

# Credit Card Lead Prediction

Yucong Chen, Haocheng Liao, Kaiqi Peng, Danlei Wang, Wenting Yang

## Contents

<b>1. Introduction</b>	<b>2</b>
<b>2. Description of the Data</b>	<b>3</b>
Histograms & Fitted Distribution & Pie Charts . . . . .	4
Correlation . . . . .	13
Statistical Summary . . . . .	14
<b>3. Data Analysis / Model</b>	<b>15</b>
OLS Model . . . . .	15
Probit Model . . . . .	17
Logit Model . . . . .	18
Probit Training . . . . .	22
Logit Training . . . . .	23
Tests for Logit Model . . . . .	25
Multinomial Logit Model . . . . .	28
Instrumental Model . . . . .	29
Identify Potential Customers . . . . .	34
<b>4. Conclusion</b>	<b>45</b>
<b>5. Future Work</b>	<b>46</b>
<b>6. References</b>	<b>47</b>

# 1. Introduction

Happy Customer Bank is a mid-sized private bank that operates a variety of banking products, such as savings accounts, current accounts, investment products, credit products, etc. It also cross sells products to existing customers through different means of communication. In this case, the bank hopes to cross sell its credit cards to existing customers.

Our goal of this project is to help Happy Customer Bank to conduct a credit card lead prediction - to identify customers who have higher intention to recommended credit cards from a qualified group identified by the bank - through data analysis and models.

## 2. Description of the Data

Happy Customer Bank has collected information from 105,312 customers and put their information into a dataset named "test.csv". For each observation in the dataset, there are eight variables that are the most relevant for our models: "Gender", "Age", "Occupation", "Channel\_Code", "Vintage", "Credit\_Product", "Avg\_Account\_Balance", and "Is\_Active". The variable "ID" is irrelevant because it only serves as an identification for customers and shows us that every observation is unique. For "Region\_Code", we figured that as there are too many values, it might be distractive for our models. Thus, these two variables will be deleted during our data analysis and modeling, but will be used after we move on to identify potential customers.

In this section, we will give an overview of our variables, visualize them, and encode the categorical variables into numerical ones to prepare for our models.

For the eight relevant variables, "Age", "Vintage", and "Avg\_Account\_Balance" are numerical variables. "Age" refers to the age of the customer. "Vintage" refers to the vintage for the customer in months. "Avg\_Account\_Balance" refers to the average account balance in the last 12 months.

"Gender", "Occupation", "Channel\_Code", "Credit\_Product", and "Is\_Active" are categorical variables, and we will encode them into numbers. Among them, "Gender" and "Is\_Active" take just two values, so they are dummy variables. "Gender" refers to the gender of the customer and it takes the value of "Female" or "Male". We will let "Female"=1 and "Male"=0. "Is\_Active" describes if the customer is active in the last three months and it takes the value of "Yes" or "No". We will let "Yes"=1 and "No"=0. We will take "Is\_Active" as our dependent variable.

"Occupation" refers to the job of the customer. There are four different kinds of jobs: "Salaried", "Self\_Employed", "Entrepreneur", and "Other". We will let "Other"=1, "Salaried"=2, "Self\_Employed"=3, and "Entrepreneur"=4.

"Channel\_Code" is the encoded acquisition channel code for the customer. There are also four different kinds: "X1", "X2", "X3", and "X4". We will model "X1"=1, "X2"=2, "X3"=3, and "X4"=4.

"Credit\_Product" refers to if the customer has any active credit product, such as home loan, personal loan, credit card, etc. It takes the value of "Yes", "No", or "N/A". We will model "No"=0, "Yes"=1, "N/A"=2.

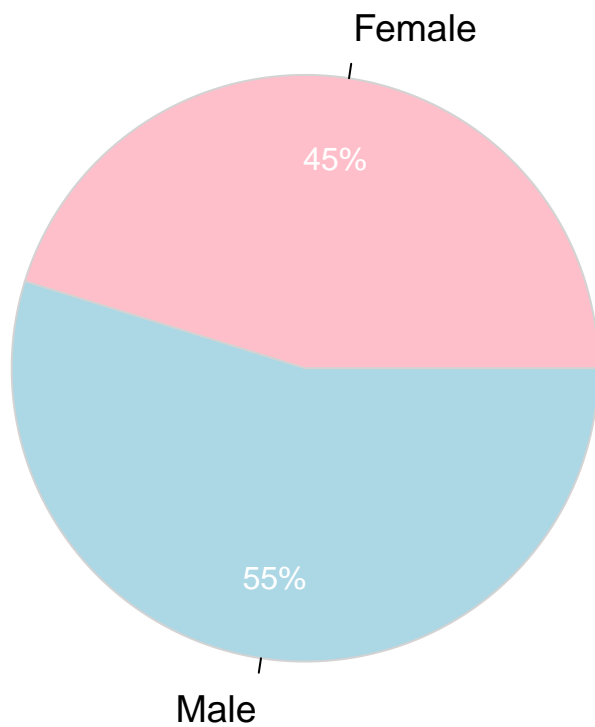
```
# Load the data test.csv
data <- read_csv("test.csv")

## Parsed with column specification:
## cols(
##   ID = col_character(),
##   Gender = col_character(),
##   Age = col_double(),
##   Region_Code = col_character(),
##   Occupation = col_character(),
##   Channel_Code = col_character(),
##   Vintage = col_double(),
##   Credit_Product = col_character(),
##   Avg_Account_Balance = col_double(),
##   Is_Active = col_character()
## )
```

## Histograms & Fitted Distribution & Pie Charts

```
# Gender
attach(data)
PieChart(Gender, hole = 0, values = "%", data = data,
         fill = c("pink", "lightblue"), main = "Gender Pie Chart", quiet = TRUE)
```

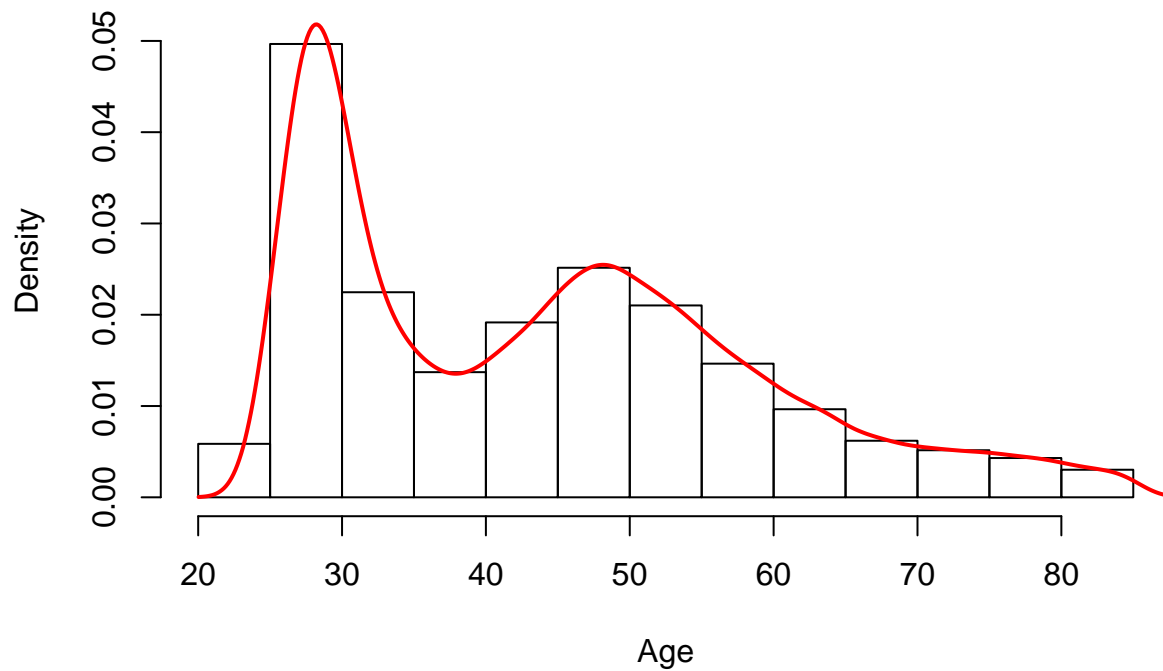
### Gender Pie Chart



The pie chart of "Gender" shows that 55% of the customers are male while 45% of them are female.

```
# Age
hist(Age, prob = TRUE, ylim = c(0, max(density(Age)$y)))
lines(density(Age), lwd = 2, col = "red")
```

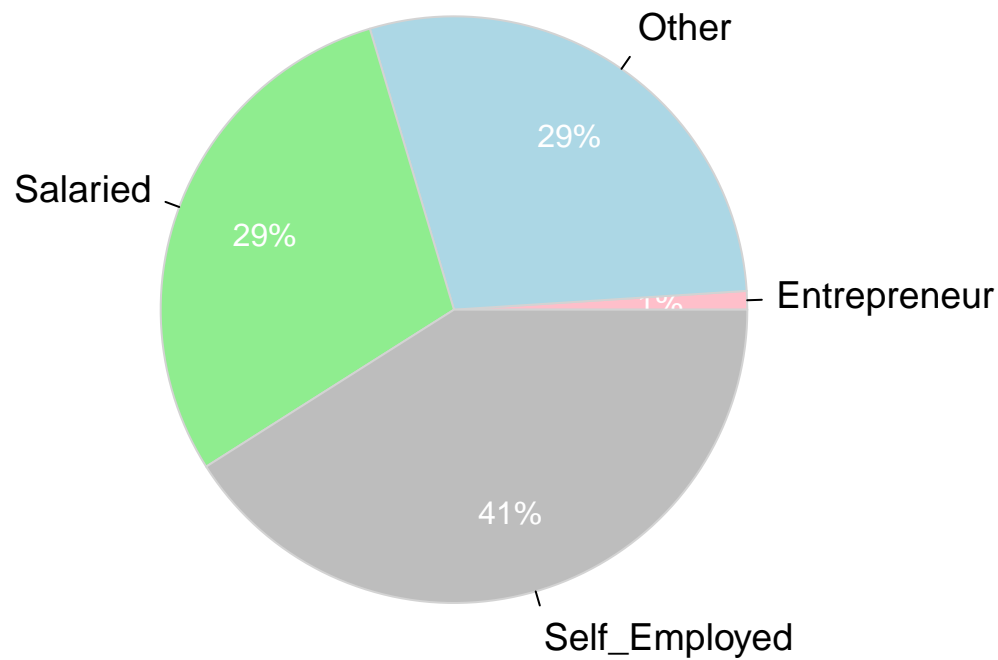
**Histogram of Age**



The histogram of “Age” shows that most people’s age falls within the interval [20, 80]. However, people around 25 to 30 years old are the largest client crowd in this dataset.

```
# Occupation
PieChart(Occupation, hole = 0, values = "%", data = data,
         fill = c("pink", "lightblue", "lightgreen", "grey"),
         main = "Occupation Pie Chart", quiet = TRUE)
```

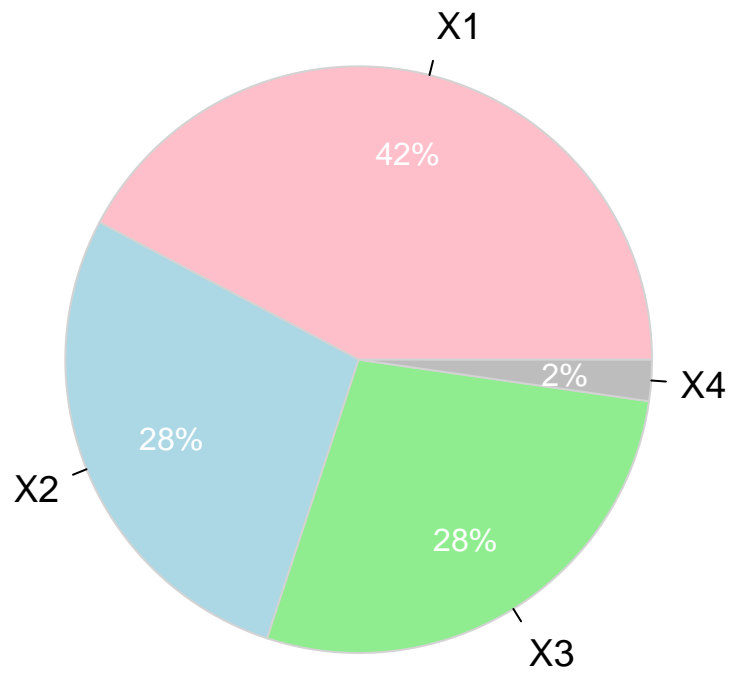
## Occupation Pie Chart



The pie chart of "Occupation" shows that 41% of customers are self-employed, 29% are salaried, 1% are entrepreneurs, and 29% have other occupations.

```
# Channel_Code
PieChart(Channel_Code, hole = 0, values = "%", data = data,
         fill = c("pink", "lightblue", "lightgreen", "grey"),
         main = "Channel Code Pie Chart", quiet = TRUE)
```

