

DATA_621_HW4

Chi Pong, Euclid Zhang, Jie Zou, Joseph Connolly, LeTicia Cancel

3/14/2022

```
# train_df <- read.csv("https://raw.githubusercontent.com/ezaccountz/DATA_621/main/HW4/insurance_trainin
# test_df <- read.csv("https://raw.githubusercontent.com/ezaccountz/DATA_621/main/HW4/insurance-evaluat

train_df <- read.csv("insurance_training_data.csv")
test_df <- read.csv("insurance-evaluation-data.csv")
```

DATA EXPLORATION

Data Summary

First, let's correct the formats/values of the data

```
train_df$INDEX <- NULL
train_df$INCOME <- as.numeric(gsub('[,$]', ' ', train_df$INCOME))
train_df$HOME_VAL <- as.numeric(gsub('[,$]', ' ', train_df$HOME_VAL))
train_df$BLUEBOOK <- as.numeric(gsub('[,$]', ' ', train_df$BLUEBOOK))
train_df$OLDCLAIM <- as.numeric(gsub('[,$]', ' ', train_df$OLDCLAIM))
train_df$PARENT1 <- gsub("z_", "", train_df$PARENT1)
train_df$MSTATUS <- gsub("z_", "", train_df$MSTATUS)
train_df$SEX <- gsub("z_", "", train_df$SEX)
train_df$EDUCATION <- gsub("z_", "", train_df$EDUCATION)
train_df$EDUCATION <- gsub("<", "Less Than", train_df$EDUCATION)
train_df$JOB <- gsub("z_", "", train_df$JOB)
train_df$CAR_TYPE <- gsub("z_", "", train_df$CAR_TYPE)
train_df$URBANICITY <- ifelse(train_df$URBANICITY == "Highly Urban/ Urban", "Urban", "Rural")

train_df[c("TARGET_FLAG", "PARENT1", "MSTATUS", "SEX", "EDUCATION", "JOB", "CAR_TYPE",
          "RED_CAR", "URBANICITY", "CAR_USE", "REVOKED")] <-
  lapply(train_df[c("TARGET_FLAG", "PARENT1", "MSTATUS", "SEX",
                  "EDUCATION", "JOB", "CAR_TYPE", "RED_CAR",
                  "URBANICITY", "CAR_USE", "REVOKED")], factor)
```

Below is the summary of the cleaned up data.

```
summary(train_df)
```

```
##   TARGET_FLAG  TARGET_AMT        KIDSDRV          AGE          HOMEKIDS
## 0:6008      Min.    : 0     Min.   :0.0000  Min.   :16.00  Min.   :0.0000
## 1:2153      1st Qu.: 0     1st Qu.:0.0000  1st Qu.:39.00  1st Qu.:0.0000
```

```

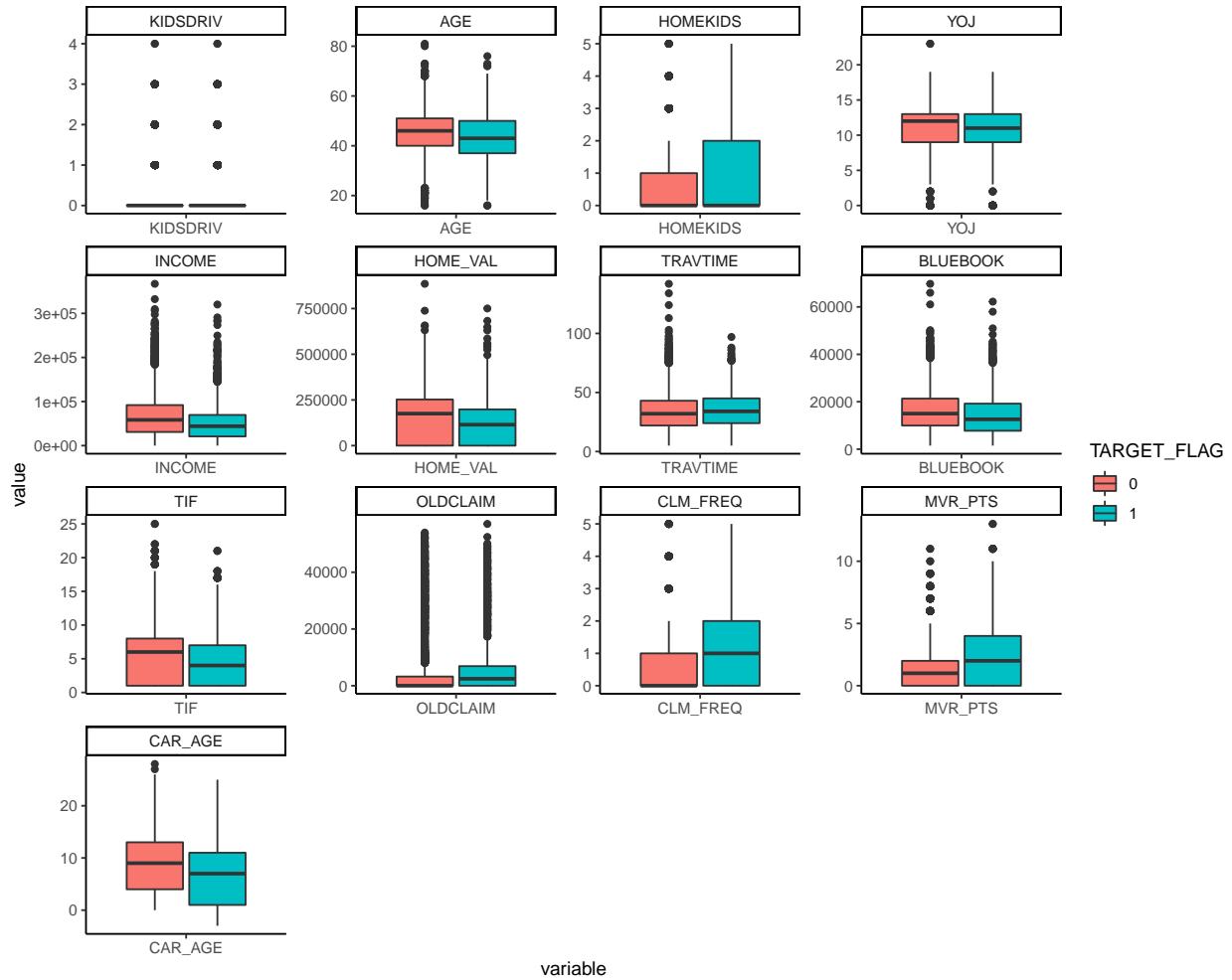
##          Median :      0    Median :0.0000    Median :45.00    Median :0.0000
##          Mean   : 1504    Mean   :0.1711    Mean   :44.79    Mean   :0.7212
##          3rd Qu.: 1036    3rd Qu.:0.0000    3rd Qu.:51.00    3rd Qu.:1.0000
##          Max.   :107586    Max.   :4.0000    Max.   :81.00    Max.   :5.0000
##                               NA's    :6
##          YOJ           INCOME        PARENT1      HOME_VAL      MSTATUS
##          Min.   : 0.0    Min.   : 0    No :7084    Min.   : 0    No :3267
##          1st Qu.: 9.0    1st Qu.: 28097 Yes:1077    1st Qu.: 0    Yes:4894
##          Median :11.0    Median : 54028
##          Mean   :10.5    Mean   : 61898
##          3rd Qu.:13.0    3rd Qu.: 85986
##          Max.   :23.0    Max.   :367030   Max.   :885282
##          NA's    :454     NA's    :445      NA's    :464
##          SEX           EDUCATION      JOB          TRAVTIME
##          F:4375      Bachelors     :2242    Blue Collar :1825    Min.   : 5.00
##          M:3786      High School   :2330    Clerical   :1271    1st Qu.: 22.00
##          Less ThanHigh School:1203  Professional:1117  Median   : 33.00
##          Masters       :1658    Manager    : 988    Mean   : 33.49
##          PhD          : 728     Lawyer     : 835    3rd Qu.: 44.00
##          Student       : 712     Student   : 712    Max.   :142.00
##          (Other)       :1413
##          CAR_USE       BLUEBOOK      TIF          CAR_TYPE
##          Commercial:3029  Min.   : 1500  Min.   : 1.000  Minivan   :2145
##          Private   :5132    1st Qu.: 9280  1st Qu.: 1.000  Panel Truck: 676
##          Median   :14440   Median   : 4.000  Pickup    :1389
##          Mean     :15710   Mean     : 5.351  Sports Car : 907
##          3rd Qu.:20850   3rd Qu.: 7.000  SUV       :2294
##          Max.    :69740   Max.    :25.000  Van      : 750
##          RED_CAR       OLDCLAIM      CLM_FREQ     REVOKED      MVR_PTS
##          no :5783      Min.   : 0    Min.   :0.0000  No :7161    Min.   : 0.000
##          yes:2378     1st Qu.: 0    1st Qu.:0.0000  Yes:1000   1st Qu.: 0.000
##          Median   : 0    Median   :0.0000
##          Mean     : 4037   Mean     :0.7986
##          3rd Qu.: 4636   3rd Qu.:2.0000
##          Max.    :57037   Max.    :5.0000
##          CAR_AGE       URBANICITY
##          Min.   : -3.000 Rural:1669
##          1st Qu.:  1.000 Urban:6492
##          Median   :  8.000
##          Mean     :  8.328
##          3rd Qu.:12.000
##          Max.    :28.000
##          NA's    :510

```

YOJ, INCOME, HOME_VAL, CAR_AGE have a lot of missing values, we will perform multiple imputations to fill in the missing values. **CAR_AGE** also has an incorrect value of -3. We will also replace it by imputation.

Box Plots

```
data.m <- melt(train_df[c("TARGET_FLAG", "KIDSDRV", "AGE", "HOMEKIDS", "YOJ", "INCOME",
                         "HOME_VAL", "TRAVTIME", "BLUEBOOK", "TIF", "OLDCLAIM", "CLM_FREQ",
                         "MVR PTS", "CAR AGE")], id.vars = 'TARGET_FLAG')
ggplot(data.m, aes(x = variable, y = value, fill = TARGET_FLAG)) + geom_boxplot() +
  facet_wrap(~ variable, scales = 'free') + theme_classic()
```



The box plots show that a lot of numeric variables are right skewed, we will transform the variables to reduce outliers.

Distribution plots

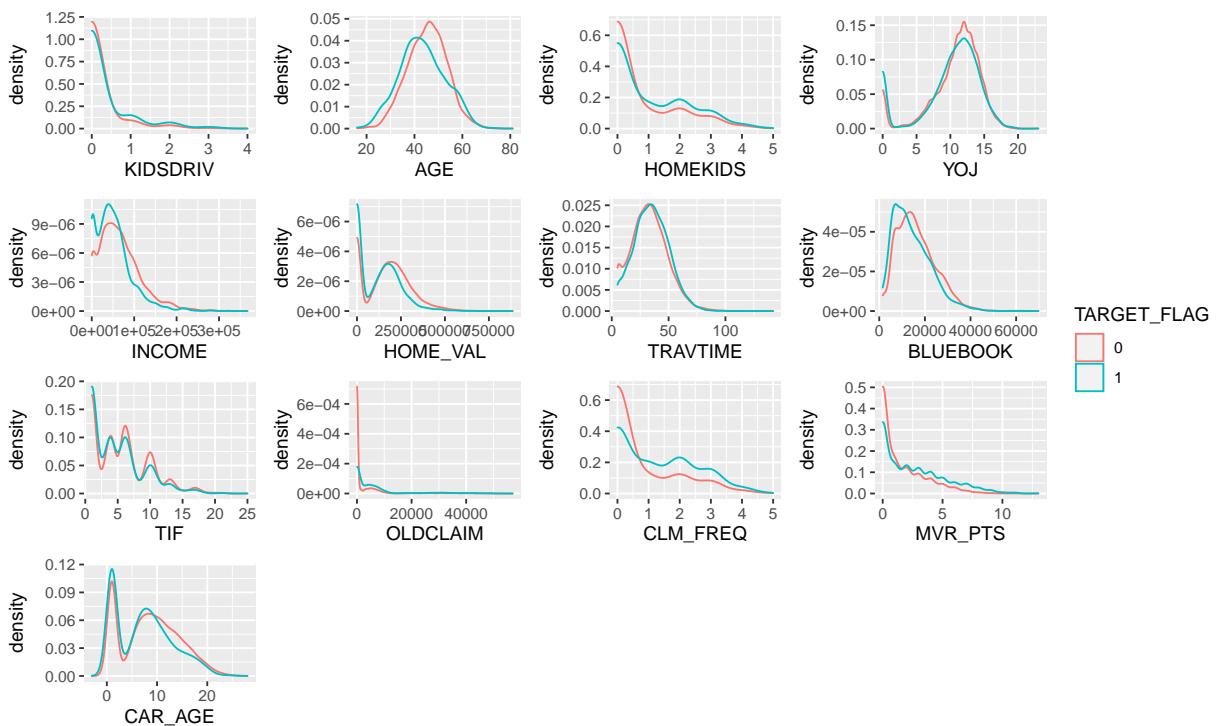
```
plot_KIDSDRV <- ggplot(train_df, aes(x=KIDSDRV, color=TARGET_FLAG)) + geom_density(na.rm =TRUE, bw=0.1)
plot_AGE <- ggplot(train_df, aes(x=AGE, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_HOMEKIDS <- ggplot(train_df, aes(x=HOMEKIDS, color=TARGET_FLAG)) + geom_density(na.rm =TRUE, bw=0.1)
plot_YOJ <- ggplot(train_df, aes(x=YOJ, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_INCOME <- ggplot(train_df, aes(x=INCOME, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_HOME_VAL <- ggplot(train_df, aes(x=HOME_VAL, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
```

```

plot_TRAVTIME <- ggplot(train_df, aes(x=TRAVTIME, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_BLUEBOOK <- ggplot(train_df, aes(x=BLUEBOOK, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_TIF <- ggplot(train_df, aes(x=TIF, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_OLDCLAIM <- ggplot(train_df, aes(x=OLDCLAIM, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_CLM_FREQ <- ggplot(train_df, aes(x=CLM_FREQ, color=TARGET_FLAG)) + geom_density(na.rm =TRUE, bw=0.4)
plot_MVR PTS <- ggplot(train_df, aes(x=MVR PTS, color=TARGET_FLAG)) + geom_density(na.rm =TRUE, bw=0.4)
plots_CAR AGE <- ggplot(train_df, aes(x=CAR AGE, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)

plot_KIDSDRV+plot_AGE+plot_HOMEKIDS+plot_YOJ+plot_INCOME+plot_HOME_VAL+
  plot_TRAVTIME+plot_BLUEBOOK+plot_TIF+plot_OLDCLAIM+plot_CLM_FREQ+
  plot_MVR PTS+plots_CAR AGE+
  plot_layout(ncol = 4, guides = "collect")

```



Most of the distributions are similar for target = 0 and target = 1. **OLDCLAIM** and **CLM_FREQ** are good candidates predicting whether there is a crash.

