

Report on Default of Credit Card Clients Dataset

Masayoshi Sato

2021/6/26

Introduction

Finance is one of fields where machine learning is commonly used. It deals with a huge amount of data and is also surrounded by a lot of uncertainties. To predict outcome using pre-existing data is a crucial part of finance. In this paper, we will deal with a problem many credit card companies have been facing. Can we predict whether a credit user will pay their debt or fail based on objective data? Traditionally, finding a credible borrower have been a kind of know-how, or skill and experience nurtured by financial institutions. Instead, we try to build machine learning models. using a dataset which is open to public.

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. Each data was anonymously collected and labeled with individual ID.

Our goal is to find a classification model which predicts the most accurate outcome, default or not. We need to bear it in mind that its distribution of these outcomes 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. Our 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, “tidyverse¹”, “DataExplorer²”, “gridExtra³”, “rpart⁴”, “caret⁵”, and “ranger⁶”.

We use a dataset stored in Kaggle⁷ 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.

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

²<https://cran.r-project.org/web/packages/DataExplorer/index.html>

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

⁴<https://cran.r-project.org/web/packages/rpart/index.html>

⁵<https://cran.r-project.org/web/packages/caret/index.html>

⁶<https://cran.r-project.org/web/packages/ranger/index.html>

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

Kaggle requires registration to download the data. For the sake of convenience, the data file is stored in my GitHub repository⁸.

Data Exploration

First, we need to check the downloaded dataset. Columns are as follows.

```
## spec_tbl_df [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(),
```

⁸https://github.com/masa951125/Final_project/raw/main/UCI_Credit_Card.csv

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

It has 30000 rows and 25 columns. “Default.payment.next.month” is an outcome . Other features seem to be either numerical or categorical data. “SEX”, “EDUCATION”, “MARRIAGE”, “PAY_0” -“PAY_6”, and “default.payment.next.month” look like categorical data, as their values are limited number of integers. Other features seem to be numerical.

We know there are no NAs, Nulls in the dataset.

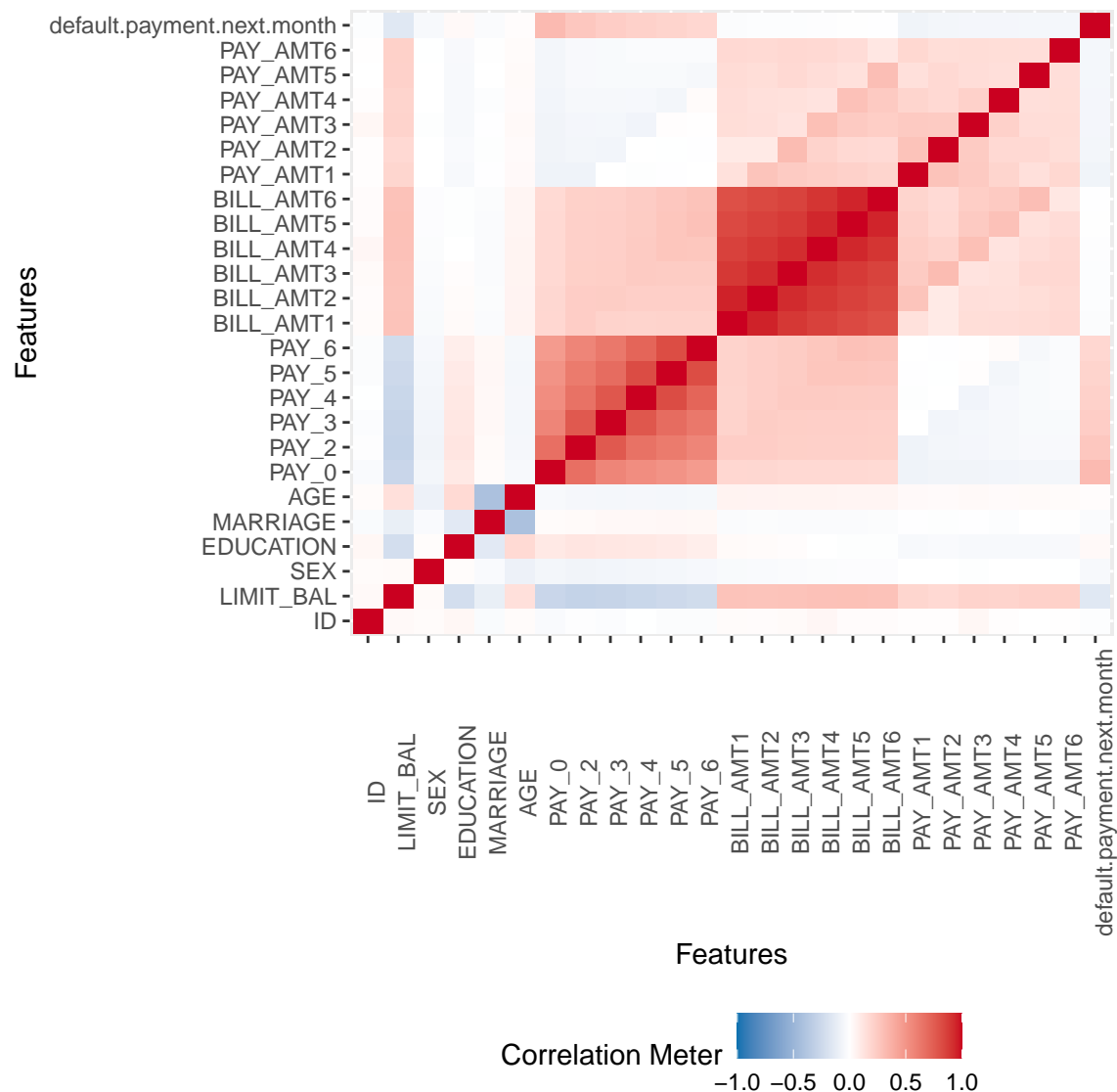
```
#Number of NAs
sum(is.na(original_default))
```

```
## [1] 0
```

```
#Number of Nulls
sum(is.null(original_default))
```

```
## [1] 0
```

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.”⁹ We show the proportion of “0”, “1”.

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

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

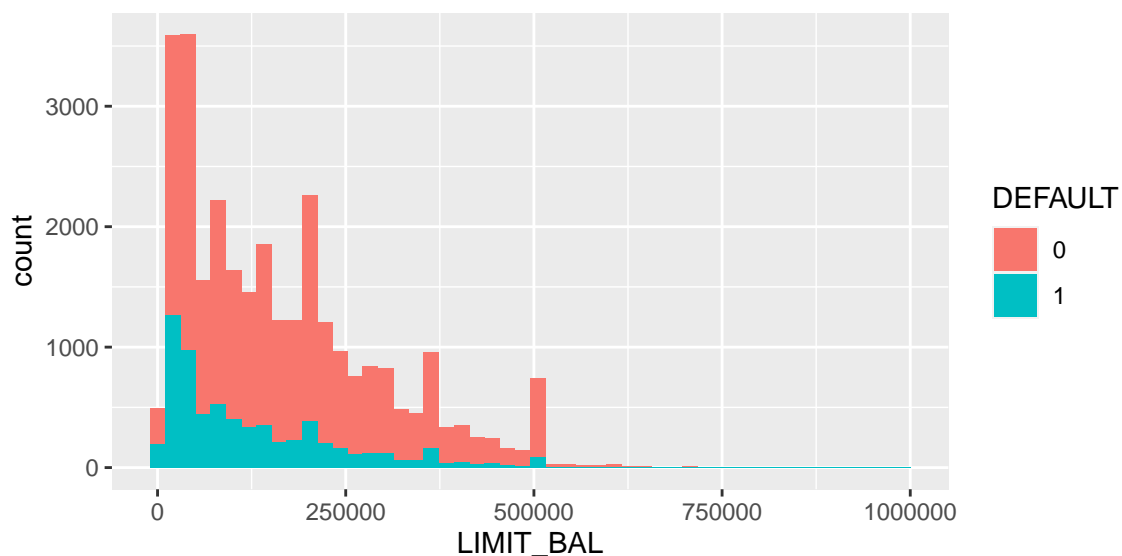
This means that if we predict all the outcome as “0”, we will get 77.9% accuracy. We need to take into account this fact. We change the name, “default.payment.next.month”, to “DEFAULT” for the sake of convenience. Also, we change this numeric variable into factor.

2 “LIMIT_BAL”

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

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

We draw its distribution filling the proportion of default.



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

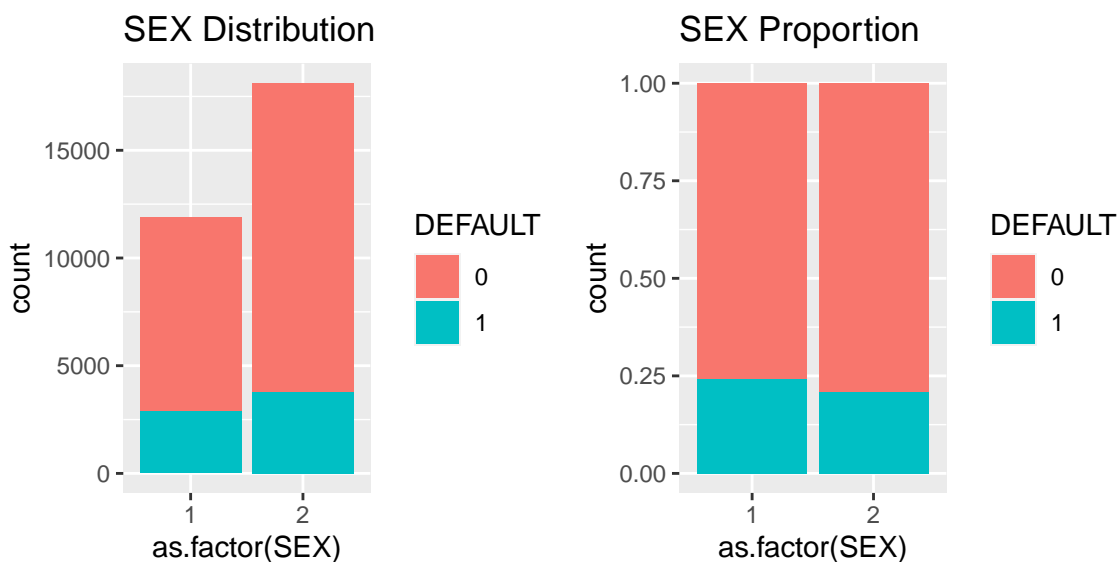
3 “SEX”

The values “1”, “2” correspond to male and female respectively¹¹. Male is 40% and female is 60%.