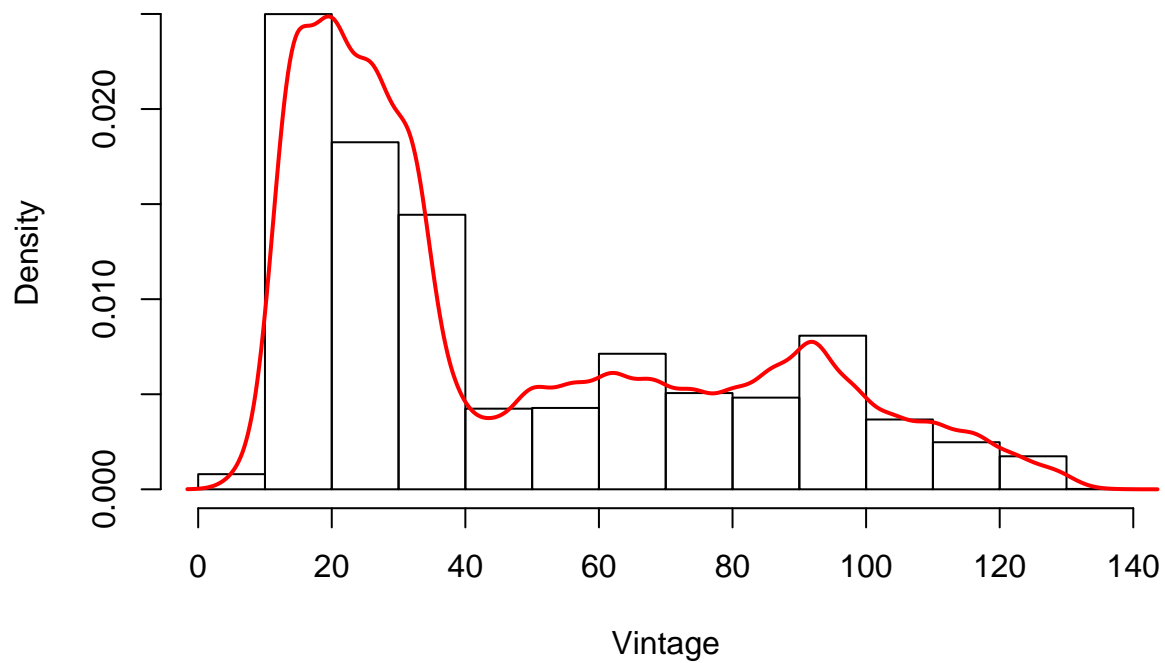
## Channel Code Pie Chart



The pie chart of "Channel\_Code" shows that 42% of channel codes are X1, 28% are X2, 28% are X3, and 2% are X4.

```
# Vintage
hist(Vintage, prob = TRUE, ylim = c(0, max(density(Vintage)$y)))
lines(density(Vintage), lwd = 2, col = "red")
```

**Histogram of Vintage**

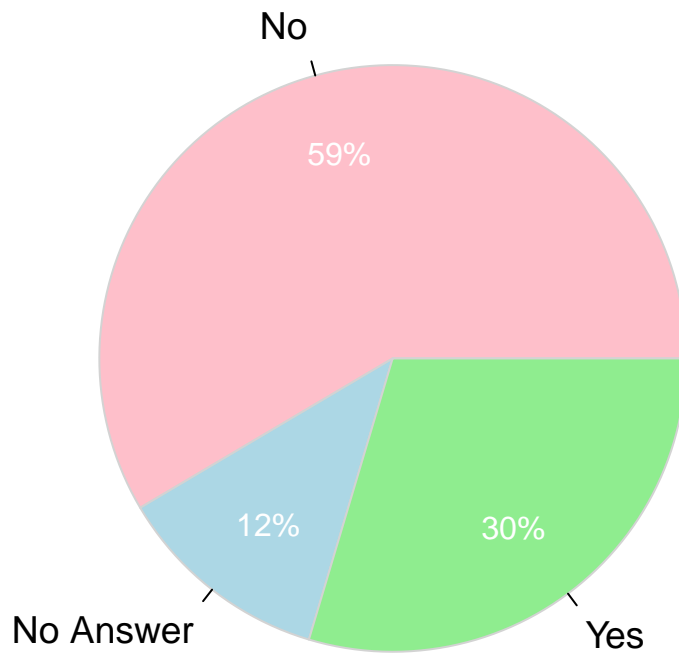


The histogram of “Vintage” indicates that the vintage for most clients falls within the interval [0, 140].



```
# Credit_Product
# We change "N/A" to "No Answer" because N/As are not counted in pie charts
data$Credit_Product[is.na(data$Credit_Product)] <- "No Answer"
PieChart(Credit_Product, hole = 0, values = "%", data = data,
         fill = c("pink", "lightblue", "lightgreen"),
         main = "Credit Product Pie Chart", quiet = TRUE)
```

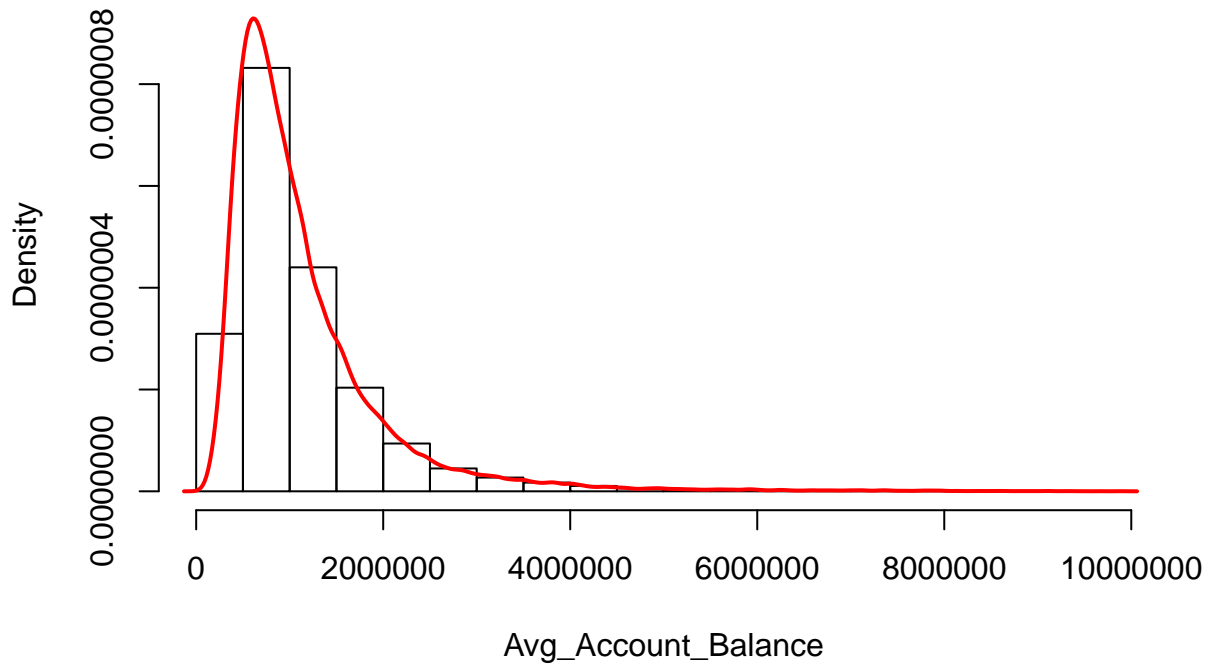
## Credit Product Pie Chart



The pie chart of "Credit\_Product" shows that 30% of the customers answered "Yes", 59% answered "No", and 12% preferred not to answer.

```
# Avg_Account_Balance
hist(Avg_Account_Balance, prob = TRUE, ylim = c(0, max(density(Avg_Account_Balance)$y)))
lines(density(Avg_Account_Balance), lwd = 2, col = "red")
```

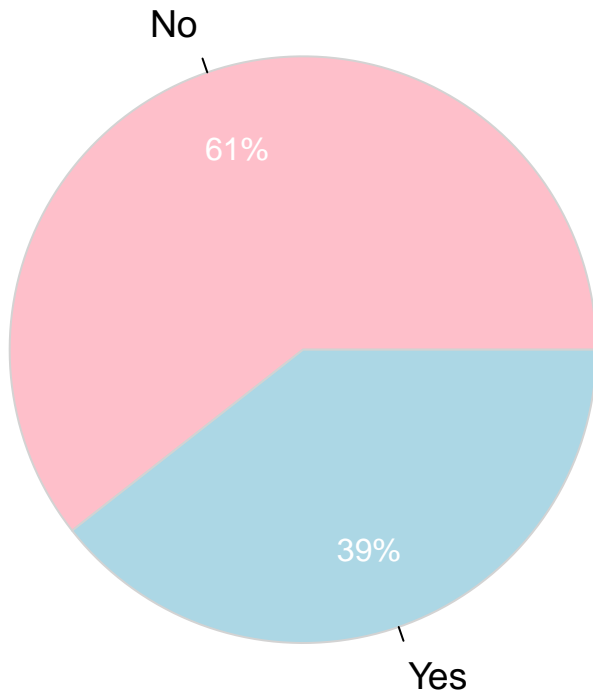
### Histogram of Avg\_Account\_Balance



The histogram of “Avg\_Account\_Balance” shows that the average account balance of clients falls within the interval [0, 4000000], while only a few clients have a balance that is above 4,000,000. The financial resources of the clients vary greatly.

```
# Is_Active
PieChart(Is_Active, hole = 0, values = "%", data = data,
        fill = c("pink", "lightblue"),
        main = "Active Pie Chart", quiet = TRUE)
```

## Active Pie Chart



The pie chart of "Is\_Active" shows that 39% of the customers answered "Yes" while 61% of them answered "No".

Now we have visualized histograms and pie charts, we can encode the categorical variables into numerical ones for a correlation plot and models.

```
# Encode the categorical variables
data$Gender[data$Gender == "Female"] <- 1
data$Gender[data$Gender == "Male"] <- 0
data$Gender <- gsub(",", "", data$Gender)
data$Gender <- as.numeric(as.character(data$Gender))

data$Channel_Code[data$Channel_Code == "X1"] <- 1
data$Channel_Code[data$Channel_Code == "X2"] <- 2
data$Channel_Code[data$Channel_Code == "X3"] <- 3
data$Channel_Code[data$Channel_Code == "X4"] <- 4
data$Channel_Code <- gsub(",", "", data$Channel_Code)
data$Channel_Code <- as.numeric(as.character(data$Channel_Code))

data$Occupation[data$Occupation == "Other"] <- 1
data$Occupation[data$Occupation == "Salaried"] <- 2
data$Occupation[data$Occupation == "Self_Employed"] <- 3
data$Occupation[data$Occupation == "Entrepreneur"] <- 4
data$Occupation <- gsub(",", "", data$Occupation)
data$Occupation <- as.numeric(as.character(data$Occupation))

data$Credit_Product[data$Credit_Product == "No Answer"] <- 1
data$Credit_Product[data$Credit_Product == "Yes"] <- 2
data$Credit_Product[data$Credit_Product == "No"] <- 3
data$Credit_Product <- gsub(",", "", data$Credit_Product)
data$Credit_Product <- as.numeric(as.character(data$Credit_Product))

data$Is_Active[data$Is_Active == "Yes"] <- 1
data$Is_Active[data$Is_Active == "No"] <- 0
data$Is_Active <- gsub(",", "", data$Is_Active)
data$Is_Active <- as.numeric(as.character(data$Is_Active))

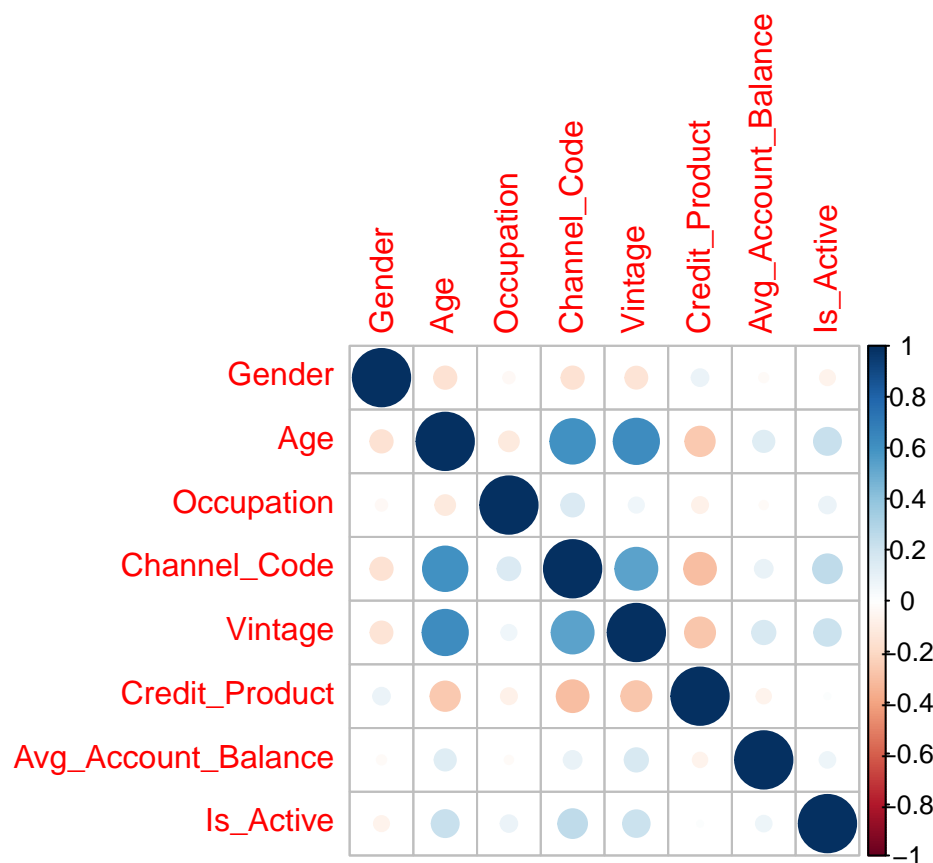
data <- data[, -1] # Delete column "ID"
data <- data[, -3] # Delete column "Region_Code"

sapply(data, class) # Check all columns are numeric now
```

##	Gender	Age	Occupation	Channel_Code
##	"numeric"	"numeric"	"numeric"	"numeric"
##	Vintage	Credit_Product	Avg_Account_Balance	Is_Active
##	"numeric"	"numeric"	"numeric"	"numeric"

## Correlation

```
corrplot(cor(data))
```



We have a correlation matrix here that shows the correlation coefficients between our variables. Each cell in the matrix indicates the correlation between two specific variables. The color index on the right measures the level of association between the two variables from -1 to 1. Notice that the correlation coefficients along the diagonal of the matrix are all dark blue (=1) because each variable has a perfect uphill linear relationship with itself.

