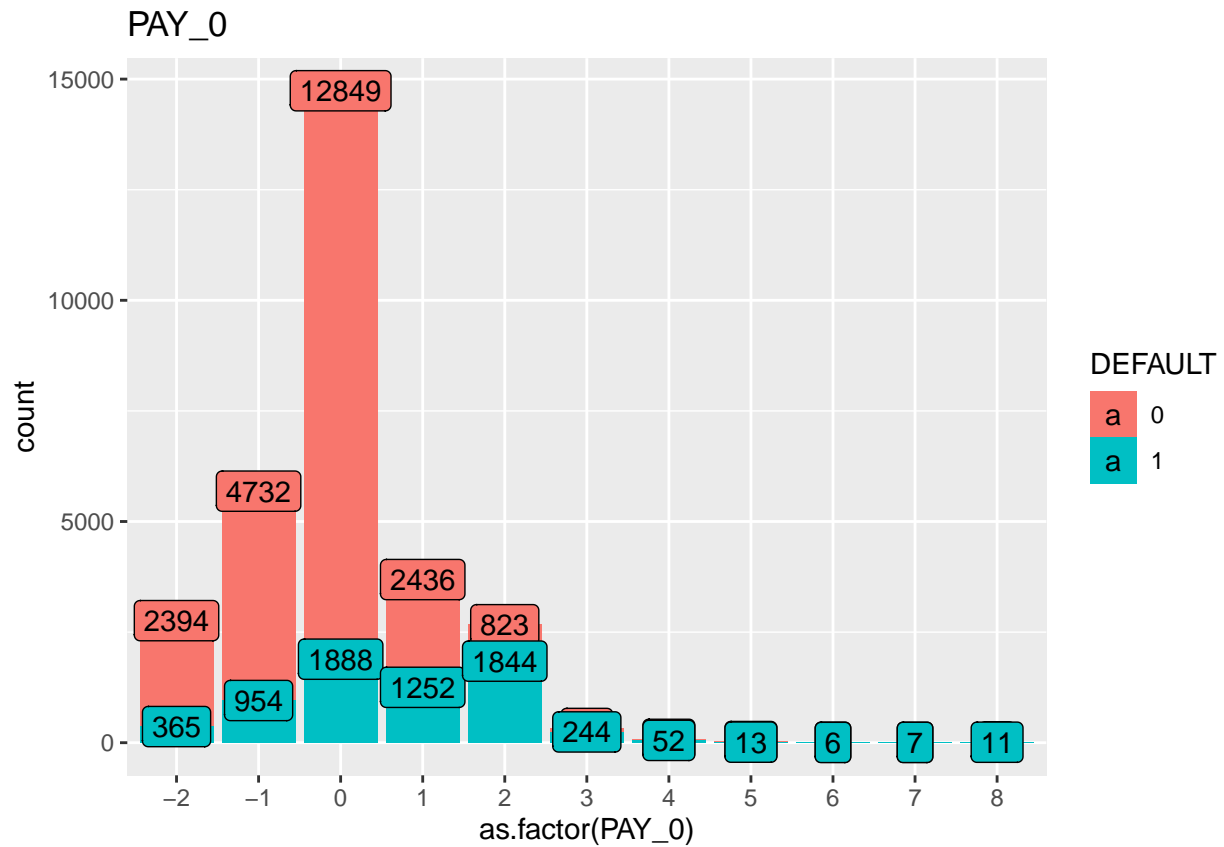
```
unique(original_default$PAY_0)
```

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

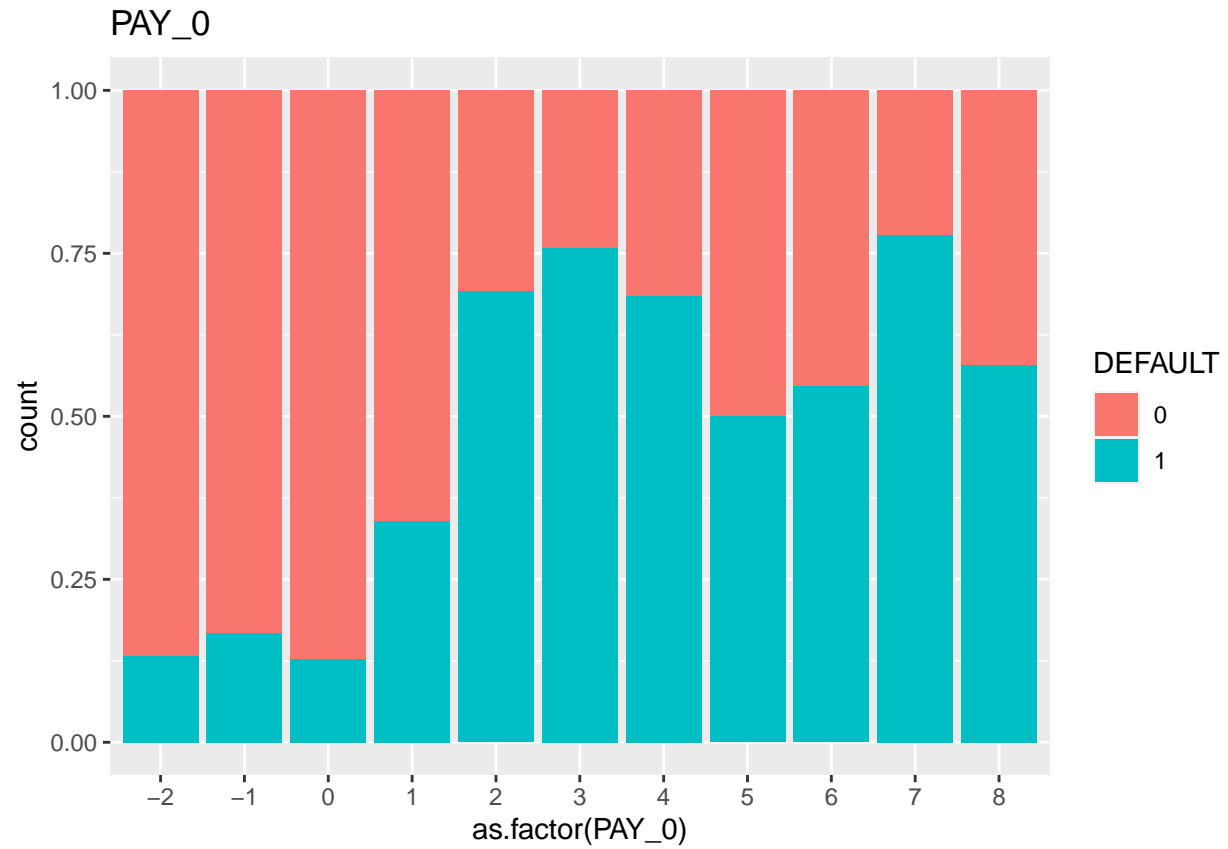
They are categorical data.