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

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

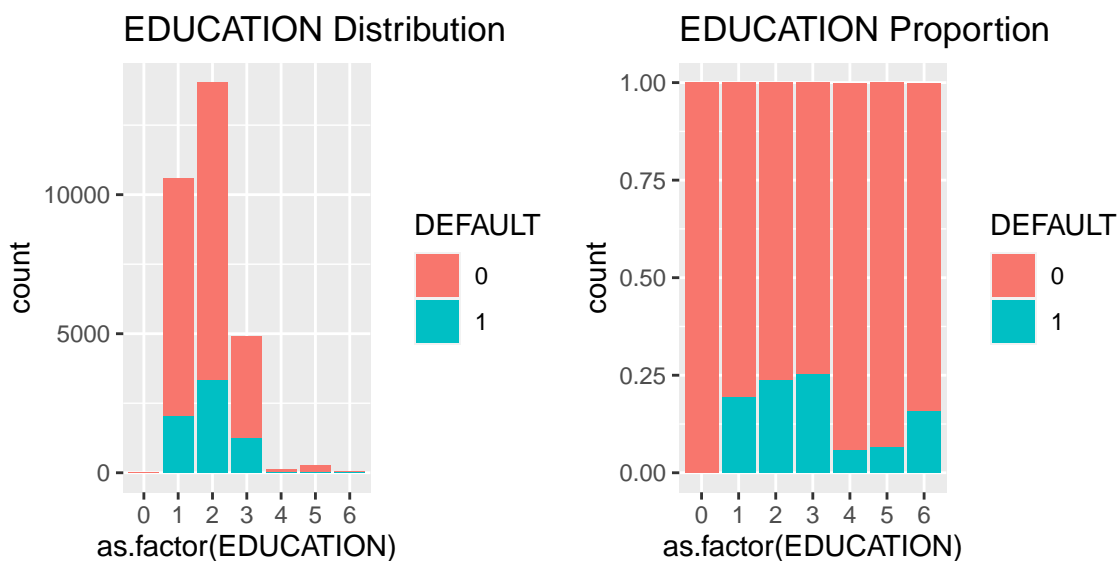
¹¹<https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>



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=unknown¹². We plot its distribution and stacked bar graph.



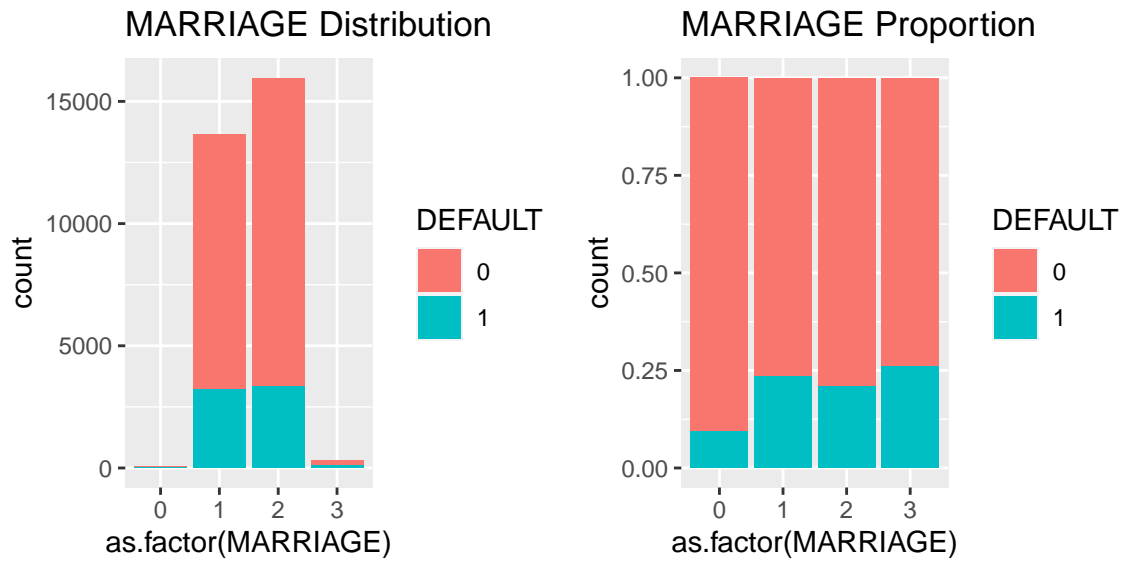
People whose final education is high school have relatively high default rate. On the other hand, people whose final education is graduate school have low default rate.

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

5 “MARRIAGE”

This is also categorical data. The value number means clients’ marital status¹³. 1=married, 2=single, 3=others. There are 54 individuals whose values are 0.

We draw its distribution graph and stacked bar graph.



Comparing married and single, married people are a little likely to become default.

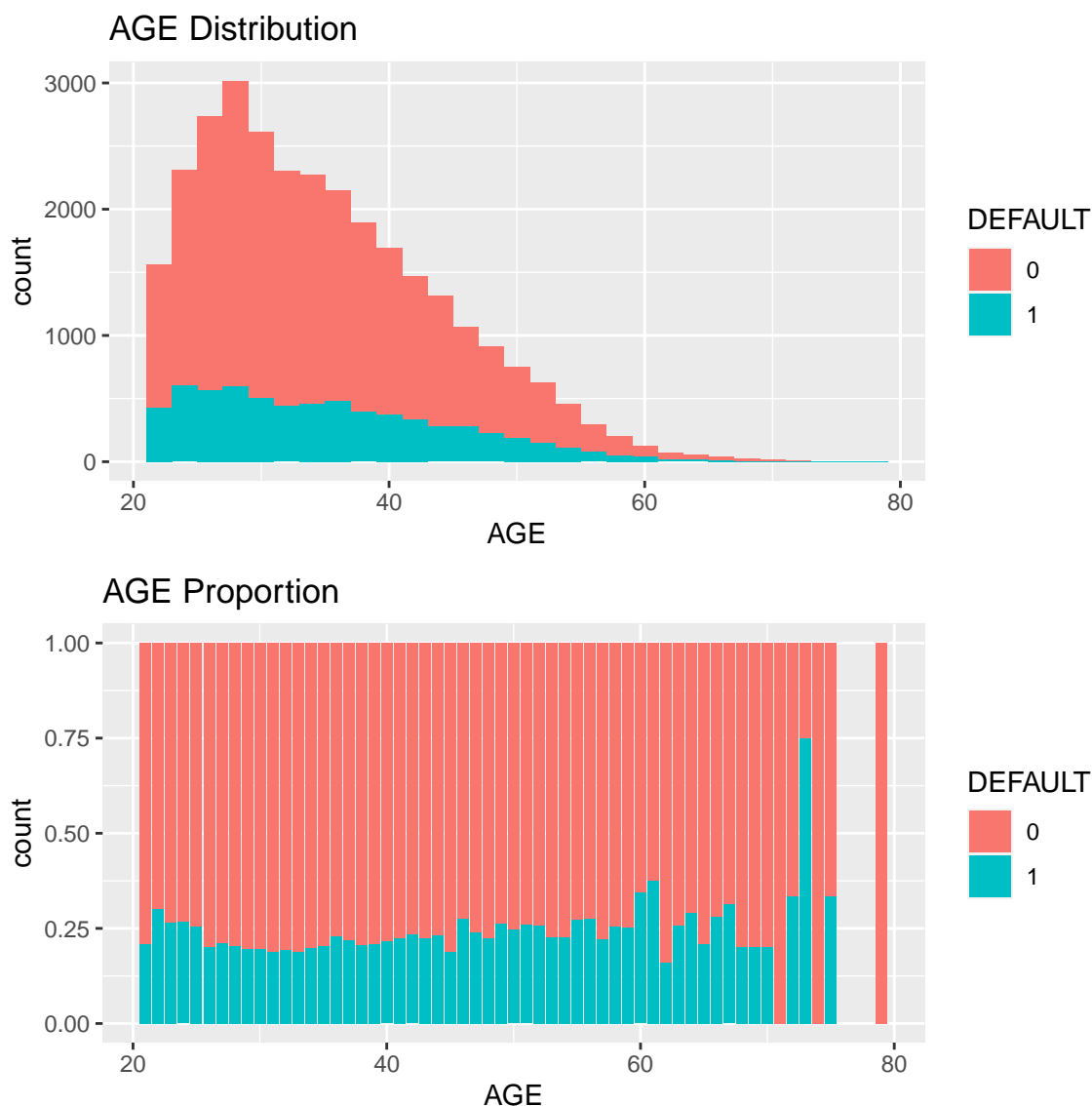
6 “AGE”

As its name indicates, it is numerical data.

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

We plots its distribution as well as its default rate.

¹³<https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>



The number of clients are decreasing as they get older. Regarding default clients, younger people in their early 20s and older age groups around 60s show higher proportions of default than other age groups.

