

Homework 4

Team 1

November 8, 2020

1. Data Exploration

The auto insurance training dataset has 26 variables and 8161 observations. Of the variables, 24 of them are predictors for two responses: TARGET_FLAG and TARGET_AMT is numerical.

To explore the training data: - used the summary function to see means, medians, and quartiles of predictors - used str function to see the data type of each predictor - explored TARGET_FLAG in relation to some other variables such as AGE and CAR_AGE - looked at distribution of some numerical variables such as AGE and MVR_PTS

From the summary function, the TARGET_FLAG is binary and 26% of the 8161 records were accidents.

2. Data Preparation

This data was prepared to build both a binary logistic model and a multiple linear regression model. The binary logistic model was used to predict the TARGET_FLAG response variable and the multiple linear regression model was used to predict the TARGET_AMT variable.

We want to train the multiple linear regression model on records that actually have a valid TARGET_AMT variable, so its training dataset is a subset of the full dataset where TARGET_FLAG is 1.

We made dummy variable columns for all variables that had NA (AGE, YOJ, CAR_AGE) and then filled those columns with their median values.

The training dataset for the binary logistic regression model was labeled train_df. The training dataset for the multiple linear regression model was titled train_amt_df.

3. Build Models

First, we built two models using most predictors as numerics. Then we used the step AIC function to find the best variables for each model.

One model was a Binary Logistic Regression model for the TARGET_FLAG response titled step_BLR. The second model was a Multiple Linear Regression for the TARGET_AMT response titled MLR_all_vars.

4. Select Models

To finally select a model, we used Stepwise AIC (both backward and forward) to do model selection and ended with a Binary Logistic 7661.4.

Appendix

Import Libraries and Data

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

## corrrplot 0.84 loaded

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##   select

## Loading required package: lattice

## Type 'citation("pROC")' for a citation.

##
## Attaching package: 'pROC'

## The following objects are masked from 'package:stats':
##
##   cov, smooth, var

## -- Attaching packages ----- tidyverse 1.3.0 --

## v tibble 3.0.3    v purrr 0.3.4
## v tidyr  1.1.2    v stringr 1.4.0
## v readr  1.4.0    v forcats 0.5.0

## -- Conflicts ----- tidyverse_conflicts() --
## x tidyr::complete() masks Rcurl::complete()
## x dplyr::filter()   masks stats::filter()
## x dplyr::lag()      masks stats::lag()
## x purrr::lift()     masks caret::lift()
## x MASS::select()   masks dplyr::select()

##
## Attaching package: 'kableExtra'
```

```
## The following object is masked from 'package:dplyr':
##
##   group_rows
```

```
# Loading the data
git_dir <- 'https://raw.githubusercontent.com/odonnell31/DATA621-HW4/main/data'
#class_data = read.csv(paste(git_dir, "/classification-output-data.csv", sep=""))
train_df = read.csv(paste(git_dir, "/insurance_training_data.csv", sep=""))
test_df = read.csv(paste(git_dir, "/insurance-evaluation-data.csv", sep = ""))
head(train_df, 2)
```

```
##   INDEX TARGET_FLAG TARGET_AMT KIDSDRIV AGE HOMEKIDS YOJ INCOME PARENT1
## 1     1           0           0         0  60         0  11 $67,349      No
## 2     2           0           0         0  43         0  11 $91,449      No
##   HOME_VAL MSTATUS SEX      EDUCATION      JOB TRAVTIME  CAR_USE BLUEBOOK
## 1         $0   z_No  M          PhD  Professional      14   Private $14,230
## 2 $257,252   z_No  M z_High School z_Blue Collar      22 Commercial $14,940
##   TIF CAR_TYPE RED_CAR OLDCLAIM CLM_FREQ REVOKED MVR_PTS CAR_AGE
## 1  11 Minivan    yes   $4,461         2      No        3      18
## 2   1 Minivan    yes      $0         0      No        0       1
##
##   URBANICITY
## 1 Highly Urban/ Urban
## 2 Highly Urban/ Urban
```

Data Exploration & Preparation

See a summary of each column in the train_df set