Plot.

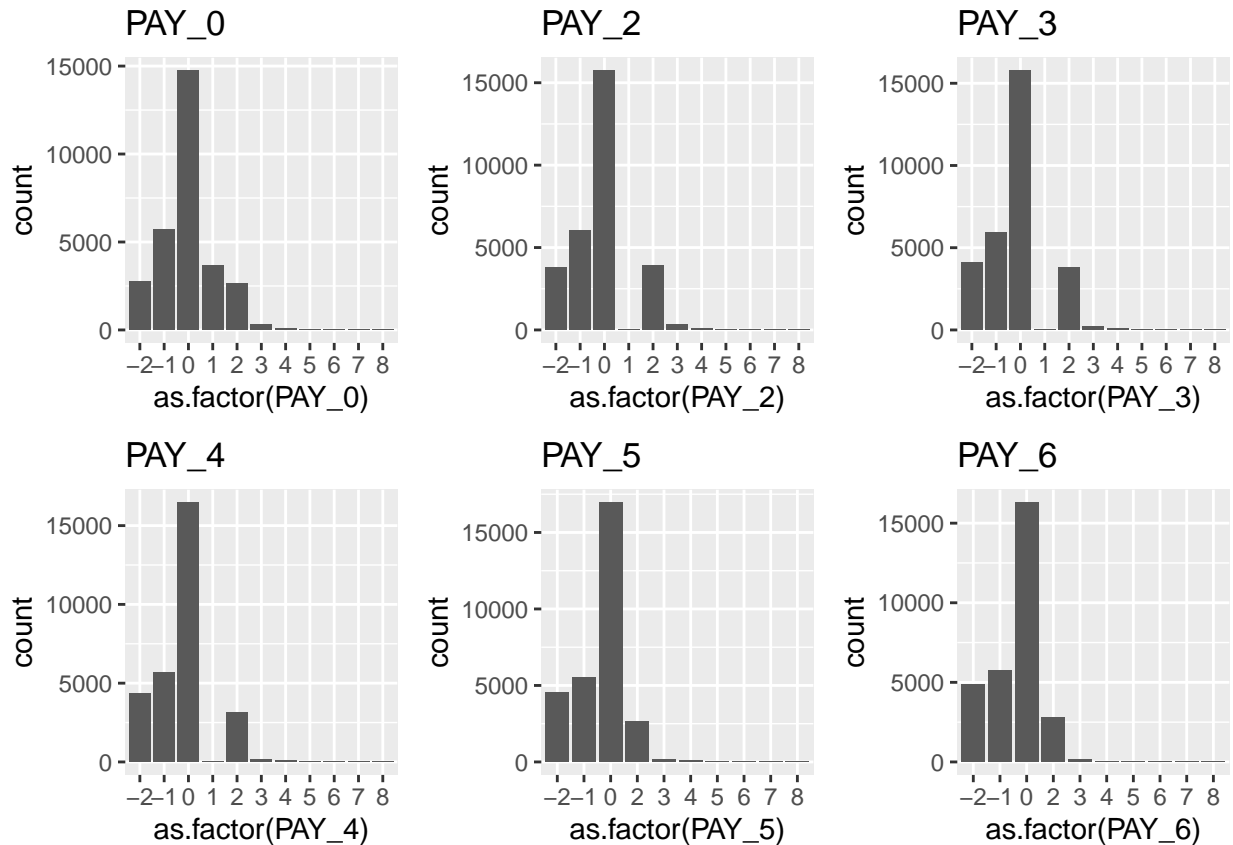
```
original_default %>% ggplot(aes(x=as.factor(PAY_0), fill= DEFAULT)) +
  geom_bar() +
  ggtitle("PAY_0")+
  stat_count(aes(label = ..count..), geom = "label")# illustrate numbers
```



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)
```

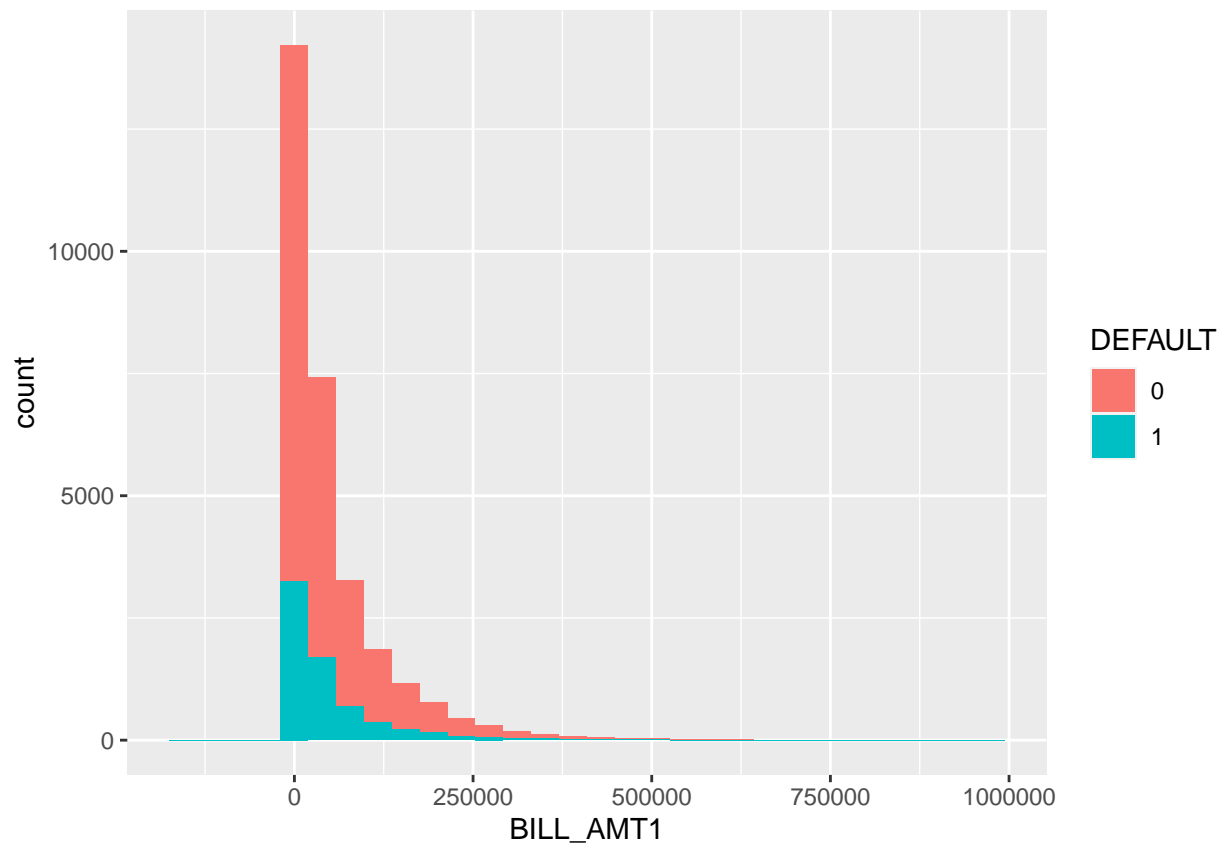
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -165580   3559    22382   51223   67091   964511
```

These are numerical data.