7 “PAY”

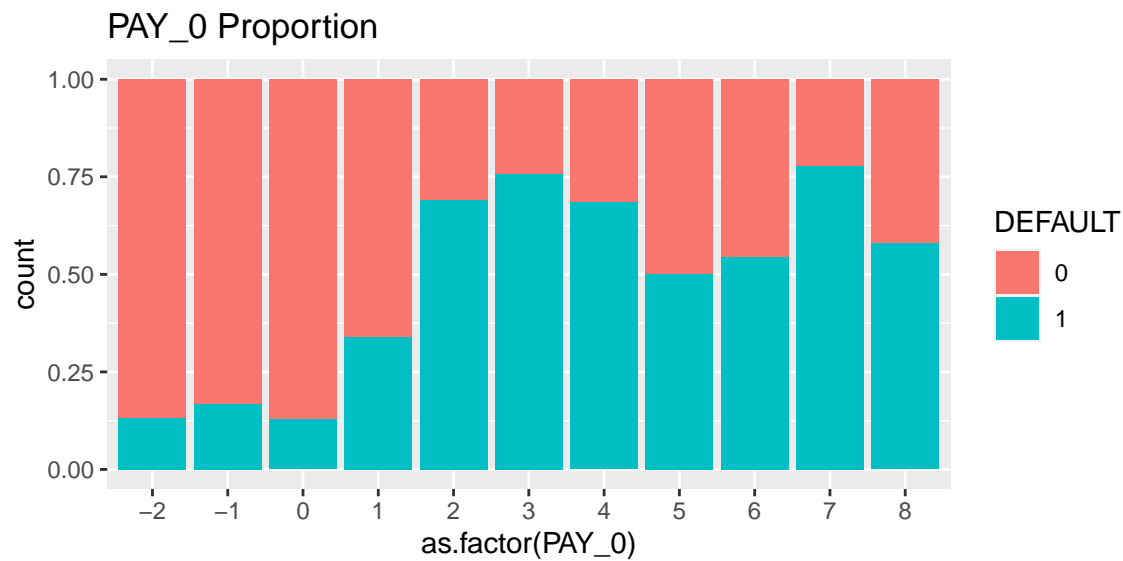
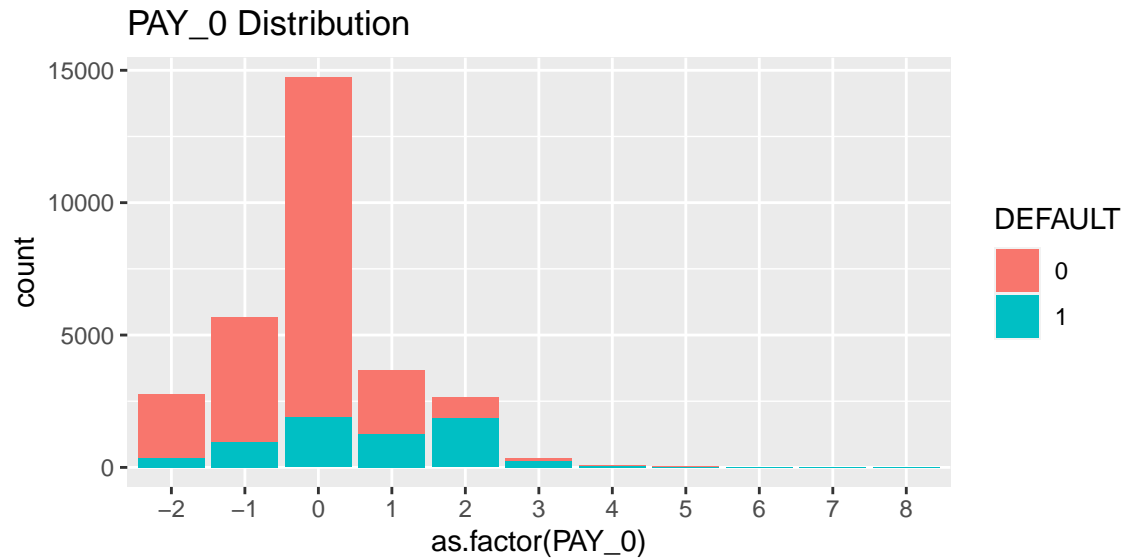
Variable from PAY_0, PAY_2 ~PAY_6 have the same values, as are explained in the description¹⁴. They are categorical data. PAY_0 means repayment status in September, 2005. Then go back in time by a month until April, 2005. Values are ;

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

They mean; “-1”= pay duly, “1”= payment delay for 1 month, “2” = payment delay for 2 months, ... 9 = payment delay for 9 months and above. “0” and “-2” are not defined.

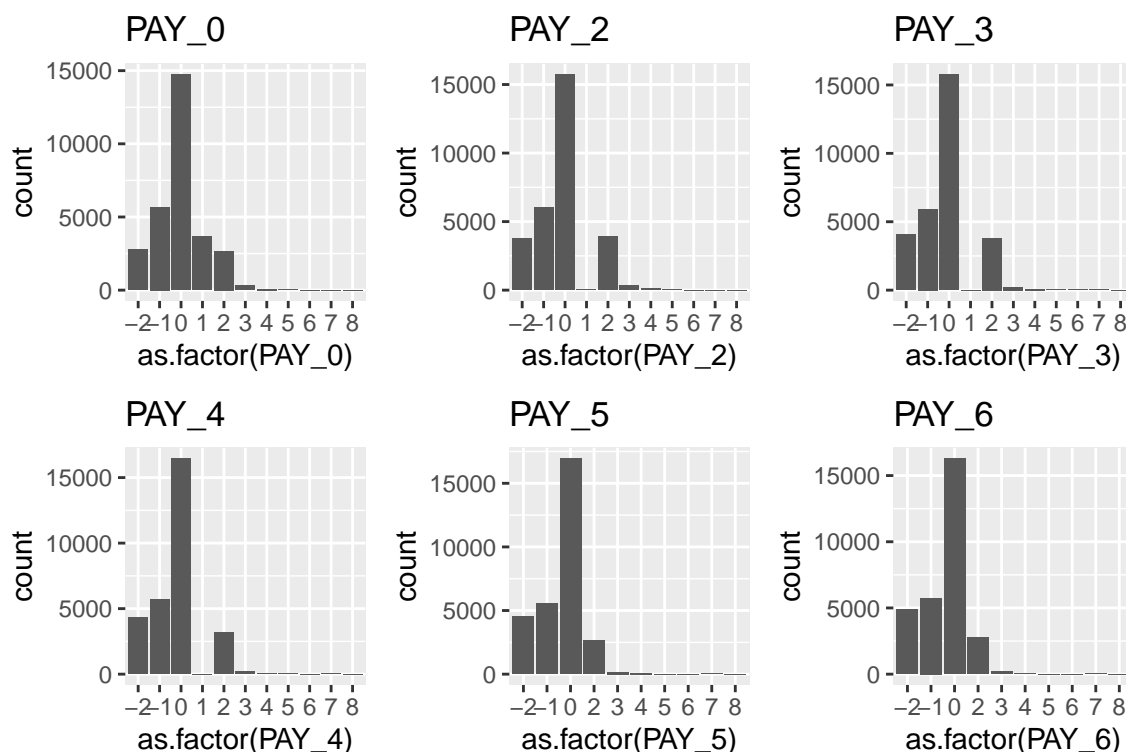
We draw PAY_0 distribution and stacked bar graph.hey are categorical data.

¹⁴<https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>



We are not sure what “-2” and “0” mean. But from these two graphs we understand that as payment delay becomes longer, clients are more likely to come to a default.

PAY_2 ~ PAY_6 ’s structures are almost as the same as PAY_0. We show their distribution.



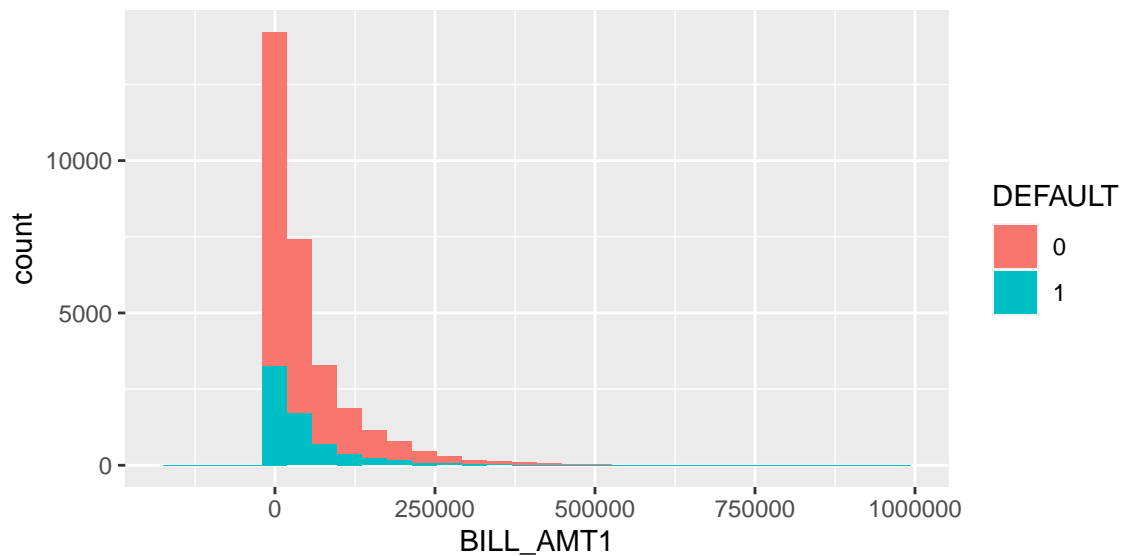
8 “BILL_AMT”

This variable means an amount of bill statement. This is numerical data.