```
# view a summary of all columns
summary(train_df)
```

```
##      INDEX      TARGET_FLAG      TARGET_AMT      KIDSDRIV
##  Min.   :    1  Min.   :0.0000  Min.   :    0  Min.   :0.0000
## 1st Qu.: 2559 1st Qu.:0.0000 1st Qu.:    0 1st Qu.:0.0000
## Median : 5133 Median :0.0000 Median :    0 Median :0.0000
## Mean   : 5152 Mean   :0.2638 Mean   : 1504 Mean   :0.1711
## 3rd Qu.: 7745 3rd Qu.:1.0000 3rd Qu.: 1036 3rd Qu.:0.0000
## Max.   :10302 Max.   :1.0000 Max.   :107586 Max.   :4.0000
##
##      AGE      HOMEKIDS      YOJ      INCOME
##  Min.   :16.00  Min.   :0.0000  Min.   : 0.0  Length:8161
## 1st Qu.:39.00 1st Qu.:0.0000 1st Qu.: 9.0  Class :character
## Median :45.00 Median :0.0000 Median :11.0  Mode  :character
## Mean   :44.79 Mean   :0.7212 Mean   :10.5
## 3rd Qu.:51.00 3rd Qu.:1.0000 3rd Qu.:13.0
## Max.   :81.00 Max.   :5.0000 Max.   :23.0
## NA's   :6      NA's   :454
##   PARENT1      HOME_VAL      MSTATUS      SEX
## Length:8161 Length:8161 Length:8161 Length:8161
## Class :character Class :character Class :character Class :character
## Mode  :character Mode  :character Mode  :character Mode  :character
##
```

```
##
##
##
## EDUCATION          JOB          TRAVTIME          CAR_USE
## Length:8161        Length:8161    Min.   : 5.00    Length:8161
## Class :character    Class :character 1st Qu.: 22.00    Class :character
## Mode  :character    Mode  :character Median : 33.00    Mode  :character
##                                     Mean  : 33.49
##                                     3rd Qu.: 44.00
##                                     Max.   :142.00
##
## BLUEBOOK           TIF           CAR_TYPE          RED_CAR
## Length:8161        Min.   : 1.000    Length:8161    Length:8161
## Class :character    1st Qu.: 1.000    Class :character  Class :character
## Mode  :character    Median : 4.000    Mode  :character  Mode  :character
##                                     Mean  : 5.351
##                                     3rd Qu.: 7.000
##                                     Max.   :25.000
##
## OLDCLAIM           CLM_FREQ        REVOKED           MVR_PTS
## Length:8161        Min.   :0.0000    Length:8161    Min.   : 0.000
## Class :character    1st Qu.:0.0000    Class :character 1st Qu.: 0.000
## Mode  :character    Median :0.0000    Mode  :character  Median : 1.000
##                                     Mean  :0.7986
##                                     3rd Qu.:2.0000
##                                     Max.   :5.0000
##                                     Mean  : 1.696
##                                     3rd Qu.: 3.000
##                                     Max.   :13.000
##
## CAR_AGE            URBANICITY
## Min.   : -3.000    Length:8161
## 1st Qu.: 1.000    Class :character
## Median : 8.000    Mode  :character
## Mean    : 8.328
## 3rd Qu.:12.000
## Max.    :28.000
## NA's    :510
```

Look at the data type of each variable

```
# data type of predictors
str(train_df)
```

```
## 'data.frame': 8161 obs. of 26 variables:
## $ INDEX : int 1 2 4 5 6 7 8 11 12 13 ...
## $ TARGET_FLAG: int 0 0 0 0 0 1 0 1 1 0 ...
## $ TARGET_AMT : num 0 0 0 0 0 ...
## $ KIDSDRIV : int 0 0 0 0 0 0 0 1 0 0 ...
## $ AGE : int 60 43 35 51 50 34 54 37 34 50 ...
## $ HOMEKIDS : int 0 0 1 0 0 1 0 2 0 0 ...
## $ YOJ : int 11 11 10 14 NA 12 NA NA 10 7 ...
## $ INCOME : chr "$67,349" "$91,449" "$16,039" "" ...
## $ PARENT1 : chr "No" "No" "No" "No" ...
## $ HOME_VAL : chr "$0" "$257,252" "$124,191" "$306,251" ...
## $ MSTATUS : chr "z_No" "z_No" "Yes" "Yes" ...
```