Here is BILL_AMT1’s plot.

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

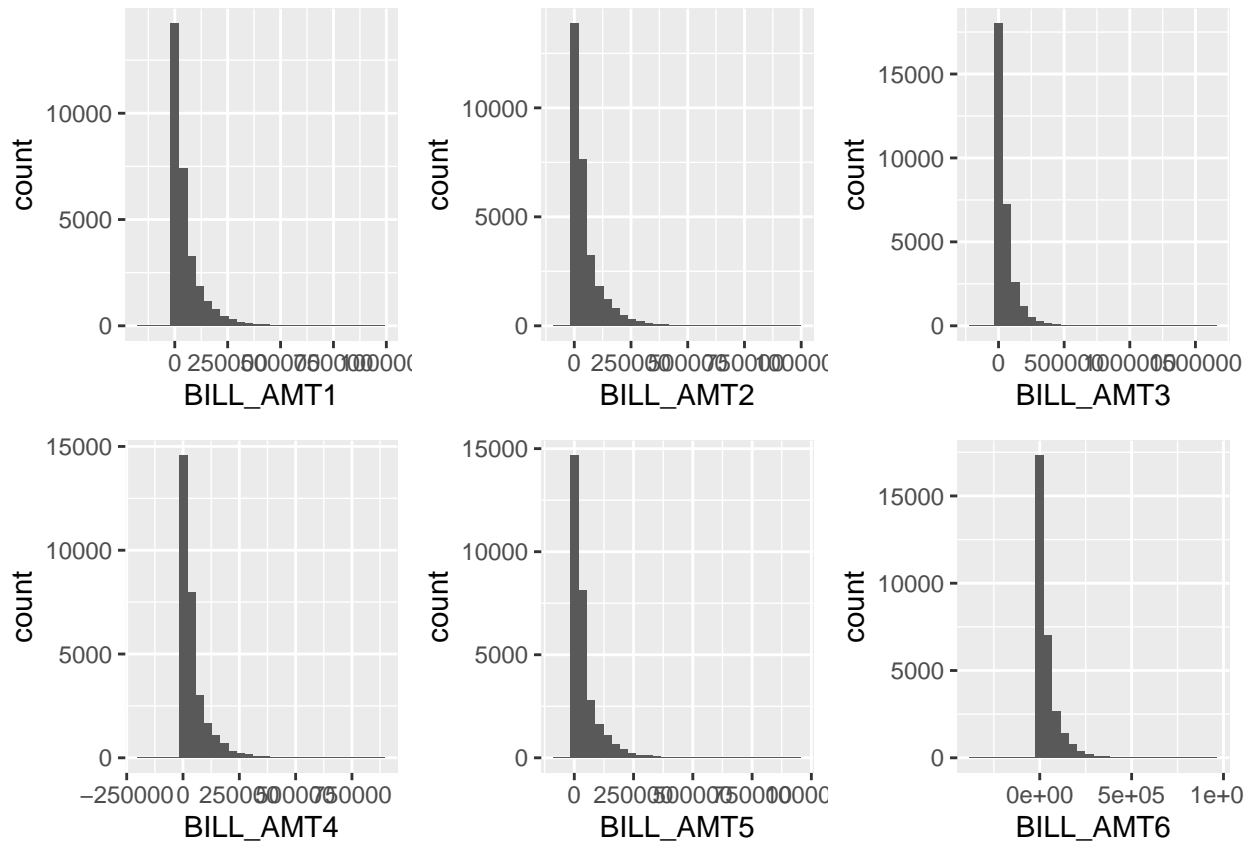
```
## ‘stat_bin()’ using ‘bins = 30’. Pick better value with ‘binwidth’.
```

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)
```

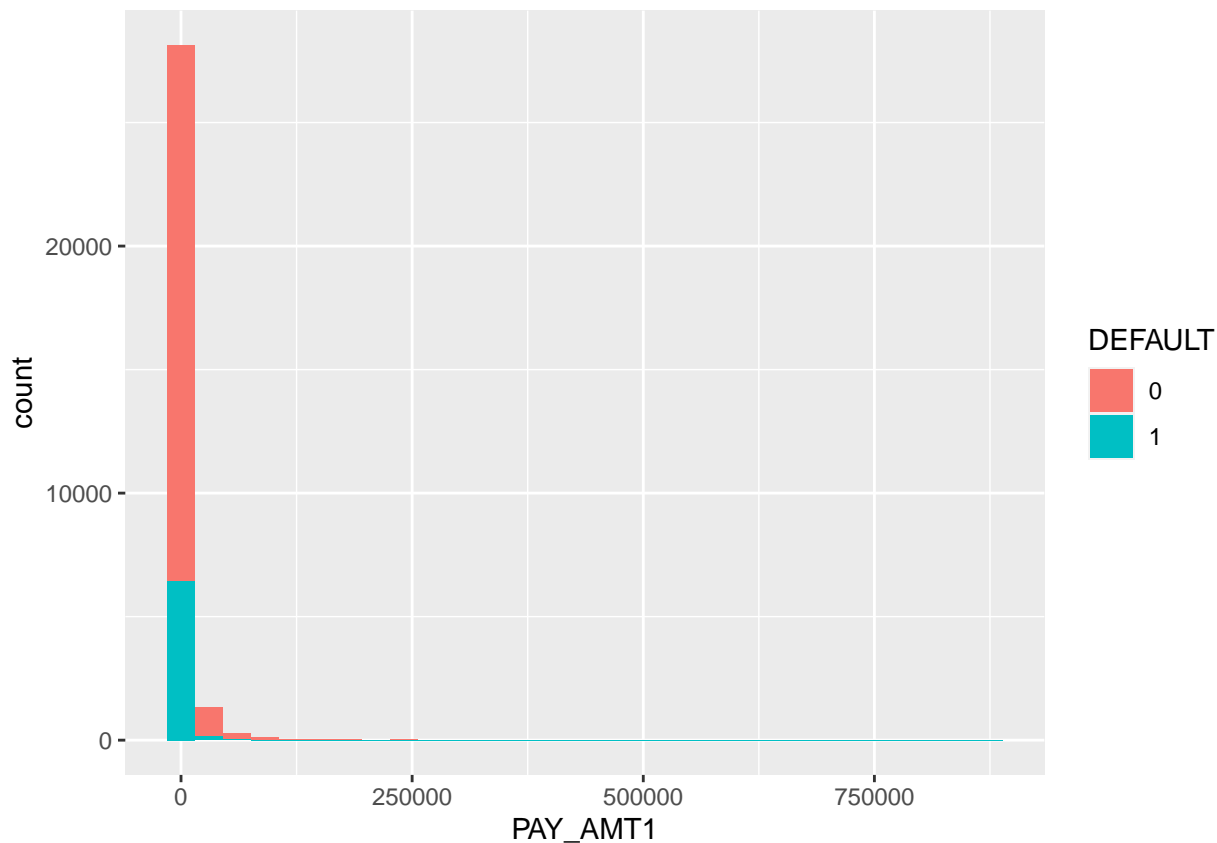
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         0    1000    2100   5664   5006 873552
```