It is quite clear that some variables are very strongly and positively correlated. For example, "Channel\_Code" & "Age", "Vintage" & "Age". Their values of correlation are all larger or close to 0.8 which means there is a very strong positive linear relationship between each two.

Furthermore, we have some other variables that are strongly and positively correlated, but not as much as above. For instance, "Channel\_Code" & "Vintage". Their values of correlation are close to 0.7 which means there is a strong uphill linear relationship between each two.

The matrix also shows us some variables that are basically not correlated, or not correlated at all with another. For examples, "Gender" & "Age", "Gender" & "Channel\_Code", "Gender" & "Vintage", "Gender" & "Occupation", "Gender" & "Credit\_Product", "Gender" & "Avg\_Account\_Balance", "Gender" & "Is\_Active", "Age" & "Occupation", "Age" & "Avg\_Account\_Balance", "Occupation" & "Channel\_Code", "Occupation" & "Vintage", "Occupation" & "Credit\_Product", "Occupation" & "Avg\_Account\_Balance", "Occupation" & "Is\_Active", "Channel\_Code" & "Avg\_Account\_Balance", "Vintage" & "Occupation", "Channel\_Code" & "Avg\_Account\_Balance". The correlations between them are close or equal to zero which indicates very small or no linear correlation at all.

Last but not least, there are variables that have a weak negative linear relationship with another variable: "Age" & "Credit\_Product", "Age" & "Is\_Active", "Channel\_Code" & "Credit\_Product", "Channel\_Code" &

“Is\_Active”, “Vintage” & “Credit\_Product”, “Vintage” & “Is\_Active”. The correlations are close to -0.3 which prove that they are weakly negatively correlated with the other.

## Statistical Summary

```
summary(data)
```

##	Gender	Age	Occupation	Channel_Code
##	Min. :0.0000	Min. :24.00	Min. :1.000	Min. :1.000
##	1st Qu.:0.0000	1st Qu.:30.00	1st Qu.:1.000	1st Qu.:1.000
##	Median :0.0000	Median :43.00	Median :2.000	Median :2.000
##	Mean :0.4521	Mean :43.87	Mean :2.144	Mean :1.901
##	3rd Qu.:1.0000	3rd Qu.:54.00	3rd Qu.:3.000	3rd Qu.:3.000
##	Max. :1.0000	Max. :85.00	Max. :4.000	Max. :4.000
##	Vintage	Credit_Product	Avg_Account_Balance	Is_Active
##	Min. : 7.00	Min. :1.000	Min. : 22597	Min. :0.0000
##	1st Qu.: 20.00	1st Qu.:2.000	1st Qu.: 603982	1st Qu.:0.0000
##	Median : 32.00	Median :3.000	Median : 896634	Median :0.0000
##	Mean : 46.84	Mean :2.466	Mean :1134195	Mean :0.3942
##	3rd Qu.: 73.00	3rd Qu.:3.000	3rd Qu.:1371598	3rd Qu.:1.0000
##	Max. :135.00	Max. :3.000	Max. :9908858	Max. :1.0000

From our statistical summary, we notice that the minimum of “Age” is 24 while the maximum is 85, so there are customers of all ages [24, 85]. The median is 43.00 and the mean is 43.87 (i.e., there is not much difference), which means that our data is unbiased. The minimum of “Vintage” is 7 while the maximum is 135, and the mean is 46.84. The range is very wide in this case, which means that the bank has investigated both new and regular customers. The minimum of “Avg\_Account\_Balance” is 22,597 while the maximum is 9,908,858, and the mean is 1,134,195. The range for average account balance is wide as well.

For our dummy variables “Gender” and “Is\_Active”, their maximum is 1 and minimum is 0, and their means are 0.4521 and 0.3942 respectively, which indicates the distribution between females and males, and active and inactive customers are unbiased.

The minimum of “Occupation” is 1 while the maximum is 4, and the mean of it is 2.144. Note that we have four possible numbers for “Occupation” with respect to salaried, self-employed, entrepreneur, and other (“Other”=1, “Salaried”=2, “Self\_Employed”=3, “Entrepreneur”=4). This indicates that our occupation distribution is close to even.

The minimum of “Channel\_Code” is 1 while the maximum is 4, and the mean of it is 1.901. Note that we have four possible numbers for this variable with respect to X1, X2, X3, and X4 (“X1”=1, “X2”=2, “X3”=3, “X4”=4).

The minimum of “Credit\_Product” is 1 while the maximum is 3, and the mean of it is 2.466. Note that we have three possible numbers for this variable with respect to “Yes”, “No”, and “NA” (“NA”=1, “Yes”=2, “NO”=3).

### 3. Data Analysis / Model

In this section, we will predict customers' interests in recommended credit cards. Since this involves individuals making "either-or" choices, represented by the binary variable "Is\_Active", we will start with the OLS model, the probit model, and the logit model. Then, we will divide the data into training and testing for evaluation. In addition, multinomial logit model and instrumental variables will be used for possible improvements. We will also provide a series of tests to find the best model and use that model to identify potential customers for Happy Customer Bank.

#### OLS Model

```
OLS.Model <- lm(Is_Active ~ Gender + Age + Occupation + Channel_Code + Vintage
               + Credit_Product + Avg_Account_Balance, data = data)
summary(OLS.Model)
```

```
##
## Call:
## lm(formula = Is_Active ~ Gender + Age + Occupation + Channel_Code +
##     Vintage + Credit_Product + Avg_Account_Balance, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.8090 -0.3663 -0.2310  0.4848  0.9968
##
## Coefficients:
##              Estimate      Std. Error t value
## (Intercept)  -0.299549257025  0.009612541207 -31.162
## Gender       -0.022986164668  0.002928938041  -7.848
## Age           0.003414044624  0.000141362426  24.151
## Occupation    0.042592111664  0.001798125383  23.687
## Channel_Code  0.091150174537  0.002196421548  41.499
## Vintage       0.001197522106  0.000059795931  20.027
## Credit_Product 0.084426992521  0.002185858181  38.624
## Avg_Account_Balance 0.000000022464  0.000000001681  13.366
##
##              Pr(>|t|)
## (Intercept)  < 0.0000000000000002
## Gender       0.000000000000000427
## Age          < 0.0000000000000002
## Occupation   < 0.0000000000000002
## Channel_Code < 0.0000000000000002
## Vintage      < 0.0000000000000002
## Credit_Product < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## Residual standard error: 0.4653 on 105304 degrees of freedom
## Multiple R-squared:  0.09339,    Adjusted R-squared:  0.09333
## F-statistic: 1550 on 7 and 105304 DF,  p-value: < 0.00000000000000022
confint(OLS.Model) # Confidence interval

##              2.5 %              97.5 %
## (Intercept)  -0.31838970814349 -0.28070880590589
## Gender       -0.02872684372470 -0.01724548561119
## Age           0.00313697617657  0.00369111307178
```

```
## Occupation          0.03906781016438  0.04611641316309
## Channel_Code        0.08684521792636  0.09545513114680
## Vintage             0.00108032288839  0.00131472132315
## Credit_Product      0.08014273996761  0.08871124507449
## Avg_Account_Balance 0.00000001916964  0.00000002575765
```

All variables are statistically significant in the OLS model. The 95% confidence interval for each variable shows us that all estimators are statistically significant since all intervals do not contain 0. We will use this model as a benchmark.

From the model, we notice that most estimators are logical: one year increase in age will increase the probability of taking credit cards by 0.3414%; one month increase in vintage will increase the probability by 0.1198%; and, 1,000,000 dollars increase in average account balance will increase the probability by 0.02246%.

Different occupations and channel codes lead to different levels of interest. A customer who is salaried has a higher probability of taking credit cards than one with other occupations by 4.259%, a self-employed customer has a higher probability than a salaried customer by 4.259%, and an entrepreneur has a higher probability than a self-employed customer by 4.259%.

For the four channel codes, a customer who chooses “X2” has a higher probability of taking credit cards than a customer who chooses “X1” by 9.115%, a customer who chooses “X3” has a higher probability than the one who chooses “X2” by 9.115%. The customers who choose “X4” have a higher probability than those who choose “X3” by 9.115%.

However, we also notice that females have a lower probability of 2.299% to take credit cards than males. This shows that gender plays an important role in customers’ interest in credit cards. In addition, a customer who has an active credit product has a lower probability of taking a credit card than a customer who does not by 8.443%.

```
# Accuracy
ols.pred.classes <- ifelse(fitted(OLS.Model) > 0.5, 1, 0)
table(ols.pred.classes, data$Is_Active)

##
## ols.pred.classes      0      1
##                0 50837 24597
##                1 12960 16918
mean(ols.pred.classes == data$Is_Active)

## [1] 0.643374
```

The result shows that the accuracy is 64.3374%.

```
fit = lm(Is_Active ~ Gender + Age + Occupation + Channel_Code + Vintage
          + Credit_Product + Avg_Account_Balance, data = data, x = TRUE, y = TRUE)
cv.lm(fit, k = 5) # Perform a 5-fold cv

## Mean absolute error      : 0.4330399
## Sample standard deviation : 0.001216585
##
## Mean squared error       : 0.2165393
## Sample standard deviation : 0.001296877
##
## Root mean squared error   : 0.4653362
## Sample standard deviation : 0.001393213
```

Root mean squared error is  $0.4653289 < 0.5$ , so this is an okay model to use.



## Probit Model

```
Probit.Model <- glm(Is_Active ~ Gender + Age+ Occupation + Channel_Code
+ Vintage + Credit_Product + Avg_Account_Balance,
family = binomial(link = "probit"), data = data)
summary(Probit.Model)
```

```
##
## Call:
## glm(formula = Is_Active ~ Gender + Age + Occupation + Channel_Code +
##     Vintage + Credit_Product + Avg_Account_Balance, family = binomial(link = "probit"),
##     data = data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8032  -0.9430  -0.7242   1.1509   2.2053
##
## Coefficients:
##              Estimate      Std. Error z value
## (Intercept)  -2.180335017947  0.027576324477 -79.065
## Gender       -0.064381114529  0.008242539477  -7.811
## Age          0.009629843417  0.000391666488  24.587
## Occupation   0.119468005030  0.005014942079  23.822
## Channel_Code  0.245926846656  0.006053012264  40.629
## Vintage      0.003081554713  0.000164247726  18.762
## Credit_Product 0.226043718193  0.006165347867  36.664
## Avg_Account_Balance 0.000000063303  0.000000004686  13.508
##
##              Pr(>|z|)
## (Intercept)  < 0.0000000000000002
## Gender       0.00000000000000568
## Age          < 0.0000000000000002
## Occupation   < 0.0000000000000002
## Channel_Code  < 0.0000000000000002
## Vintage      < 0.0000000000000002
## Credit_Product < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 141243  on 105311  degrees of freedom
## Residual deviance: 131140  on 105304  degrees of freedom
## AIC: 131156
##
## Number of Fisher Scoring iterations: 4
confint(Probit.Model) # Confidence interval
```

```
## Waiting for profiling to be done...
##
##              2.5 %          97.5 %
## (Intercept)  -2.23417029886423 -2.12655226692457
## Gender       -0.08054099303620 -0.04822090305290
## Age          0.00886264219527  0.01039699211410
## Occupation   0.10966840911315  0.12926925405325
## Channel_Code  0.23405590725779  0.25779949948908
```

```
## Vintage          0.00276085765945  0.00340234543331
## Credit_Product   0.21408563878054  0.23800909019379
## Avg_Account_Balance 0.00000005417661  0.00000007242973
```

All variables are statistically significant in the probit model. The 95% confidence interval for each variable shows us that all estimators are statistically significant since all intervals do not contain 0.

We notice all estimators have the same signs and significant codes as they are in the OLS model. Many estimators are also close in value, such as the estimators of “Age”, “Vintage”, and “Avg\_Account\_Balance”. So, one year increase in age will increase the probability of taking credit cards by 0.963%; one month’s increase in Vintage will increase the probability by 0.3082%; 1,000,000 dollars increase in average account balance will increase the probability by 0.0633%.