```
## $ SEX      : chr  "M" "M" "z_F" "M" ...
## $ EDUCATION : chr  "PhD" "z_High School" "z_High School" "<High School" ...
## $ JOB       : chr  "Professional" "z_Blue Collar" "Clerical" "z_Blue Collar" ...
## $ TRAVTIME  : int   14 22 5 32 36 46 33 44 34 48 ...
## $ CAR_USE   : chr  "Private" "Commercial" "Private" "Private" ...
## $ BLUEBOOK  : chr  "$14,230" "$14,940" "$4,010" "$15,440" ...
## $ TIF       : int   11 1 4 7 1 1 1 1 7 ...
## $ CAR_TYPE  : chr  "Minivan" "Minivan" "z_SUV" "Minivan" ...
## $ RED_CAR   : chr  "yes" "yes" "no" "yes" ...
## $ OLDCLAIM  : chr  "$4,461" "$0" "$38,690" "$0" ...
## $ CLM_FREQ  : int   2 0 2 0 2 0 0 1 0 0 ...
## $ REVOKED   : chr  "No" "No" "No" "No" ...
## $ MVR_PTS   : int   3 0 3 0 3 0 0 10 0 1 ...
## $ CAR_AGE   : int   18 1 10 6 17 7 1 7 1 17 ...
## $ URBANICITY : chr  "Highly Urban/ Urban" "Highly Urban/ Urban" "Highly Urban/ Urban" "Highly Urban/ Urban"
```

Corr analysis

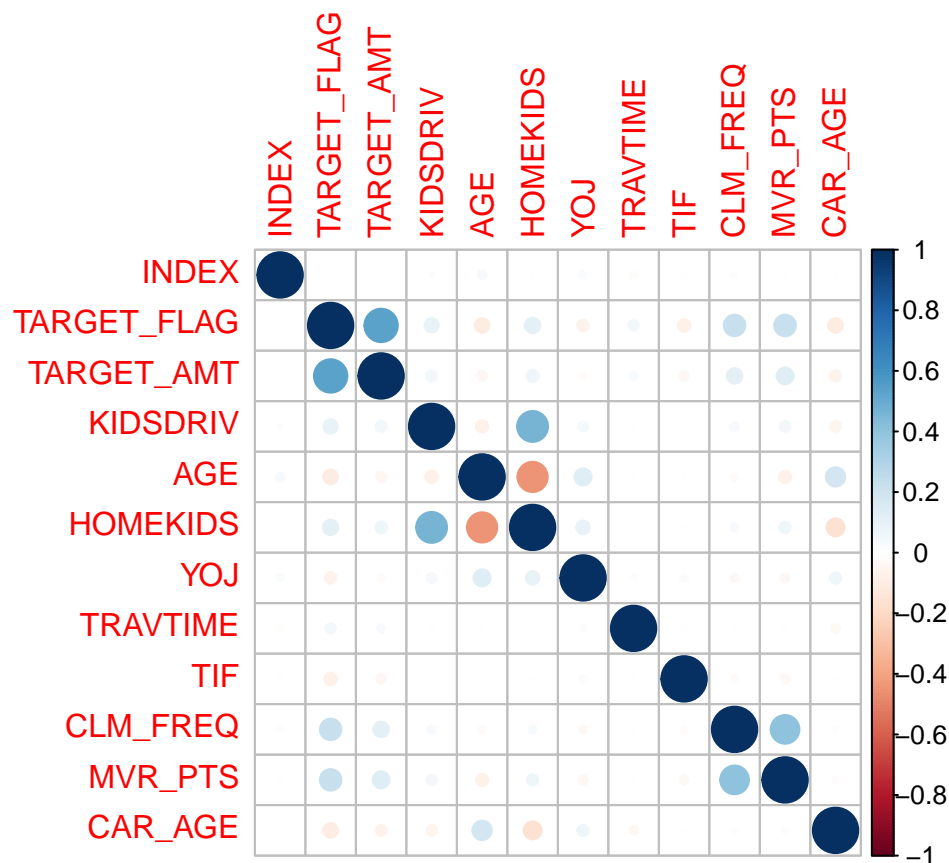
```
# Correlations
#cor_train <- cor(train_df)
#corrplot(cor_train)