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’.
```

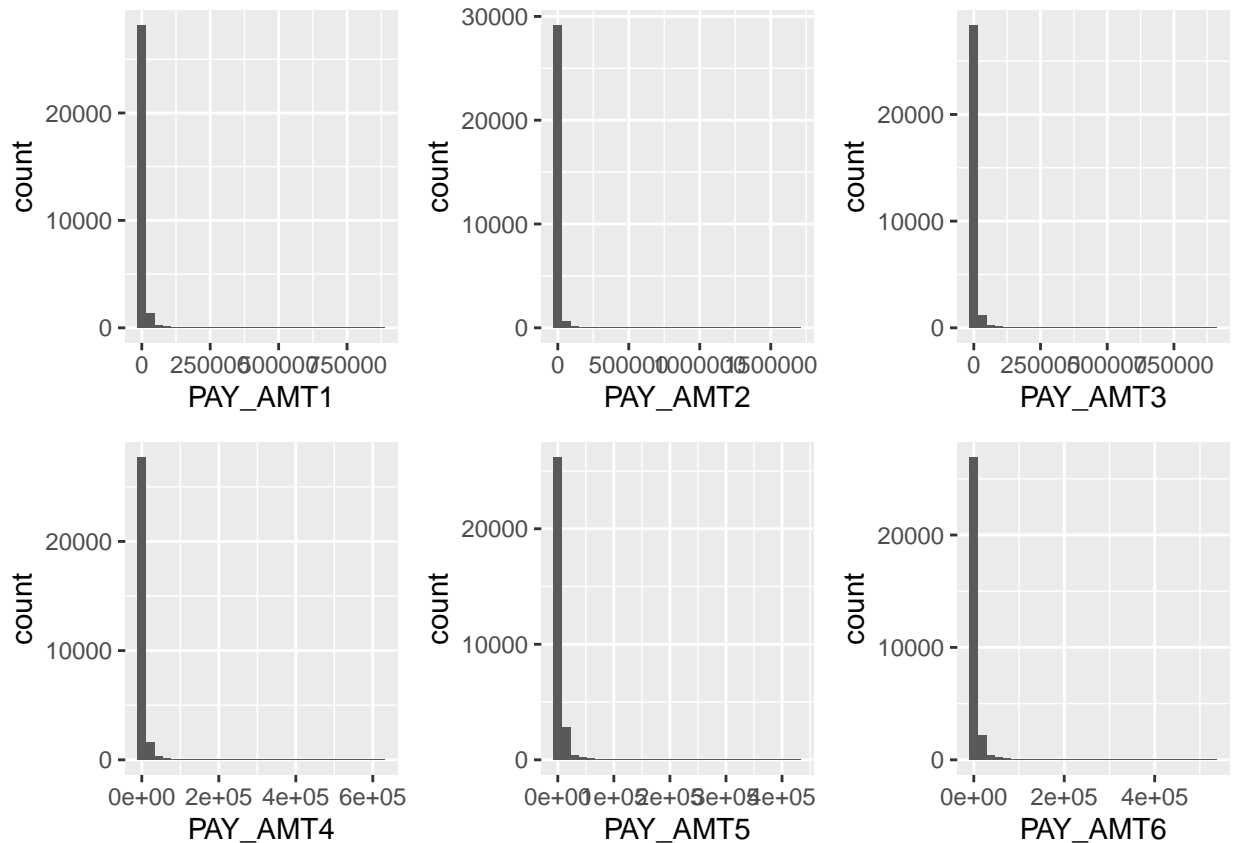


From PAY_AMT1 to PAY_AMT6, their structures are almost the same as are shown in following plots.

```
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/ 4 levels "1","2","3","4": 2 2 2 2 2 1 1 2 3 3 ...
## $ MARRIAGE  : Factor w/ 3 levels "1","2","3": 1 2 2 1 1 2 2 2 1 2 ...
## $ 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   1:10585   1:13659 Min.   : -1.5715
## 1st Qu.: -0.9055 2:18112   2:14030   2:15964 1st Qu.: -0.8121
## Median : -0.2118          3: 4917    3:  377 Median : -0.1612
## Mean   :  0.0000          4:  468      Mean   :  0.0000
## 3rd Qu.:  0.5589                        3rd Qu.:  0.5982
## Max.   :  6.4164                        Max.   :  4.7207
##

```