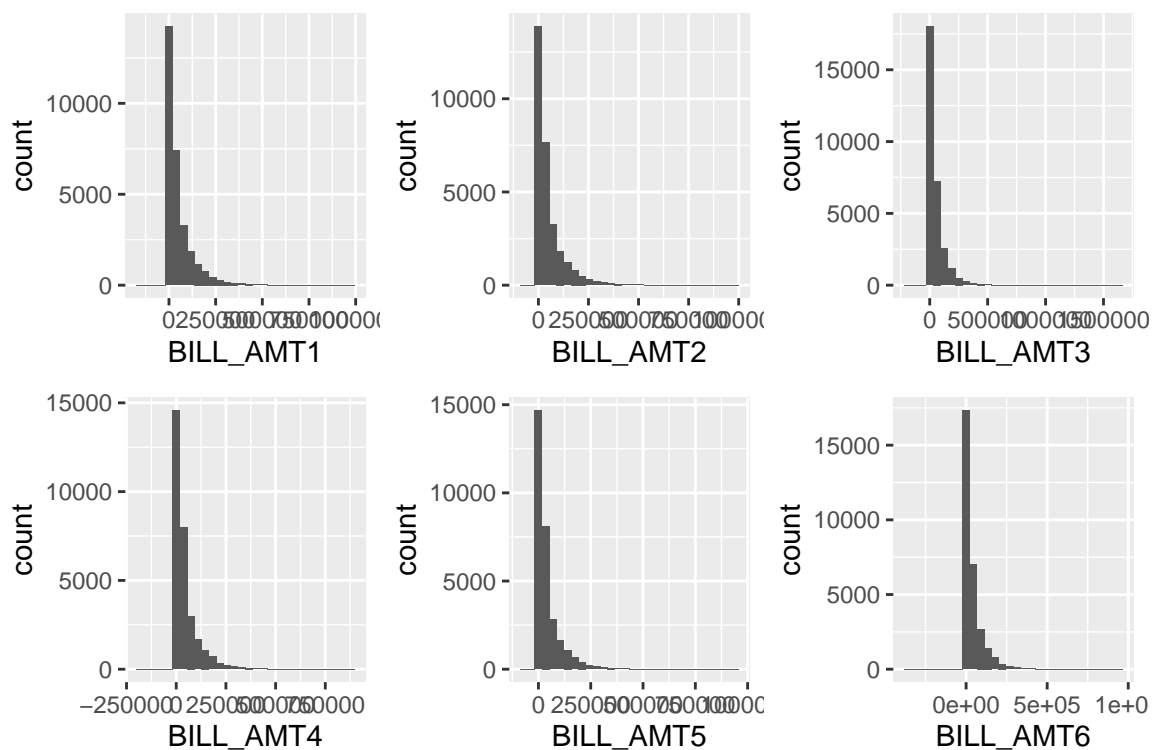
##	BILL_AMT1	BILL_AMT2	BILL_AMT3	BILL_AMT4
##	Min. : -165580	Min. : -69777	Min. : -157264	Min. : -170000
##	1st Qu.: 3559	1st Qu.: 2985	1st Qu.: 2666	1st Qu.: 2327
##	Median : 22382	Median : 21200	Median : 20089	Median : 19052
##	Mean : 51223	Mean : 49179	Mean : 47013	Mean : 43263
##	3rd Qu.: 67091	3rd Qu.: 64006	3rd Qu.: 60165	3rd Qu.: 54506
##	Max. : 964511	Max. : 983931	Max. : 1664089	Max. : 891586
##	BILL_AMT5	BILL_AMT6		
##	Min. : -81334	Min. : -339603		
##	1st Qu.: 1763	1st Qu.: 1256		
##	Median : 18105	Median : 17071		
##	Mean : 40311	Mean : 38872		
##	3rd Qu.: 50191	3rd Qu.: 49198		
##	Max. : 927171	Max. : 961664		

BILL_AMT1 is a record in September, 2005¹⁵. Then go back in time by a month until April, 2005. Likewise previous variables, we plot BILL_AMT distribution.

¹⁵<https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>



From BILL_AMT1 to BILL_AMT6, their structures are almost the same as are shown in following plots.



9 “PAY_AMT”

This variable means an amount of previous payment. This is numerical data.

##	PAY_AMT1	PAY_AMT2	PAY_AMT3	PAY_AMT4
##	Min. : 0	Min. : 0	Min. : 0	Min. : 0
##	1st Qu.: 1000	1st Qu.: 833	1st Qu.: 390	1st Qu.: 296
##	Median : 2100	Median : 2009	Median : 1800	Median : 1500

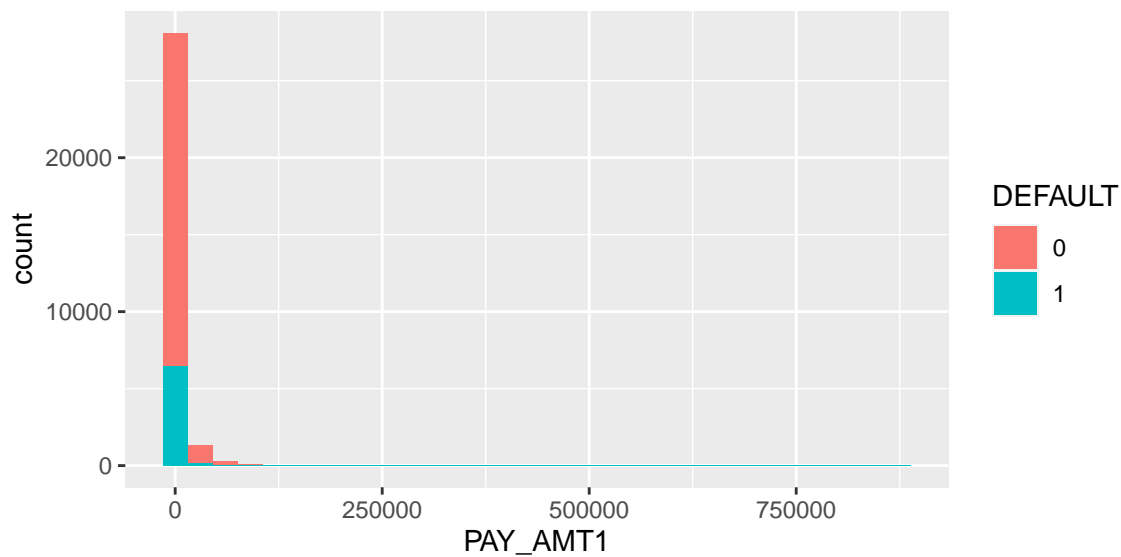
```

## Mean   : 5664   Mean   : 5921   Mean   : 5226   Mean   : 4826
## 3rd Qu.: 5006   3rd Qu.: 5000   3rd Qu.: 4505   3rd Qu.: 4013
## Max.   :873552   Max.   :1684259   Max.   :896040   Max.   :621000
## PAY_AMT5      PAY_AMT6
## Min.    : 0.0   Min.    : 0.0
## 1st Qu.: 252.5   1st Qu.: 117.8
## Median : 1500.0   Median : 1500.0
## Mean    : 4799.4   Mean    : 5215.5
## 3rd Qu.: 4031.5   3rd Qu.: 4000.0
## Max.    :426529.0   Max.    :528666.0

```

PAY_AMT1 is an amount of previous payment in September, 2005¹⁶. Likewise BILL_AMT, PAY_AMT goes back in time by a month from August to April, 2005 which is PAY_AMT6.

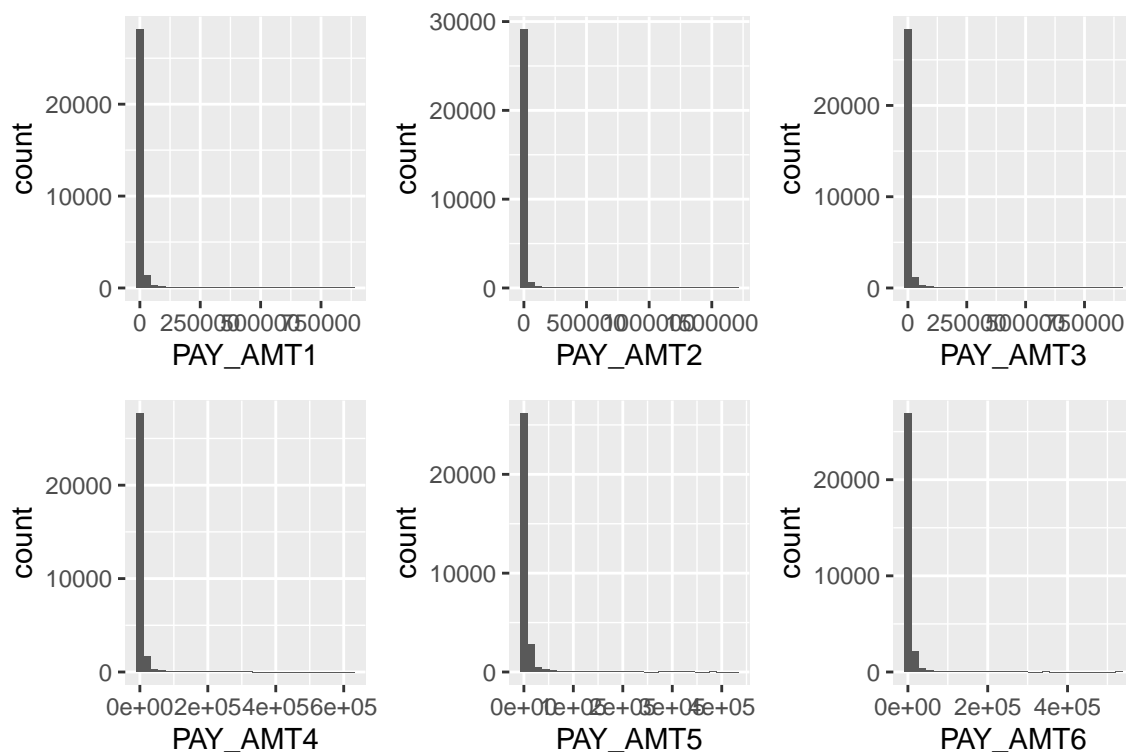
Here is PAY_AMT1's plot.



As we have seen a distribution summary, it is right skewed significantly. Maximum amount is more than 850,000 NT dollars, but its mean and median are 5664 and 2100 respectively.