Some estimators are larger, such as “Gender”, “Credit\_Product”, “Occupation”, and “Channel\_Code”. Males have a higher probability of 6.438% to take credit cards than females. A customer who has an active credit product has a lower probability of taking a credit card than a customer who does not by 22.6%.

A customer who is salaried has a higher probability of taking credit cards than one with other occupations by 11.95%, a self-employed customer has a higher probability than a salaried customer by 11.95%, and an entrepreneur has a higher probability than a self-employed customer by 11.95%.

A customer who chooses “X2” has a higher probability of taking credit cards than a customer who chooses “X1” by 24.59%, a customer who chooses “X3” has a higher probability than the one who chooses “X2” by 24.59%. The customers who choose “X4” have a higher probability than those who choose “X3” by 24.59%.

Here the estimators of “Gender”, “Credit\_Product”, “Occupation”, and “Channel\_Code” are all more statistically significant, and differences in those variables cause a larger change in determining if the customer will be interested in taking a recommended credit card.

```
# Accuracy
probit.pred.classes <- ifelse(fitted(Probit.Model) > 0.5, 1, 0)
table(probit.pred.classes, data$Is_Active)
```

```
##
## probit.pred.classes      0      1
##                0 50877 24639
##                1 12920 16876
mean(probit.pred.classes == data$Is_Active)
```

```
## [1] 0.643355
```

64.3355% accuracy, which is a little bit worse than the OLS model but almost the same.

## Logit Model

```
Logit.Model <- glm(Is_Active ~ Gender + Age + Occupation + Channel_Code
                  + Vintage + Credit_Product + Avg_Account_Balance,
                  family = binomial(link = "logit"), data = data)
summary(Logit.Model)
```

```
##
## Call:
## glm(formula = Is_Active ~ Gender + Age + Occupation + Channel_Code +
##      Vintage + Credit_Product + Avg_Account_Balance, family = binomial(link = "logit"),
##      data = data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -1.7883 -0.9402 -0.7270 1.1498 2.1520
##
## Coefficients:
##              Estimate      Std. Error z value
## (Intercept)  -3.559385787506  0.045834397349 -77.658
## Gender       -0.107130460020  0.013537195808  -7.914
## Age          0.015755659588  0.000637563722  24.712
## Occupation   0.197108705206  0.008191904106  24.061
## Channel_Code  0.395956528697  0.009820731328  40.318
## Vintage      0.005048973287  0.000265879353  18.990
## Credit_Product 0.372790725550  0.010134072578  36.786
## Avg_Account_Balance 0.000000101389 0.000000007657  13.241
##
##              Pr(>|z|)
## (Intercept)  < 0.0000000000000002
## Gender       0.00000000000000025
## Age          < 0.0000000000000002
## Occupation   < 0.0000000000000002
## Channel_Code  < 0.0000000000000002
## Vintage      < 0.0000000000000002
## Credit_Product < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 141243 on 105311 degrees of freedom
## Residual deviance: 131169 on 105304 degrees of freedom
## AIC: 131185
##
## Number of Fisher Scoring iterations: 4
```

```
confint(Logit.Model) # Confidence interval
```

```
## Waiting for profiling to be done...
```

```
##              2.5 %          97.5 %
## (Intercept)  -3.64932783252265 -3.4696589123572
## Gender       -0.13366275639513 -0.0805976281131
## Age          0.01450593334352  0.0170051689946
## Occupation   0.18105683099962  0.2131687281415
## Channel_Code  0.37671474687297  0.4152119138685
## Vintage      0.00452800993617  0.0055702487087
## Credit_Product 0.35294709007543  0.3926722703695
## Avg_Account_Balance 0.00000008638431 0.0000001164014
```

All variables are statistically significant in the logit model. The 95% confidence interval for each variable also shows us that all estimators are statistically significant since all intervals do not contain 0.

Similarities are noticed. One year increase in age will increase the probability of taking credit cards by 1.576%. One unit increase in vintage will increase the probability by 0.5049%. 1,000,000 dollars increase in average account balance will increase the probability by 0.1041%. These results are not too different from the previous two models.

However, we notice that males have a higher probability of 10.71% of taking credit cards than females. Also, a customer who has an active credit product has a lower probability than a customer who does not by 37.28%. The logit model suggests that a customer's gender and experiences in having an active credit product may have a stronger influence on his or her choice of credit cards.

Other estimators, such as “Occupation”, suggests that a customer who is salaried has a higher probability of taking credit cards than one with other occupations by 19.71%, a self-employed customer has a higher probability than a salaried customer by 19.71%, and an entrepreneur has a higher probability than a self-employed customer by 19.71%.

For the four channel codes, a customer who chooses “X2” has a higher probability of taking credit cards than a customer who chooses “X1” by 39.6%, a customer who chooses “X3” has a higher probability than the one who chooses “X2” by 39.6%. The customers choose “X4” have a higher probability of than those who choose “X3” by 39.6%. These two results are similar to the probit model but very different from the OLS model.

```
# Accuracy
logit.pred.classes <- ifelse(fitted(Logit.Model) > 0.5, 1, 0)
table(logit.pred.classes, data$Is_Active)
```

```
##
## logit.pred.classes      0      1
##                0 50901 24608
##                1 12896 16907
mean(logit.pred.classes == data$Is_Active)
```

```
## [1] 0.6438772
```

64.38% accuracy, so it seems like the logit model is better than the OLS model and the probit model.

```
# OLS average marginal effect
sum_phi <- mean(dnorm(predict(OLS.Model, type = "response")))
ame.ols = sum_phi*coef(OLS.Model)
```

```
# Probit average marginal effect
sum_phi <- mean(dnorm(predict(Probit.Model, type = "link")))
ame.probit = sum_phi*coef(Probit.Model)
```

```
# Logit average marginal effect
sum_phi <- mean(dnorm(predict(Logit.Model, type = "link")))
ame.logit = sum_phi*coef(Logit.Model)
```

```
stargazer(ame.ols, ame.probit, ame.logit, type = "text", column.labels = c("OLS", "Probit", "Logit"))
```

```
##
## =====
## (Intercept) Gender Age Occupation Channel_Code Vintage Credit_Product Avg_Account_Balance
## -----
## -0.110 -0.008 0.001 0.016 0.033 0.0004 0.031 0
## -----
##
## =====
## (Intercept) Gender Age Occupation Channel_Code Vintage Credit_Product Avg_Account_Balance
## -----
## -0.774 -0.023 0.003 0.042 0.087 0.001 0.080 0.00000
## -----
##
## =====
## (Intercept) Gender Age Occupation Channel_Code Vintage Credit_Product Avg_Account_Balance
## -----
## -1.076 -0.032 0.005 0.060 0.120 0.002 0.113 0.00000
## -----
```

The average marginal effect has an influence on probability.

From the results of the OLS average marginal effect, only “Gender” has a negative effect, that is, males have a higher probability of taking credit cards than females by 0.8%. Both “Channel\_Code” and “Credit\_Product” have dominated effect, which means that “Channel\_Code” and “Credit\_Product” increase the probability of “Is\_Active” by about 3%.

From the average marginal effect results of the probit and the logit models, the direction of marginal effects is consistent with the OLS model. However, compared with the OLS model, the probit and the logit model are more influential. In particular, “Channel\_Code” and “Credit\_Product” in the logit model have a positive impact on customers’ credit cards purchase intention by about 12%.

```
# Linear model prediction
(tab <- table(predict(OLS.Model) > 0.5, data$Is_Active))

##
##           0      1
##  FALSE 50837 24597
##   TRUE  12960 16918

c(tab[1, 1]/sum(tab[, 1]), tab[2, 2]/sum(tab[, 2]))

## [1] 0.7968557 0.4075154

# Probit model prediction
(tab <- table(predict(Probit.Model, type = "response") > 0.5, data$Is_Active))

##
##           0      1
##  FALSE 50877 24639
##   TRUE  12920 16876

c(tab[1, 1]/sum(tab[, 1]), tab[2, 2]/sum(tab[, 2]))

## [1] 0.7974826 0.4065037

# Logit model prediction
(tab <- table(predict(Logit.Model, type = "response") > 0.5, data$Is_Active))

##
##           0      1
##  FALSE 50901 24608
##   TRUE  12896 16907

c(tab[1, 1]/sum(tab[, 1]), tab[2, 2]/sum(tab[, 2]))

## [1] 0.7978588 0.4072504
```

The prediction accuracy of the OLS model for inactive customers is 79.69% and that for active customers is 40.75%. The prediction accuracy of the probit model is 79.75% for inactive customers and 40.65% for active customers. The prediction accuracy of the logit model is 79.79% for inactive customers and 40.73% for active customers. It can be seen from the statistical data that the accuracy of these three models in predicting whether customers are active is very close, and the logit model is slightly higher than the other two models.

From all results above, the logit model is the best fitting model among all three models. In order to evaluate the logit model and find if there is a more suitable one, we will divide the data into testing and training data, and compare the sensitivity, specificity, and RMSE.

## Probit Training

```
intraining <- createDataPartition(data$Is_Active, p = 0.75, list = FALSE)
training <- data[intraining,]
testing <- data[-intraining,]
train_control <- trainControl(method = "cv", number = 5)
probit.mod <- train(as.factor(Is_Active)~.,
                    data = training,
                    method = "glm",
                    family = "binomial"(link = "probit"),
                    trControl = train_control)
summary(probit.mod)
```

```
##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.7710  -0.9448  -0.7249   1.1517   2.2068
##
## Coefficients:
##              Estimate      Std. Error z value
## (Intercept)  -2.179548393680  0.031821294853 -68.493
## Gender       -0.068258677749  0.009517345538  -7.172
## Age          0.009528314019  0.000451488403  21.104
## Occupation   0.123004456938  0.005792013006  21.237
## Channel_Code  0.248275487374  0.006988048172  35.529
## Vintage      0.002940853178  0.000188962222  15.563
## Credit_Product 0.227192452193  0.007104783872  31.977
## Avg_Account_Balance 0.000000059972 0.000000005403  11.099
##
##              Pr(>|z|)
## (Intercept) < 0.0000000000000002
## Gender      0.0000000000000739
## Age         < 0.0000000000000002
## Occupation  < 0.0000000000000002
## Channel_Code < 0.0000000000000002
## Vintage     < 0.0000000000000002
## Credit_Product < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 105930  on 78983  degrees of freedom
## Residual deviance:  98412  on 78976  degrees of freedom
## AIC: 98428
##
## Number of Fisher Scoring iterations: 4
```

Compare to our previous probit model, the estimators of the variables do not change much and the signs are consistent.

```
# Predict using the testing data
pred_is_active <- predict(probit.mod, newdata = testing)
# Evaluate performance
```

```
confusionMatrix(data = pred_is_active, reference = as.factor(testing$Is_Active))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction      0      1
##           0 12811  6193
##           1  3136  4188
##
##           Accuracy : 0.6457
##           95% CI : (0.6398, 0.6514)
##           No Information Rate : 0.6057
##           P-Value [Acc > NIR] : < 0.00000000000000022
##
##           Kappa : 0.218
##
## Mcnemar's Test P-Value : < 0.00000000000000022
##
##           Sensitivity : 0.8033
##           Specificity : 0.4034
##           Pos Pred Value : 0.6741
##           Neg Pred Value : 0.5718
##           Prevalence : 0.6057
##           Detection Rate : 0.4866
##           Detection Prevalence : 0.7218
##           Balanced Accuracy : 0.6034
##
##           'Positive' Class : 0
##
```

```
# Training RMSE
```

```
sqrt(mean((training$Is_Active - predict(Probit.Model, training)) ^ 2))
```

```
## [1] 0.8657024
```

```
# Testing RMSE
```

```
sqrt(mean((testing$Is_Active - predict(Probit.Model, testing)) ^ 2))
```

```
## [1] 0.8672898
```

From the result, we see that accuracy is about 0.64. Sensitivity is about 0.79, which is close to 1.0 and shows that the number of correct positive predictions is high. Specificity is about 0.40, which is a little bit low and shows that the number of correct negative predictions is a bit low.

Both RMSE for training and testing are about 0.86. This is a large number that indicates probit model may not be a suitable model for our data.