```
##      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_mdl <- 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	3	20174	20182
##	+ PAY_AMT6	1	20220	20224
##	+ SEX	1	20257	20261
##	+ MARRIAGE	2	20277	20283
##	+ 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
##	+ EDUCATION	3	17324	17352
##	+ PAY_AMT5	1	17329	17353

```
## + PAY_AMT4      1      17338 17362
## + PAY_AMT6      1      17352 17376
## + SEX            1      17364 17388
## + BILL_AMT5     1      17377 17401
## + BILL_AMT6     1      17377 17401
## + MARRIAGE      2      17375 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
## + EDUCATION  3      17062 17110
## + PAY_AMT3   1      17067 17111
## + PAY_2      9      17052 17112
## + PAY_AMT4   1      17075 17119
## + PAY_AMT6   1      17082 17126
## + SEX        1      17091 17135
## + BILL_AMT6  1      17095 17139
## + BILL_AMT5  1      17096 17140
## + BILL_AMT4  1      17099 17143
## + MARRIAGE   2      17099 17145
## + 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
## + EDUCATION  3      16929 16979
```

```

## + PAY_AMT1    1    16934 16980
## + PAY_2       9    16919 16981
## + BILL_AMT2   1    16936 16982
## + BILL_AMT1   1    16937 16983
## + PAY_AMT5    1    16941 16987
## + SEX         1    16943 16989
## + MARRIAGE    2    16941 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
## + EDUCATION    3    16827 16895
## + PAY_AMT1    1    16834 16898
## + BILL_AMT3   1    16837 16901
## + PAY_2       9    16823 16903
## + BILL_AMT4   1    16840 16904
## + MARRIAGE    2    16839 16905
## + 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
##	+ EDUCATION	3	16792	16862
##	+ BILL_AMT5	1	16800	16866
##	+ BILL_AMT6	1	16805	16871
##	+ MARRIAGE	2	16803	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
##	+ EDUCATION	3	16753	16825
##	+ PAY_AMT1	1	16758	16826
##	+ PAY_3	10	16744	16830
##	+ SEX	1	16767	16835
##	+ MARRIAGE	2	16765	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=16824.78
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## EDUCATION
```

```
## 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 16732 16806
## + PAY_3 10 16718 16810
## + SEX 1 16740 16814
## + PAY_AMT5 1 16741 16815
## + MARRIAGE 2 16739 16815
## + PAY_AMT3 1 16744 16818
## + PAY_AMT4 1 16745 16819
## + AGE 1 16746 16820
## + PAY_5 9 16730 16820
## + BILL_AMT6 1 16747 16821
## + PAY_2 9 16731 16821
## + BILL_AMT5 1 16748 16822
## + PAY_AMT6 1 16749 16823
## <none> 16753 16825
## + BILL_AMT4 1 16751 16825
## + BILL_AMT2 1 16753 16827
## + BILL_AMT1 1 16753 16827
## - EDUCATION 3 16780 16846
## - BILL_AMT3 1 16792 16862
## - PAY_AMT2 1 16809 16879
## - PAY_6 9 16855 16909
## - PAY_4 10 16866 16918
## - LIMIT_BAL 1 16872 16942
## - PAY_0 10 18314 18366
##
## Step: AIC=16806.31
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## EDUCATION + PAY_AMT1
```

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

```
##           Df Deviance   AIC
## + PAY_3 10 16694 16788
## + SEX 1 16719 16795
## + MARRIAGE 2 16719 16797
## + PAY_AMT5 1 16723 16799
```

```

## + AGE      1      16726 16802
## + BILL_AMT6 1      16726 16802
## + PAY_5     9      16710 16802
## + PAY_AMT3  1      16726 16802
## + BILL_AMT5 1      16727 16803
## + PAY_AMT4  1      16727 16803
## + PAY_AMT6  1      16730 16806
## <none>      16732 16806
## + BILL_AMT4 1      16731 16807
## + PAY_2     9      16716 16808
## + BILL_AMT1 1      16732 16808
## + BILL_AMT2 1      16732 16808
## - PAY_AMT1  1      16753 16825
## - EDUCATION 3      16758 16826
## - PAY_AMT2  1      16780 16852
## - BILL_AMT3 1      16781 16853
## - PAY_6     9      16833 16889
## - PAY_4     10     16844 16898
## - LIMIT_BAL 1      16843 16915
## - PAY_0     10     18268 18322
##
## Step:  AIC=16788.3
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          EDUCATION + PAY_AMT1 + PAY_3
##
##           Df Deviance   AIC
## + SEX      1      16682 16778
## + MARRIAGE  2      16681 16779
## + PAY_AMT5  1      16685 16781
## + AGE      1      16687 16783
## + BILL_AMT6 1      16688 16784
## + PAY_5     9      16672 16784
## + PAY_AMT3  1      16688 16784
## + BILL_AMT5 1      16689 16785
## + PAY_AMT4  1      16689 16785
## + PAY_AMT6  1      16691 16787
## <none>      16694 16788
## + BILL_AMT4 1      16693 16789
## + BILL_AMT1 1      16694 16790
## + BILL_AMT2 1      16694 16790
## + PAY_2     9      16686 16798
## - PAY_3     10     16732 16806
## - EDUCATION 3      16720 16808
## - PAY_AMT1  1      16718 16810
## - PAY_AMT2  1      16730 16822
## - PAY_4     10     16752 16826
## - BILL_AMT3 1      16737 16829
## - PAY_6     9      16787 16863
## - LIMIT_BAL 1      16796 16888
## - PAY_0     10     18019 18093
##
## Step:  AIC=16777.92
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          EDUCATION + PAY_AMT1 + PAY_3 + SEX

```

```

##
##           Df Deviance   AIC
## + MARRIAGE  2    16668 16768
## + PAY_AMT5  1    16672 16770
## + PAY_AMT3  1    16675 16773
## + PAY_5     9    16660 16774
## + BILL_AMT6 1    16676 16774
## + PAY_AMT4  1    16676 16774
## + BILL_AMT5 1    16677 16775
## + AGE       1    16677 16775
## + PAY_AMT6  1    16679 16777
## <none>      16682 16778
## + BILL_AMT4 1    16680 16778
## + BILL_AMT1 1    16681 16779
## + BILL_AMT2 1    16682 16780
## + PAY_2     9    16674 16788
## - SEX       1    16694 16788
## - PAY_3     10   16719 16795
## - EDUCATION 3    16707 16797
## - PAY_AMT1  1    16706 16800
## - PAY_AMT2  1    16718 16812
## - PAY_4     10   16740 16816
## - BILL_AMT3 1    16725 16819
## - PAY_6     9    16773 16851
## - LIMIT_BAL 1    16782 16876
## - PAY_0     10   18007 18083
##
## Step:  AIC=16768.08
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##           EDUCATION + PAY_AMT1 + PAY_3 + SEX + MARRIAGE
##
##           Df Deviance   AIC
## + PAY_AMT5  1    16659 16761
## + PAY_5     9    16645 16763
## + PAY_AMT3  1    16662 16764
## + BILL_AMT6 1    16662 16764
## + BILL_AMT5 1    16663 16765
## + PAY_AMT4  1    16663 16765
## + PAY_AMT6  1    16665 16767
## <none>      16668 16768
## + BILL_AMT4 1    16666 16768
## + BILL_AMT1 1    16667 16769
## + AGE       1    16668 16770
## + BILL_AMT2 1    16668 16770
## + PAY_2     9    16660 16778
## - MARRIAGE  2    16682 16778
## - SEX       1    16681 16779
## - PAY_3     10   16705 16785
## - EDUCATION 3    16693 16787
## - PAY_AMT1  1    16692 16790
## - PAY_AMT2  1    16704 16802
## - PAY_4     10   16725 16805
## - BILL_AMT3 1    16710 16808
## - PAY_6     9    16759 16841