# Correlation
train_df_2 <- train_df[,!names(train_df) %in% c("PARENT1", "MSTATUS", "EDUCATION", "JOB", "CAR_USE", "SEX", "AGE")]
a <- cor(train_df_2, use="complete.obs")
a
```

	INDEX	TARGET_FLAG	TARGET_AMT	KIDSDRIV	AGE
## INDEX	1.000000000	0.001329397	-0.001396512	0.010076793	0.039465360
## TARGET_FLAG	0.001329397	1.000000000	0.539929292	0.097295404	-0.103649398
## TARGET_AMT	-0.001396512	0.539929292	1.000000000	0.056857228	-0.046745810
## KIDSDRIV	0.010076793	0.097295404	0.056857228	1.000000000	-0.072361631
## AGE	0.039465360	-0.103649398	-0.046745810	-0.072361631	1.000000000
## HOMEKIDS	-0.002644576	0.115481537	0.068857861	0.463046635	-0.442383841
## YOJ	0.026014951	-0.066586612	-0.021446311	0.048112090	0.139566052
## TRAVTIME	-0.018846768	0.051491295	0.032708168	0.008979590	0.004555303
## TIF	-0.009216514	-0.077186438	-0.045259254	-0.003423442	0.002871951
## CLM_FREQ	0.015389421	0.223381685	0.115156936	0.035087170	-0.026312189
## MVR_PTS	0.007192153	0.225262361	0.137708292	0.055019621	-0.073523273
## CAR_AGE	-0.002148739	-0.104357704	-0.062833451	-0.055877063	0.182184524
##	HOMEKIDS	YOJ	TRAVTIME	TIF	CLM_FREQ
## INDEX	-0.002644576	0.02601495	-0.018846768	-0.009216514	0.015389421
## TARGET_FLAG	0.115481537	-0.06658661	0.051491295	-0.077186438	0.223381685
## TARGET_AMT	0.068857861	-0.02144631	0.032708168	-0.045259254	0.115156936
## KIDSDRIV	0.463046635	0.04811209	0.008979590	-0.003423442	0.035087170
## AGE	-0.442383841	0.13956605	0.004555303	0.002871951	-0.026312189
## HOMEKIDS	1.000000000	0.09041645	-0.007787772	0.004673246	0.030695809
## YOJ	0.090416449	1.000000000	-0.015762889	0.029302946	-0.030658029
## TRAVTIME	-0.007787772	-0.01576289	1.000000000	-0.009343232	0.009306981
## TIF	0.004673246	0.02930295	-0.009343232	1.000000000	-0.024972898
## CLM_FREQ	0.030695809	-0.03065803	0.009306981	-0.024972898	1.000000000