We can also look at the categorical variables:

```

plot_PARENT1 <- ggplot(train_df,aes(x=PARENT1,fill=TARGET_FLAG))+geom_bar(position = position_dodge())
plot_MSTATUS <- ggplot(train_df,aes(x=MSTATUS,fill=TARGET_FLAG))+geom_bar(position = position_dodge())
plot_SEX <- ggplot(train_df,aes(x=SEX,fill=TARGET_FLAG))+geom_bar(position = position_dodge())
plot_EDUCATION <- ggplot(train_df,aes(x=substring(train_df$EDUCATION,1,5),fill=TARGET_FLAG))+
  geom_bar(position = position_dodge())+xlab("EDUCATION")
plot_JOB <- ggplot(train_df,aes(x=substring(train_df$JOB,1,2),fill=TARGET_FLAG))+
  geom_bar(position = position_dodge())+xlab("JOB")
plot_CAR_TYPE <- ggplot(train_df,aes(x=substring(train_df$CAR_TYPE,1,4),fill=TARGET_FLAG))+
  geom_bar(position = position_dodge())+xlab("CAR_TYPE")
plot_RED_CAR <- ggplot(train_df,aes(x=RED_CAR,fill=TARGET_FLAG))+geom_bar(position = position_dodge())
plot_URBANICITY <- ggplot(train_df,aes(x=URBANICITY,fill=TARGET_FLAG))+geom_bar(position = position_dodge())

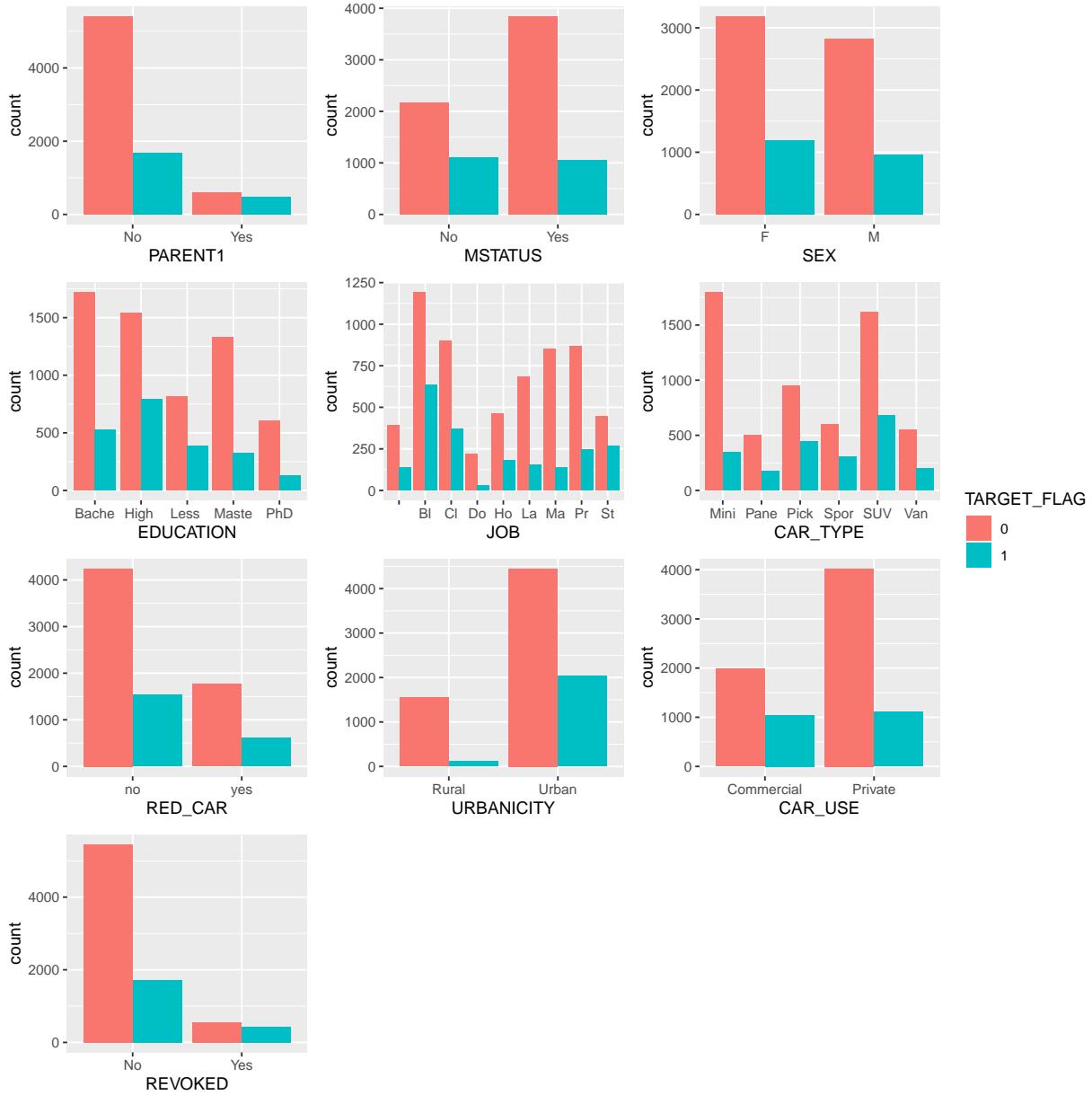
```

```

plot_CAR_USE <- ggplot(train_df,aes(x=CAR_USE,fill=TARGET_FLAG))+geom_bar(position = position_dodge())
plot_REVOKED <- ggplot(train_df,aes(x=REVOKED,fill=TARGET_FLAG))+geom_bar(position = position_dodge())

plot_PARENT1+plot_MSTATUS+plot_SEX+plot_EDUCATION+plot_JOB+plot_CAR_TYPE+plot_RED_CAR+
  plot_URBANICITY+plot_CAR_USE+plot_REVOKED+plot_layout(ncol = 3, guides = "collect")

```

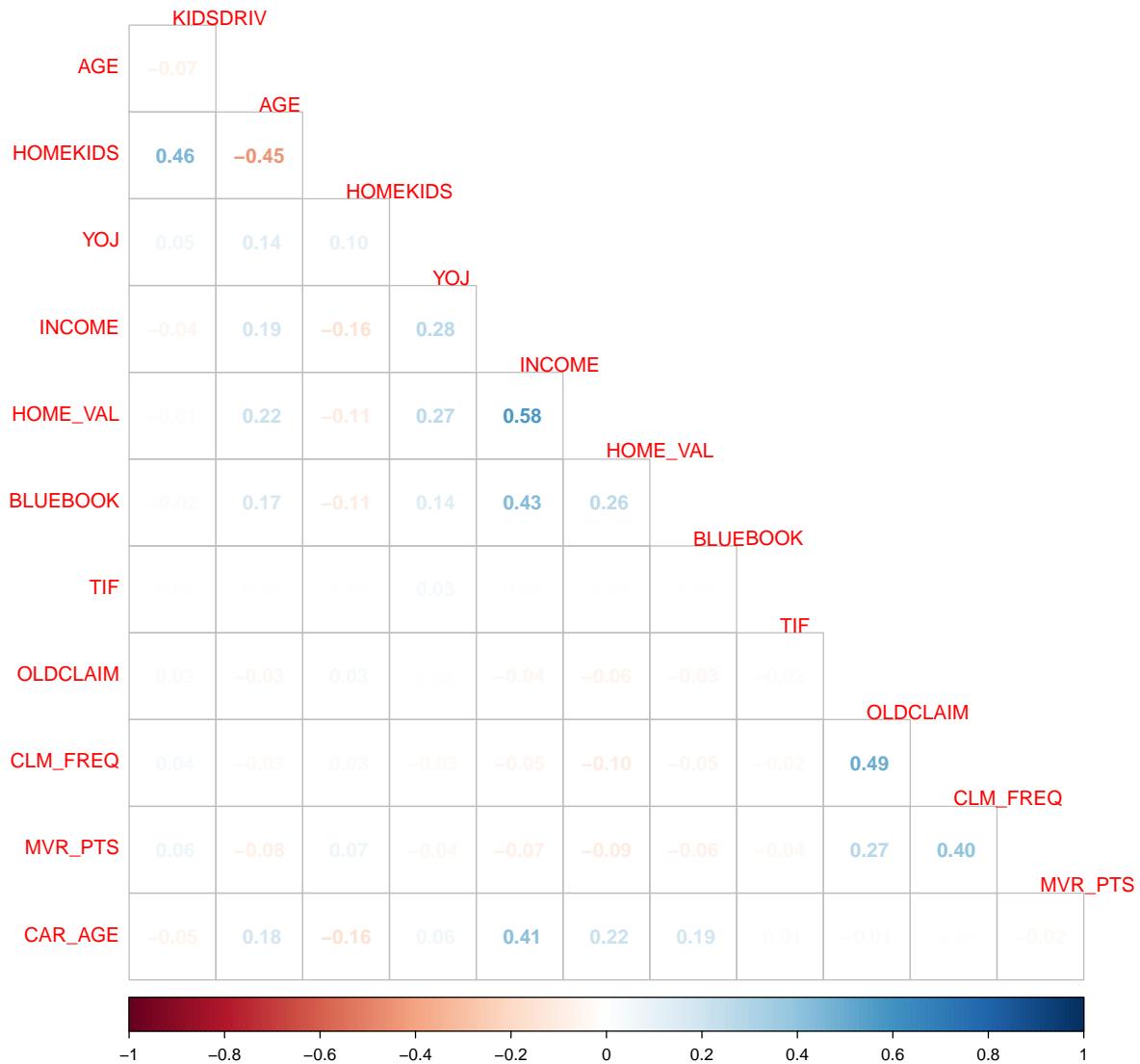


PARENT1, MSTATUS, URBANICITY, CAR_USE AND REVOKED seem to have notable difference in the distributions between target = 0 and target = 1

Correlations

Now let's look at the correlations between the variables

```
corrplot::corrplot(cor(train_df[c("KIDSDRV", "AGE", "HOMEKIDS", "YOJ", "INCOME",
                                "HOME_VAL", "BLUEBOOK", "TIF", "OLDCLAIM", "CLM_FREQ",
                                "MVR PTS", "CAR AGE")], use = "na.or.complete"),
                   method = 'number', type = 'lower', diag = FALSE, tl.srt = 0.1)
```



None of the variables have very strong correlations.

DATA PREPARATION

Data Imputation

```
#save the indicators of missing values. It will be used to verify the distributions
#of the imputed values
YOJ_NA <- is.na(train_df$YOJ)
INCOME_NA <- is.na(train_df$INCOME)
HOME_VAL_NA <- is.na(train_df$HOME_VAL)
CAR_AGE_NA <- is.na(train_df$CAR_AGE)

#remove incorrect CAR_AGE value for imputation
train_df$CAR_AGE[train_df$CAR_AGE < 0] <- NA

#temporary exclude TARGET_FLAG and TARGET_AMT in our imputation
TARGET_FLAG <- train_df$TARGET_FLAG
TARGET_AMT <- train_df$TARGET_AMT
train_df$TARGET_FLAG <- NULL
train_df$TARGET_AMT <- NULL

#save the imputation models to impute the test data set later
mickey <- parlmice(train_df, maxit = 5, m = 1, printFlag = FALSE, seed = 2022)

#save the imputation result
train_df <- complete(mickey,1)

#Add TARGET_FLAG and TARGET_AMT back to our dataframe
train_df$TARGET_FLAG <- TARGET_FLAG
train_df$TARGET_AMT <- TARGET_AMT
TARGET_FLAG <- NULL
TARGET_AMT <- NULL

#write.csv(train_df,"train_df.csv", row.names = FALSE)

# train_df <- read.csv("train_df.csv", stringsAsFactors = TRUE)
# train_df$TARGET_FLAG <- as.factor(train_df$TARGET_FLAG)
```

The plots on the top row below show the distributions of the values from the original data

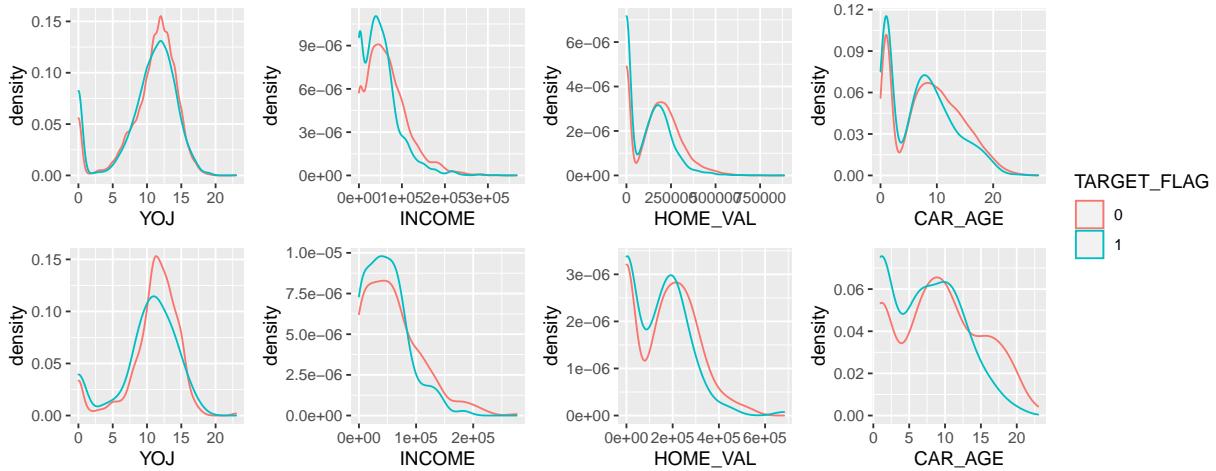
The plots on the bottom row below show the distributions of the imputed values

```
plot_YOJ <- ggplot(train_df[!YOJ_NA], aes(x=YOJ, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_INCOME <- ggplot(train_df[!INCOME_NA], aes(x=INCOME, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_HOME_VAL <- ggplot(train_df[!HOME_VAL_NA], aes(x=HOME_VAL, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_CAR_AGE <- ggplot(train_df[!CAR_AGE_NA], aes(x=CAR_AGE, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)

plot_YOJ2 <- ggplot(train_df[YOJ_NA], aes(x=YOJ, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_INCOME2 <- ggplot(train_df[INCOME_NA], aes(x=INCOME, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_HOME_VAL2 <- ggplot(train_df[HOME_VAL_NA], aes(x=HOME_VAL, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)
plot_CAR_AGE2 <- ggplot(train_df[CAR_AGE_NA], aes(x=CAR_AGE, color=TARGET_FLAG)) + geom_density(na.rm =TRUE)

plot_YOJ+plot_INCOME+plot_HOME_VAL+plot_CAR_AGE+
```

```
plot_YOJ2+plot_INCOME2+plot_HOME_VAL2+plot_CAR_AGE2+
plot_layout(ncol = 4, guides = "collect")
```



The distributions look similar and so the imputed values are plausible

Data Transformation

Since **YOJ** and **HOME_VAL** are zero-inflated. We would add a dummy variable for each of them indicating whether the variable is 0. The effect of variables

YOJ: The density plot shows the variable is zero-inflated. The coefficient for YOJ=0 and the coefficient for YOJ>0 may be significantly different. Therefore, we would add a dummy variable indicating whether the variable is 0. **HOME_VAL:** The variable is also zero-inflated, we would add a dummy variable indicating whether the person has a house. **INCOME:** We would add a dummy variable indicating whether the person has a job. Practically, it is a key factor in insurance pricing. **OLDCLAIM:** We would add a dummy variable indicating whether the person had an old claim. The coefficient for OLDCLAIM=0 and the coefficient for OLDCLAIM>0 may be significantly different.

INCOME, HOME_VAL, BLUEBOOK, OLDCLAIM: We will log transform all monetary variables as they are right-skewed.

```
train_df$YOJ_Y <- as.factor(ifelse(train_df$YOJ == 0, 0, 1))
train_df$INCOME_Y <- as.factor(ifelse(train_df$INCOME == 0, 0, 1))
train_df$HOME_VAL_Y <- as.factor(ifelse(train_df$HOME_VAL == 0, 0, 1))
train_df$OLDCLAIM_Y <- as.factor(ifelse(train_df$OLDCLAIM == 0, 0, 1))

train_df$INCOME_LOG <- log(train_df$INCOME+1)
train_df$HOME_VAL_LOG <- log(train_df$HOME_VAL+1)
train_df$BLUEBOOK_LOG <- log(train_df$BLUEBOOK)
train_df$OLDCLAIM_LOG <- log(train_df$OLDCLAIM+1)

train_df$INCOME <- NULL
train_df$HOME_VAL <- NULL
train_df$BLUEBOOK <- NULL
train_df$OLDCLAIM <- NULL

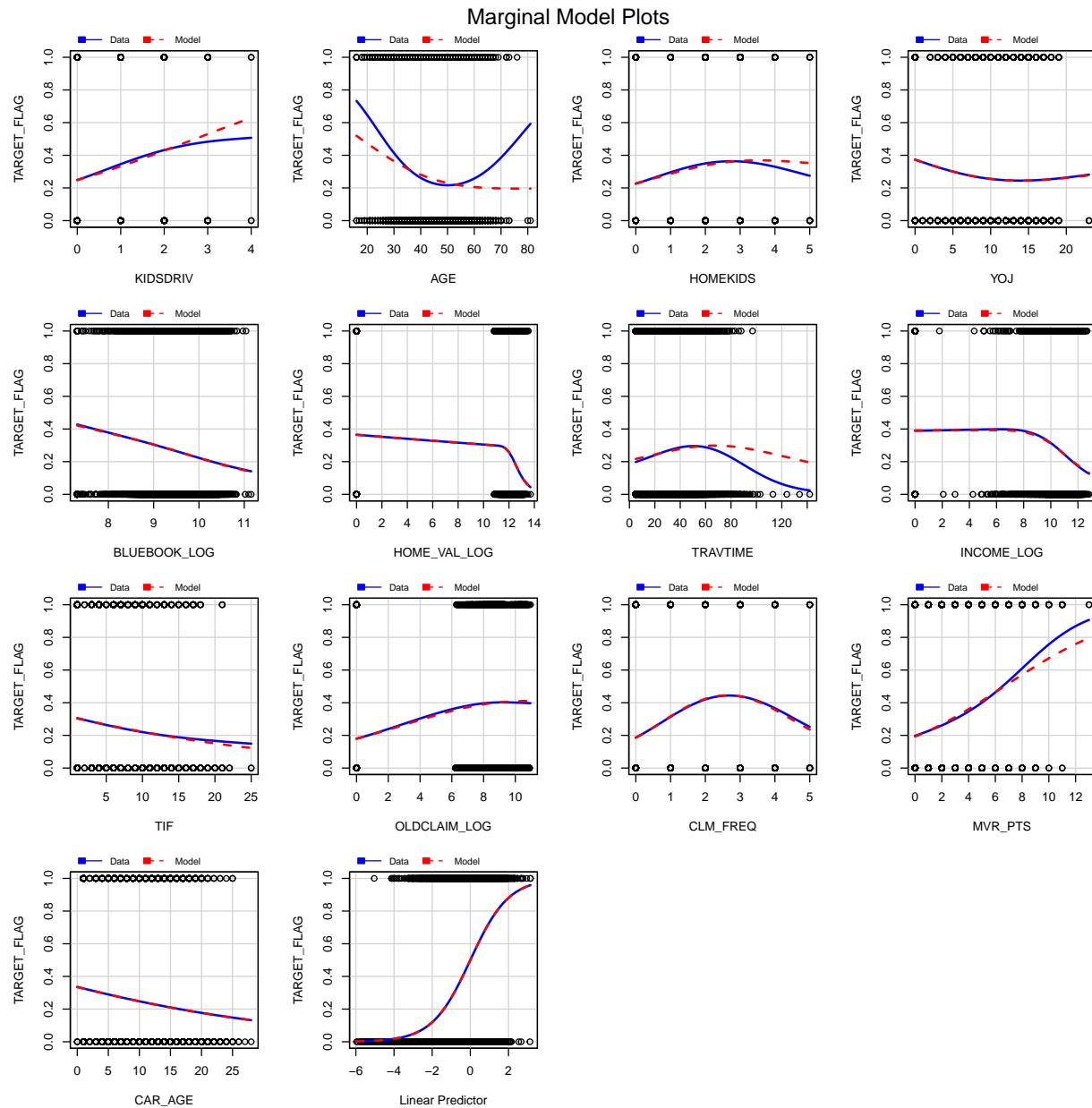
logistic_train_df <- train_df
```