```

```

## - LIMIT_BAL 1 16776 16874
## - PAY_0 10 17988 18068
##
## Step: AIC=16760.67
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## EDUCATION + PAY_AMT1 + PAY_3 + SEX + MARRIAGE + PAY_AMT5
##
## Df Deviance AIC
## + PAY_5 9 16636 16756
## + BILL_AMT5 1 16654 16758
## + PAY_AMT3 1 16654 16758
## + PAY_AMT4 1 16655 16759
## + BILL_AMT6 1 16656 16760
## + PAY_AMT6 1 16657 16761
## <none> 16659 16761
## + BILL_AMT4 1 16657 16761
## + AGE 1 16658 16762
## + BILL_AMT1 1 16658 16762
## + BILL_AMT2 1 16658 16762
## - PAY_AMT5 1 16668 16768
## + PAY_2 9 16650 16770
## - MARRIAGE 2 16672 16770
## - SEX 1 16672 16772
## - PAY_3 10 16696 16778
## - EDUCATION 3 16683 16779
## - PAY_AMT1 1 16680 16780
## - PAY_AMT2 1 16690 16790
## - PAY_4 10 16717 16799
## - BILL_AMT3 1 16704 16804
## - PAY_6 9 16747 16831
## - LIMIT_BAL 1 16759 16859
## - PAY_0 10 17978 18060
##
## Step: AIC=16755.58
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
## EDUCATION + PAY_AMT1 + PAY_3 + SEX + MARRIAGE + PAY_AMT5 +
## PAY_5

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

## Df Deviance AIC
## + PAY_AMT3 1 16629 16751
## + BILL_AMT5 1 16630 16752
## + BILL_AMT6 1 16633 16755
## + PAY_AMT4 1 16633 16755
## + PAY_AMT6 1 16633 16755
## + BILL_AMT4 1 16634 16756
## <none> 16636 16756
## + AGE 1 16635 16757
## + BILL_AMT1 1 16635 16757
## + BILL_AMT2 1 16635 16757
## - PAY_5 9 16659 16761
## - PAY_AMT5 1 16645 16763
## + PAY_2 9 16627 16765

```



```

## - MARRIAGE      2      16649 16765
## - SEX           1      16649 16767
## - PAY_4         10      16673 16773
## - PAY_3         10      16673 16773
## - PAY_AMT1      1      16656 16774
## - EDUCATION     3      16660 16774
## - PAY_6         9      16682 16784
## - PAY_AMT2      1      16667 16785
## - BILL_AMT3     1      16680 16798
## - LIMIT_BAL     1      16734 16852
## - PAY_0         10      17929 18029
##
## Step:  AIC=16751.29
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          EDUCATION + PAY_AMT1 + PAY_3 + SEX + MARRIAGE + PAY_AMT5 +
##          PAY_5 + PAY_AMT3

```

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

```

##          Df Deviance   AIC
## + BILL_AMT5  1      16626 16750
## + PAY_AMT4   1      16627 16751
## <none>              16629 16751
## + PAY_AMT6   1      16628 16752
## + BILL_AMT6  1      16628 16752
## + BILL_AMT2  1      16629 16753
## + AGE        1      16629 16753
## + BILL_AMT1  1      16629 16753
## + BILL_AMT4  1      16629 16753
## - PAY_AMT3   1      16636 16756
## - PAY_5      9      16654 16758
## - PAY_AMT5   1      16638 16758
## - MARRIAGE   2      16643 16761
## + PAY_2      9      16621 16761
## - SEX        1      16643 16763
## - PAY_4      10      16663 16765
## - PAY_AMT1   1      16647 16767
## - EDUCATION  3      16653 16769
## - PAY_3      10      16667 16769
## - PAY_AMT2   1      16657 16777
## - PAY_6      9      16675 16779
## - BILL_AMT3  1      16675 16795
## - LIMIT_BAL  1      16721 16841
## - PAY_0      10      17920 18022
##
## Step:  AIC=16750.12
## DEFAULT ~ PAY_0 + PAY_4 + LIMIT_BAL + PAY_6 + PAY_AMT2 + BILL_AMT3 +
##          EDUCATION + PAY_AMT1 + PAY_3 + SEX + MARRIAGE + PAY_AMT5 +
##          PAY_5 + PAY_AMT3 + BILL_AMT5

```

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

```
##          Df Deviance   AIC
```

```
## <none>          16626 16750
## + PAY_AMT6      1    16624 16750
## + BILL_AMT2     1    16625 16751
## + BILL_AMT4     1    16625 16751
## - BILL_AMT5     1    16629 16751
## + PAY_AMT4      1    16625 16751
## + AGE           1    16626 16752
## + BILL_AMT1     1    16626 16752
## + BILL_AMT6     1    16626 16752
## - PAY_AMT3      1    16630 16752
## - PAY_AMT5      1    16635 16757
## - PAY_5         9    16651 16757
## - MARRIAGE      2    16640 16760
## + PAY_2         9    16618 16760
## - SEX           1    16639 16761
## - PAY_4         10   16661 16765
## - PAY_AMT1      1    16644 16766
## - PAY_3         10   16664 16768
## - EDUCATION     3    16650 16768
## - BILL_AMT3     1    16650 16772
## - PAY_6         9    16671 16777
## - PAY_AMT2      1    16657 16779
## - LIMIT_BAL     1    16714 16836
## - PAY_0         10   17919 18023
```

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

```
step_prob <- predict(step_md1, 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 3568  715
##           1  171  347
##
##           Accuracy : 0.8155
##           95% CI : (0.8042, 0.8263)
##           No Information Rate : 0.7788
##           P-Value [Acc > NIR] : 2.329e-10
##
##           Kappa : 0.3441
##
##           Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9543
##           Specificity : 0.3267
##           Pos Pred Value : 0.8331
```

```
##          Neg Pred Value : 0.6699
##          Prevalence : 0.7788
##          Detection Rate : 0.7432
##          Detection Prevalence : 0.8921
##          Balanced Accuracy : 0.6405
##
##          '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	0.9542658	0.326742	0.6405039

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.

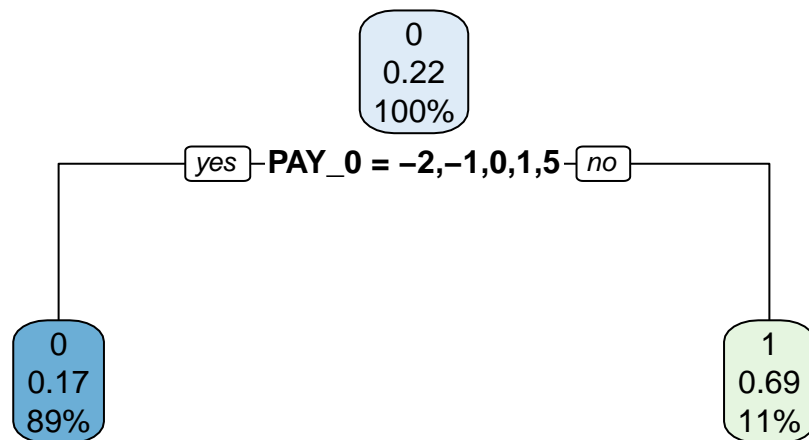
```
confusionMatrix(rpart_pred, validation_set$DEFAULT)
```