```
## MVR_PTS      0.062776101 -0.03917262  0.009937566 -0.037174513  0.400121265
## CAR_AGE      -0.156534495  0.06122969 -0.037055196  0.009125709 -0.011538390
##              MVR_PTS      CAR_AGE
## INDEX        0.007192153 -0.002148739
## TARGET_FLAG  0.225262361 -0.104357704
## TARGET_AMT   0.137708292 -0.062833451
## KIDSDRIV     0.055019621 -0.055877063
## AGE          -0.073523273  0.182184524
## HOMEKIDS     0.062776101 -0.156534495
## YOJ          -0.039172617  0.061229694
## TRAVTIME     0.009937566 -0.037055196
## TIF          -0.037174513  0.009125709
## CLM_FREQ     0.400121265 -0.011538390
## MVR_PTS      1.000000000 -0.019363647
## CAR_AGE      -0.019363647  1.000000000
```

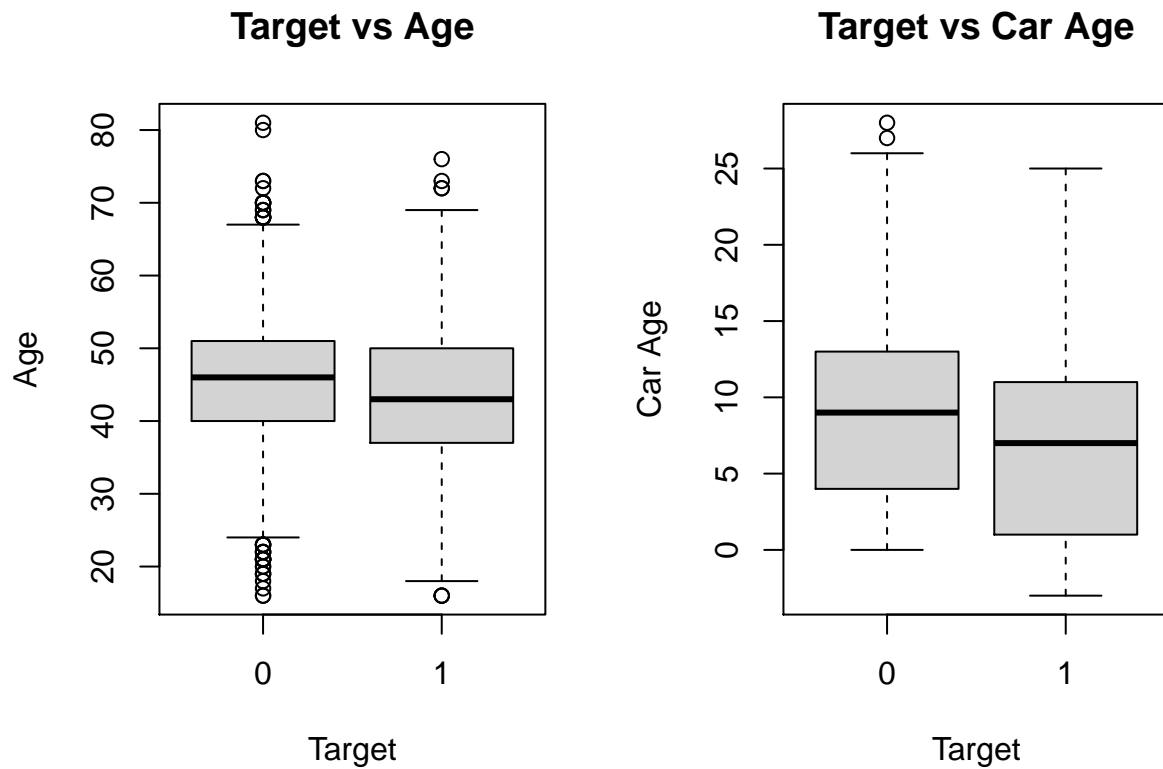
```
corrplot(a, method="circle")
```



Look at the relationship between TARGET_FLAG and some of the numerical variables.