Logistic Models

First, let build a test model to see if any additional transformations are needed to fit our logistic models

```
test_model <- glm(TARGET_FLAG~.-TARGET_AMT, family = binomial, logistic_train_df)
```

```
marginalModelPlots(test_model, ~KIDSDRV + AGE + HOMEKIDS + YOJ + BLUEBOOK_LOG +
HOME_VAL_LOG + TRAVTIME+ INCOME_LOG + TIF + OLDCLAIM_LOG +
CLM_FREQ + MVR PTS + CAR_AGE, layout =c(4,4))
```



Additonal Transformations

Additional transformation are needed for **KIDSDRV**, **AGE**, **HOMEKIDS**, **TRAVTIME**, and **MVR_PTS**

From the density plots above, the see that **AGE** is approximately normal for both target = 0 and target = 1. From the text book *A Modern Approach To Regression With R*, if the variance of the variable is different for the two response value, then a squared term should be added.

```
data.frame(Variance_of_AGE_TARGET0 = c(var(logistic_train_df$AGE[logistic_train_df$TARGET_FLAG == 0])),  
           Variance_of_AGE_TARGET1 = c(var(logistic_train_df$AGE[logistic_train_df$TARGET_FLAG == 1])))  
  
##  Variance_of_AGE_TARGET0 Variance_of_AGE_TARGET1  
## 1                 67.26915                 92.02521  
  
var.test(AGE ~ TARGET_FLAG, logistic_train_df, alternative = "two.sided")  
  
##  
##  F test to compare two variances  
##  
## data: AGE by TARGET_FLAG  
## F = 0.73099, num df = 6007, denom df = 2152, p-value < 2.2e-16  
## alternative hypothesis: true ratio of variances is not equal to 1  
## 95 percent confidence interval:  
##  0.6814095 0.7832593  
## sample estimates:  
## ratio of variances  
##                 0.7309861
```

The variance is apparently different. We will add a squared term and check if that fits the model.

```
logistic_train_df$AGE_Squared <- logistic_train_df$AGE^2
```

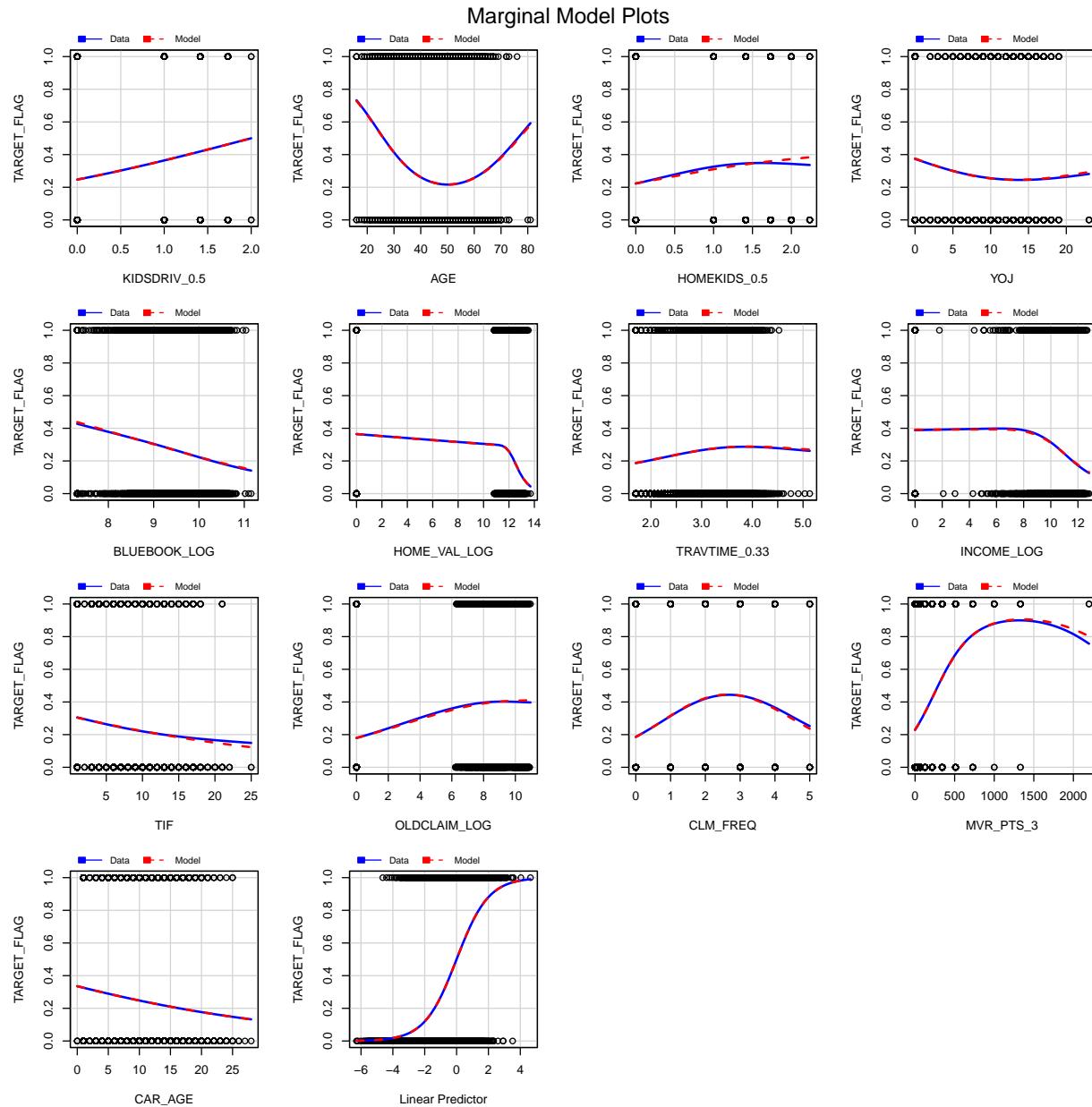
For **KIDSDRV**, **HOMEKIDS**, **TRAVTIME**, and **MVR_PTS**, power transformations are used and the powers are determined by trial and error

```
logistic_train_df$KIDSDRV_0.5 <- (logistic_train_df$KIDSDRV)^0.5  
logistic_train_df$HOMEKIDS_0.5 <- (logistic_train_df$HOMEKIDS)^0.5  
logistic_train_df$MVR PTS_3 <- (logistic_train_df$MVR PTS)^3  
logistic_train_df$TRAVTIME_0.33 <- (logistic_train_df$TRAVTIME)^0.33  
  
logistic_train_df$KIDSDRV <- NULL  
logistic_train_df$HOMEKIDS <- NULL  
logistic_train_df$MVR PTS <- NULL  
logistic_train_df$TRAVTIME <- NULL
```

After all the transformations, the test model now fits our data well

```
test_model <- glm(TARGET_FLAG~.-TARGET_AMT, family = binomial, logistic_train_df)
```

```
marginalModelPlots(test_model, ~KIDSDRV_0.5 + AGE + HOMEKIDS_0.5 + YOJ + BLUEBOOK_LOG +
HOME_VAL_LOG + TRAVTIME_0.33 + INCOME_LOG + TIF + OLDCLAIM_LOG +
CLM_FREQ + MVR PTS_3 + CAR_AGE, layout =c(4,4))
```



Building Models

Full Model

First we build a full model with all predictors

```
logi_full <- glm(TARGET_FLAG~.-TARGET_AMT, family = binomial, logistic_train_df)
```

```

summary(logi_full)

##
## Call:
## glm(formula = TARGET_FLAG ~ . - TARGET_AMT, family = binomial,
##      data = logistic_train_df)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -2.6640 -0.7034 -0.3886  0.5867  3.0479
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)                3.6444309  0.9039895  4.031 5.54e-05 ***
## AGE                      -0.2071383  0.0261547 -7.920 2.38e-15 ***
## YOJ                       0.0104534  0.0121841  0.858 0.390922
## PARENT1Yes                 0.2684085  0.1179489  2.276 0.022868 *
## MSTATUSYes                -0.5272217  0.0909780 -5.795 6.83e-09 ***
## SEXM                      0.0333886  0.1087867  0.307 0.758906
## EDUCATIONHigh School       0.4134866  0.0917760  4.505 6.63e-06 ***
## EDUCATIONLess ThanHigh School 0.3609941  0.1193708  3.024 0.002493 **
## EDUCATIONMasters            0.1041391  0.1412517  0.737 0.460965
## EDUCATIONPhD                 0.0891467  0.1784974  0.499 0.617477
## JOBBlue Collar              0.4057605  0.1859322  2.182 0.029087 *
## JOBClerical                  0.4679805  0.1971209  2.374 0.017593 *
## JOBDoctor                     -0.4484738  0.2678414 -1.674 0.094052 .
## JOBHome Maker                 0.0047536  0.2311896 -0.021 0.983596
## JOBLawyer                     0.1455225  0.1699818  0.856 0.391939
## JOBManager                   -0.5224542  0.1710772 -3.054 0.002259 **
## JOBProfessional               0.2274060  0.1785379  1.274 0.202765
## JOBStudent                    -0.0840792  0.2405653 -0.350 0.726709
## CAR_USEPrivate                0.7514255  0.0926888 -8.107 5.19e-16 ***
## TIF                         -0.0559303  0.0074030 -7.555 4.19e-14 ***
## CAR_TYPEPanel Truck           0.5586394  0.1519431  3.677 0.000236 ***
## CAR_TYPEPickup                 0.5955788  0.1014877  5.868 4.40e-09 ***
## CAR_TYPESports Car            0.8816775  0.1295854  6.804 1.02e-11 ***
## CAR_TYPESUV                   0.7163950  0.1080258  6.632 3.32e-11 ***
## CAR_TYPEVan                   0.6841286  0.1266153  5.403 6.55e-08 ***
## RED_CARyes                   -0.0590874  0.0874078 -0.676 0.499042
## CLM_FREQ                      0.0490468  0.0449383  1.091 0.275086
## REVOKEDYes                   0.8615942  0.0893347  9.645 < 2e-16 ***
## CAR_AGE                      -0.0044851  0.0076074 -0.590 0.555480
## URBANICITYUrban                2.3737654  0.1135661 20.902 < 2e-16 ***
## YOJ_Y1                        -0.4426728  0.3129938 -1.414 0.157269
## INCOME_Y1                      0.5974304  0.5749808  1.039 0.298784
## HOME_VAL_Y1                   2.9135134  1.4037424  2.076 0.037937 *
## OLDCLAIM_Y1                   1.7559812  0.4180983  4.200 2.67e-05 ***
## INCOME_LOG                     -0.0947848  0.0552928 -1.714 0.086486 .
## HOME_VAL_LOG                  -0.2662495  0.1155750 -2.304 0.021240 *
## BLUEBOOK_LOG                  -0.3268810  0.0597510 -5.471 4.48e-08 ***
## OLDCLAIM_LOG                  -0.1588576  0.0459550 -3.457 0.000547 ***
## AGE_Squared                    0.0022856  0.0002860  7.993 1.32e-15 ***
## KIDSDRIV_0.5                  0.6893959  0.0861530  8.002 1.22e-15 ***

```

```

## HOMEKIDS_0.5          0.0036184  0.0695918  0.052 0.958533
## MVR PTS_3            0.0016580  0.0002609  6.355 2.08e-10 ***
## TRAVTIME_0.33         0.4464462  0.0557077  8.014 1.11e-15 ***
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 9418.0  on 8160  degrees of freedom
## Residual deviance: 7187.4  on 8118  degrees of freedom
## AIC: 7273.4
##
## Number of Fisher Scoring iterations: 5

```

Backward Elimination by AIC

```

logi_AIC <- step(logi_full,trace=0)

summary(logi_AIC)

##
## Call:
## glm(formula = TARGET_FLAG ~ AGE + PARENT1 + MSTATUS + EDUCATION +
##     JOB + CAR_USE + TIF + CAR_TYPE + REVOKED + URBANICITY + HOME_VAL_Y +
##     OLDCALL_Y + INCOME_LOG + HOME_VAL_LOG + BLUEBOOK_LOG + OLDCALL_LOG +
##     AGE_Squared + KIDSDRIV_0.5 + MVR PTS_3 + TRAVTIME_0.33, family = binomial,
##     data = logistic_train_df)
##
## Deviance Residuals:
##    Min      1Q   Median      3Q      Max
## -2.6745 -0.7037 -0.3901  0.5913  3.0388
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)                3.5777438  0.8248926  4.337 1.44e-05 ***
## AGE                      -0.2097624  0.0251070 -8.355 < 2e-16 ***
## PARENT1Yes                 0.2824317  0.1018929  2.772 0.005574 **
## MSTATUSYes                -0.5123869  0.0860036 -5.958 2.56e-09 ***
## EDUCATIONHigh School       0.4380366  0.0846295  5.176 2.27e-07 ***
## EDUCATIONLess ThanHigh School 0.3957984  0.1101180  3.594 0.000325 ***
## EDUCATIONMasters            0.0849668  0.1351185  0.629 0.529459
## EDUCATIONPhD                0.0656773  0.1734382  0.379 0.704927
## JOBBlue Collar              0.4159810  0.1857392  2.240 0.025117 *
## JOBClerical                  0.4863252  0.1957699  2.484 0.012985 *
## JOBDoctor                   -0.4460820  0.2675188 -1.667 0.095419 .
## JOBHome Maker                0.0597830  0.2219439  0.269 0.787652
## JOBLawyer                     0.1432302  0.1698509  0.843 0.399078
## JOBManager                  -0.5220403  0.1709412 -3.054 0.002259 **
## JOBProfessional               0.2315211  0.1784232  1.298 0.194426
## JOBStudent                   -0.0106369  0.2261358 -0.047 0.962483
## CAR USEPrivate                0.7486815  0.0925857 -8.086 6.15e-16 ***
## TIF                         -0.0558830  0.0073949 -7.557 4.13e-14 ***

```

```

## CAR_TYPEPanel Truck          0.5594484  0.1455788  3.843 0.000122 ***
## CAR_TYPEPickup            0.5966744  0.1013620  5.887 3.94e-09 ***
## CAR_TYPESports Car        0.8818399  0.1093565  8.064 7.39e-16 ***
## CAR_TYPESUV               0.7172636  0.0865353  8.289 < 2e-16 ***
## CAR_TYPEVan               0.6828012  0.1231073  5.546 2.92e-08 ***
## REVOKEDYes                0.8604889  0.0892496  9.641 < 2e-16 ***
## URBANICITYUrban           2.3708057  0.1134411 20.899 < 2e-16 ***
## HOME_VAL_Y1                3.1808717  1.2921465  2.462 0.013828 *
## OLDCLAIM_Y1                1.8515523  0.4065296  4.555 5.25e-06 ***
## INCOME_LOG                 -0.0621858  0.0144766 -4.296 1.74e-05 ***
## HOME_VAL_LOG                -0.2879553  0.1063213 -2.708 0.006762 **
## BLUEBOOK_LOG                -0.3309724  0.0556910 -5.943 2.80e-09 ***
## OLDCLAIM_LOG                -0.1580488  0.0459130 -3.442 0.000577 ***
## AGE_Squared                  0.0023221  0.0002771  8.379 < 2e-16 ***
## KIDSDRIV_0.5                 0.6985398  0.0748798  9.329 < 2e-16 ***
## MVR PTS_3                    0.0016498  0.0002604  6.336 2.36e-10 ***
## TRAVTIME_0.33                 0.4469181  0.0556669  8.028 9.87e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 9418  on 8160  degrees of freedom
## Residual deviance: 7192  on 8126  degrees of freedom
## AIC: 7262
##
## Number of Fisher Scoring iterations: 5