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

¹⁶<https://www.kaggle.com/uciml/default-of-credit-card-clients-dataset>



Data Preparation

Before going further, we will arrange our dataset to make it easier to investigate. First, we remove ID column, as it is irrelevant to the outcome.

Currently our data is composed of numerical values. Categorical values need to be changed to factors. As we saw in data exploration, SEX, EDUCATION, MARRIAGE, PAY_0~PAY_6 are categorical values.

Regarding numerical values, “LIMIT_BAL” ranges from 10000 to 1000000, but “AGE” from 21 to 79 as we saw in the previous section. When doing regression, these wide varieties of ranges of variables cause inaccuracy. Thus we standardize these variables, using following formula.

$$Transformed.Values = \frac{Values - Mean}{Standard.Deviation}$$

Then we get variables whose means are 0, and standard deviations are 1 using function “scale”.

Then, we split the data into two. As we have relatively a large amount of data, 30000, the proportion we use is 80% and 20%. The smaller dataset will be used when evaluating a model at the final stage. We call it test_set.

We are going to use decision tree, and random forest later. They can be tuned to produce improved results by finding hyperparameters. If we tune the model using hyperparameters again and again, this might cause overfitting. To avoid this, we split the larger dataset again and produce two datasets, “train_set” and “validation_set”. Split proportion is the same as before, 80% and 20% respectively.

The three datasets have the almost same outcome ratio.

```
## # A tibble: 2 x 4
##   outcome train_set validation_set test_set
```

```
##      <dbl> <table>      <table>      <table>
## 1      0 0.7788311 0.7787961      0.7787035
## 2      1 0.2211689 0.2212039      0.2212965
```

Model analysis

1 Evaluation metrics

Generally speaking, overall accuracy is used to evaluate a model. But this dataset has imbalanced proportion of outcomes. Guessing all responses are 0, we calculate its confusion matrix.

```
##      Reference
## Prediction    0    1
##           0 3739 1062
##           1    0    0
```

From this confusion matrix, we know even such a guess produces fairly good results.

- Accuracy (True positive + True Negative / Total) 0.7788
- Sensitivity (True Positive / True Positive+ False Negative) 1.0

In contrast,

- Specificity (True Negative / True Negative + False Positive) 0
- Balanced accuracy (arithmetic mean of sensitivity and specificity) 0.5

As such, the credit company falsely give credit to a lot of clients who will fail to repay a debt. The loss for the company would be huge. Therefore, our goal is to pursue more accurate accuracy than guessing all responses are 0, taking account of improving specificity. We will use both *accuracy* and *balanced accuracy* to evaluate each model.

2 Logistic regression

Firstly, we choose logistic regression. Logistic regression is commonly used in binary classification problems (outcome is 0 and 1, or True and False). As it uses logistic transformed odds, it produces outcome ranging from 0 to 1.

We have 24 features in the dataset. There are many ways to deal with such many predictors. One of them is step-wise regression. This method is to repeatedly examine statistically significance of each variable in a regression model. There are three ways; forward selection, backward elimination, forward and backward elimination. Forward selection starts from no variable and add a variable step by step. Backward elimination starts from full model and subtract variables one by one. Forward and backward combines the two. One of the most convenient features of this method is we can get logistic regression results without knowing details of variables in the dataset. Firstly, we use forward and backward step wise regression, as it is commonly used. Afterwards, we will use logistic regression with fewer predictors based on our data exploration.

In the step-wise regression, we make two models, null model with no predictors and full model with all predictors. Then we use “step” function to conduct step-wise regression. The result;

```
##
## Call:
## glm(formula = 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, family = binomial(link = "logit"),
##      data = train_set)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.3707  -0.6010  -0.5041  -0.3009   3.4687
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -15.98503   286.37745  -0.056  0.955487
## PAY_0-1         0.39057    0.10989   3.554  0.000379 ***
## PAY_00        -0.24892    0.10957  -2.272  0.023098 *
## PAY_01         0.77280    0.10051   7.689 1.48e-14 ***
## PAY_02         1.99325    0.11628  17.142 < 2e-16 ***
## PAY_03         1.90918    0.20210   9.447 < 2e-16 ***
## PAY_04         1.60586    0.33553   4.786 1.70e-06 ***
## PAY_05         0.26727    0.67139   0.398 0.690564
## PAY_06         1.59286    0.84034   1.895 0.058027 .
## PAY_07         0.96314    1.25174   0.769 0.441634
## PAY_08        -12.20222   371.73655  -0.033 0.973814
## PAY_4-1        -0.21537    0.13474  -1.598 0.109960
## PAY_40        -0.14822    0.15003  -0.988 0.323175
## PAY_41         0.92374  1076.95006   0.001 0.999316
## PAY_42         0.12951    0.16021   0.808 0.418855
## PAY_43         0.13295    0.30917   0.430 0.667174
## PAY_44         0.50734    0.53958   0.940 0.347087
## PAY_45        -2.37446    1.04205  -2.279 0.022689 *
## PAY_46       -29.76168   497.36889  -0.060 0.952284
## PAY_47        -7.03294  6762.09451  -0.001 0.999170
## PAY_48       -35.13222  6819.44800  -0.005 0.995890
## LIMIT_BAL     -0.25510    0.02807  -9.088 < 2e-16 ***
## PAY_6-1       -0.21039    0.10068  -2.090 0.036646 *
## PAY_60        -0.42701    0.10790  -3.957 7.58e-05 ***
## PAY_62        -0.07930    0.12554  -0.632 0.527636
## PAY_63         0.77912    0.30280   2.573 0.010081 *
## PAY_64         0.06776    0.55390   0.122 0.902637
## PAY_65        -0.17413    0.90140  -0.193 0.846818
## PAY_66         0.50944    0.95489   0.534 0.593684
## PAY_67       -11.83717   348.38903  -0.034 0.972896
## PAY_68        24.98387  6829.59209   0.004 0.997081
## PAY_AMT2       -0.27007    0.06028  -4.480 7.47e-06 ***
## BILL_AMT3       0.28795    0.05843   4.928 8.29e-07 ***
## PAY_AMT1      -0.14662    0.03934  -3.727 0.000194 ***
## PAY_3-1         0.09293    0.11794   0.788 0.430741
## PAY_30         0.15585    0.12940   1.204 0.228411
## PAY_31       -13.48725   616.91616  -0.022 0.982558
## PAY_32         0.50002    0.13127   3.809 0.000140 ***
## PAY_33         0.27342    0.23206   1.178 0.238701
## PAY_34        -0.05606    0.46813  -0.120 0.904670
## PAY_35         0.49636    0.98075   0.506 0.612785
```

```

## PAY_36      15.31057  371.73748   0.041 0.967147
## PAY_37      -0.61752   1.30540  -0.473 0.636180
## PAY_38      -5.77319 6819.44810  -0.001 0.999325
## MARRIAGE1    1.75261   0.66165   2.649 0.008077 **
## MARRIAGE2    1.60347   0.66191   2.422 0.015415 *
## MARRIAGE3    1.89062   0.68302   2.768 0.005639 **
## EDUCATION1   12.78941  286.37668   0.045 0.964379
## EDUCATION2   12.83888  286.37668   0.045 0.964241
## EDUCATION3   12.79347  286.37668   0.045 0.964368
## EDUCATION4   12.08256  286.37706   0.042 0.966346
## EDUCATION5   11.54511  286.37687   0.040 0.967842
## EDUCATION6   12.45988  286.37710   0.044 0.965296
## SEX2         -0.14270   0.04010  -3.558 0.000373 ***
## PAY_AMT5     -0.08804   0.03253  -2.706 0.006809 **
## PAY_5-1      -0.03317   0.12973  -0.256 0.798219
## PAY_50        0.06647   0.14304   0.465 0.642155
## PAY_52        0.41577   0.15993   2.600 0.009329 **
## PAY_53        0.08700   0.30026   0.290 0.772008
## PAY_54       -0.26936   0.58798  -0.458 0.646874
## PAY_55        1.20656   1.07849   1.119 0.263249
## PAY_56       26.86203  655.12680   0.041 0.967294
## PAY_57       20.58158 6762.07341   0.003 0.997571
## PAY_58       29.61922 7061.38645   0.004 0.996653
## PAY_AMT3     -0.08144   0.04312  -1.889 0.058898 .
## BILL_AMT5    -0.09897   0.05642  -1.754 0.079381 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 20288  on 19197  degrees of freedom
## Residual deviance: 16612  on 19132  degrees of freedom
## AIC: 16744
##
## Number of Fisher Scoring iterations: 13

```

From this, we understand that;

- Used variables:
LIMIT_BAL, MARRIAGE, EDUCATION, SEX, PAY_0, PAY_3, PAY_4, PAY_5, PAY_6, PAY_AMT1, PAY_AMT2, PAY_AMT3, PAY_AMT5, BILL_AMT3, BILL_AMT5
- Important variables :
PAY_0(whether its values are -1, 1, 2, 3, 4), PAY_3(whether its values are 2), PAY_6(whether its values are 0), LIMIT_BAL, BILL_AMT3, PAY_AMT1, PAY_AMT2, SEX(whether a client is a female or not)

Then predict a result and make a confusion matrix.

```

##           Reference
## Prediction    0    1
##           0 3567  714
##           1  172  348

```


Other statistic metric are;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
glm step wise	0.8154551	0.9539984	0.3276836	0.640841

Can we improve the results? In the data exploration section, we have some insights which features are important to predict the outcome. As we have seen in the correlation matrix graph, we understand predictors from PAY_0 to PAY_6 seems to be important. Regarding BILL_AMT and PAY_AMT, the structure of each feature is very similar to one another.