```
par(mfrow=c(1,2))
# plot response variable "target" against predictor variable "age" and "car_age"
boxplot(AGE ~ TARGET_FLAG, train_df,
        main="Target vs Age",
        xlab="Target",
        ylab="Age")
```

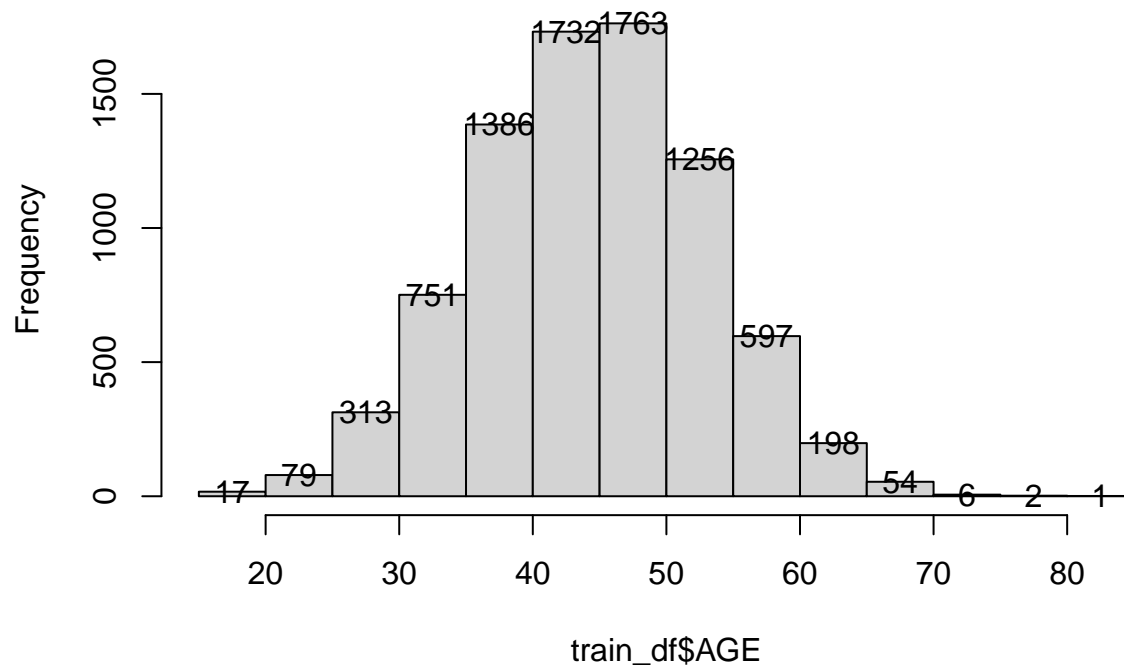
```
boxplot(CAR_AGE ~ TARGET_FLAG, train_df,
        main="Target vs Car Age",
        xlab="Target",
        ylab="Car Age")
```



Look at the distribution of some numerical variables.

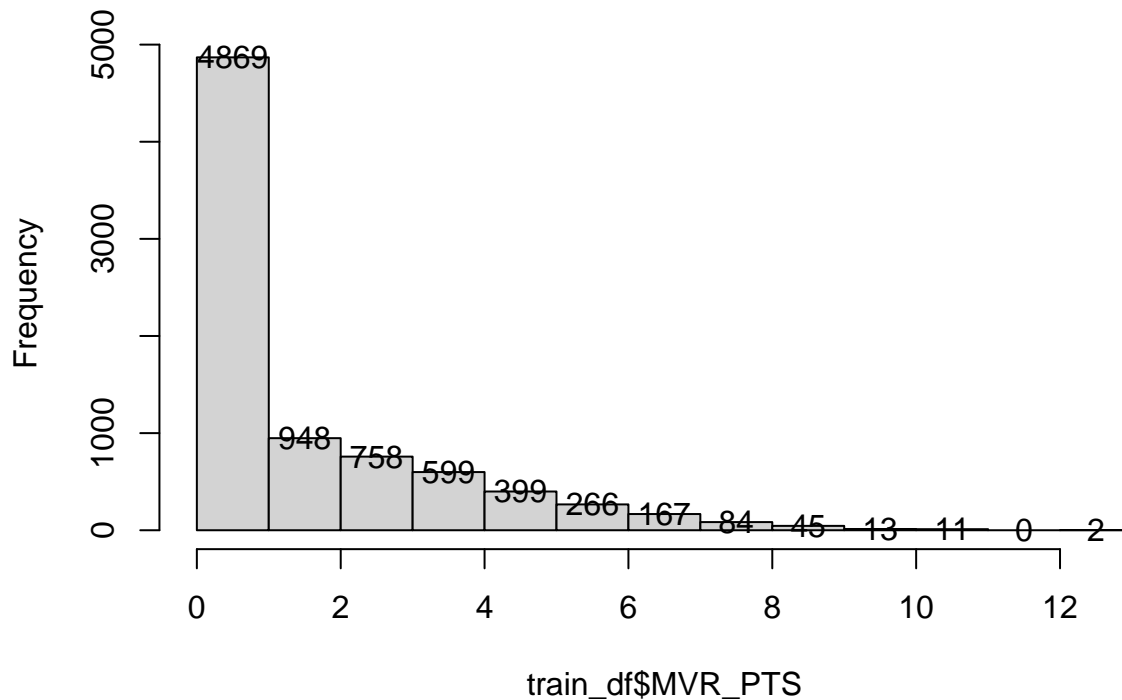
```
h <- hist(train_df$AGE)
text(h$mids, h$counts, labels=h$counts)
```

Histogram of train_df\$AGE