```

Backward Elimination by BIC

```

logi_BIC <- step(logi_full,trace=0,k=log(nrow(logistic_train_df)))

summary(logi_BIC)

##
## Call:
## glm(formula = TARGET_FLAG ~ AGE + MSTATUS + JOB + CAR_USE + TIF +
##     CAR_TYPE + REVOKED + URBANICITY + HOME_VAL_Y + OLDCLAIM_Y +
##     INCOME_LOG + HOME_VAL_LOG + BLUEBOOK_LOG + OLDCLAIM_LOG +
##     AGE_Squared + KIDSDRIV_0.5 + MVR PTS_3 + TRAVTIME_0.33, family = binomial,
##     data = logistic_train_df)
##
## Deviance Residuals:
##      Min        1Q    Median        3Q       Max
## -2.7313   -0.7095   -0.3954    0.5993    3.0082
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)              4.2740556  0.7967312  5.364 8.12e-08 ***
## AGE                     -0.2231719  0.0245256 -9.100 < 2e-16 ***
## MSTATUSYes              -0.5901268  0.0769684 -7.667 1.76e-14 ***
## JOBBlue Collar          0.6200658  0.1344731  4.611 4.01e-06 ***

```

```

## JOBCLerical      0.7134799  0.1503809  4.744 2.09e-06 ***
## JOBDoctor        -0.4918996  0.2500827 -1.967 0.049189 *
## JOBHome Maker    0.1223997  0.1998509  0.612 0.540236
## JOBLawyer         0.1234993  0.1638639  0.754 0.451047
## JOBManager       -0.5099810  0.1541476 -3.308 0.000938 ***
## JOBProfessional   0.2346458  0.1436282  1.634 0.102321
## JOBStudent        0.1919219  0.1862413  1.031 0.302775
## CAR_USEPrivate    -0.7094092  0.0872395 -8.132 4.23e-16 ***
## TIF               -0.0557254  0.0073626 -7.569 3.77e-14 ***
## CAR_TYPEPanel Truck 0.5831294  0.1433708  4.067 4.76e-05 ***
## CAR_TYPEPickup    0.6169158  0.1001922  6.157 7.40e-10 ***
## CAR_TYPESports Car 0.8875481  0.1089506  8.146 3.75e-16 ***
## CAR_TYPESUV        0.7228083  0.0862012  8.385 < 2e-16 ***
## CAR_TYPEVan        0.6882928  0.1221074  5.637 1.73e-08 ***
## REVOKEDYes         0.8579156  0.0889535  9.645 < 2e-16 ***
## URBANICITYUrban   2.3407195  0.1128161 20.748 < 2e-16 ***
## HOME_VAL_Y1        3.9797854  1.2593085  3.160 0.001576 **
## OLDCLAIM_Y1        1.8132072  0.4048088  4.479 7.49e-06 ***
## INCOME_LOG          -0.0662383  0.0143727 -4.609 4.05e-06 ***
## HOME_VAL_LOG        -0.3537823  0.1036117 -3.415 0.000639 ***
## BLUEBOOK_LOG        -0.3379876  0.0554692 -6.093 1.11e-09 ***
## OLDCLAIM_LOG        -0.1539811  0.0457257 -3.367 0.000759 ***
## AGE_Squared          0.0024381  0.0002729  8.935 < 2e-16 ***
## KIDSDRIV_0.5        0.7588245  0.0712270 10.654 < 2e-16 ***
## MVR PTS_3            0.0016533  0.0002610  6.334 2.40e-10 ***
## TRAVTIME_0.33        0.4309331  0.0553343  7.788 6.82e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 9418 on 8160 degrees of freedom
## Residual deviance: 7228 on 8131 degrees of freedom
## AIC: 7288
##
## Number of Fisher Scoring iterations: 5

```

Backward Elimination with Chi-square test

Starting with our full model, perform backward elimination with Chi-square test.

```
#Define a function to perform backward elimination with Chi-square test
#using the signficancy / alpha as one of the parameters
```

```
backward_chi <- function (train_df, signficancy) {
  glm_string <- "TARGET_FLAG~.-TARGET_AMT"
  glm_formula <- as.formula(glm_string)

  repeat{
    drop1_chi <- drop1(glm(glm_formula, family=binomial, train_df), test="Chi")

    chi_result <- data.frame(preditors = rownames(drop1_chi)[-1],
                               p_value = drop1_chi[-1,5])
```

```

chi_result <- chi_result[order(chi_result$p_value, decreasing=TRUE),]

if(chi_result[1,2] < significancy){
  break
}
else {
  glm_string <- paste0(glm_string, "-", chi_result[1,1])
  glm_formula <- as.formula(glm_string)
}

return(glm_formula)
}

```

model with alpha 0.001 (based on Chi-square test)**

```

logi_chi_0.001 <- backward_chi(logistic_train_df, 0.001)
logi_chi_0.001 <- glm(logi_chi_0.001, family=binomial, logistic_train_df)
summary(logi_chi_0.001)

```

```

##
## Call:
## glm(formula = logi_chi_0.001, family = binomial, data = logistic_train_df)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max 
## -2.6333 -0.7031 -0.3936  0.5994  3.0483 
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)    
## (Intercept)               4.2076516  0.8043810  5.231 1.69e-07 ***
## AGE                     -0.2261158  0.0245924 -9.195 < 2e-16 ***
## MSTATUSYes              -0.6162293  0.0771227 -7.990 1.35e-15 ***
## EDUCATIONHigh School     0.4620511  0.0839410  5.504 3.70e-08 ***
## EDUCATIONLess ThanHigh School 0.4382835  0.1086397  4.034 5.48e-05 ***
## EDUCATIONMasters          0.0699021  0.1348271  0.518 0.604139  
## EDUCATIONPhD              0.0217516  0.1723172  0.126 0.899550  
## JOBBlue Collar            0.4268301  0.1853849  2.302 0.021313 *  
## JOBCLerical                0.5421744  0.1943967  2.789 0.005287 ** 
## JOBDoctor                  -0.4426228  0.2664048 -1.661 0.096620 .  
## JOBHome Maker              0.1563657  0.2180316  0.717 0.473269  
## JOBLawyer                  0.1632191  0.1693621  0.964 0.335182  
## JOBManager                 -0.5033168  0.1703919 -2.954 0.003138 ** 
## JOBProfessional             0.2392954  0.1779847  1.344 0.178796  
## JOBStudent                  -0.0257806  0.2259859 -0.114 0.909174  
## CAR_USEPrivate              -0.7533441  0.0925503 -8.140 3.96e-16 ***
## TIF                        -0.0558940  0.0073870 -7.567 3.83e-14 *** 
## CAR_TYPEPanel Truck         0.5354901  0.1452482  3.687 0.000227 *** 
## CAR_TYPEPickup              0.5968769  0.1012925  5.893 3.80e-09 *** 
## CAR_TYPESports Car          0.8861276  0.1092253  8.113 4.94e-16 *** 
## CAR_TYPESUV                 0.7182018  0.0864262  8.310 < 2e-16 *** 
## CAR_TYPEVan                 0.6624746  0.1228336  5.393 6.92e-08 *** 
## REVOKEDYes                  0.8605099  0.0891430  9.653 < 2e-16 *** 

```

```

## URBANICITYUrban          2.3679319  0.1133415  20.892 < 2e-16 ***
## OLDCLAIM_Y1              1.8659739  0.4062384   4.593 4.36e-06 ***
## INCOME_LOG                -0.0699675  0.0142213  -4.920 8.66e-07 ***
## HOME_VAL_LOG               -0.0264858  0.0069905  -3.789 0.000151 ***
## BLUEBOOK_LOG                -0.3338818  0.0555400  -6.012 1.84e-09 ***
## OLDCLAIM_LOG               -0.1593016  0.0458842  -3.472 0.000517 ***
## AGE_Squared                 0.0024626  0.0002736   9.000 < 2e-16 ***
## KIDSDRIV_0.5                  0.7561342  0.0714399  10.584 < 2e-16 ***
## MVR_PTS_3                     0.0016539  0.0002605   6.349 2.17e-10 ***
## TRAVTIME_0.33                  0.4402636  0.0555662   7.923 2.31e-15 ***
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 9418.0  on 8160  degrees of freedom
## Residual deviance: 7204.9  on 8128  degrees of freedom
## AIC: 7270.9
##
## Number of Fisher Scoring iterations: 5

```

Model Selection

Since the data is imbalanced, we would not use the threshold for our model predictions.

In business, we don't want to misclassify a person with high risk to be low risk. We also don't want to lose customers by charging low risk people at a high-risk rate. Practically, we should use a cost matrix to determine the threshold for our classification. Since we don't know the cost here, we will weight the Sensitivity and Specificity equally. We will find our optimal threshold that maximize the sum of Sensitivity and Specificity

```

logi_models <- data.frame(model=c(""),
                            DF=c(0), AIC=c(0.0000), AUC=c(0.0000),
                            Optimal_Threshold=c(0.0000), Sensitivity=c(0.0000),
                            Specificity=c(0.0000), Sum_Sens_Spec=c(0.0000))

models <- list(logi_full, logi_AIC, logi_BIC, logi_chi_0.001)
model_names <- c("logi_full", "logi_AIC", "logi_BIC", "logi_chi_0.001")
for (i in 1:length(models)) {
  logi_models[i,"model"] <- model_names[i]
  logi_models[i,"DF"] <- models[[i]]$df.residual
  logi_models[i,"AIC"] <- round(models[[i]]$aic,4)
  rocCurve <- roc(logistic_train_df$TARGET_FLAG, models[[i]]$fitted.values)
  logi_models[i,"AUC"] <- round(rocCurve$auc,4)
  roc_df <- data.frame(Sensitivity = rocCurve$sensitivities, Specificity = rocCurve$specificities,
                        Sum_Sens_Spec = rocCurve$sensitivities+rocCurve$specificities,
                        Thresholds = rocCurve$thresholds)
  roc_df <- roc_df[which.max(roc_df$Sum_Sens_Spec),]
  logi_models[i,"Optimal_Threshold"] <- roc_df$Thresholds
  logi_models[i,"Sensitivity"] <- roc_df$Sensitivity
  logi_models[i,"Specificity"] <- roc_df$Specificity
  logi_models[i,"Sum_Sens_Spec"] <- roc_df$Sum_Sens_Spec
}
logi_models

```

```

##          model   DF      AIC      AUC Optimal_Threshold Sensitivity Specificity
## 1    logi_full 8118 7273.439  0.8203        0.2928494   0.7259638   0.7651465
## 2    logi_AIC  8126 7261.968  0.8201        0.2799662   0.7436136   0.7501664
## 3    logi_BIC  8131 7287.980  0.8178        0.2978043   0.7162099   0.7689747
## 4 logi_chi_0.001 8128 7270.901  0.8194        0.2832591   0.7380399   0.7524967
##   Sum_Sens_Spec
## 1    1.491110
## 2    1.493780
## 3    1.485185
## 4    1.490537

```

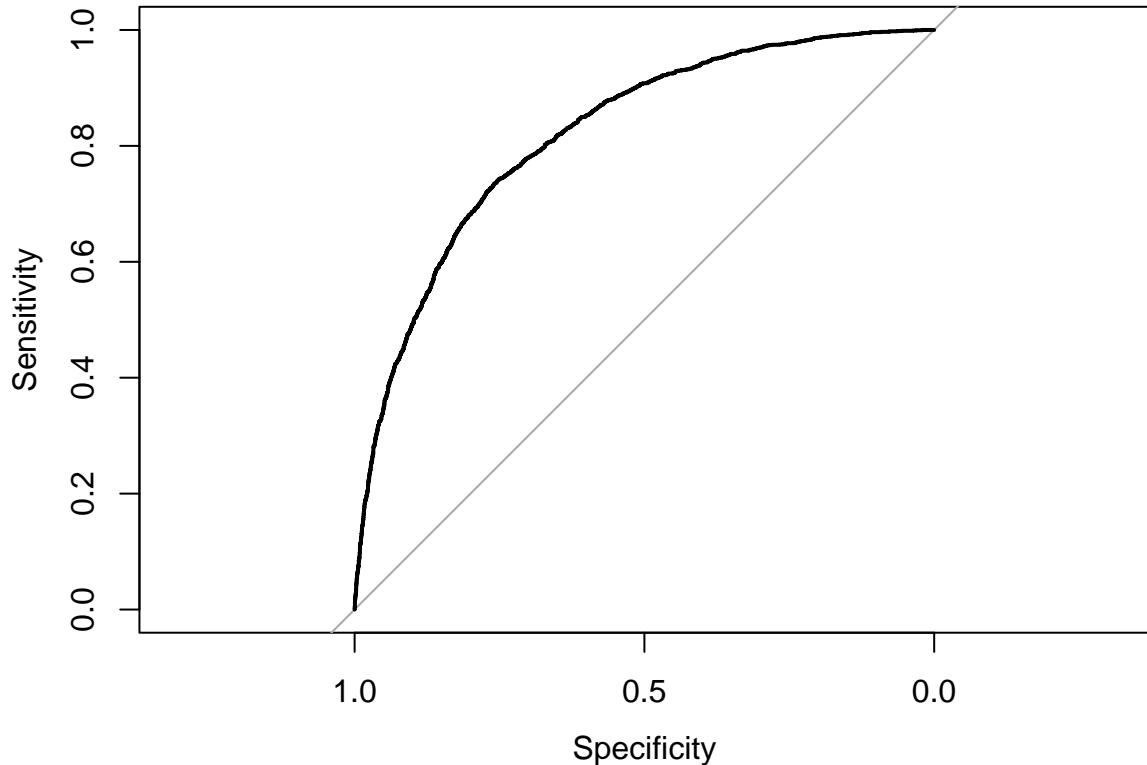
By comparing the AUC and the sum of Sensitivity and Specificity, the best model is logi_AIC

The following is the ROC of the logi_AIC model

```

rocCurve <- roc(logistic_train_df$TARGET_FLAG, logi_AIC$fitted.values)
plot(rocCurve)

```



The following is the confusion matrix of the logi_AIC model

```

predicted_class <- ifelse(logi_AIC$fitted.values>logi_models[2,"Optimal_Threshold"],1,0)
confusion_matrix <- confusionMatrix(as.factor(predicted_class),
                                     as.factor(train_df$TARGET_FLAG),positive = "1")
confusion_matrix

## Confusion Matrix and Statistics

```

```

##          Reference
## Prediction 0 1
##          0 4507 552
##          1 1501 1601
##
##          Accuracy : 0.7484
##          95% CI : (0.7389, 0.7578)
##          No Information Rate : 0.7362
##          P-Value [Acc > NIR] : 0.006038
##
##          Kappa : 0.4326
##
##          McNemar's Test P-Value : < 2.2e-16
##
##          Sensitivity : 0.7436
##          Specificity : 0.7502
##          Pos Pred Value : 0.5161
##          Neg Pred Value : 0.8909
##          Prevalence : 0.2638
##          Detection Rate : 0.1962
##          Detection Prevalence : 0.3801
##          Balanced Accuracy : 0.7469
##
##          'Positive' Class : 1
##