Bearing these in mind, we narrow down features and leave out BILL_AMT2 to BILL_AMT6 and PAY_AMT2 to PAY_AMT6. Then we predict responses using glm function. The result;

```
##
## Call:
## glm(formula = DEFAULT ~ LIMIT_BAL + SEX + EDUCATION + MARRIAGE +
##      AGE + PAY_0 + PAY_2 + PAY_3 + PAY_4 + PAY_5 + PAY_6 + BILL_AMT1 +
##      PAY_AMT1, family = binomial(link = "logit"), data = train_set)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.3794  -0.5949  -0.5034  -0.3384   3.1966
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -15.12719   174.88974  -0.086 0.931072
## LIMIT_BAL    -0.28490    0.02788 -10.218 < 2e-16 ***
## SEX2         -0.13453    0.04044  -3.326 0.000880 ***
## EDUCATION1    11.92104   174.88847   0.068 0.945655
## EDUCATION2    11.97025   174.88847   0.068 0.945431
## EDUCATION3    11.91683   174.88848   0.068 0.945675
## EDUCATION4    11.19400   174.88909   0.064 0.948965
## EDUCATION5    10.64538   174.88878   0.061 0.951463
## EDUCATION6    11.56403   174.88917   0.066 0.947281
## MARRIAGE1     1.76158    0.66388   2.653 0.007967 **
## MARRIAGE2     1.61949    0.66410   2.439 0.014744 *
## MARRIAGE3     1.88473    0.68539   2.750 0.005962 **
## AGE           0.01870    0.02287   0.818 0.413503
## PAY_0-1       0.51387    0.13320   3.858 0.000114 ***
## PAY_00        -0.21310    0.14421  -1.478 0.139493
## PAY_01        0.86547    0.10539   8.212 < 2e-16 ***
## PAY_02        2.07649    0.13116  15.832 < 2e-16 ***
## PAY_03        2.00203    0.20917   9.571 < 2e-16 ***
## PAY_04        1.68415    0.37076   4.542 5.56e-06 ***
## PAY_05        0.79507    0.71780   1.108 0.268016
## PAY_06        1.15118    1.03858   1.108 0.267680
## PAY_07       -11.76836   535.41184  -0.022 0.982464
## PAY_08       -11.30452   265.22132  -0.043 0.966002
## PAY_2-1       -0.16685    0.13846  -1.205 0.228200
## PAY_20        -0.00238    0.16886  -0.014 0.988756
## PAY_21        -0.69540    0.66605  -1.044 0.296460
## PAY_22        -0.07342    0.14308  -0.513 0.607856
## PAY_23        0.06025    0.22145   0.272 0.785574
```

```

## PAY_24      -0.52935    0.40251   -1.315  0.188469
## PAY_25      0.38857    0.86617    0.449  0.653711
## PAY_26     13.14858   535.41217    0.025  0.980408
## PAY_27      NA         NA         NA     NA
## PAY_28     14.35828   640.93480    0.022  0.982127
## PAY_3-1     0.12046    0.13251    0.909  0.363342
## PAY_30      0.20535    0.15279    1.344  0.178946
## PAY_31     -11.66599   375.43547   -0.031  0.975211
## PAY_32      0.59556    0.15522    3.837  0.000125 ***
## PAY_33      0.53110    0.27264    1.948  0.051419 .
## PAY_34      0.10310    0.54044    0.191  0.848703
## PAY_35      0.36194    1.07229    0.338  0.735710
## PAY_36     14.54758   265.22257    0.055  0.956258
## PAY_37     -0.50584    1.30903   -0.386  0.699182
## PAY_38     -5.24311  3387.97494   -0.002  0.998765
## PAY_4-1     -0.22651    0.13443   -1.685  0.091987 .
## PAY_40     -0.13770    0.14837   -0.928  0.353339
## PAY_41      0.69352   653.92392    0.001  0.999154
## PAY_42      0.12390    0.15874    0.781  0.435086
## PAY_43      0.08418    0.30941    0.272  0.785578
## PAY_44      0.49217    0.54245    0.907  0.364245
## PAY_45     -2.28211    1.03500   -2.205  0.027458 *
## PAY_46     -28.12876   352.32456   -0.080  0.936366
## PAY_47     -6.63472  3345.42735   -0.002  0.998418
## PAY_48     -33.30716  3387.97481   -0.010  0.992156
## PAY_5-1     -0.01575    0.12968   -0.121  0.903333
## PAY_50      0.03493    0.14258    0.245  0.806452
## PAY_52      0.38294    0.15975    2.397  0.016523 *
## PAY_53      0.02295    0.30108    0.076  0.939227
## PAY_54     -0.30786    0.59119   -0.521  0.602549
## PAY_55      1.08232    1.08928    0.994  0.320415
## PAY_56     25.12583   417.01832    0.060  0.951956
## PAY_57     19.18636  3345.40123    0.006  0.995424
## PAY_58     27.66820  3569.10792    0.008  0.993815
## PAY_6-1     -0.25977    0.09898   -2.624  0.008680 **
## PAY_60     -0.46644    0.10605   -4.399  1.09e-05 ***
## PAY_62     -0.11216    0.12354   -0.908  0.363942
## PAY_63      0.74573    0.30314    2.460  0.013892 *
## PAY_64      0.07752    0.55252    0.140  0.888414
## PAY_65     -0.30207    0.91535   -0.330  0.741396
## PAY_66      0.44401    0.96504    0.460  0.645446
## PAY_67     -10.90764   211.76602   -0.052  0.958921
## PAY_68     23.59520  3398.37009    0.007  0.994460
## BILL_AMT1    0.13283    0.02741    4.845  1.26e-06 ***
## PAY_AMT1     -0.16739    0.04070   -4.112  3.92e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 20288 on 19197 degrees of freedom
## Residual deviance: 16657 on 19126 degrees of freedom
## AIC: 16801
##

```

```
## Number of Fisher Scoring iterations: 12
```

From this, we understand that;

- Important variables :
LIMIT_BAL, SEX(whether a client is a female or not), PAY_0(whether its values are -1, 1, 2, 3, 4),
PAY_3(whether its values are 2), PAY_6(whether its values are 0), BILL_AMT1, PAY_AMT1

Then fit the model to the validation set and show its confusion matrix.

```
##           Reference
## Prediction    0    1
##           0 3567  713
##           1  172  349
```

Other statistic metric are;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
glm fewer features	0.8156634	0.9539984	0.3286252	0.6413118

Comparing the two logistic regression models, the later model in which some seemingly irrelevant features are left out performed better in terms of both accuracy and balanced accuracy.

3 Decision tree

A decision tree is a familiar and intuitive method. For example, if you feel sick, a physician may ask you questions following a tree chart. You are asked whether it is yes or no at each node of the chart. After following the nodes, the doctor gets your health outcome.

Decision tree algorithm recursively divides the dataset into partitions with similar values for the outcome. A metric is used to choose the partitions. One of decision tree R function “rpart” uses Gini index.

$$\text{Gini}(j) = \sum_{k=1}^K \hat{p}_{j,k}(1 - \hat{p}_{j,k})$$

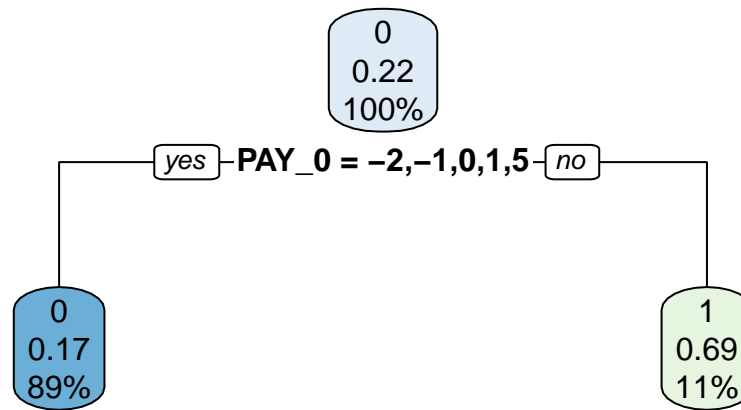
$\hat{p}_{j,k}$ as the proportion of observations in partition j that are of class k .¹⁷ If there are lots of classes in the partition, Gini index, also called Gini impurity, increases up to 1. On the other hand, if there are few classes in the partition it decreases down to 0. If you predict perfectly, the index is 0.

We use “rpart” function to make a model, then predict the outcome.

```
## n= 19198
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
## 1) root 19198 4246 0 (0.7788311 0.2211689)
##    2) PAY_0=-2,-1,0,1,5 17180 2849 0 (0.8341676 0.1658324) *
##    3) PAY_0=2,3,4,6,7,8 2018  621 1 (0.3077304 0.6922696) *
```

¹⁷Irizarry, Rafael A, “Introduction to Data Science: Data Analysis and Prediction Algorithms with R, 31.10.4 Classification (decision) trees”. *Internet archive*, <https://rafalab.github.io/dsbook/examples-of-algorithms.html#classification-decision-trees>

We draw a decision tree based on this model.



This model uses only PAY_0 to make partitions. This means the credit company need to check a client's PAY_0. If the value is -2, -1, 0, 1, 5, they are likely to repay, otherwise, they default on the debt.

Then predict the outcome, and show its confusion matrix.

```
##           Reference
## Prediction    0    1
##           0 3597  736
##           1  142  326
```

Other statistic metrics;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
rpart	0.8171214	0.9620219	0.306968	0.634495

Compared with the previous logistic regression models, its accuracy has improved but its balanced accuracy become worse.

We try to improve this result using other function “train” in caret package. When we use “train”, it conducts “cross validation” and finds optimal parameters in the decision tree algorithm.

One of the cross validation methods is k-fold cross validation. It splits the dataset into k “folds” randomly. For each group, it takes the group as a test dataset, and takes other groups as a training dataset. Then it fits a model on the training set and evaluates it on the test set. This procedure is repeated until every K-fold serve as the test set. Finally it summarizes the performance of the model as the average of these procedures.