```
## 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	0.9542658	0.326742	0.6405039
CART default	0.8171214	0.9620219	0.306968	0.6344950

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

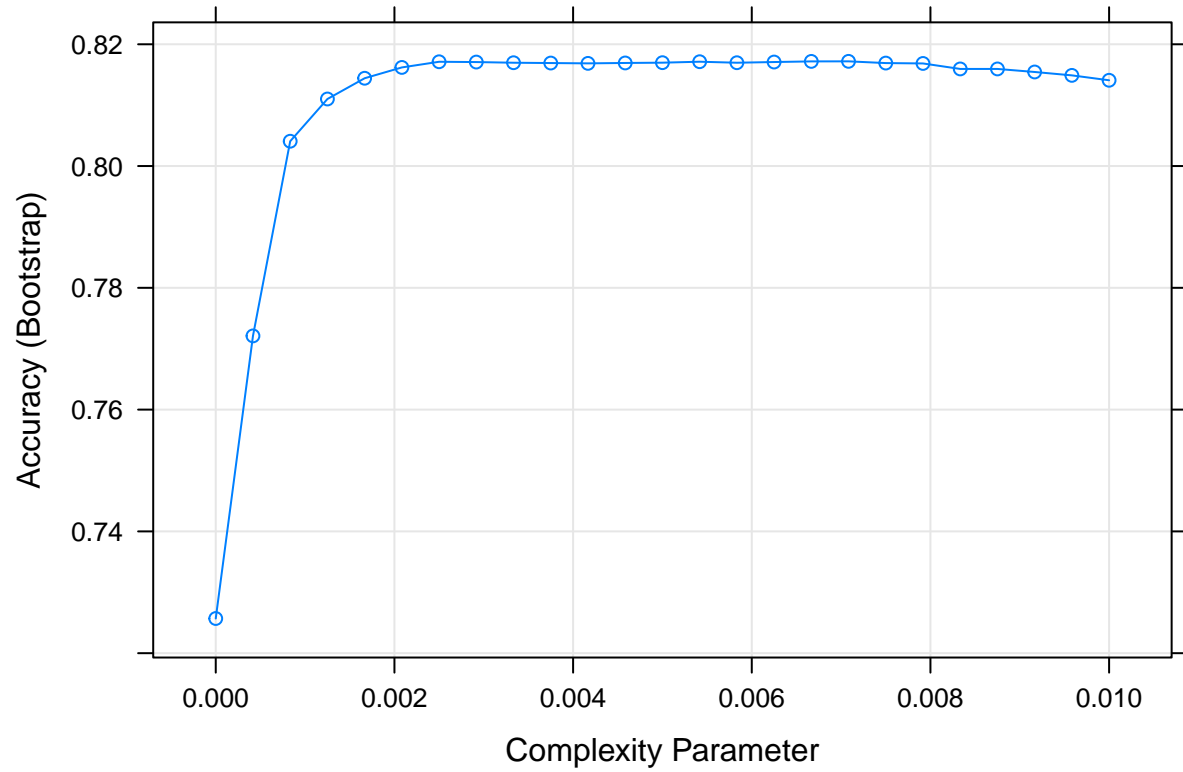
```
set.seed(2021, sample.kind = "Rounding")
```

```
## 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.

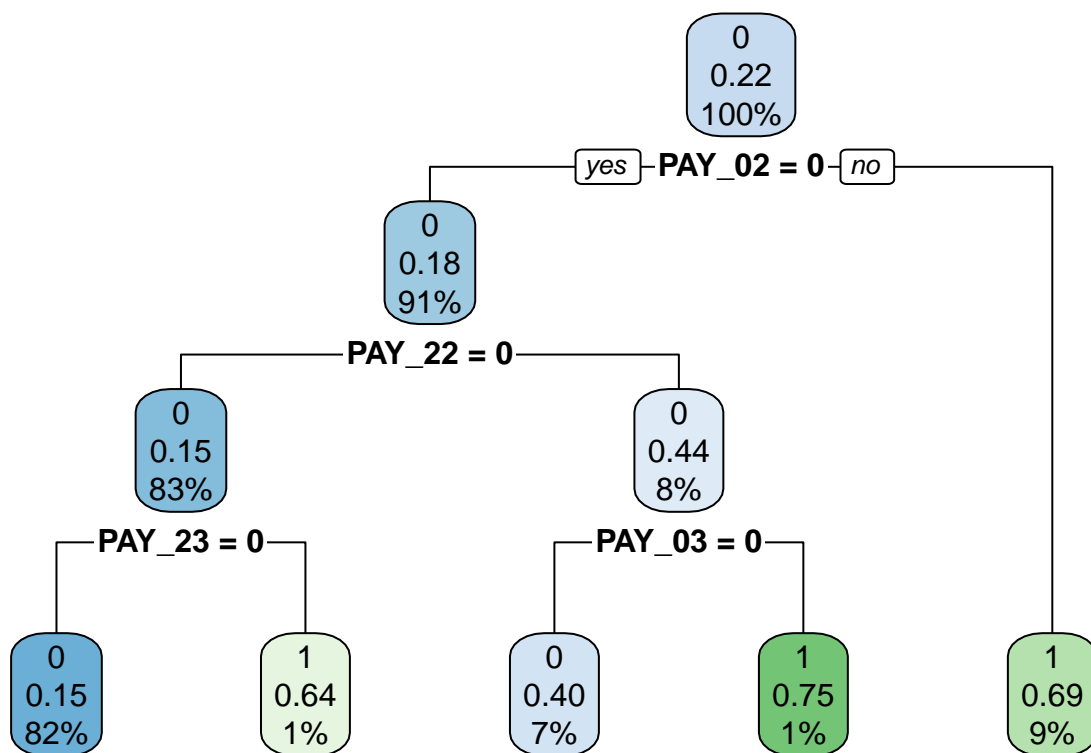
```
plot(rpart_tuned_md1)
```