```

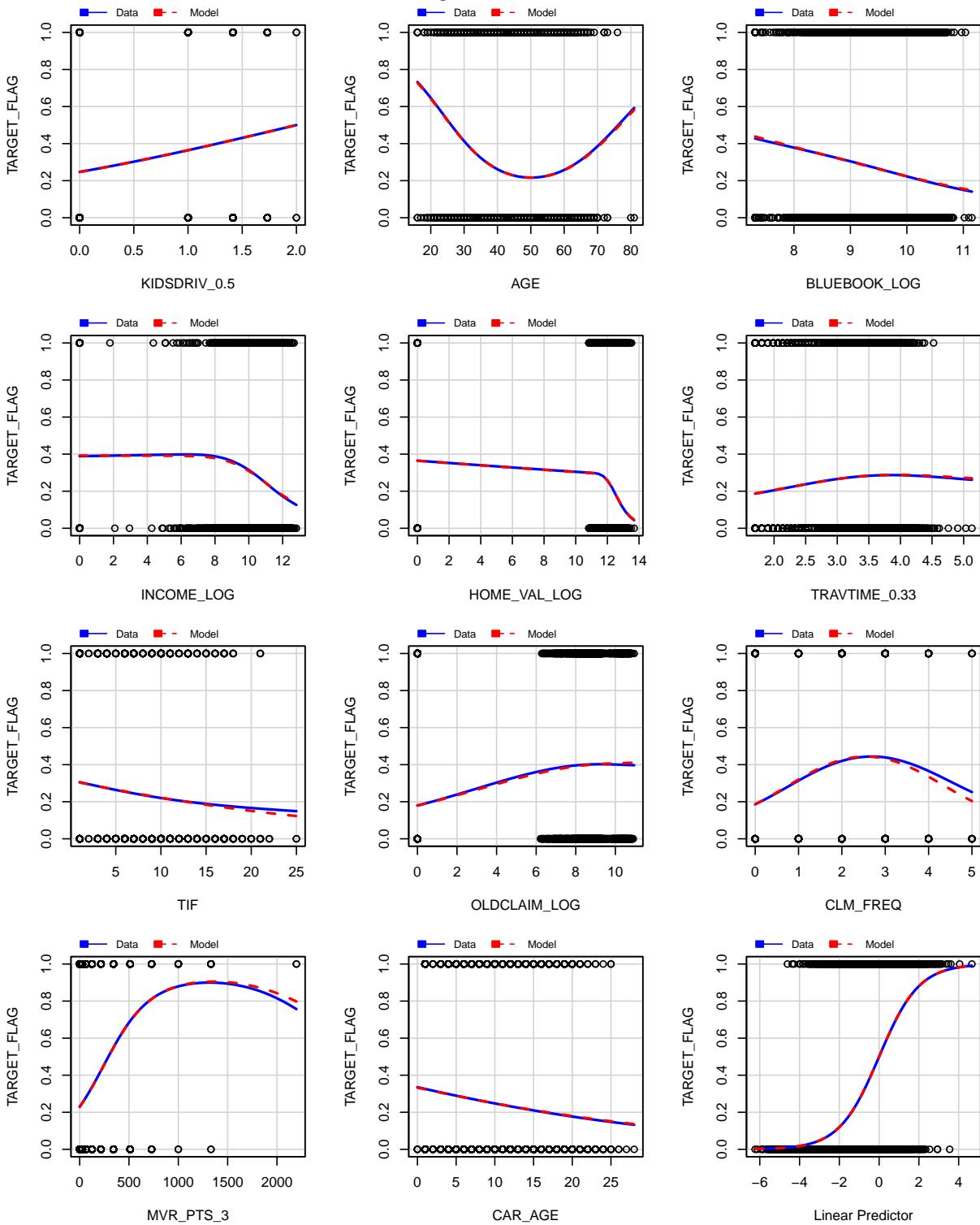
From the below marginal plots, we see no lack of fit of our model

```

marginalModelPlots(logi_AIC,~KIDSDRV_0.5 + AGE + BLUEBOOK_LOG + INCOME_LOG +
HOME_VAL_LOG + TRAVTIME_0.33 + TIF + OLDCLAIM_LOG + CLM_FREQ +
MVR PTS_3 + CAR AGE, layout =c(4,3))

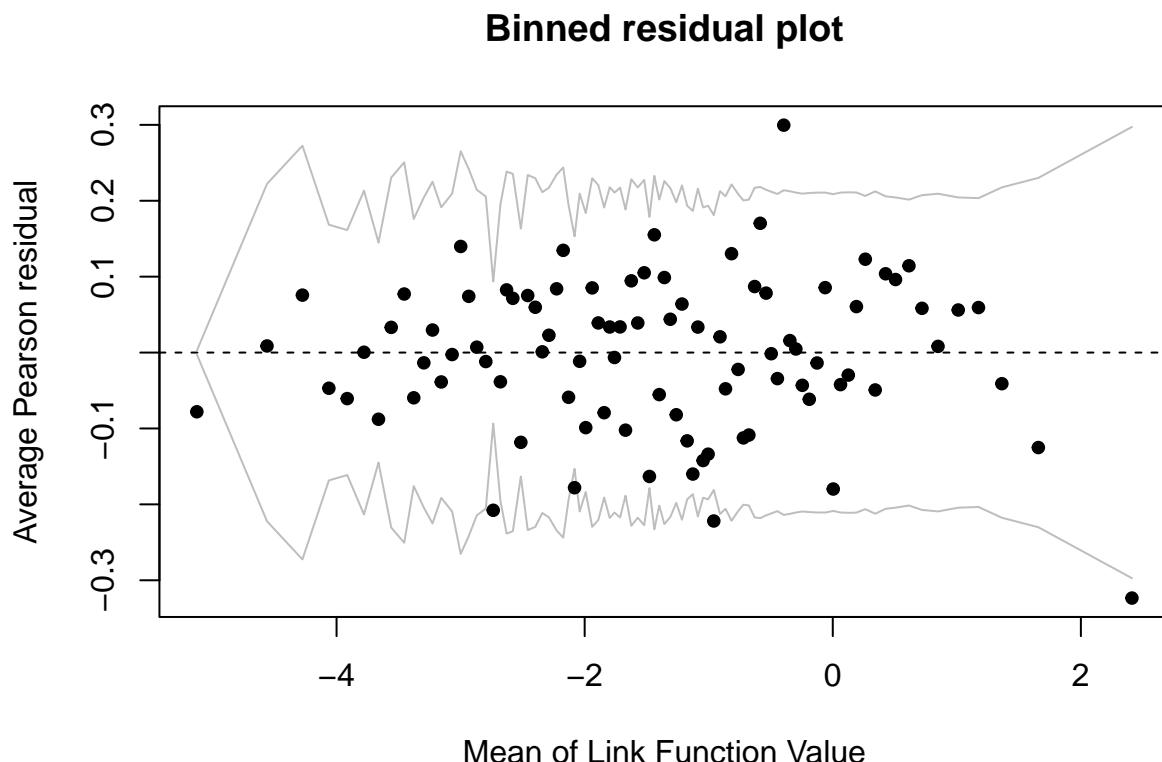
```

Marginal Model Plots



The residual plot below also shows that the pearson residuals are independent with approximately constant variance, with only a few outliers.

```
#arm::binnedplot(x = fitted(logi_AIC), y = residuals(logi_AIC, type="pearson"),
arm::binnedplot(x = predict(logi_AIC, type="link"), y = residuals(logi_AIC, type="pearson"),
  nclass = NULL,
  xlab = "Mean of Link Function Value",
  ylab = "Average Pearson residual",
  main = "Binned residual plot",
  cex.pts = 0.8,
  col.pts = 1,
  col.int = "gray")
```



We conclude that our optimal logistic model logi_AIC is valid

Linear Model

First, let's build a test model to check if any additional transformations are needed to build a valid model

```
lm_train_df <- train_df[train_df$TARGET_FLAG==1,]
lm_train_df$TARGET_FLAG <- NULL
#lm_train_df$TARGET_AMT <- log(lm_train_df$TARGET_AMT)
# lm_train_df$HOME_VAL <- log(lm_train_df$HOME_VAL+1)
# lm_train_df$BLUEBOOK <- log(lm_train_df$BLUEBOOK)
# lm_train_df$INCOME <- log(lm_train_df$INCOME+1)
```

```
lm_full <- lm(TARGET_AMT ~ ., lm_train_df)
summary(lm_full)
```

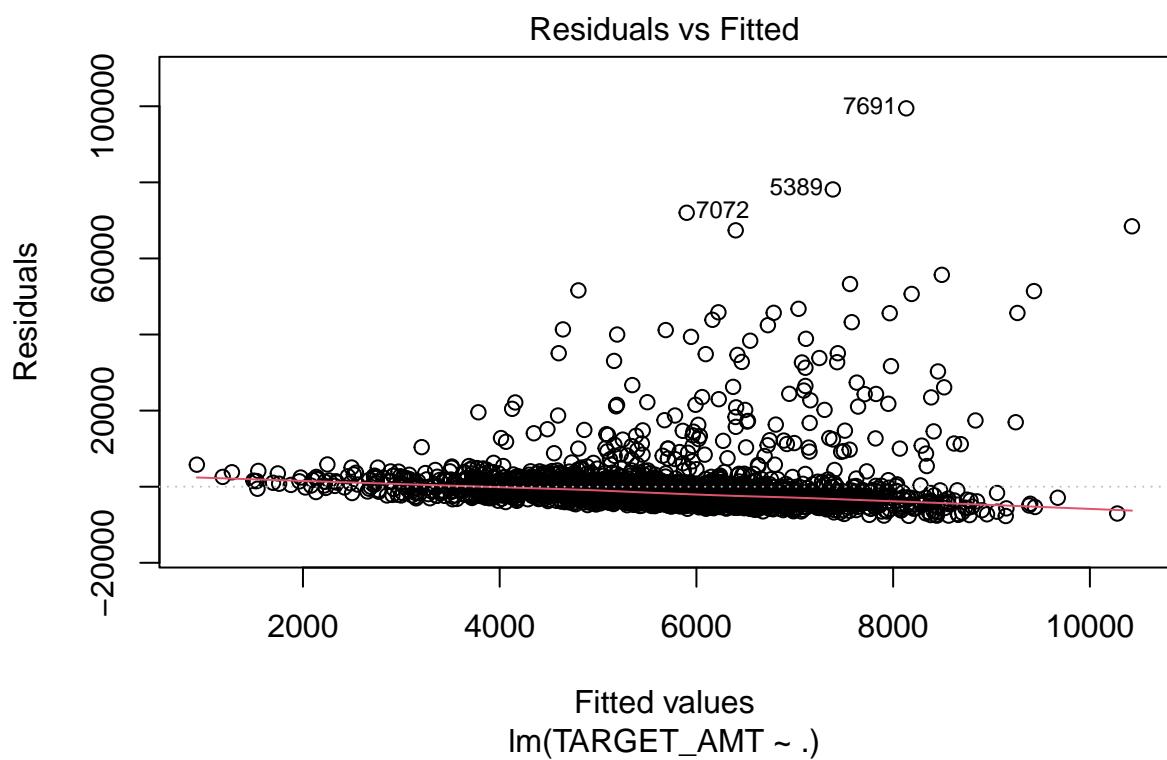
```
##
## Call:
## lm(formula = TARGET_AMT ~ ., data = lm_train_df)
##
## Residuals:
##    Min      1Q  Median      3Q     Max 
## -7704   -3177   -1528     485  99453 
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)    
## (Intercept)             -8518.3048  3664.2401 -2.325   0.0202 *  
## KIDSDRV                 -211.1151   318.4659 -0.663   0.5075    
## AGE                      18.7472    21.8723  0.857   0.3915    
## HOMEKIDS                236.3588   214.4963  1.102   0.2706    
## YOJ                      7.3643    72.5665  0.101   0.9192    
## PARENT1Yes               248.9335   590.7024  0.421   0.6735    
## MSTATUSYes              -905.1346   520.5397 -1.739   0.0822 .  
## SEXM                     1145.0450   634.2852  1.805   0.0712 .  
## EDUCATIONHigh School    -637.4099   516.2843 -1.235   0.2171    
## EDUCATIONLess ThanHigh School -210.6289   655.9132 -0.321   0.7481    
## EDUCATIONMasters          852.1937   907.3188  0.939   0.3477    
## EDUCATIONPhD              1995.6300  1122.6481  1.778   0.0756 .  
## JOBBlue Collar            583.0483  1141.9899  0.511   0.6097    
## JOBClerical                432.4797  1194.3151  0.362   0.7173    
## JOBDoctor                -2281.4659  1769.2616 -1.290   0.1974    
## JOBHome Maker              86.3057  1354.4037  0.064   0.9492    
## JOBLawyer                  299.9609  1029.5824  0.291   0.7708    
## JOBManager                 -732.8605  1067.3386 -0.687   0.4924    
## JOBProfessional            1072.9424  1127.1422  0.952   0.3412    
## JOBStudent                 252.0187  1406.7726  0.179   0.8578    
## TRAVTIME                   -0.5861   11.0883 -0.053   0.9578    
## CAR_USEPrivate              -384.8151  523.7535 -0.735   0.4626    
## TIF                         -15.1556  42.6168 -0.356   0.7222    
## CAR_TYPEPanel Truck          0.5552  883.0394  0.001   0.9995    
## CAR_TYPEPickup              -121.8783  597.2913 -0.204   0.8383    
## CAR_TYPESports Car           943.0478  736.8979  1.280   0.2008    
## CAR_TYPESUV                  641.1761  644.9835  0.994   0.3203    
## CAR_TYPEVan                  119.1634  763.1033  0.156   0.8759    
## RED_CARyes                  -195.0095  498.1558 -0.391   0.6955    
## CLM_FREQ                     -59.8093  238.5123 -0.251   0.8020    
## REVOKEDYes                  -885.8752  493.2057 -1.796   0.0726 .  
## MVR PTS                      125.1910  70.2469  1.782   0.0749 .  
## CAR_AGE                      -87.6572  44.1817 -1.984   0.0474 *  
## URBANICITYUrban                39.6688  758.5402  0.052   0.9583    
## YOJ_Y1                        705.7939  1658.0675  0.426   0.6704    
## INCOME_Y1                      2356.5642  3156.8260  0.746   0.4554    
## HOME_VAL_Y1                  -7506.1470  7876.3331 -0.953   0.3407    
## OLDCLAIM_Y1                  -789.6424  2442.0389 -0.323   0.7465    
## INCOME_LOG                      -327.2885  312.1346 -1.049   0.2945    
## HOME_VAL_LOG                  669.4137  652.5504  1.026   0.3051
```

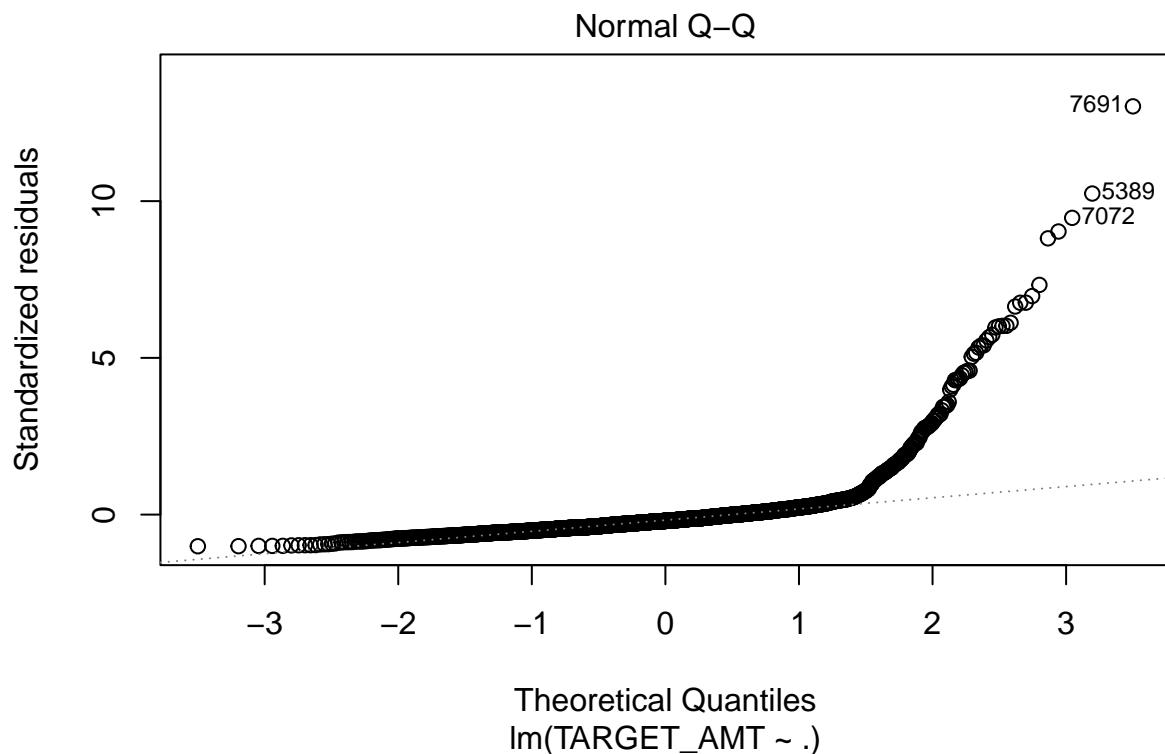
```

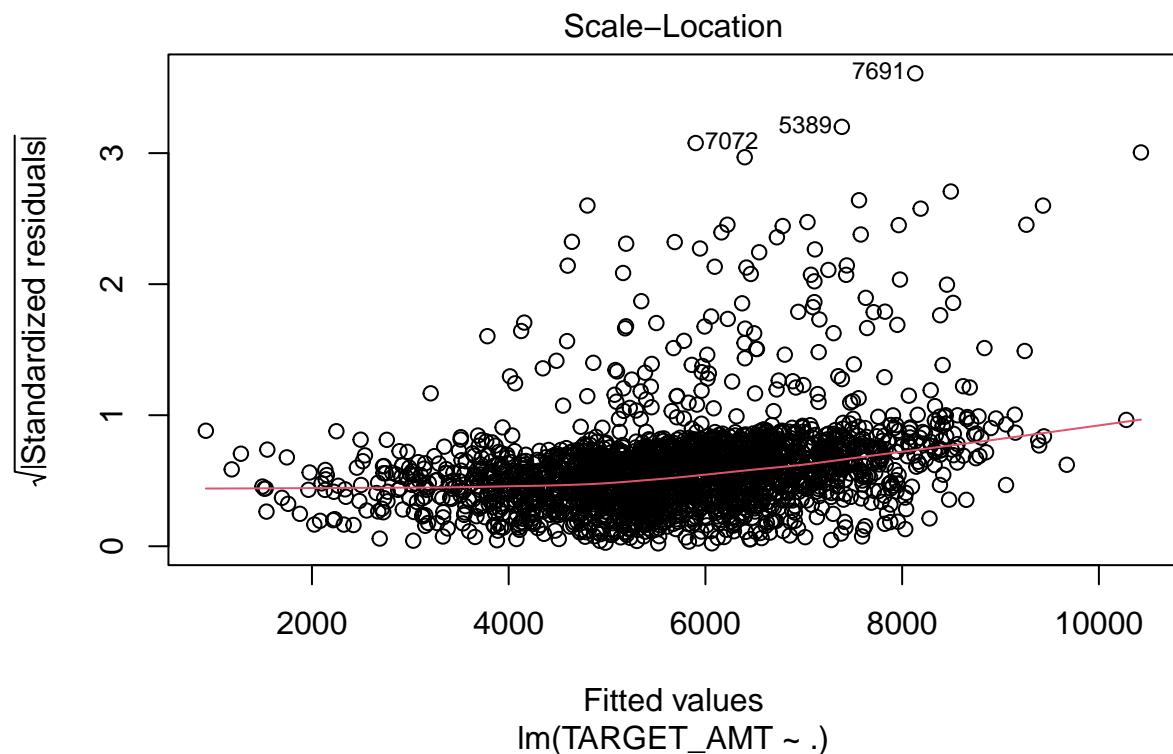
## BLUEBOOK_LOG           1418.5878   330.7068   4.290  1.87e-05 ***
## OLDCLAIM_LOG            94.1095    268.6559   0.350    0.7262
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7697 on 2111 degrees of freedom
## Multiple R-squared:  0.03081, Adjusted R-squared:  0.01199
## F-statistic: 1.637 on 41 and 2111 DF, p-value: 0.006836

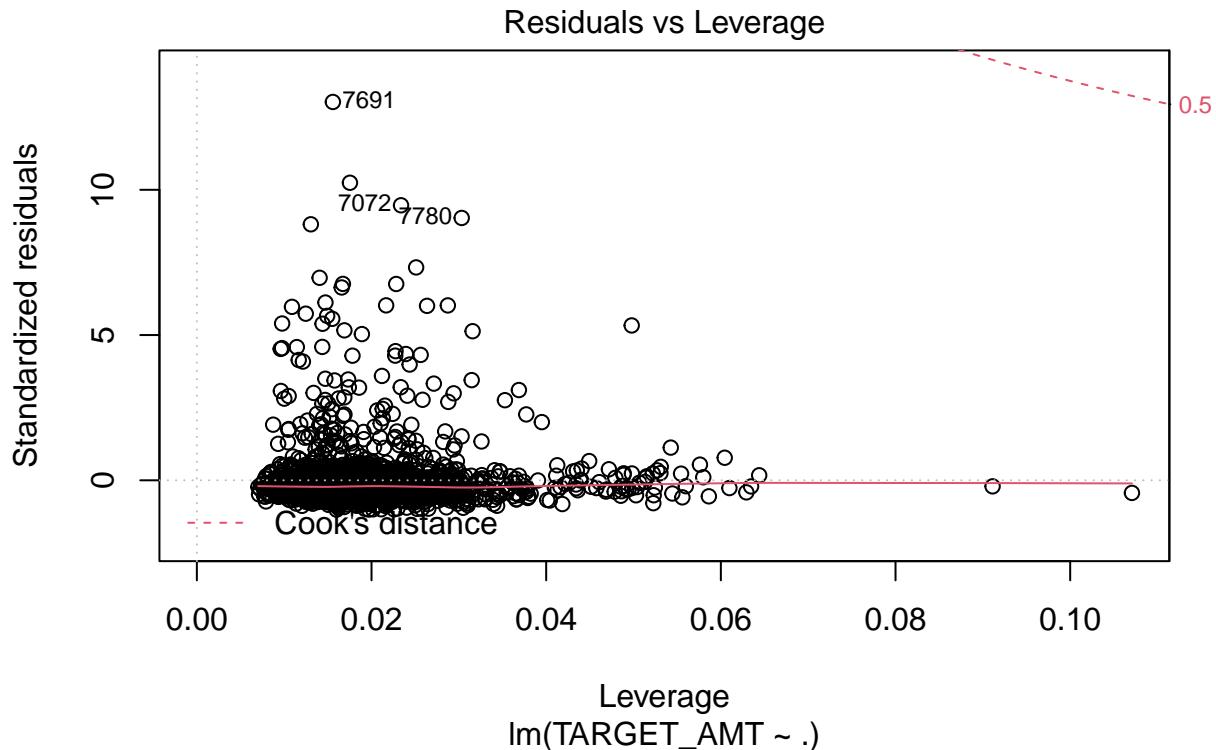
plot(lm_full)