```
h <- hist(train_df$MVR_PTS)
text(h$mids,h$counts,labels=h$counts)
```


Histogram of train_df\$MVR_PTS



Check for NA's

```
has_NA = names(which(sapply(train_df, anyNA)))
has_NA
```

```
## [1] "AGE"      "YOJ"      "CAR_AGE"
```

Check test_df for NA's

```
has_NA_test = names(which(sapply(test_df, anyNA)))
has_NA_test
```

```
## [1] "TARGET_FLAG" "TARGET_AMT"  "AGE"         "YOJ"         "CAR_AGE"
```

Since we see our test_df has NAs for the same variables as test, we need to come up with a way to handle making predictions on records that have these values as NA. We will create an "_NA" columns as dummy variables for AGE, YOJ, and CAR_AGE, 1 marking them as NA and 0 if they have a value.

```
for (col in has_NA)
{
  new_col = (paste(col, "_NA", sep=""))
  train_df[,new_col] = as.numeric(is.na(train_df[,col]))
  test_df[,new_col] = as.numeric(is.na(test_df[,col]))
  # fill missing numerics with median value
}
```

```

train_df[,col][is.na(train_df[,col])] = median(train_df[,col], na.rm=TRUE)
test_df[,col][is.na(test_df[,col])] = median(test_df[,col], na.rm=TRUE)
}

```

Create train_amt_df dataframe for multiple linear regression model

```

train_amt_df <- subset(train_df, TARGET_AMT > 0)
summary(train_amt_df$TARGET_FLAG)

```

```

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         1         1         1         1         1         1

```

Modeling

1) Binary Logistic Regression

```

# preliminary exploration with one predictor
modell1 <- glm(formula = TARGET_FLAG ~ AGE, family = binomial(), data = train_df)
summary(modell1)

```

```

##
## Call:
## glm(formula = TARGET_FLAG ~ AGE, family = binomial(), data = train_df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0728  -0.8042  -0.7403   1.4313   2.0168
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.186818   0.131990   1.415    0.157
## AGE         -0.027373   0.002954  -9.265 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 9418.0  on 8160  degrees of freedom
## Residual deviance: 9330.8  on 8159  degrees of freedom
## AIC: 9334.8
##
## Number of Fisher Scoring iterations: 4

```

Binary Logistic Regression Model with more variables

```

BLR_all_vars = glm(TARGET_FLAG ~ AGE +
                    CAR_AGE +
                    MVR_PTS +
                    YOJ +
                    CLM_FREQ +

```