## Logit Training

```
logit.mod <- train(as.factor(Is_Active)~.,
  data = training,
  method = "glm",
  family = "binomial"(link = "logit"),
  trControl = train_control)
summary(logit.mod)
```

```
##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.7614  -0.9419  -0.7276   1.1509   2.1531
##
## Coefficients:
##              Estimate      Std. Error z value
## (Intercept)   -3.557303507278    0.052886019353  -67.264
## Gender        -0.113640194507    0.015628738726   -7.271
## Age           0.015598068033    0.000735004521   21.222
## Occupation    0.202996597341    0.009460583487   21.457
## Channel_Code  0.399618010054    0.011338606320   35.244
## Vintage       0.004816358842    0.000305854968   15.747
## Credit_Product 0.374292840627    0.011677381939   32.053
## Avg_Account_Balance 0.000000095988 0.000000008829   10.872
##              Pr(>|z|)
## (Intercept)    < 0.0000000000000002
## Gender         0.0000000000000356
## Age            < 0.0000000000000002
## Occupation     < 0.0000000000000002
## Channel_Code   < 0.0000000000000002
## Vintage        < 0.0000000000000002
## Credit_Product < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 105930  on 78983  degrees of freedom
## Residual deviance:  98434  on 78976  degrees of freedom
## AIC: 98450
##
## Number of Fisher Scoring iterations: 4
```

Compare to our previous logit model, the estimators of the variables do not change much and the signs are consistent.

```
# Predict using the testing data
pred_is_active <- predict(logit.mod, newdata = testing)
# Evaluate performance
confusionMatrix(data = pred_is_active, reference = as.factor(testing$Is_Active))
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction      0      1
##              0 12819  6183
##              1  3128  4198
##
##              Accuracy : 0.6463
##              95% CI : (0.6405, 0.6521)
##      No Information Rate : 0.6057
##      P-Value [Acc > NIR] : < 0.00000000000000022
```



```
##
##           Kappa : 0.2195
##
## Mcnemar's Test P-Value : < 0.00000000000000022
##
##           Sensitivity : 0.8039
##           Specificity : 0.4044
##           Pos Pred Value : 0.6746
##           Neg Pred Value : 0.5730
##           Prevalence : 0.6057
##           Detection Rate : 0.4869
##           Detection Prevalence : 0.7217
##           Balanced Accuracy : 0.6041
##
##           'Positive' Class : 0
##
```

From the sensitivity and specificity report, we can tell it is almost the same as probit model.

```
logit.mod <- train((Is_Active)~.,
                  data = training,
                  method = "glm",
                  family = "binomial"(link = "logit"),
                  trControl = train_control)

# Predict using the testing data
pred_is_active <- predict(logit.mod, newdata = testing)
# Evaluate performance
postResample(pred = pred_is_active, obs = testing$Is_Active)

##           RMSE   Rsquared          MAE
## 0.46499540 0.09467693 0.43265514
```

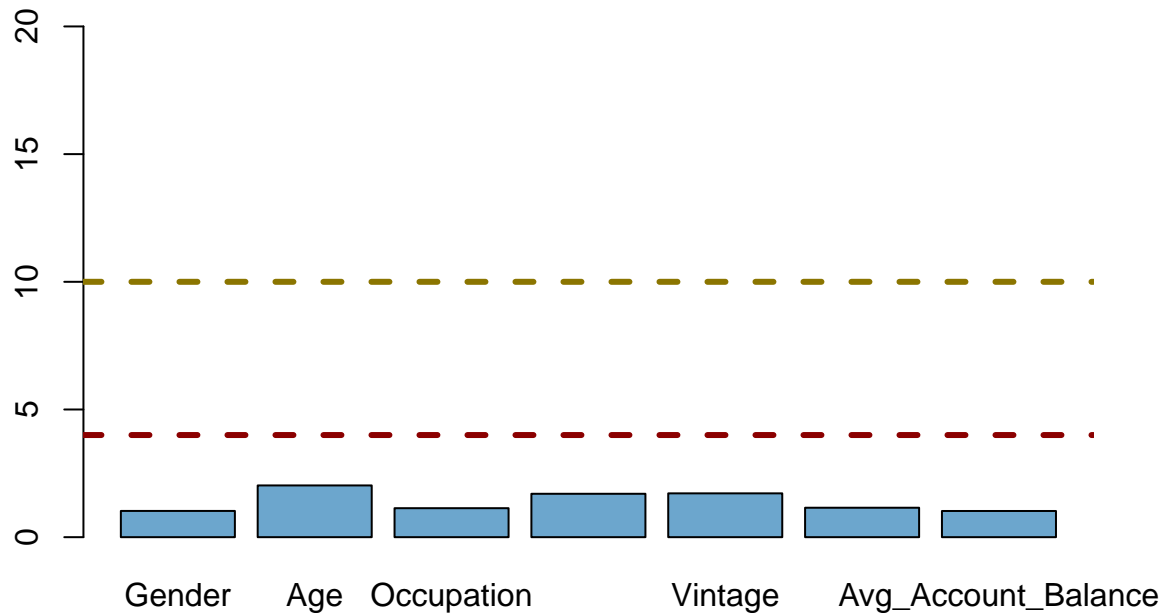
RMSE is about  $0.46 < 0.5$ , which is much smaller than what we had in the probit model. This means that the logit model fits the dataset better.

For now, we still consider logit model as our best model. We are going to do some tests on it to prove or disprove our perspective.

## Tests for Logit Model

```
# VIF
barplot(vif(Logit.Model), main = "VIF Values", col = "skyblue3", ylim = c(0.0,20.0))
abline(h = 4, lwd = 3, lty = 2, col = "red4")
abline(h = 10, lwd = 3, lty = 2, col = "gold4")
```

## VIF Values



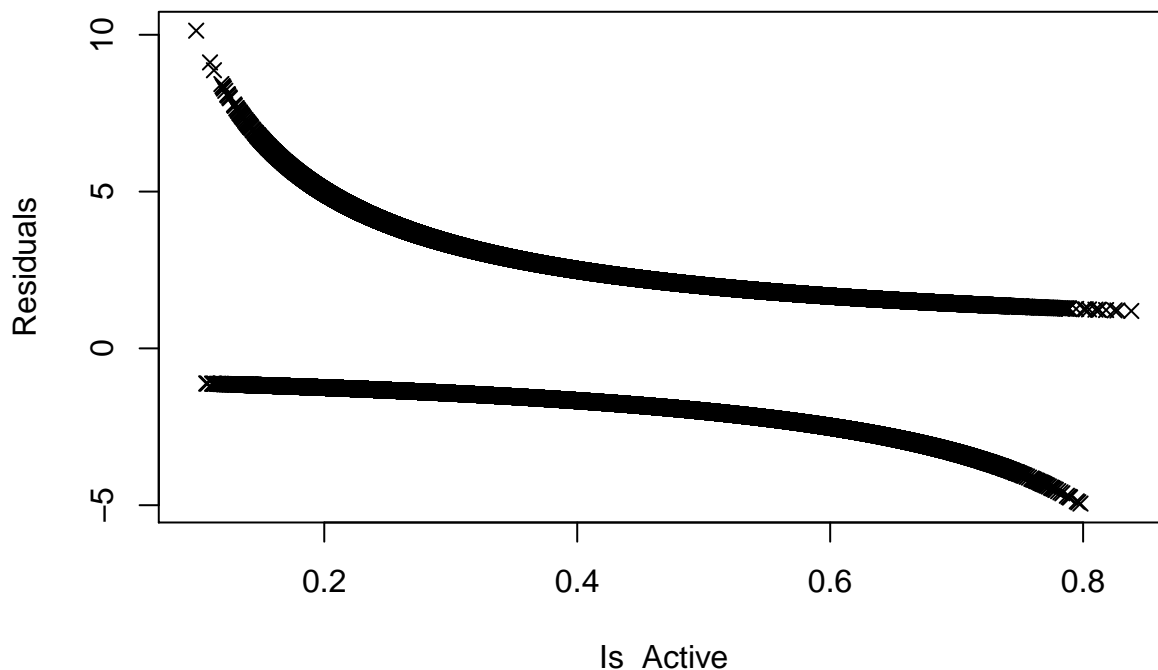
The result indicates that all of the VIF values are smaller than 4. Therefore, it is appropriate that we keep all the variables for the logit model.

```
# Reset
resettest(Logit.Model, power = 2, type = "regressor")
```

```
##
## RESET test
##
## data: Logit.Model
## RESET = 622.23, df1 = 7, df2 = 105297, p-value < 0.00000000000000022
```

p-value is smaller than 0.05, which means we may need to include a higher power term or interaction term to improve our logit model.

```
# Plot the respective vs. y hat
plot(Logit.Model$fitted.values, Logit.Model$residuals,
     xlab = "Is_Active", ylab = "Residuals", pch = 4)
```



From the graph, we can tell there is heteroskedasticity, and the residuals are not constant. We can use `bptest` to test the heteroskedasticity.

```
# Heteroskedasticity
bptest(Logit.Model)

##
## studentized Breusch-Pagan test
##
## data: Logit.Model
## BP = 5004.4, df = 7, p-value < 0.00000000000000022
```

Since the p-value is smaller than 0.05, we reject the null that homoscedasticity is presented and conclude that heteroskedasticity exists.

We convert the standard errors into robust standard errors:

```
# Robust standard errors
cov1 <- hccm(Logit.Model, type = "hc1")
Logit.Model.heter <- coeftest(Logit.Model, vcov. = cov1)
tidy(Logit.Model.heter)

## # A tibble: 8 x 5
##   term                estimate  std.error statistic  p.value
##   <chr>                <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)        -3.56      0.113     -31.6 1.13e-219
## 2 Gender             -0.107     0.0324     -3.30 9.58e- 4
## 3 Age                0.0158     0.00152    10.3 4.78e- 25
## 4 Occupation         0.197     0.0198      9.96 2.25e- 23
## 5 Channel_Code        0.396     0.0233     17.0 1.13e- 64
## 6 Vintage            0.00505    0.000640      7.89 3.05e- 15
## 7 Credit_Product      0.373     0.0252     14.8 1.60e- 49
## 8 Avg_Account_Balance 0.000000101 0.0000000187    5.42 6.01e- 8
```

## Multinomial Logit Model

In our logit model, we noticed that the variable “Credit\_Product” has a significant effect on the model. “Credit\_Product” measures if the consumers have any active credit products, so intuitively, the variables such as “Gender” and “Age” may affect “Credit\_Product” in some way, though this has not been seen in our model or tests. Then, we note that for customers, there may be multinomial choice situations, because they actually face more than two choices. In other words, if we ignore the choice of “Credit\_Product”=“N/A”, there are at least three combinations of “Is\_Active” and “Credit\_Product”, indicating that customers can make at least three different choices. We can put the customers into three groups: inactive customers, active customers who have active credit product, and active customers who does not have active credit product.

In order to help Happy Customer Bank to cross sell, it is necessary for us to analyze customers’ behavior in different groups.

So, we create a new variable “ia\_cp” based on “Is\_Active” and “Credit\_Product”. We ignore the observations where “Credit\_Product” = “N/A” as they could be confusing. We let the variable “ia\_cp” ranges from 1 to 3 for three possible combinations of “Is\_Active” and “Credit\_Product”, so we can use a multinomial logit model to see the similarities and differences among the four groups.

```
# Create a new variable
data2 = data [data$Credit_Product != 1, ]
data2$ia_cp [data2$Is_Active == 0] <- 1
data2$ia_cp [data2$Is_Active == 1 & data2$Credit_Product == 3] <- 2
data2$ia_cp [data2$Is_Active == 1 & data2$Credit_Product == 2] <- 3

table(data2$ia_cp)
```

```
##
##      1      2      3
## 57387 25855  9548
```

From the table, it seems like our dataset is a little bit biased.

```
attach(data)

## The following objects are masked from data (pos = 3):
##
##      Age, Avg_Account_Balance, Channel_Code, Credit_Product, Gender,
##      Is_Active, Occupation, Vintage

multi <- multinom(ia_cp ~ Gender + Age + Occupation + Channel_Code + Vintage
                  + Avg_Account_Balance, data = data2)

## # weights:  24 (14 variable)
## initial value 101940.234265
## iter  10 value 83544.812683
## iter  20 value 77162.953771
## final value 77087.447339
## converged

summary(multi)

## Call:
## multinom(formula = ia_cp ~ Gender + Age + Occupation + Channel_Code +
##      Vintage + Avg_Account_Balance, data = data2)
##
## Coefficients:
##      (Intercept)      Gender      Age Occupation Channel_Code      Vintage
## 2      -2.537160 -0.0631628 0.01776665  0.1853642    0.2895855 -0.0005045724
```

```
## 3    -4.767244 -0.2393539 0.01388448  0.2224651    0.5839420  0.0119428055
##      Avg_Account_Balance
## 2      0.00000008731495
## 3      0.00000012917576
##
## Std. Errors:
##              (Intercept)                Gender                Age
## 2 0.0000000000000003034117 0.000000000000001411385 0.0000000000001340968
## 3 0.0000000000000003752311 0.000000000000001412151 0.0000000000001902959
##              Occupation                Channel_Code                Vintage
## 2 0.0000000000000006446172 0.000000000000005657823 0.0000000000001377464
## 3 0.0000000000000008335946 0.000000000000008821456 0.0000000000002457466
##      Avg_Account_Balance
## 2      0.000000005352584
## 3      0.000000007447146
##
## Residual Deviance: 154174.9
## AIC: 154202.9
```

From the summary, we notice that most estimators are similar to what we had in our logit model. Almost all estimators are consistent in values and signs. “ia\_cp”=1 is missing, so it is the baseline. “Age” and “Avg\_Account\_Balance” increase the probability of taking credit cards, while the probability decreases in “Gender” (females have a lower probability). The estimators for “Occupation” and “Channel\_Code” show that for customers with different occupations and different interests in channels, their responses towards credit cards are different. However, the estimators for “Vintage” have different signs. This indicates that vintage may affect different groups of customers differently. We will come back to this after a few tests to decide if this model is worth analyzing.

```
# Tests for multinomial logit model
multi$AIC
```

```
## [1] 154202.9
```

```
multi$edf
```

```
## [1] 14
```

AIC is large indicates that the model does not fit the data very well, and edf is small implies that the (effective) number of degrees of freedom used by the model is low. Therefore, the multinomial logit model is questionable. Hence, we will not choose multinomial logit model as our best model.