```







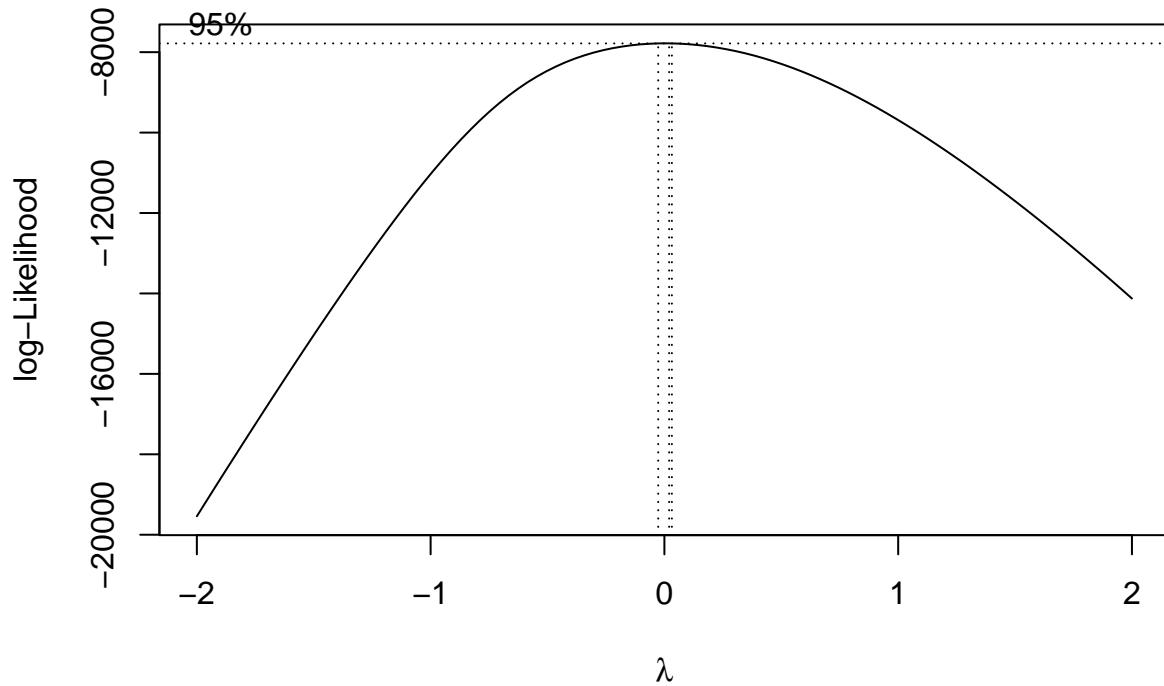


Additioanl Transformation

The plots show that there is a non-linear relationship between the response variable and the predictors.

Let's see what transformation Box-Cox would suggest for the response variable.

```
bc <- boxcox(lm_full)
```



```
lambda <- bc$x[which.max(bc$y)]
lambda
```

```
## [1] 0.02020202
```

It result indicates a log-transformation is appropriate.

```
lm_train_df$TARGET_AMT_LOG <- log(lm_train_df$TARGET_AMT)
lm_train_df$TARGET_AMT <- NULL
```

Buidling Models

Full Model

```
lm_full <- lm(TARGET_AMT_LOG ~ ., lm_train_df)
summary(lm_full)
```

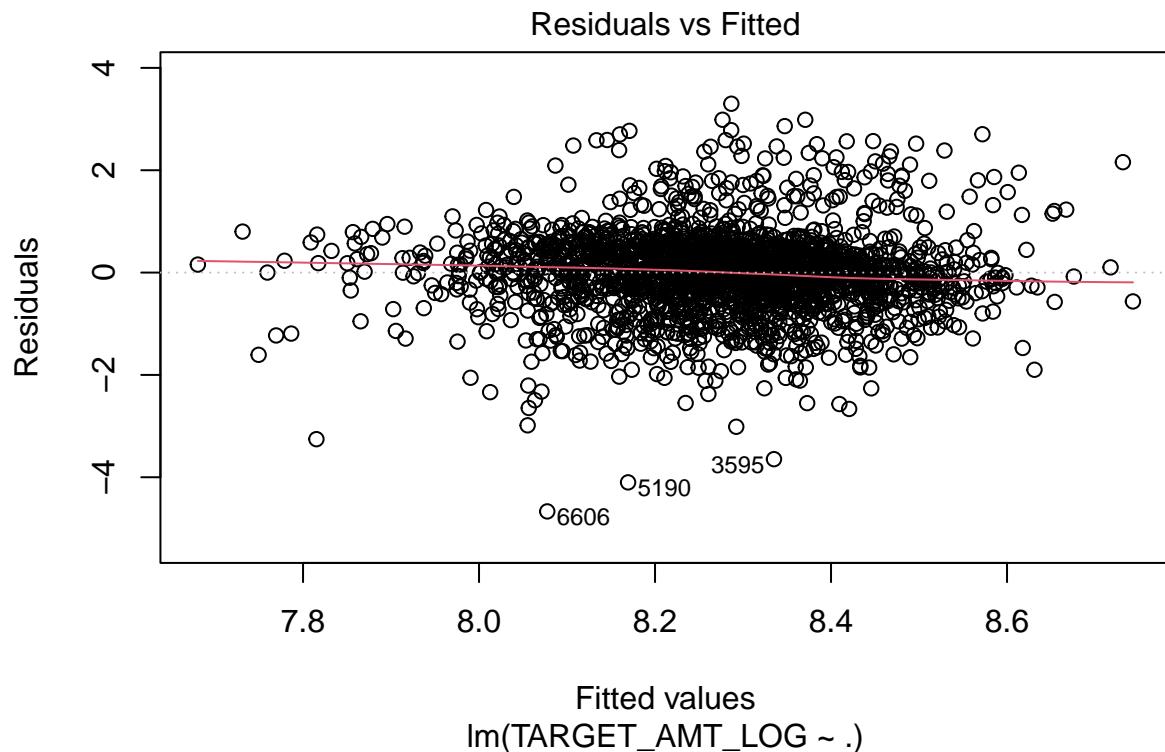
```
##
## Call:
## lm(formula = TARGET_AMT_LOG ~ ., data = lm_train_df)
##
## Residuals:
```

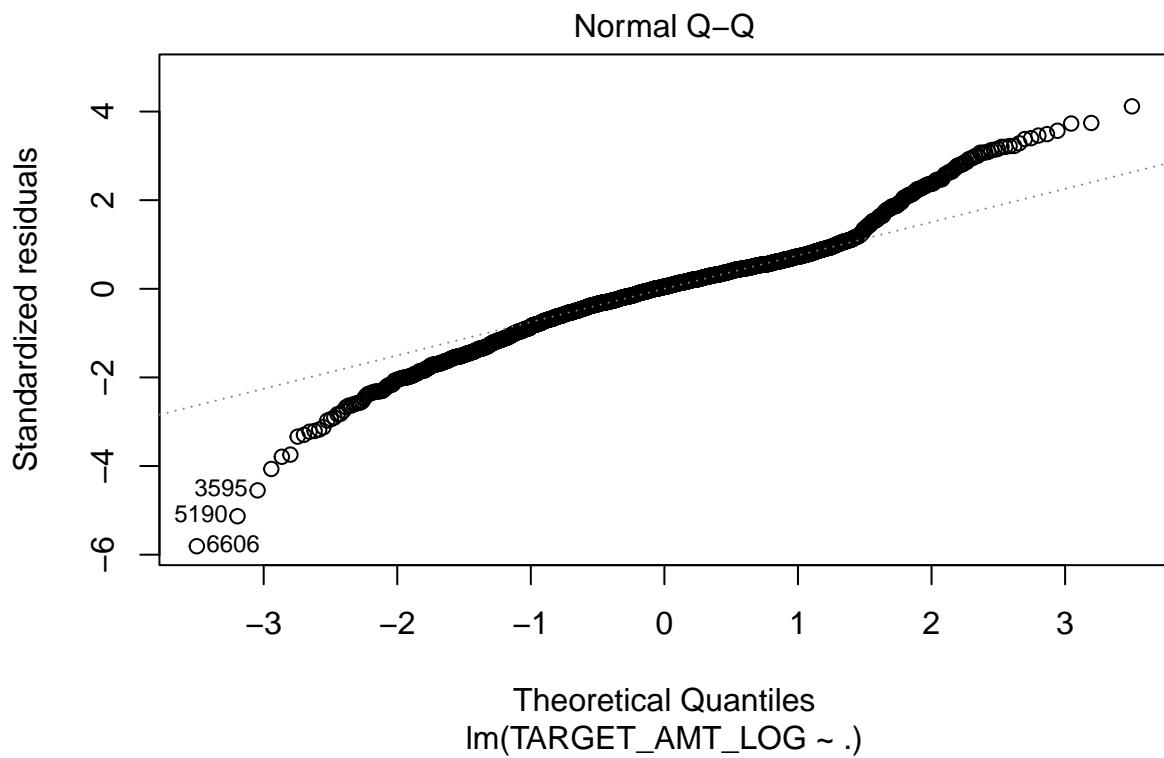
```

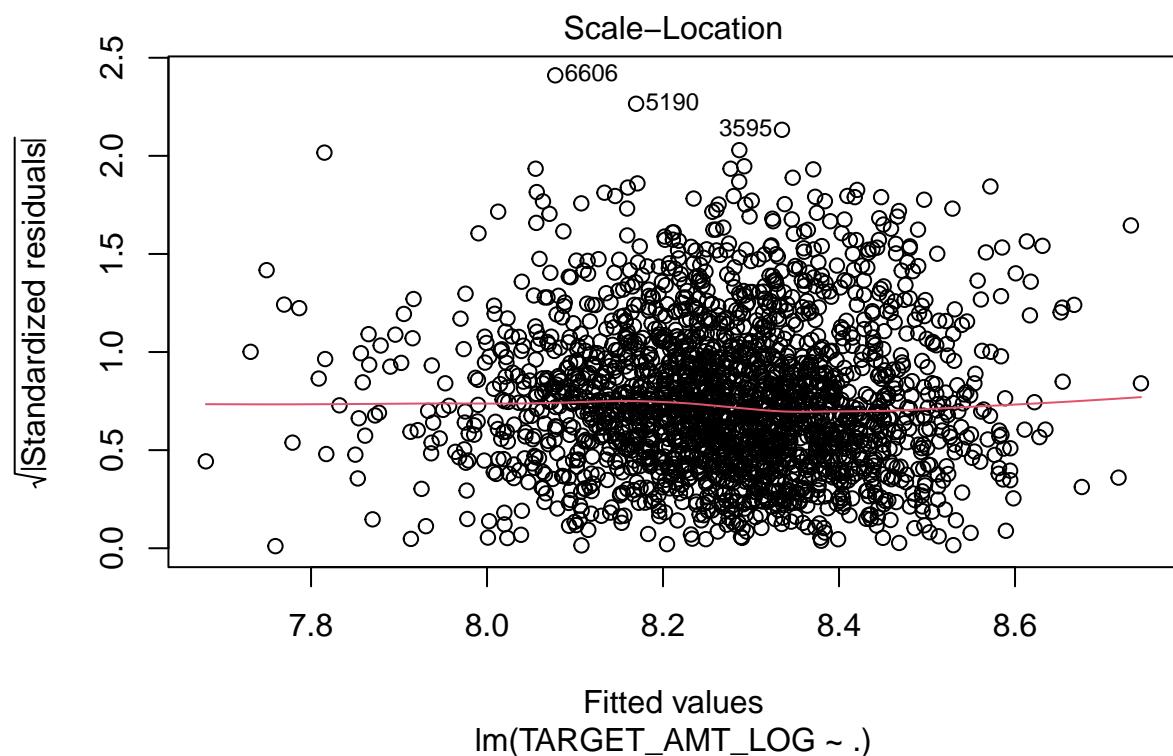
##      Min     1Q   Median     3Q    Max
## -4.6671 -0.4053  0.0383  0.4072  3.2993
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)                6.3403818  0.3844869 16.490 < 2e-16 ***
## KIDSDRV                  -0.0374447  0.0334165 -1.121  0.2626
## AGE                      0.0019603  0.0022951  0.854  0.3931
## HOMEKIDS                 0.0265044  0.0225070  1.178  0.2391
## YOJ                      -0.0052966  0.0076144 -0.696  0.4868
## PARENT1Yes                0.0286128  0.0619821  0.462  0.6444
## MSTATUSYes                -0.0955843  0.0546200 -1.750  0.0803 .
## SEXM                     0.0933562  0.0665552  1.403  0.1609
## EDUCATIONHigh School      0.0466162  0.0541735  0.860  0.3896
## EDUCATIONLess ThanHigh School 0.0433291  0.0688247  0.630  0.5291
## EDUCATIONMasters          0.1704169  0.0952045  1.790  0.0736 .
## EDUCATIONPhD              0.2526967  0.1177989  2.145  0.0321 *
## JOBBlue Collar            0.0720705  0.1198285  0.601  0.5476
## JOBClerical                0.0674620  0.1253189  0.538  0.5904
## JOBDoctor                 -0.0465039  0.1856478 -0.250  0.8022
## JOBHome Maker              -0.0198850  0.1421169 -0.140  0.8887
## JOBLawyer                  -0.0207339  0.1080336 -0.192  0.8478
## JOBManager                 0.0217136  0.1119953  0.194  0.8463
## JOBProfessional            0.1086867  0.1182705  0.919  0.3582
## JOBStudent                 0.0673101  0.1476120  0.456  0.6484
## TRAVTIME                   -0.0004643  0.0011635 -0.399  0.6899
## CAR_USEPrivate              -0.0022455  0.0549572 -0.041  0.9674
## TIF                         -0.0018590  0.0044718 -0.416  0.6777
## CAR_TYPEPanel Truck         0.0283144  0.0926569  0.306  0.7600
## CAR_TYPEPickup              0.0286601  0.0626735  0.457  0.6475
## CAR_TYPESports Car          0.0752778  0.0773223  0.974  0.3304
## CAR_TYPESUV                 0.0914787  0.0676778  1.352  0.1766
## CAR_TYPEVan                 -0.0275008  0.0800721 -0.343  0.7313
## RED_CARYes                  0.0256810  0.0522712  0.491  0.6233
## CLM_FREQ                     -0.0481653  0.0250270 -1.925  0.0544 .
## REVOKEDYes                  -0.0606707  0.0517518 -1.172  0.2412
## MVR PTS                     0.0148261  0.0073710  2.011  0.0444 *
## CAR_AGE                      -0.0008229  0.0046360 -0.177  0.8591
## URBANICITYUrban              0.0457258  0.0795933  0.574  0.5657
## YOJ_Y1                       0.0448365  0.1739802  0.258  0.7967
## INCOME_Y1                    0.3442065  0.3312442  1.039  0.2989
## HOME_VAL_Y1                  0.2408393  0.8264598  0.291  0.7708
## OLDCLAIM_Y1                  -0.1380212  0.2562420 -0.539  0.5902
## INCOME_LOG                   -0.0343748  0.0327521 -1.050  0.2940
## HOME_VAL_LOG                 -0.0167809  0.0684718 -0.245  0.8064
## BLUEBOOK_LOG                  0.1762226  0.0347009  5.078 4.14e-07 ***
## OLDCLAIM_LOG                  0.0249890  0.0281899  0.886  0.3755
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8076 on 2111 degrees of freedom
## Multiple R-squared:  0.03099,    Adjusted R-squared:  0.01217
## F-statistic: 1.646 on 41 and 2111 DF,  p-value: 0.006271

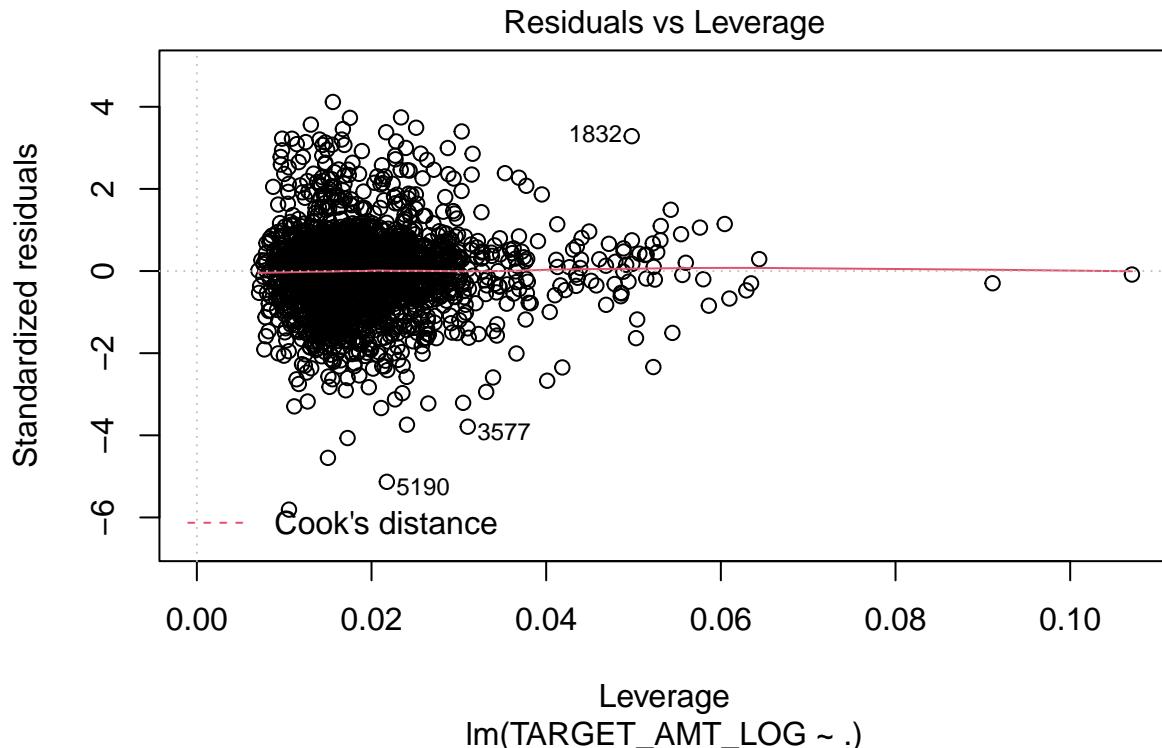
```

```
plot(lm_full)
```