“ModelLookup” function in caret package tells us what parameters in rpart we can tune.

```
modelLookup("rpart")
```

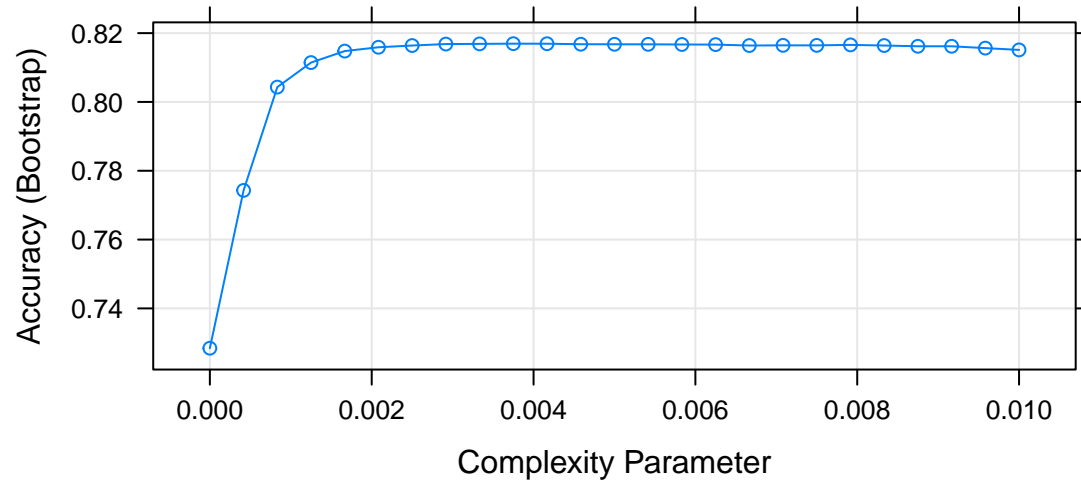
```
##   model parameter          label forReg forClass probModel
## 1 rpart          cp Complexity Parameter    TRUE    TRUE    TRUE
```

Cp stands for complexity parameter. It determines the complexity of the model. As its value gets smaller, the algorithm produces more trees. In the previous decision tree model we fit, cp is 0.01. We want to know whether cp values less than 0.01 will improve the performance.

We train the dataset using “train” function. This function conducts 10 folds cross validation as a default.

```
## CART
##
## 19198 samples
##    23 predictor
##    2 classes: '0', '1'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 19198, 19198, 19198, 19198, 19198, 19198, ...
## Resampling results across tuning parameters:
##
##    cp            Accuracy    Kappa
## 0.0000000000    0.7284228    0.2235331
## 0.0004166667    0.7742936    0.2859917
## 0.0008333333    0.8043130    0.3325194
## 0.0012500000    0.8114391    0.3445758
## 0.0016666667    0.8147804    0.3520765
## 0.0020833333    0.8158908    0.3523277
## 0.0025000000    0.8163979    0.3484179
## 0.0029166667    0.8168234    0.3458840
## 0.0033333333    0.8168740    0.3431768
## 0.0037500000    0.8169363    0.3425090
## 0.0041666667    0.8169238    0.3412075
## 0.0045833333    0.8167714    0.3405336
## 0.0050000000    0.8167598    0.3402173
## 0.0054166667    0.8167432    0.3408068
## 0.0058333333    0.8166985    0.3413150
## 0.0062500000    0.8166645    0.3414413
## 0.0066666667    0.8163935    0.3386679
## 0.0070833333    0.8164501    0.3381807
## 0.0075000000    0.8164443    0.3386962
## 0.0079166667    0.8165811    0.3373467
## 0.0083333333    0.8163949    0.3357974
## 0.0087500000    0.8161803    0.3344378
## 0.0091666667    0.8161803    0.3344378
## 0.0095833333    0.8156400    0.3303031
## 0.0100000000    0.8151237    0.3273365
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was cp = 0.00375.
```

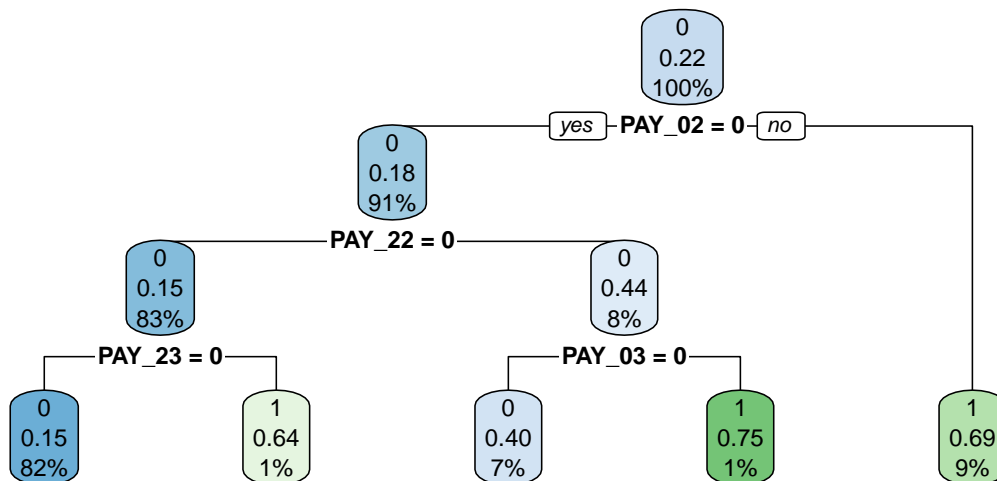
Plot cp.



Optimal cp value is

```
##          cp
## 10 0.00375
```

Then draw a decision tree.



As PAY variables are factored and turned into a dummy variable, a node “PAY_02”, for example, is asking whether a client’s PAY_0 is “-1” or not. “-1” is the second level in PAY_0 original column.

Fit the model to the validation set and show confusion matrix.

```
##          Reference
## Prediction    0    1
##           0 3587  730
##           1  152  332
```

Other statistic metrics;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
rpart tuned	0.8162883	0.9593474	0.3126177	0.6359826

Comparing two decision models, the tuned model has improved in terms of specificity and balanced accuracy. But both performed worse than logistic regression using fewer predictors in terms of balanced accuracy.

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
rpart	0.8171214	0.9620219	0.3069680	0.6344950
rpart tuned	0.8162883	0.9593474	0.3126177	0.6359826

4 Random forest

Decision tree algorithm is easy to understand as it follow human decision making process. However, it “can easily over-train” and “is not very flexible.”¹⁸ To improve this, we introduce random forest algorithm.

Random forest is one of the most popular machine learning algorithms. As its name implies, it makes multiple trees “**randomly** different” from a dataset, then produce a final prediction based on the average prediction of “combination of trees”, namely “**forest**”.¹⁹

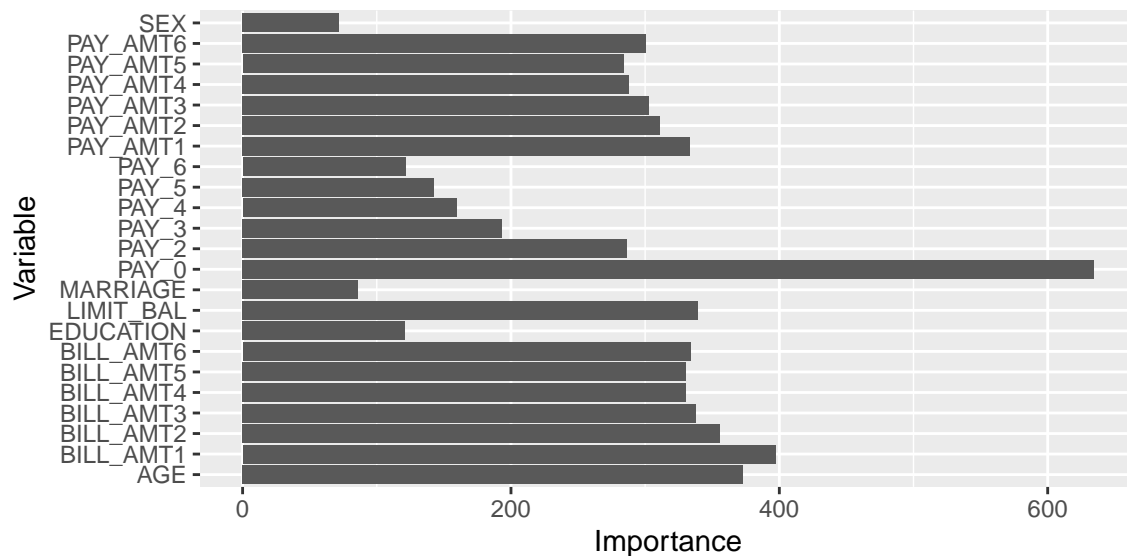
In this paper, we use “ranger” function. And show the model details;

```
## Ranger result
##
## Call:
## ranger(formula = DEFAULT ~ ., data = train_set, importance = "impurity",      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:              impurity
## Splitrule:                             gini
## OOB prediction error:                  18.30 %
```

Plot the variables’ importance;

¹⁸Irizarry, Rafael A, “Introduction to Data Science: Data Analysis and Prediction Algorithms with R, 31.10.4 Classification (decision) trees”. *Internet archive*, <<https://rafalab.github.io/dsbook/examples-of-algorithms.html#classification-decision-trees>>

¹⁹Irizarry, Rafael A, “Introduction to Data Science: Data Analysis and Prediction Algorithms with R, 31.11 Random forests”. *Internet archive*, <https://rafalab.github.io/dsbook/examples-of-algorithms.html#random-forests>



As we saw in the decision tree model, PAY_0 is considerably important compared with other features. As for other features, BILL_AMT1, PAY_AMT1, and LIMIT_BAL, and AGE are important.