## Instrumental Model

After plotting correlation plot, we noticed that the variables “Channel\_Code” and “Vintage” have high correlations with “Age”. Therefore, we have two potential instrumental variables: “Channel\_Code” & “Vintage”, with “Age” as the endogenous variable, and “Gender”, “Occupation”, “Credit\_Product”, and “Avg\_Account\_Balance” as exogenous variables. We will establish a new model here and conduct a range of tests for our new model.

```
# Model with two instrumental variables Channel_Code & Vintage
iv.mod <- ivreg(Is_Active ~ Gender + Age + Occupation + Credit_Product
               + Avg_Account_Balance | Gender + Occupation + Credit_Product
               + Avg_Account_Balance + Channel_Code + Vintage, data = data)
summary(iv.mod)
```

```
##
```

```
## Call:
```

```
## ivreg(formula = Is_Active ~ Gender + Age + Occupation + Credit_Product +
##       Avg_Account_Balance | Gender + Occupation + Credit_Product +
##       Avg_Account_Balance + Channel_Code + Vintage, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9908 -0.3723 -0.1742  0.4569  1.1186
##
## Coefficients:
##              Estimate      Std. Error t value
## (Intercept)   -0.608840694249   0.011968922434 -50.868
## Gender        -0.012312605340   0.003024377414  -4.071
## Age           0.013373474581   0.000153001911  87.407
## Occupation    0.082096256617   0.001782966241  46.045
## Credit_Product 0.091987247448   0.002278836936  40.366
## Avg_Account_Balance 0.000000016758 0.000000001726  9.709
##
##              Pr(>|t|)
## (Intercept)    < 0.0000000000000002
## Gender         0.0000468
## Age            < 0.0000000000000002
## Occupation     < 0.0000000000000002
## Credit_Product < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## Residual standard error: 0.4766 on 105306 degrees of freedom
## Multiple R-Squared: 0.04889, Adjusted R-squared: 0.04885
## Wald test: 1918 on 5 and 105306 DF, p-value: < 0.00000000000000022
```

From the summary, females have a lower probability of 1.23% to take credit cards than males. One year increase in age will increase the probability by 1.34%. 1,000,000 dollars increase in average account balance will increase the probability by 0.0167%. A customer who has an active credit product has a lower probability of taking a credit card than a customer who has not by 9.19%. In brief, we find the signs of the estimators are the same between the IV model and the logit model, and almost all estimators are smaller in value.

```
# Hausman Test
summary(iv.mod, diagnostics = TRUE)
```

```
##
## Call:
## ivreg(formula = Is_Active ~ Gender + Age + Occupation + Credit_Product +
##       Avg_Account_Balance | Gender + Occupation + Credit_Product +
##       Avg_Account_Balance + Channel_Code + Vintage, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9908 -0.3723 -0.1742  0.4569  1.1186
##
## Coefficients:
##              Estimate      Std. Error t value
## (Intercept)   -0.608840694249   0.011968922434 -50.868
## Gender        -0.012312605340   0.003024377414  -4.071
## Age           0.013373474581   0.000153001911  87.407
## Occupation    0.082096256617   0.001782966241  46.045
## Credit_Product 0.091987247448   0.002278836936  40.366
## Avg_Account_Balance 0.000000016758 0.000000001726  9.709
```

```
##                                Pr(>|t|)
## (Intercept)                  < 0.0000000000000002
## Gender                       0.0000468
## Age                          < 0.0000000000000002
## Occupation                   < 0.0000000000000002
## Credit_Product               < 0.0000000000000002
## Avg_Account_Balance < 0.0000000000000002
##
## Diagnostic tests:
##                                df1      df2 statistic      p-value
## Weak instruments             2 105305   47151.4 <0.0000000000000002
## Wu-Hausman                   1 105305    2340.5 <0.0000000000000002
## Sargan                       1      NA     195.3 <0.0000000000000002
##
## Residual standard error: 0.4766 on 105306 degrees of freedom
## Multiple R-Squared: 0.04889, Adjusted R-squared: 0.04885
## Wald test: 1918 on 5 and 105306 DF, p-value: < 0.00000000000000022
```

In the Diagnostic test, the result indicates that our tested instruments are not weak, which is good. We reject the null which essentially means that the OLS is not consistent, and endogeneity is present. The least squares estimator is a better and more efficient estimator. Furthermore, all of the coefficients are statistically significantly different from 0 and both instruments ("Channel\_Code" & "Vintage") are valid.  $p\text{-value} < 2e-16$ , so we conclude that both the least squares estimator and the instrumental variables estimator are consistent. The least squares estimator is a better and more efficient estimator.

```
# Test the validity of instrument
linearHypothesis(OLS.Model, c("Channel_Code = 0", "Vintage = 0"),
  vcov = vcovHC, type = "HC1")
```

```
## Linear hypothesis test
##
## Hypothesis:
## Channel_Code = 0
## Vintage = 0
##
## Model 1: restricted model
## Model 2: Is_Active ~ Gender + Age + Occupation + Channel_Code + Vintage +
## Credit_Product + Avg_Account_Balance
##
## Note: Coefficient covariance matrix supplied.
##
##   Res.Df Df      F      Pr(>F)
## 1 105306
## 2 105304  2 1199.9 < 0.00000000000000022
```

Since  $p\text{-value} < 0.05$ , we reject the null hypothesis and conclude that the estimators of "Channel\_Code" and "Vintage" are significantly different from 0, and  $f > 10$  suggests they are strong instruments.

```
# t-test of coefficients
coefTest(iv.mod, vcov = vcovHC, type = "HC1")
```

```
##
## t test of coefficients:
##
##              Estimate      Std. Error t value
## (Intercept) -0.608840694249  0.011836959890 -51.4356
```

```
## Gender          -0.012312605340  0.003017367535  -4.0806
## Age             0.013373474581  0.000153410391  87.1745
## Occupation      0.082096256617  0.001843416159  44.5348
## Credit_Product  0.091987247448  0.002370877352  38.7988
## Avg_Account_Balance 0.000000016758  0.000000001814   9.2383
##
##                      Pr(>|t|)
## (Intercept)      < 0.00000000000000022
## Gender           0.00004496
## Age              < 0.00000000000000022
## Occupation       < 0.00000000000000022
## Credit_Product   < 0.00000000000000022
## Avg_Account_Balance < 0.00000000000000022
```

We conclude that all coefficients are statistically significantly different from 0.

```
# Test the validity of the instruments (J-statistic) via
# the overidentifying restrictions test
iv_OR <- lm(residuals(iv.mod) ~ Gender + Occupation + Credit_Product
            + Avg_Account_Balance + Channel_Code + Vintage, data = data)
linearHypothesis(iv_OR, c("Channel_Code = 0", "Vintage = 0"), test = "Chisq")
```

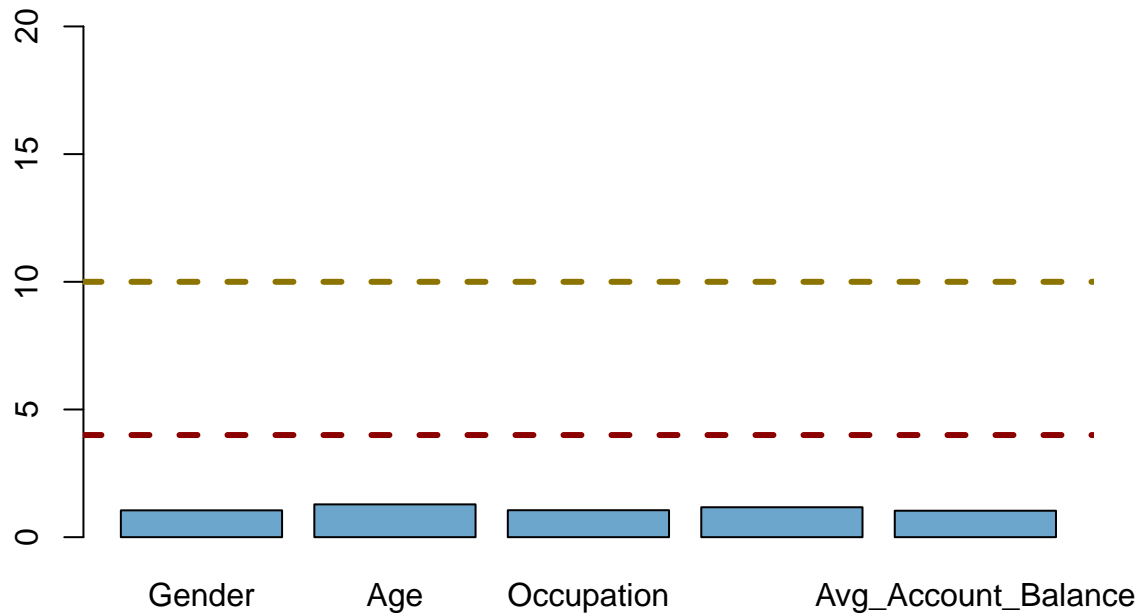
```
## Linear hypothesis test
##
## Hypothesis:
## Channel_Code = 0
## Vintage = 0
##
## Model 1: restricted model
## Model 2: residuals(iv.mod) ~ Gender + Occupation + Credit_Product + Avg_Account_Balance +
##          Channel_Code + Vintage
##
##   Res.Df  RSS Df Sum of Sq  Chisq          Pr(>Chisq)
## 1 105307 23920
## 2 105305 23875   2    44.365 195.68 < 0.00000000000000022
```

$\text{Pr} < 2.2\text{e-}16$ , so we conclude that both instruments are valid.

```
# VIF for IV
barplot(vif(iv.mod), main = "VIF Values", col = "skyblue3", ylim = c(0.0,20.0))
abline(h = 4, lwd = 3, lty = 2, col = "red4")
abline(h = 10, lwd = 3, lty = 2, col = "gold4")
```



## VIF Values



The result indicates that all of the vif values are smaller than 4, so it is appropriate that we keep all variables for the IV model.

```
# RESET for IV
resettest(iv.mod, power = 2, type = "regressor")
```

```
##
## RESET test
##
## data: iv.mod
## RESET = 689.49, df1 = 7, df2 = 105299, p-value < 0.00000000000000022
```

With a statistic of 689.49 and a p-value of ~0.000, the RESET test suggests that the IV model is NOT correctly specified. So we reject H0.

```
# Heteroskedasticity for IV
bptest(iv.mod)
```

```
##
## studentized Breusch-Pagan test
##
## data: iv.mod
## BP = 4770.7, df = 5, p-value < 0.00000000000000022
```

Since the p-value is smaller than 0.05, we reject the null that homoscedasticity is presented and conclude that heteroskedasticity exists.

Therefore, our IV model is not so useful according to the tests. The logit model remains to be our best model.

## Identify Potential Customers

So, after going through different models and running different tests, we decide that the logit model is the most suitable model for our data. Now we will use the logit model to identify potential customers who will take credit cards.

```
# load test.csv again because we need the customers' IDs
new_data <- read_csv("test.csv")
```

```
## Parsed with column specification:
## cols(
##   ID = col_character(),
##   Gender = col_character(),
##   Age = col_double(),
##   Region_Code = col_character(),
##   Occupation = col_character(),
##   Channel_Code = col_character(),
##   Vintage = col_double(),
##   Credit_Product = col_character(),
##   Avg_Account_Balance = col_double(),
##   Is_Active = col_character()
## )
```

```
# Store the coefficients
b1 = Logit.Model$coefficients[1]
b2 = Logit.Model$coefficients[2]
b3 = Logit.Model$coefficients[3]
b4 = Logit.Model$coefficients[4]
b5 = Logit.Model$coefficients[5]
b6 = Logit.Model$coefficients[6]
b7 = Logit.Model$coefficients[7]
b8 = Logit.Model$coefficients[8]
```

We identify a list of most potential customers by calculating the probability (or we can call it the level) of customers' interest in taking recommended credit cards, and store the value as "Is\_Lead".

```
# data$Is_Lead <- b1 + b2*data$Gender + b3*data$Age + b4*data$Occupation
# + b5*data$Channel_Code + b6*data$Vintage + b7*data$Credit_Product
# + b8*data$Avg_Account_Balance
```

```
data$Is_Lead <- b1 + b2*data$Gender + b3*data$Age + b4*data$Occupation + b5*data$Channel_Code + b6*data$
```

```
# Store "Is_Lead" to the new dataset with IDs
new_data$Is_Lead = data$Is_Lead
```