The residual plots show that the relationship is now linear. The only problem is that the distribution of the residuals is not normal. Since the optimal transformation suggested by Box-Cox would not fix this problem, a GLM regression would be more appropriate to fit the data in this case. As requested by this assignment, we would keep the linear models. Since the normality of the residuals is violated, we would not judge the significance of the coefficient by the t-values. We will compare the performance of different models by the adjusted R-squared and the Root of Mean Square Errors.

Backward Elimination By AIC

```

lm_AIC <- step(lm_full, trace = 0)
summary(lm_AIC)

##
## Call:
## lm(formula = TARGET_AMT_LOG ~ MSTATUS + SEX + CLM_FREQ + MVR PTS +
##     BLUEBOOK_LOG, data = lm_train_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -4.6994 -0.4043  0.0425  0.4106  3.2129 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 6.807358  0.248943 27.345 < 2e-16 ***
## MSTATUSYes -0.073271  0.034682 -2.113  0.0347 *  
## 
```

```

## SEXM          0.054122   0.034990   1.547   0.1221
## CLM_FREQ     -0.022761   0.014548  -1.565   0.1178
## MVR_PTS      0.017331   0.007044   2.460   0.0140 *
## BLUEBOOK_LOG 0.156247   0.026314   5.938  3.36e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.804 on 2147 degrees of freedom
## Multiple R-squared:  0.02314,    Adjusted R-squared:  0.02087
## F-statistic: 10.17 on 5 and 2147 DF,  p-value: 1.201e-09

```

Backward Elimination By BIC

```

lm_BIC <- step(lm_full, trace = 0, k = log(nrow(lm_train_df)))
summary(lm_BIC)

```

```

##
## Call:
## lm(formula = TARGET_AMT_LOG ~ BLUEBOOK_LOG, data = lm_train_df)
##
## Residuals:
##       Min     1Q Median     3Q    Max
## -4.7888 -0.3907  0.0425  0.3914  3.2417
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 6.77803   0.24679 27.465 < 2e-16 ***
## BLUEBOOK_LOG 0.15976   0.02626  6.085 1.38e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8058 on 2151 degrees of freedom
## Multiple R-squared:  0.01692,    Adjusted R-squared:  0.01646
## F-statistic: 37.02 on 1 and 2151 DF,  p-value: 1.377e-09

```

Model with only characteristics of the cars

```

lm_car <- lm(TARGET_AMT_LOG~CAR_USE+BLUEBOOK_LOG+CAR_TYPE+RED_CAR+CAR_AGE, data = lm_train_df)
summary(lm_car)

##
## Call:
## lm(formula = TARGET_AMT_LOG ~ CAR_USE + BLUEBOOK_LOG + CAR_TYPE +
##     RED_CAR + CAR_AGE, data = lm_train_df)
##
## Residuals:
##       Min     1Q Median     3Q    Max
## -4.7809 -0.3998  0.0401  0.3949  3.2345
## 

```

```

## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)           6.7720339  0.3032735 22.330 < 2e-16 ***
## CAR_USEPrivate       0.0062649  0.0414251  0.151   0.880
## BLUEBOOK_LOG        0.1552993  0.0315060  4.929 8.89e-07 ***
## CAR_TYPEPanel Truck  0.0703853  0.0847465  0.831   0.406
## CAR_TYPEPickup      0.0359615  0.0607888  0.592   0.554
## CAR_TYPESports Car  0.0211452  0.0667915  0.317   0.752
## CAR_TYPESUV         0.0347350  0.0570779  0.609   0.543
## CAR_TYPEVan          0.0012866  0.0759776  0.017   0.986
## RED_CARyes          0.0549875  0.0463059  1.187   0.235
## CAR_AGE              0.0002214  0.0032346  0.068   0.945
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8068 on 2143 degrees of freedom
## Multiple R-squared:  0.0183, Adjusted R-squared:  0.01417
## F-statistic: 4.438 on 9 and 2143 DF,  p-value: 8.833e-06

```

Model Selection

```

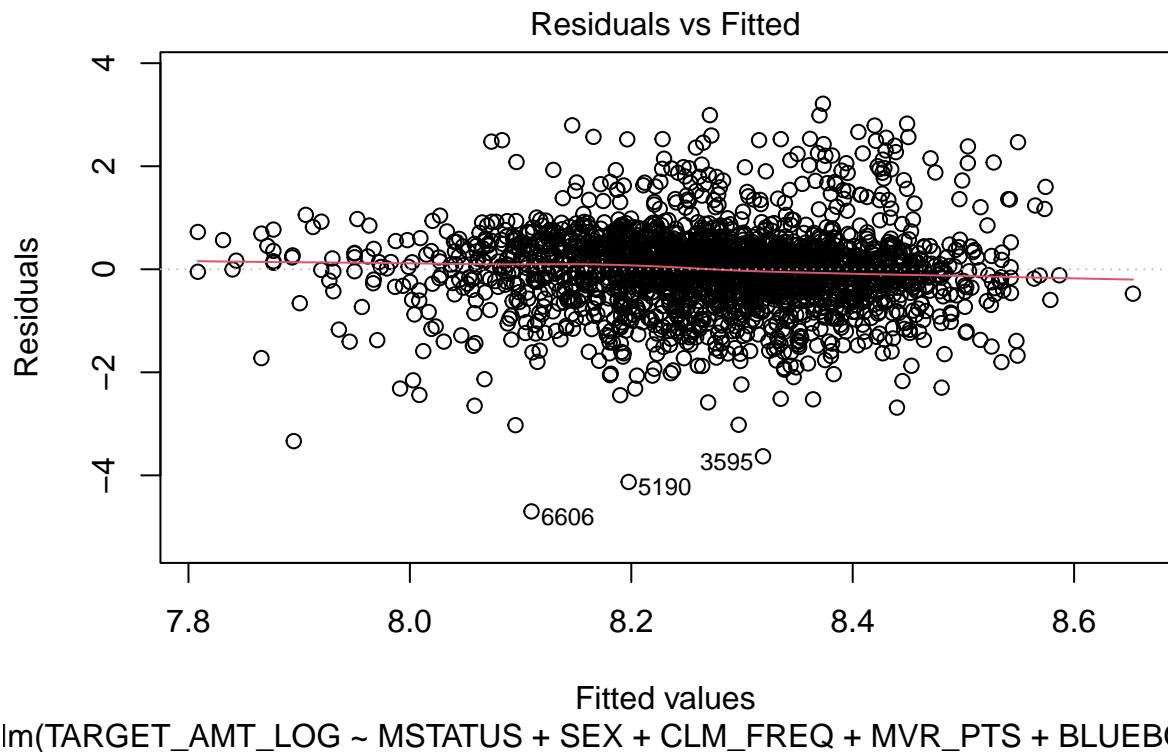
lm_models <- data.frame(model=c(""),
                         Num_of_Coefficients=c(0),
                         R_squared_adj=c(0.0000), RMSE=c(0.0000))
models <- list(lm_full, lm_AIC, lm_BIC, lm_car)
model_names <- c("lm_full", "lm_AIC", "lm_BIC", "lm_car")
for (i in c(1:length(models))) {
  lm_models[i,"model"] <- model_names[i]
  lm_models[i,"Num_of_Coefficients"] <- length(models[[i]]$coefficients) - 1
  lm_models[i,"R_squared_adj"] <- summary(models[[i]])$adj.r.squared
  lm_models[i,"RMSE"] <- sqrt(mean(models[[i]]$residuals^2))
}
lm_models

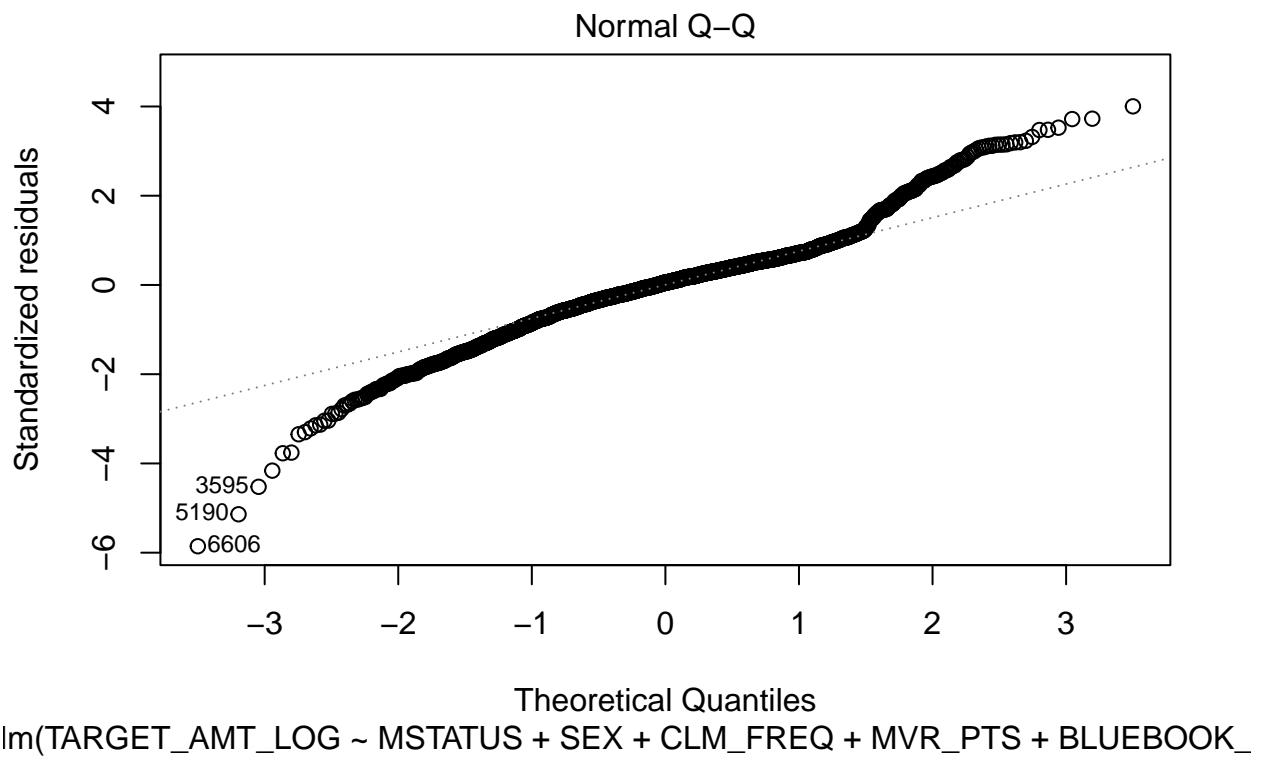
##      model Num_of_Coefficients R_squared_adj      RMSE
## 1 lm_full                  41    0.01216516 0.7996880
## 2 lm_AIC                   5    0.02086660 0.8029181
## 3 lm_BIC                   1    0.01646443 0.8054703
## 4 lm_car                   9    0.01417443 0.8049065

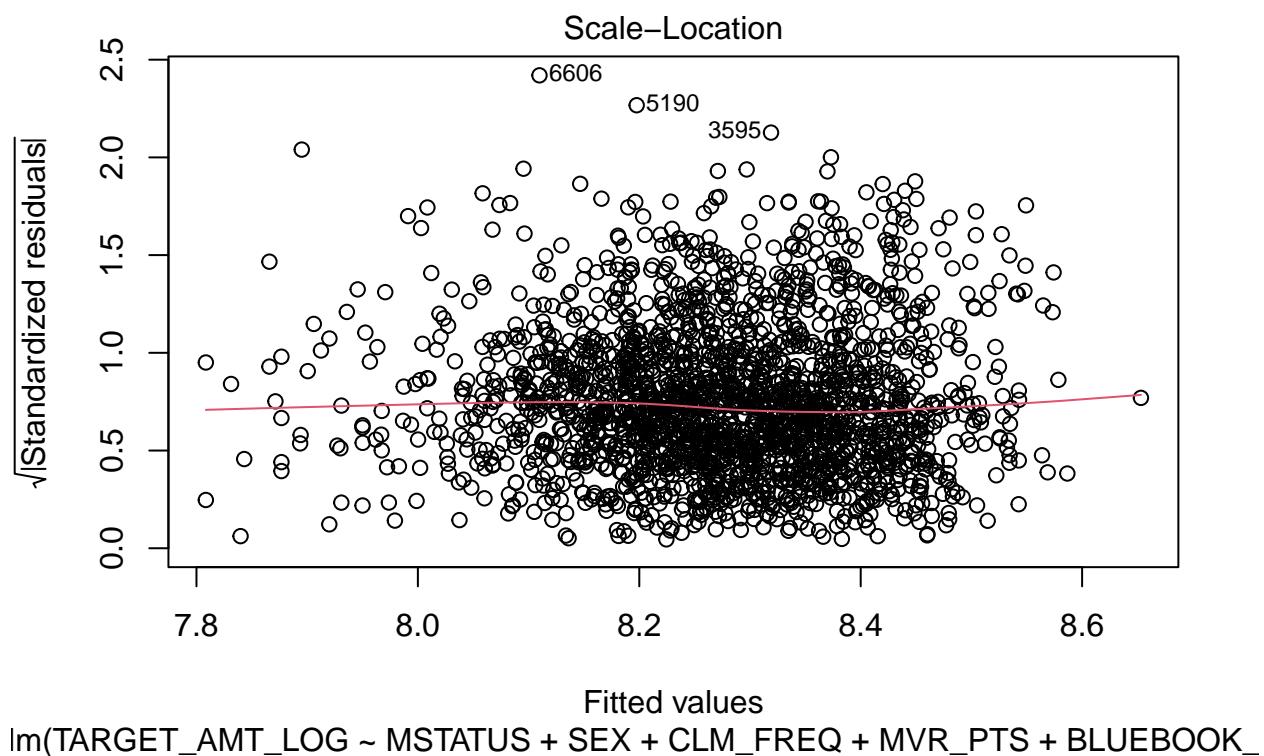
```