```
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	0.9542658	0.3267420	0.6405039
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:      18.27 %
```

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 3556  700
##           1  183  362
##
##           Accuracy : 0.8161
##           95% CI : (0.8048, 0.8269)
##       No Information Rate : 0.7788
##       P-Value [Acc > NIR] : 1.154e-10
##
##           Kappa : 0.3535
##
##  Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9511
##           Specificity : 0.3409
##           Pos Pred Value : 0.8355
##           Neg Pred Value : 0.6642
##           Prevalence : 0.7788
##           Detection Rate : 0.7407
##       Detection Prevalence : 0.8865
##           Balanced Accuracy : 0.6460
##
##           '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	0.9542658	0.3267420	0.6405039
CART default	0.8171214	0.9620219	0.3069680	0.6344950
CART tuned cp	0.8162883	0.9593474	0.3126177	0.6359826
random forest default	0.8160800	0.9510564	0.3408663	0.6459614

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))

```

```

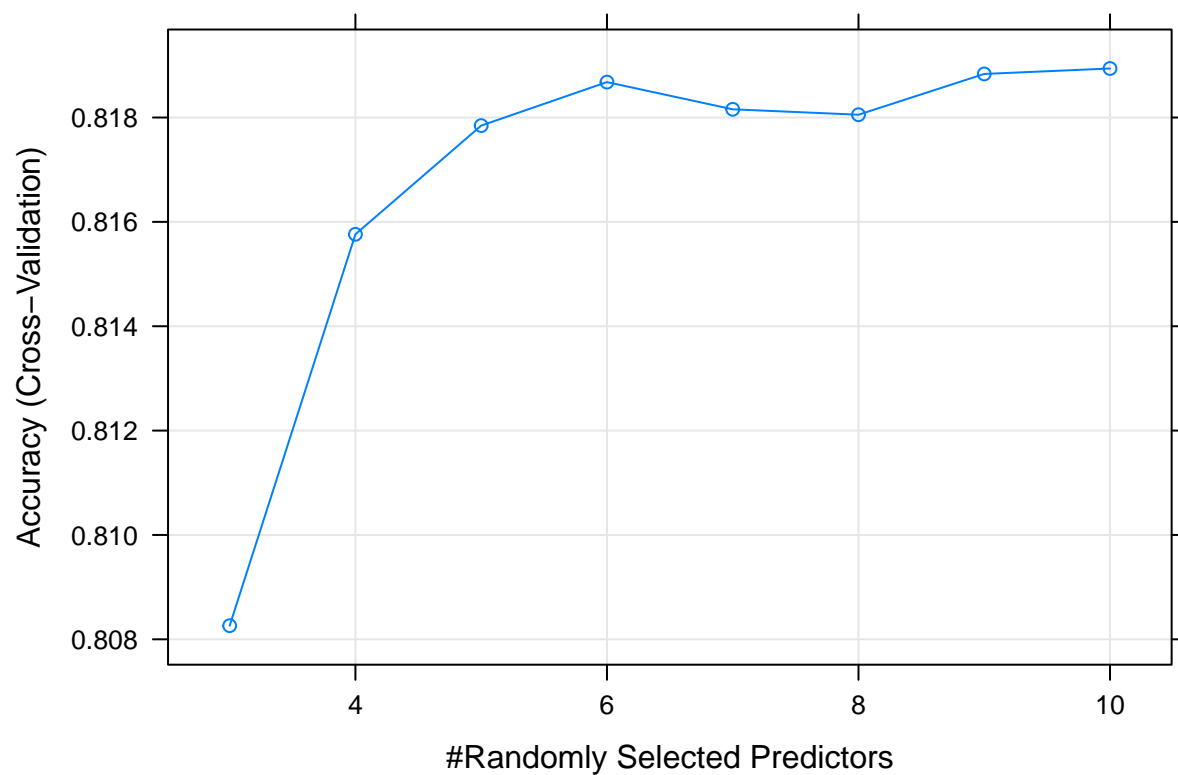
## Growing trees.. Progress: 87%. Estimated remaining time: 4 seconds.
## Growing trees.. Progress: 79%. Estimated remaining time: 8 seconds.
## Growing trees.. Progress: 86%. Estimated remaining time: 5 seconds.
## Growing trees.. Progress: 79%. Estimated remaining time: 8 seconds.
## Growing trees.. Progress: 100%. Estimated remaining time: 0 seconds.
## Growing trees.. Progress: 92%. Estimated remaining time: 2 seconds.

```

```
## Growing trees.. Progress: 78%. Estimated remaining time: 8 seconds.
## Growing trees.. Progress: 96%. Estimated remaining time: 1 seconds.
## Growing trees.. Progress: 90%. Estimated remaining time: 3 seconds.
## Growing trees.. Progress: 78%. Estimated remaining time: 8 seconds.
## Growing trees.. Progress: 100%. Estimated remaining time: 0 seconds.
## Growing trees.. Progress: 87%. Estimated remaining time: 4 seconds.
## Growing trees.. Progress: 80%. Estimated remaining time: 7 seconds.
## Growing trees.. Progress: 55%. Estimated remaining time: 25 seconds.
```

Plot.

```
plot(rf_cv_mdl)
```



Prediction.

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

Confusion Matrix

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

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
```

```
##          0 3570  710
##          1  169  352
##
##          Accuracy : 0.8169
##          95% CI : (0.8057, 0.8278)
##    No Information Rate : 0.7788
##    P-Value [Acc > NIR] : 4.443e-11
##
##          Kappa : 0.3501
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##          Sensitivity : 0.9548
##          Specificity : 0.3315
##    Pos Pred Value : 0.8341
##    Neg Pred Value : 0.6756
##          Prevalence : 0.7788
##    Detection Rate : 0.7436
##    Detection Prevalence : 0.8915
##    Balanced Accuracy : 0.6431
##
##    '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	0.9542658	0.3267420	0.6405039
CART default	0.8171214	0.9620219	0.3069680	0.6344950
CART tuned cp	0.8162883	0.9593474	0.3126177	0.6359826
random forest default	0.8160800	0.9510564	0.3408663	0.6459614
random forest tuned	0.8169131	0.9548007	0.3314501	0.6431254

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
##           0.9505671           0.3576807           0.8389046
##       Neg Pred Value           Precision           Recall
##           0.6728045           0.8389046           0.9505671
##           F1           Prevalence           Detection Rate
##           0.8912520           0.7787035           0.7402100
##   Detection Prevalence   Balanced Accuracy
##           0.8823529           0.6541239
```

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])

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])
```

```

    Balanced_Accuracy = confusionMatrix(final_pred_rf, test_set$DEFAULT)$byClass
final_results %>% knitr::kable()

```

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
CART default	0.8226962	0.9621228	0.3320783	0.6471006
Random forest default	0.8193634	0.9505671	0.3576807	0.6541239

Conclusion

###