Now we have the level of customers' interest in taking credit cards. We expect "Is\_Lead" to be between 0 and 1, but there may be numbers > 1 or < 0 due to the residuals. So we can now identify potential customers. We divide customers into three levels and follow the rule that the higher "Is\_Lead" is, the more likely the customer is to be interested in a credit card. We will print some of the IDs of the customers from the first level, and we will provide the number of customers at all levels.

```
# First level: The most potential customers
level.1 <- new_data[new_data$Is_Lead > 0.7, ]
nrow(level.1)
```

```
## [1] 3591
```

```
print(level.1$ID[1:100])
```

```
## [1] "MZZAQMPT" "MXETLUP4" "JTYNEKRQ" "5NTYJJDV" "YRLEWZfq" "NWKHMEOF"
## [7] "LXHCDSVH" "BNIUPSVK" "KLHUI2SE" "CUAUZJH8" "6EDAM4EU" "NZKCSdz4"
## [13] "ENUB4N9S" "SGJFSSJK" "FPU4ZI8H" "27F5ZEQA" "ATJAEQCJ" "UGIADHRB"
## [19] "UMVSSZ9N" "NCVH9WFE" "DIVE6ZUL" "JTCERL4D" "DHARALYV" "MPQQFX4I"
## [25] "CJHC43KV" "FTHIZFGT" "VQEMTVX" "7V6GQUOD" "QPQZPSMH" "50Y8PVUA"
## [31] "5PRTDKXK" "4WTK9RBR" "JWPYWNEU" "GYKEBWCJ" "LGMCRtXe" "SXHNHYFX"
## [37] "QADPUPRH" "SKAFGCWQ" "GY3GGMD6" "BKAUTCvF" "NX428PF5" "DAVGfGE3"
## [43] "XRDU4PBA" "DAJT6YYA" "T5TGKHOM" "R4MRAHPP" "QWT7FVWQ" "8Q4FXOQI"
## [49] "M8PSGEVT" "EIBVIG5Z" "BX8KYC5S" "KR8BxBYI" "KHUdKRPY" "XKAQ5I8F"
## [55] "NMJN8TQV" "H6SXNKC7" "W4RJBEBS" "FAHIFIFF" "EZBSYwQU" "NCWSS9X9"
## [61] "BDJWIXPB" "5S2GPQID" "ULZYVWRZ" "CYLS7DNT" "JMPMFHLR" "CHMWB8HB"
## [67] "H8IKUYS2" "J9QWFQF9" "B4PB8ZAG" "T4IKMKME" "8CHLIViy" "WECEW6V7"
## [73] "GKUDDQSV" "ZZVS5V5X" "BKMHGyUD" "5PRAGM92" "4QVZMAT4" "ITOJTXUY"
## [79] "A7K2XHTP" "QKTNQGVs" "XRZDNUSE" "XECJYFEB" "AKYV6SS7" "NR3BBJTI"
## [85] "YTTS5ZR2" "FV8BZ7IQ" "KSZDM7NB" "LFYXBxAV" "FJUQUVFK" "FGR8EWNZ"
## [91] "2TWCWRGK" "DFARFVGA" "PRM72CHV" "PFBPHDB6" "K8PVPHAP" "BCY5Q3NT"
## [97] "Q29W7QSQ" "ZUEJN22B" "HH5NJGP0" "LS8U4YWW"
```

```
# Second level: Also potential customers
level.2 <- new_data[new_data$Is_Lead <= 0.7, ]
level.2 <- level.2[level.2$Is_Lead > 0.5, ]
nrow(level.2)
```

```
## [1] 5049
```

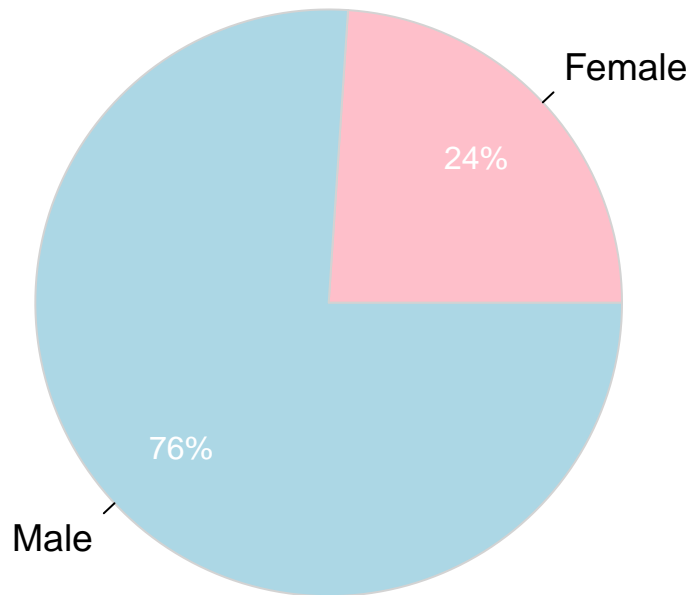
```
# Third level: Least recommended customers, but the bank can still have a try
level.3 <- new_data[new_data$Is_Lead <= 0.5, ]
nrow(level.3)
```

```
## [1] 96672
```

From here, we can see some of the IDs of the customers most likely to take credit cards. We will recommend customers to Happy Customer Bank according to their levels (i.e., the most recommended customers are from level.1, then level.2, and the least recommended customers are from level.3). Now, we are here at our last step before the conclusion. We are going to construct histograms and pie charts of the 3,591 most potential customers in level.1 to find out the patterns in their information (i.e., how are “Gender”, “Age”, “Occupation”, etc. distributed in potential customers).

```
# Gender
PieChart(Gender, hole = 0, values = "%", data = level.1,
         fill = c("pink", "lightblue"), main = "Gender Pie Chart", quiet = TRUE)
```

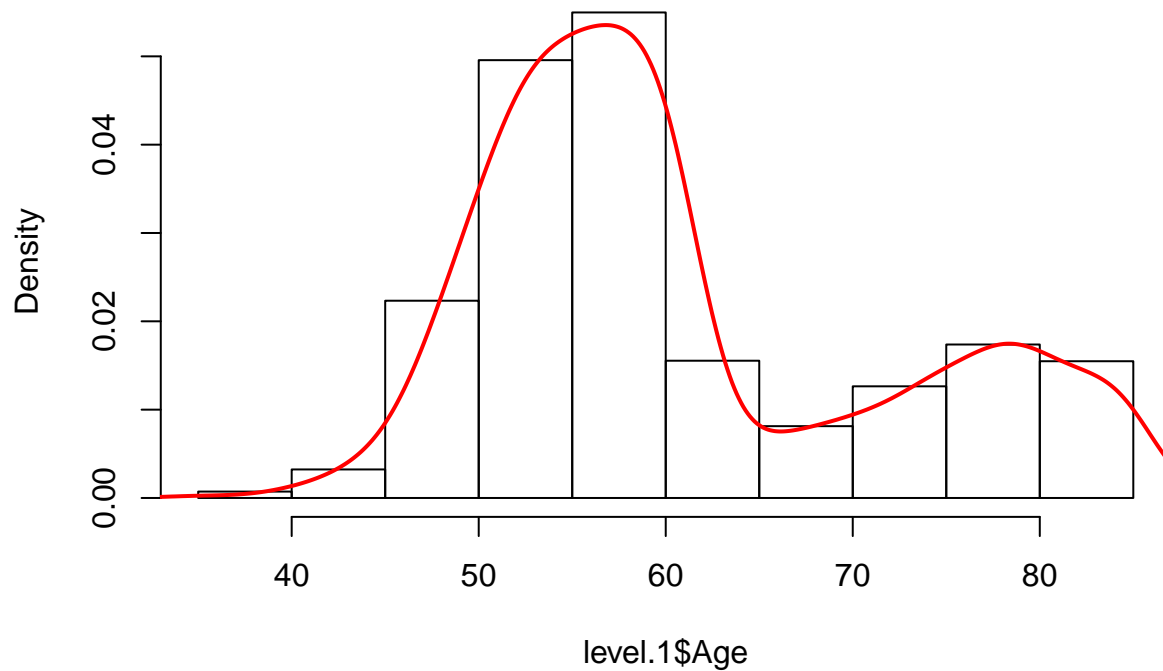
## Gender Pie Chart



The pie chart of "Gender" shows that 76% of the customers are male while only 24% of them are female.

```
# Age
hist(level.1$Age, prob = TRUE, ylim = c(0, max(density(level.1$Age)$y)))
lines(density(level.1$Age), lwd = 2, col = "red")
```

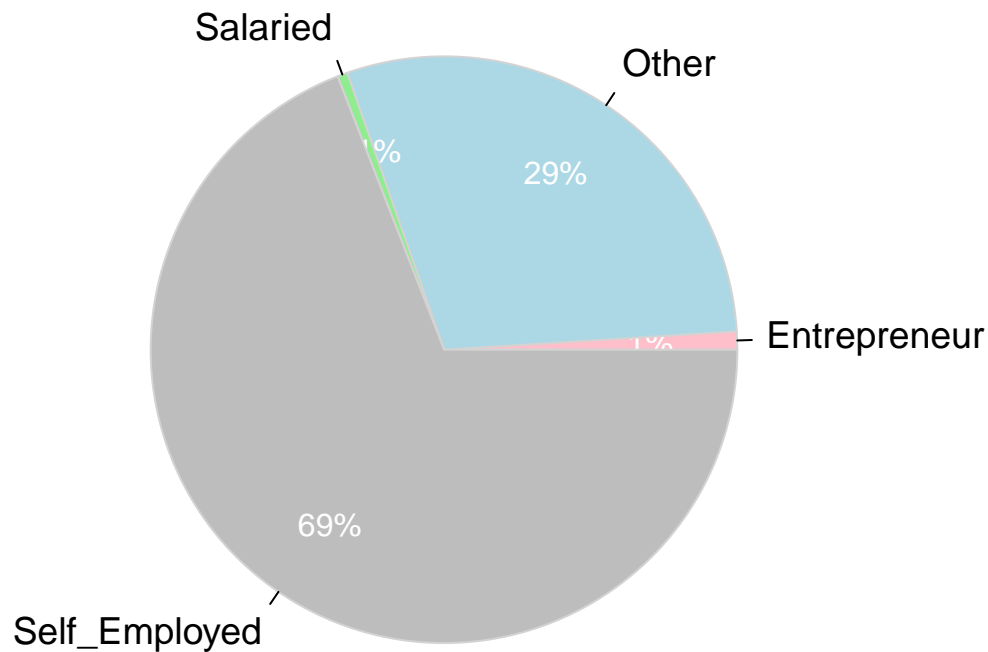
**Histogram of level.1\$Age**



The customers' age falls within the interval [40, 80]; many customers are about 50 to 60 years old.

```
# Occupation
PieChart(Occupation, hole = 0, values = "%", data = level.1,
         fill = c("pink", "lightblue", "lightgreen", "grey"),
         main = "Occupation Pie Chart", quiet = TRUE)
```

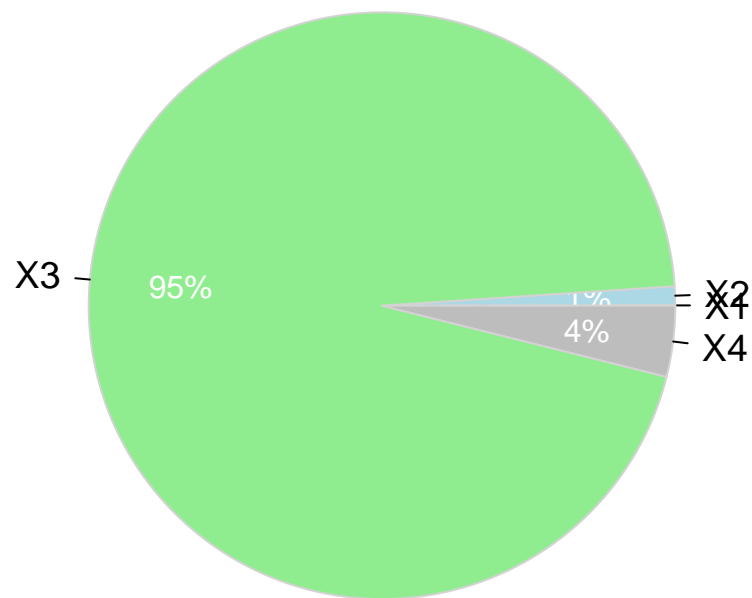
## Occupation Pie Chart



The pie chart of "Occupation" shows that most customers are self-employed or other occupations. There are very few salaried customers or entrepreneurs. This shows that different occupation does lead to different probabilities of taking credit cards (but with bias).

```
# Channel_Code
PieChart(Channel_Code, hole = 0, values = "%", data = level.1,
         fill = c("pink", "lightblue", "lightgreen", "grey"),
         main = "Channel Code Pie Chart", quiet = TRUE)
```