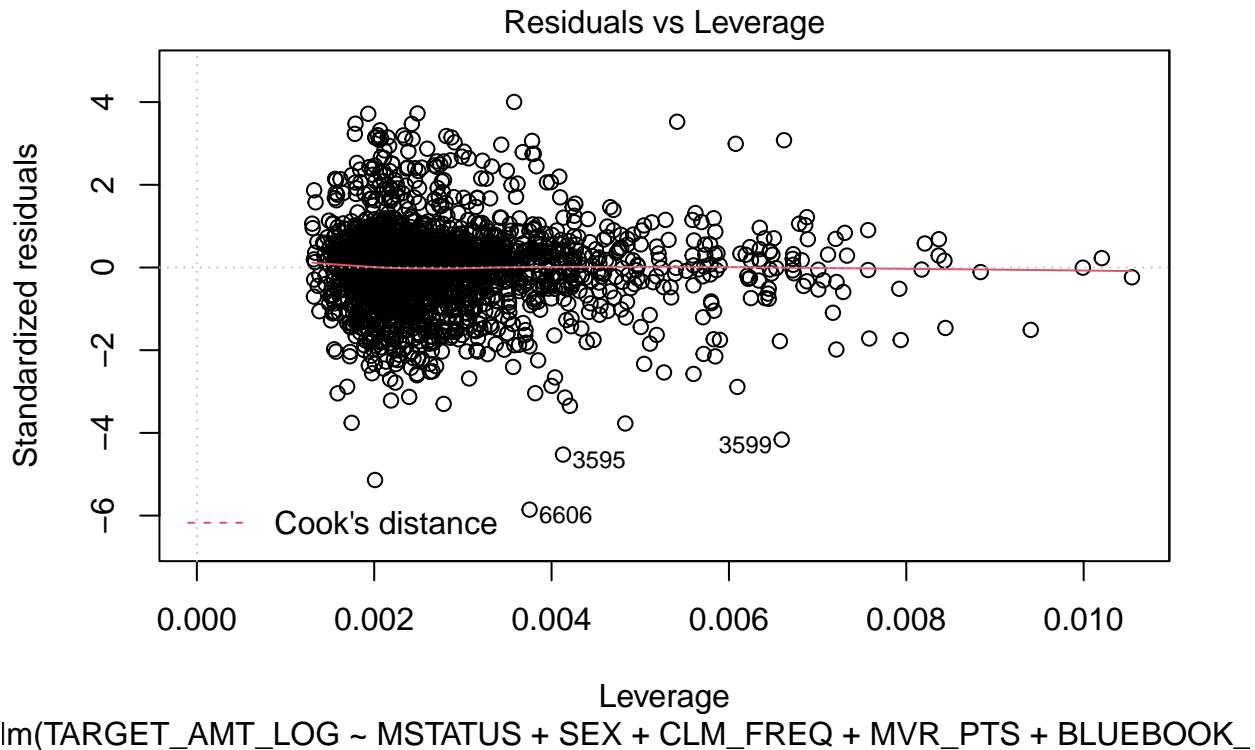
Model lm_AIC has the highest adjusted R-squared and the RMSE is very close to the full model. Our optimal linear model is lm_AIC

```
plot(lm_AIC)
```









Evaluation Data Prediction

```

test_df$INDEX <- NULL
test_df$INCOME <- as.numeric(gsub('[,]', '', test_df$INCOME))
test_df$HOME_VAL <- as.numeric(gsub('[,]', '', test_df$HOME_VAL))
test_df$BLUEBOOK <- as.numeric(gsub('[,]', '', test_df$BLUEBOOK))
test_df$OLDCLAIM <- as.numeric(gsub('[,]', '', test_df$OLDCLAIM))
test_df$PARENT1 <- gsub("z_", "", test_df$PARENT1)
test_df$MSTATUS <- gsub("z_", "", test_df$MSTATUS)
test_df$SEX <- gsub("z_", "", test_df$SEX)
test_df$EDUCATION <- gsub("z_", "", test_df$EDUCATION)
test_df$EDUCATION <- gsub("<", "Less Than", test_df$EDUCATION)
test_df$JOB <- gsub("z_", "", test_df$JOB)
test_df$CAR_TYPE <- gsub("z_", "", test_df$CAR_TYPE)
test_df$URBANICITY <- ifelse(test_df$URBANICITY == "Highly Urban/ Urban", "Urban", "Rural")

test_df[c("TARGET_FLAG", "PARENT1", "MSTATUS", "SEX", "EDUCATION", "JOB", "CAR_TYPE",
        "RED_CAR", "URBANICITY", "CAR_USE", "REVOKED")] <-
  lapply(test_df[c("TARGET_FLAG", "PARENT1", "MSTATUS", "SEX",
                 "EDUCATION", "JOB", "CAR_TYPE", "RED_CAR",
                 "URBANICITY", "CAR_USE", "REVOKED")], factor)

```

```

test_df$TARGET_FLAG <- NULL
test_df$TARGET_AMT <- NULL

test_df <- mice.reuse(mickey, test_df, maxit = 5, printFlag = FALSE, seed = 2022)[[1]]

summary(test_df)

##      KIDSDRV          AGE          HOMEKIDS          YOJ
##  Min.   :0.0000  Min.   :17.00  Min.   :0.0000  Min.   : 0.00
##  1st Qu.:0.0000  1st Qu.:39.00  1st Qu.:0.0000  1st Qu.: 9.00
##  Median :0.0000  Median :45.00  Median :0.0000  Median :11.00
##  Mean   :0.1625  Mean   :45.01  Mean   :0.7174  Mean   :10.36
##  3rd Qu.:0.0000  3rd Qu.:51.00  3rd Qu.:1.0000  3rd Qu.:13.00
##  Max.   :3.0000  Max.   :73.00  Max.   :5.0000  Max.   :19.00
##
##      INCOME          PARENT1        HOME_VAL        MSTATUS        SEX
##  Min.   :    0  No :1875  Min.   :     0  No : 847  F:1170
##  1st Qu.:25749 Yes: 266  1st Qu.:     0  Yes:1294  M: 971
##  Median :51626                   Median :159254
##  Mean   :60494                   Mean   :153751
##  3rd Qu.:86473                   3rd Qu.:237057
##  Max.   :291182                  Max.   :669271
##
##      EDUCATION          JOB          TRAVTIME
##  Bachelors       :581  Blue Collar :463  Min.   : 5.00
##  High School     :622  Clerical   :319  1st Qu.:22.00
##  Less ThanHigh School:312 Professional:291  Median :33.00
##  Masters         :420  Manager    :269  Mean   :33.15
##  PhD             :206  Home Maker :202  3rd Qu.:43.00
##                      Lawyer    :196  Max.   :105.00
##                      (Other)   :401
##
##      CAR_USE          BLUEBOOK          TIF          CAR_TYPE
##  Commercial: 760  Min.   :1500  Min.   :1.000  Minivan   :549
##  Private   :1381  1st Qu.:8870  1st Qu.:1.000  Panel Truck:177
##                      Median :14170  Median :4.000  Pickup    :383
##                      Mean   :15469  Mean   :5.245  Sports Car :272
##                      3rd Qu.:21050  3rd Qu.:7.000  SUV       :589
##                      Max.   :49940  Max.   :25.000 Van      :171
##
##      RED_CAR          OLDCLAIM          CLM_FREQ          REVOKED          MVR PTS
##  no :1543  Min.   :    0  Min.   :0.000  No :1880  Min.   : 0.000
##  yes: 598  1st Qu.:    0  1st Qu.:0.000  Yes: 261  1st Qu.: 0.000
##                      Median :    0  Median :0.000  Median : 1.000
##                      Mean   :4022  Mean   :0.809  Mean   : 1.766
##                      3rd Qu.:4718  3rd Qu.:2.000  3rd Qu.: 3.000
##                      Max.   :54399 Max.   :5.000  Max.   :12.000
##
##      CAR_AGE          URBANICITY
##  Min.   : 0.000  Rural: 403
##  1st Qu.: 1.000  Urban:1738
##  Median : 8.000
##  Mean   : 8.185
##  3rd Qu.:13.000

```

```

##  Max.    :26.000
## 

test_df$Y0J_Y <- as.factor(ifelse(test_df$Y0J == 0,0,1))
test_df$INCOME_Y <- as.factor(ifelse(test_df$INCOME == 0,0,1))
test_df$HOME_VAL_Y <- as.factor(ifelse(test_df$HOME_VAL == 0,0,1))
test_df$OLDCLAIM_Y <- as.factor(ifelse(test_df$OLDCLAIM == 0,0,1))

test_df$INCOME_LOG <- log(test_df$INCOME+1)
test_df$HOME_VAL_LOG <- log(test_df$HOME_VAL+1)
test_df$BLUEBOOK_LOG <- log(test_df$BLUEBOOK)
test_df$OLDCLAIM_LOG <- log(test_df$OLDCLAIM+1)

test_df$INCOME <- NULL
test_df$HOME_VAL <- NULL
test_df$BLUEBOOK <- NULL
test_df$OLDCLAIM <- NULL

logistic_test_df <- test_df

```

```
logistic_test_df$AGE_Squared <- logistic_test_df$AGE^2
```

```

logistic_test_df$KIDSDRV_0.5 <- (logistic_test_df$KIDSDRV)^0.5
logistic_test_df$HOMEKIDS_0.5 <- (logistic_test_df$HOMEKIDS)^0.5
logistic_test_df$MVR PTS_3 <- (logistic_test_df$MVR PTS)^3
logistic_test_df$TRAVTIME_0.33 <- (logistic_test_df$TRAVTIME)^0.33

logistic_test_df$KIDSDRV <- NULL
logistic_test_df$HOMEKIDS <- NULL
logistic_test_df$MVR PTS <- NULL
logistic_test_df$TRAVTIME <- NULL

```

```
logistic_test_df$TARGET_FLAG <- ifelse(predict(logi_AIC,logistic_test_df, type="response") > logi_models[2,"Optimal_Threshold"],1,0)
```

```

test_predict <- logistic_test_df$TARGET_FLAG
train_predict <- ifelse(logi_AIC$fitted.values>logi_models[2,"Optimal_Threshold"],1,0)

```

```

dist_df <- data.frame(rbind(
  cbind(train_predict,"train_predict"),
  cbind(test_predict,"test_predict")
))
colnames(dist_df) <- c("value","data")
dist_df <- table(dist_df)
dist_df[,1] <- dist_df[,1]/sum(dist_df[,1])
dist_df[,2] <- dist_df[,2]/sum(dist_df[,2])
dist_df

```

```

##      data
## value test_predict train_predict
##     0      0.6081270      0.6198995
##     1      0.3918730      0.3801005

```

The model produces similar result for both the training and testing data. Around 61% of the cases are classified as no crash and 39% of the cases are classified as having a crash. Our logistic model has similar performance for predicting unseen results.

```

lm_test_df <- test_df
lm_test_df$TARGET_FLAG <- NULL

lm_test_df$TARGET_AMT <- predict(lm_AIC, lm_test_df)
lm_test_df$TARGET_AMT <- exp(lm_test_df$TARGET_AMT)

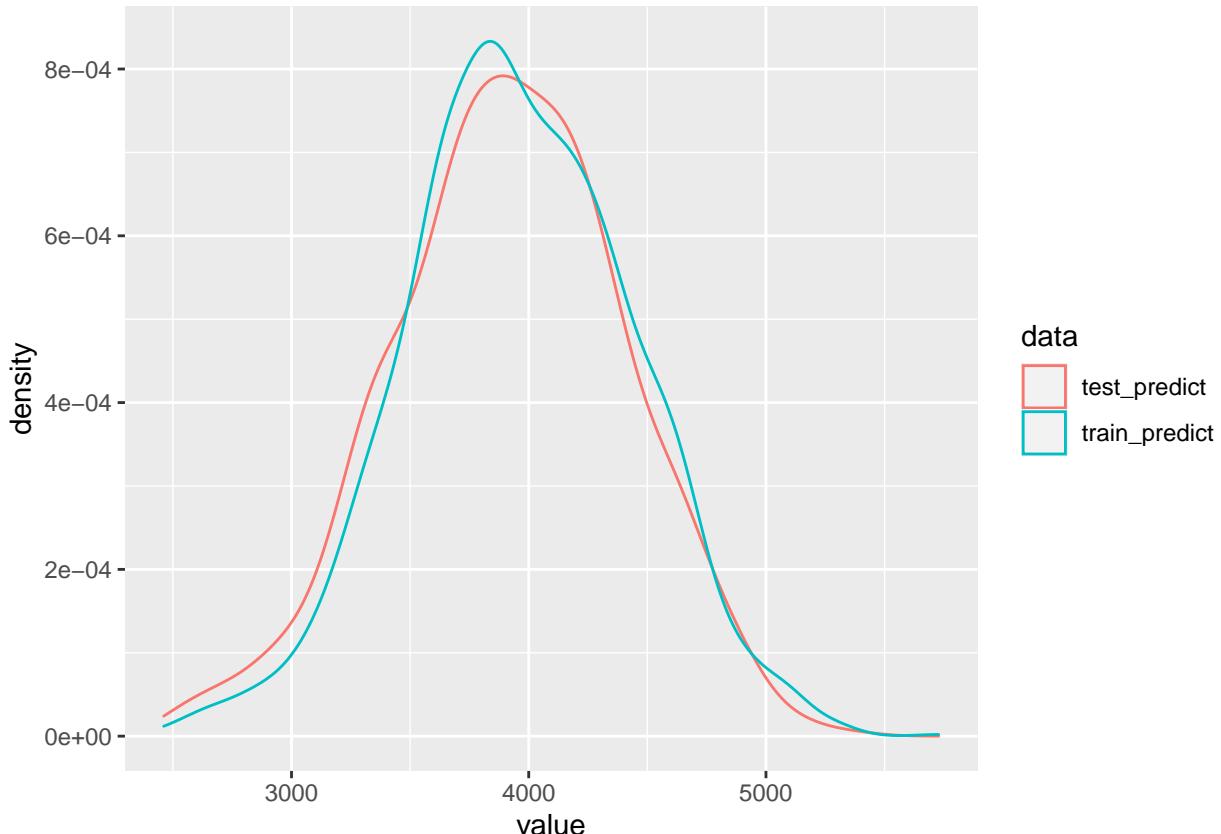
test_df$TARGET_FLAG <- logistic_test_df$TARGET_FLAG
test_df$TARGET_AMT <- lm_test_df$TARGET_AMT * logistic_test_df$TARGET_FLAG

train_predict <- exp(lm_AIC$fitted.values)
test_predict <- test_df$TARGET_AMT[test_df$TARGET_AMT > 0]

dist_df <- data.frame(rbind(
  cbind(train_predict, "train_predict"),
  cbind(test_predict, "test_predict")
), stringsAsFactors=FALSE)
colnames(dist_df) <- c("value", "data")
dist_df$value <- as.numeric(dist_df$value)

ggplot(dist_df, aes(x=value, color=data)) +
  geom_density()

```



The prediction of claim amounts have similar distributions for the training and testing data. Our linear model has stable performance in predicting unseen results.