```
TIF, family = binomial(), data = train_df)
summary(BLR_all_vars)
```

```
##
## Call:
## glm(formula = TARGET_FLAG ~ AGE + CAR_AGE + MVR_PTS + YOJ + CLM_FREQ +
##       TIF, family = binomial(), data = train_df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8021  -0.7630  -0.6108   0.9899   2.4099
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.095985   0.153832   0.624    0.533
## AGE         -0.019810   0.003107  -6.376 1.82e-10 ***
## CAR_AGE     -0.035902   0.004949  -7.254 4.04e-13 ***
## MVR_PTS      0.147989   0.012363  11.971 < 2e-16 ***
## YOJ         -0.025942   0.006464  -4.013 5.99e-05 ***
## CLM_FREQ     0.293062   0.022906  12.794 < 2e-16 ***
## TIF         -0.045555   0.006689  -6.811 9.72e-12 ***
## ---
## 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: 8704.2  on 8154  degrees of freedom
## AIC: 8718.2
##
## Number of Fisher Scoring iterations: 4
```

Step through AIC scores to find best model

```
step_BLR = stepAIC(BLR_all_vars)
```

```
## Start:  AIC=8718.2
## TARGET_FLAG ~ AGE + CAR_AGE + MVR_PTS + YOJ + CLM_FREQ + TIF
##
##           Df Deviance    AIC
## <none>          8704.2 8718.2
## - YOJ           1   8720.1 8732.1
## - AGE           1   8745.2 8757.2
## - TIF           1   8752.2 8764.2
## - CAR_AGE       1   8757.7 8769.7
## - MVR_PTS       1   8847.9 8859.9
## - CLM_FREQ      1   8864.3 8876.3
```

```
summary(step_BLR)
```

```
##
```

```
## Call:
## glm(formula = TARGET_FLAG ~ AGE + CAR_AGE + MVR_PTS + YOJ + CLM_FREQ +
##      TIF, family = binomial(), data = train_df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8021  -0.7630  -0.6108   0.9899   2.4099
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.095985   0.153832   0.624   0.533
## AGE         -0.019810   0.003107  -6.376 1.82e-10 ***
## CAR_AGE     -0.035902   0.004949  -7.254 4.04e-13 ***
## MVR_PTS      0.147989   0.012363  11.971 < 2e-16 ***
## YOJ         -0.025942   0.006464  -4.013 5.99e-05 ***
## CLM_FREQ     0.293062   0.022906  12.794 < 2e-16 ***
## TIF         -0.045555   0.006689  -6.811 9.72e-12 ***
## ---
## 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: 8704.2  on 8154  degrees of freedom
## AIC: 8718.2
##
## Number of Fisher Scoring iterations: 4
```

2) Multiple Linear Regression

Multiple Linear Regression models with many variables

```
MLR_all_vars = lm(TARGET_AMT ~ AGE +
                  CAR_AGE +
                  MVR_PTS +
                  YOJ +
                  CLM_FREQ +
                  TIF, data = train_amt_df)
summary(MLR_all_vars)
```

```
##
## Call:
## lm(formula = TARGET_AMT ~ AGE + CAR_AGE + MVR_PTS + YOJ + CLM_FREQ +
##      TIF, data = train_amt_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
##  -6311  -3111  -1579   160 101042
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4117.71     892.22   4.615 4.16e-06 ***
## AGE          22.58       17.80   1.268  0.2048
```

```
## CAR_AGE      -23.46      31.64  -0.741  0.4586
## MVR_PTS      132.33      68.03   1.945  0.0519 .
## YOJ          56.31      38.36   1.468  0.1423
## CLM_FREQ     -64.15     140.39  -0.457  0.6478
## TIF          -7.77      42.47  -0.183  0.8549
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7739 on 2146 degrees of freedom
## Multiple R-squared:  0.003804,    Adjusted R-squared:  0.001019
## F-statistic: 1.366 on 6 and 2146 DF,  p-value: 0.2248
```

```
# step_BLR prediction on test
test_preds_BLR = round(predict(step_BLR, newdata=test_df, type='response'))
test_df$TARGET_FLAG = test_preds_BLR
test_preds_MLR = predict(MLR_all_vars, newdata=test_df)
test_df$TARGET_AMT = test_preds_MLR

# write out evaluation data with predictions
write.csv(test_df, 'eval_with_preds.csv')
```