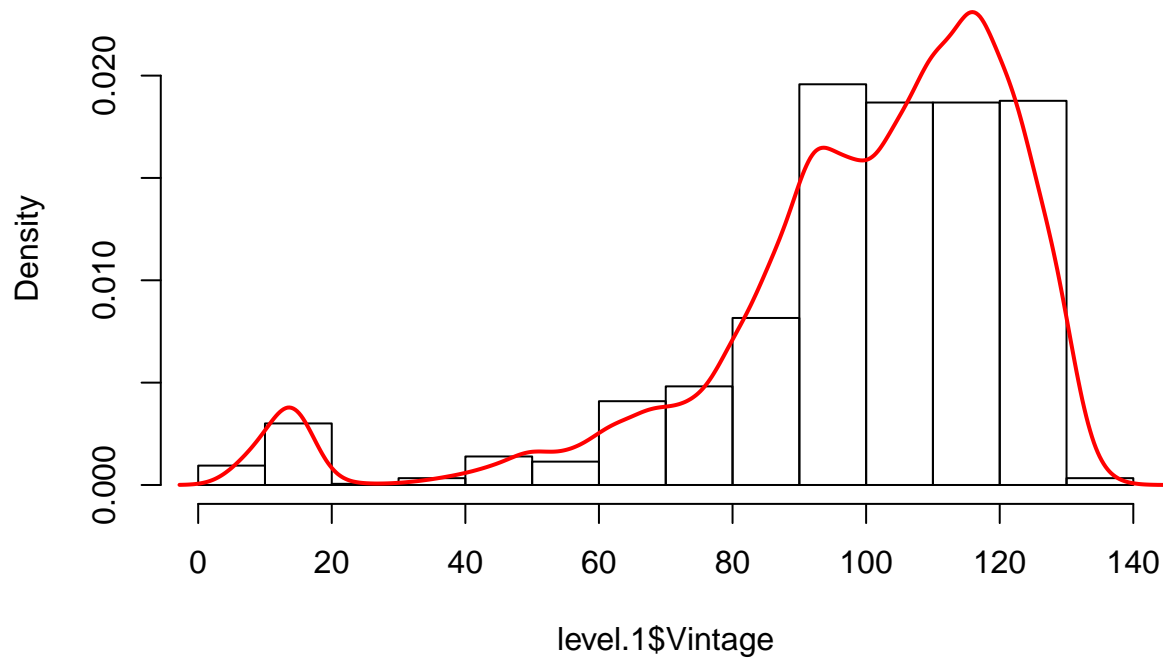
## Channel Code Pie Chart



Surprisingly, almost all customers choose channel code “X3”. This may be investigated more in the future.

```
# Vintage
hist(level.1$Vintage, prob = TRUE, ylim = c(0, max(density(level.1$Vintage)$y)))
lines(density(level.1$Vintage), lwd = 2, col = "red")
```

**Histogram of level.1\$Vintage**

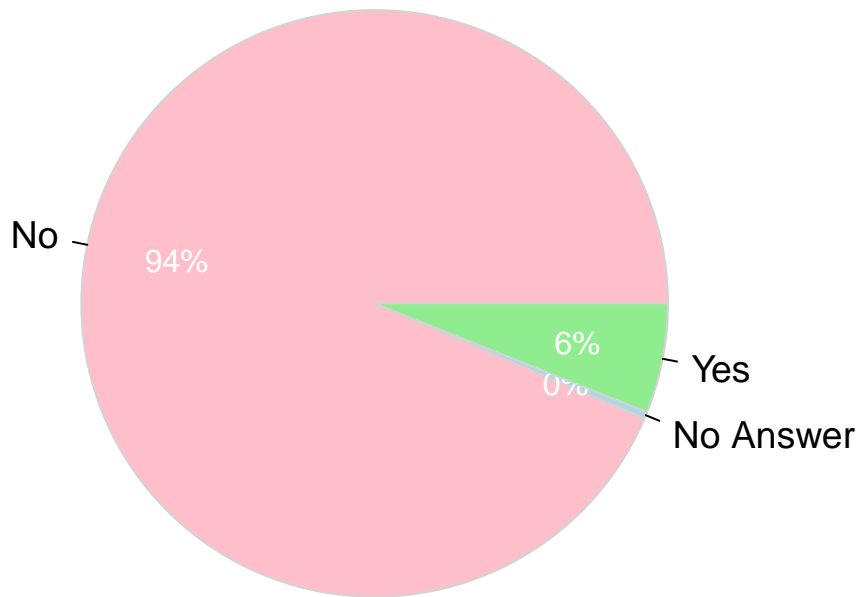


The vintage is left-skewed and is clustered around 120.



```
# Credit_Product
# We change "N/A" to "No Answer" because N/As are not counted in pie charts
level.1$Credit_Product[is.na(level.1$Credit_Product)] <- "No Answer"
PieChart(Credit_Product, hole = 0, values = "%", data = level.1,
         fill = c("pink", "lightblue", "lightgreen"),
         main = "Credit Product Pie Chart", quiet = TRUE)
```

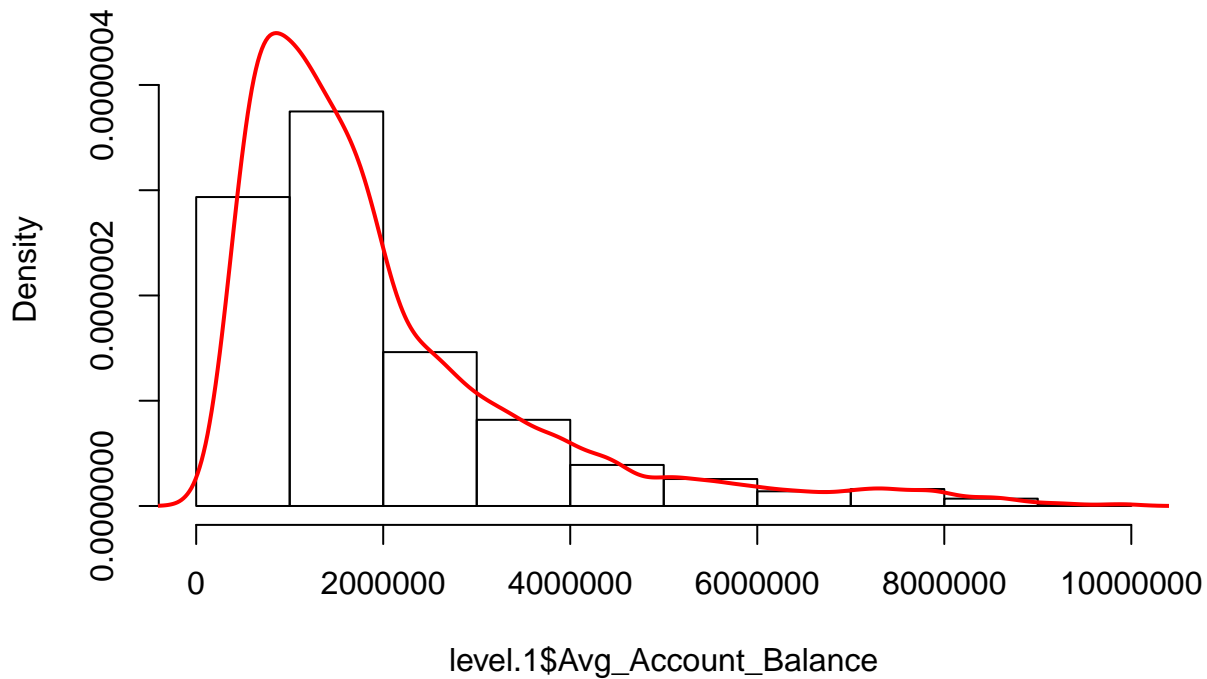
## Credit Product Pie Chart



The pie chart of "Credit\_Product" shows that 94% of the customers answered "No" while only 6% answered "Yes".

```
# Avg_Account_Balance
hist(level.1$Avg_Account_Balance, prob = TRUE, ylim = c(0, max(density(level.1$Avg_Account_Balance)$y)))
lines(density(level.1$Avg_Account_Balance), lwd = 2, col = "red")
```

### Histogram of level.1\$Avg\_Account\_Balance



The histogram is similar to what we had before. We calculate the mean to see if there is any change.

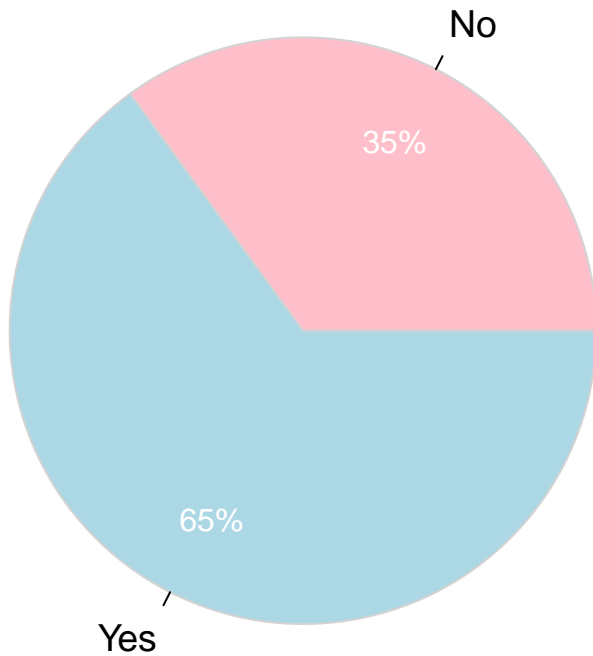
```
mean(level.1$Avg_Account_Balance)
```

```
## [1] 1985474
```

Therefore, customers with average account balance of around 1,985,474 have a higher probability of taking credit cards, much higher than the overall mean.

```
# Is_Active
PieChart(Is_Active, hole = 0, values = "%", data = level.1,
         fill = c("pink", "lightblue"),
         main = "Active Pie Chart", quiet = TRUE)
```

## Active Pie Chart



The pie chart of "Is\_Active" shows that 65% of customers answered "Yes" while 35% answered "No".

```
# Region_Code
Region <- as.factor(level.1$Region_Code)
Region <- Region[-length(Region)]
table(Region)

## Region
## RG250 RG251 RG252 RG253 RG254 RG255 RG256 RG257 RG258 RG259 RG260 RG261 RG262
##      9      72      14      41     505      11      12      35      10      17      23      39      22
## RG263 RG264 RG265 RG266 RG267 RG268 RG269 RG270 RG271 RG272 RG273 RG274 RG275
##      25      11      13       6       3    845      88      42       2      45      32      31      19
## RG276 RG277 RG278 RG279 RG280 RG281 RG282 RG283 RG284
##      21     127      21      33     137      45      57     736     441
```

The table shows that 845 most potential customers live in RG268. Happy Customer Bank could set more banking outlets for customers' convenience / attract more customers around RG268.

From these graphs, we can see the specific range of the 3,591 most recommended customers' information, specifically the share of their genders, occupations, channel codes, whether they have or have not active credit products, and whether they are active or not. We also have a more comprehensive understanding of the distribution of their ages, vintages, and average account balances. Hereto, we have chosen the best model for our data and used it to successfully identify potential customers, and have had a thorough understanding of the distributions of these customers.

Based on our analysis, we suggest that Happy Customer Bank could set more banking outlets for customers' convenience around RG268. Furthermore, medium-to-high income people have a higher probability of taking credit cards, so it is effective to set banking outlets in central city regions. Self-employed people consist of 70% of potential customers, so the Happy Customer Bank could develop some financing products or petty loans to attract them.

## 4. Conclusion

In conclusion, by using and comparing various models and analysis results, we have successfully identified the logit model as our best model, and come up with a recommended list of potential customers who are most likely to take credit cards from Happy Customer Bank.

The best model was chosen with some difficulty, as throughout our analysis, a range of possible models seemed suitable at first sight. For example, the OLS model, the probit model, and the logit model showed similar results in their summary, their model accuracy, their average marginal effect for each variable, etc. All estimators were statistically significant in these three models, and the signs of the estimators were consistent. The logit model showed slightly better model accuracy. To evaluate their performance better, we divided the data into training data. And logit model still had higher accuracy as well as lower RMSE.

However, this was far from the time when we determined the logit model as our best model. A large range of tests was provided. For example, RESET test was used to check the model specification, and BP test was used to detect heteroskedasticity. The logit model passed most of the tests. Furthermore, a multinomial logit model and an IV model were used to challenge the logit model to see if they were better, but they failed because other tests have shown that they were either questionable or not well specified.

Although we have identified the logit model as our best model, we still have many shortcomings at present, such as heteroskedasticity. Though we converted our standard errors into robust standard errors, this did not help much.

Nevertheless, applying the logit model to our data, lists of recommended customers sorted according to their level of interest in taking credit cards came out. Based on our result, 3,591 customers are most likely to take a recommended credit card. Top recommended customers turn out to be mostly males, aged between 40 and 80, self-employed, and are those who tend to choose channel code "X3", have no active credit product, have the vintage of about 120 and average account balance between 2 million and 10 million dollars. Most are active users and live in RG268. These results show the common features of the most recommended customers and we hope it will help Happy Customer Bank make future decisions about recommending credit cards and other financial products.

## 5. Future Work

For future work, to improve our model, more data collection is necessary to minimize the standard errors of each term. This could help us to have a more unbiased confidence interval and better hypothesis tests. Also, according to the law of large numbers, as more data are collected, our sample means for each variable will become closer and closer to the true expected results.

For data collection, it can be more accurate and logical. This may help our model greatly. For example, the variable "Occupation" can be numbered in some way (e.g., based on the average income) to become an ordered variable. Then we can see a more reliable relationship between "Occupation" and "Is\_Active", and may discover more correlations among variables, which may help improve our instrumental variable model.

In addition, if the variable "Credit\_Product" only contains "Yes" or "No" (there is no such option as "N/A"), we can model it as a dummy variable. This may help us to make a better analysis in the multinomial logit model, so as to truly realize the hope of cross-selling. Frequencies for each channel code can be included to make it into count data, so a Poisson regression model can be used.

Finally, we noticed that  $R^2$  in the OLS model is relatively low, which means our independent variables are poorly explained by the explanatory variables. In order to improve our model, it is necessary to detect some new variables, such as pre-retirement wealth and debt.

## 6. References

<https://www.kaggle.com/sajidhussain3/jobathon-may-2021-credit-card-lead-prediction?select=test.csv>