Then fit the model to the validation set and show its confusion matrix.

```
##           Reference
## Prediction    0    1
##           0 3564  698
##           1  175  364
```

Other statistic metrics;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
random forest	0.8181629	0.953196	0.3427495	0.6479728

The results has been much improved, especially specificity, compared to other models. Can we improve the result further? Using “modelLookup”, we can find parameters which can be tuned.

```
##   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
```

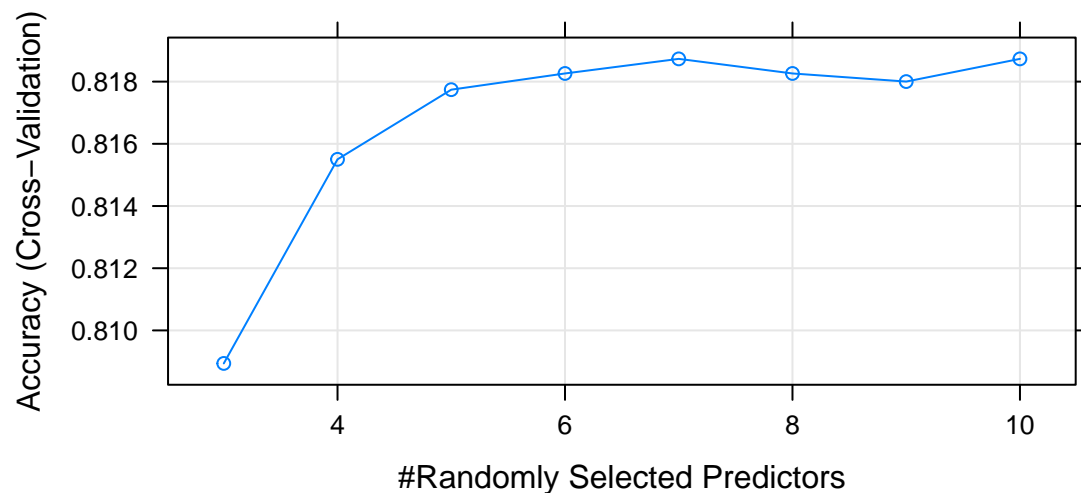
Mtry is a number of variables to possibly split at in each node. Splitrule is a splitting rule. Min.node.size is a minimal node size. In the previous random forest model, its mtry is 4, we used Gini index as a splitting rule, and minimal node size is 1. In the following model, we compare mtry ranging from 3 to 10 in terms of accuracy.

```
## Random Forest
##
## 19198 samples
##   23 predictor
##   2 classes: '0', '1'
```



```
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 17278, 17279, 17278, 17278, 17279, 17278, ...
## Resampling results across tuning parameters:
##
##   mtry  Accuracy  Kappa
##   3     0.8089386 0.2799217
##   4     0.8155016 0.3355440
##   5     0.8177415 0.3582670
##   6     0.8182623 0.3620446
##   7     0.8187315 0.3660563
##   8     0.8182623 0.3655969
##   9     0.8180020 0.3658315
##  10     0.8187312 0.3700717
##
## Tuning parameter 'splitrule' was held constant at a value of gini
##
## Tuning parameter 'min.node.size' was held constant at a value of 1
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were mtry = 7, splitrule = gini
## and min.node.size = 1.
```

Plot mtry and accuracy.

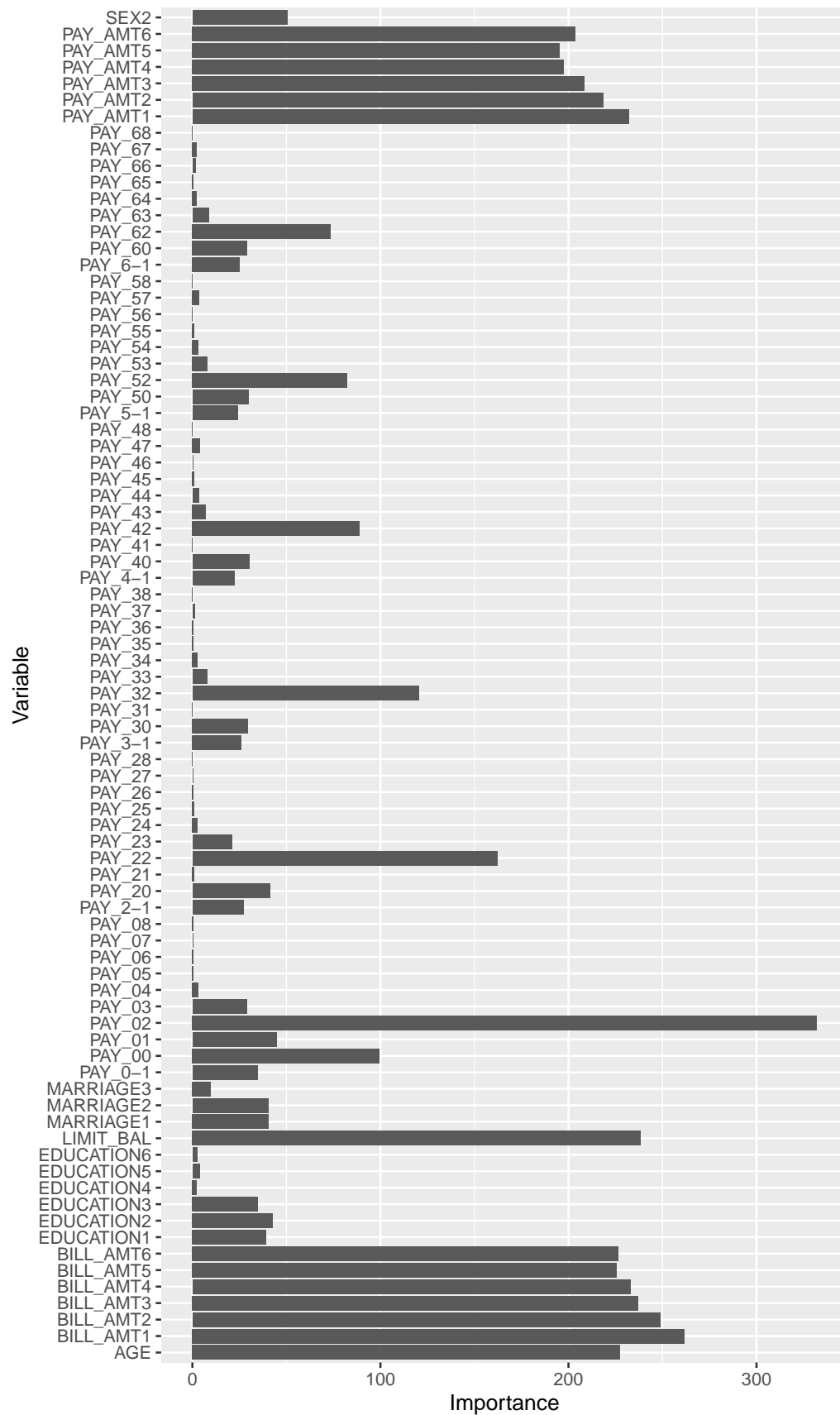


The best mtry is;

```
rf_cv_md1$bestTune
```

```
##   mtry splitrule min.node.size
## 5     7      gini             1
```

Plot the variables' importance;



Contrary to the previous model, not PAY_0 but PAY_02's importance is outstanding.

Then fit the model to the validation set and show confusion matrix;

```
##           Reference
## Prediction    0    1
##           0 3578  725
##           1  161  337
```

Other statistic metrics;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
random forest tuned	0.8154551	0.9569404	0.3173258	0.6371331

Comparing these two random forest models, the results of the tuned model gets worse than the default model in terms of both accuracy and balanced accuracy. There occurred over fitting in the tuned model.

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
random forest	0.8181629	0.9531960	0.3427495	0.6479728
random forest tuned	0.8154551	0.9569404	0.3173258	0.6371331

Evaluation

We look at the table in which the results of our 6 models are shown.

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
glm step wise	0.8154551	0.9539984	0.3276836	0.6408410
glm fewer features	0.8156634	0.9539984	0.3286252	0.6413118
rpart	0.8171214	0.9620219	0.3069680	0.6344950
rpart tuned	0.8162883	0.9593474	0.3126177	0.6359826
random forest	0.8181629	0.9531960	0.3427495	0.6479728
random forest tuned	0.8154551	0.9569404	0.3173258	0.6371331

Among them, the best performance in terms of both accuracy and balanced accuracy is produced by “**random forest default model**”. As we mentioned before, **PAY_0** is the most important variables compared with others. Then follows **BILL_AMT1**, **PAY_AMT1**, and **LIMIT_BAL**, and **AGE**.

We fit this model to the test set and extract final evaluation.

Final statistic metrics are as follows;

method	Accuracy	Sensitivity	Specificity	Balanced_Accuracy
final random forest	0.8226962	0.9621228	0.3320783	0.6471006

Fitting the model to the test set, we get **accuracy 0.8226962** and **balanced accuracy 0.6471006**. This results outperformed considerably our first guess (guessing all outcomes are 0), accuracy 0.7788 and balanced accuracy 0.5.

Conclusion

In this paper, we collected an open data, “Default of Credit Card Clients Dataset” and explored its information seeking for insights which contributed to the model making. We used three methods, logistic regression, decision tree, and random forest. In each method, we tried to improve the results by changing variables and tuning hyperparameters. As its outcomes are imbalanced, we used balanced accuracy to compare the models results.

Random forest algorithm showed the best performance. Fitting it to the test set, we got accuracy 0.8226962 and balanced accuracy 0.6471006. But when tuned, the model showed over-fitted results. The logistic regression models outperform the decision tree models. As well as the random forest models, the default decision tree model showed better performance than the tuned model.

As we use rather simple methods in the paper, we assume limitations are revealed. Firstly, we used variables almost as they are. But some variables, for example, PAY, PAY_AMT, and BILL_AMT, seem to have some relations. If we tweak these variables, this would produce more accurate results.

Moreover, these variables, PAY, PAY_AMT, and BILL_AMT, are historical data over a half year. If we can look into these historical changes further, we might be able to acquire more accurate results.

Finally, from a business point of view, overall economic condition may play a very important role in predicting people’s financial behavior. In the late 2000s, many people could not afford to pay their loans. On the other hand, when economy is booming, people repay the debt easily. Of course, our dataset did not have such information, but if we want to make a more credible model, we need to consider what kind of data is offered.