

Data 621 - Homework 4

Group 2: William Aiken, Donald Butler, Michael Ippolito, Bharani Nittala, and Leticia Salazar

November 20, 2022

Contents

Overview:	1
Objective:	1
Description:	2
Data Exploration:	4
Data Preparation:	15
	18
Model Building:	19
	26
Select Models:	26
Predictions:	35
Appendix:	36
References:	43

Overview:

In this homework assignment, you will explore, analyze and model a data set containing approximately 8,000 records representing a customer at an auto insurance company. Each record has two response variables. The first response variable, TARGET_FLAG, is a 1 or a 0. A “1” means that the person was in a car crash. A zero means that the person was not in a car crash. The second response variable is TARGET_AMT. This value is zero if the person did not crash their car. But if they did crash their car, this number will be a value greater than zero.

Objective:

Your objective is to build multiple linear regression and binary logistic regression models on the training data to predict the probability that a person will crash their car and also the amount of money it will cost if the person does crash their car. You can only use the variables given to you (or variables that you derive from the variables provided). Below is a short description of the variables of interest in the data set:

Description:

Below is a short description of the variables of interest in the data set:

VARIABLE NAME:	DEFINITION:	THEORETICAL EFFECT:
INDEX	Identification Variable (do not use)	None
TARGET_FLAG	Was Car in a crash? 1 = YES 2 = NO	None
TARGET_AMT	If car was in a crash, what was the cost	None
AGE	Age of Driver	Very young people tend to be risky. Maybe very old people also.
BLUEBOOK	Value of Vehicle	Unknown effect on probability of collision, but probably effect the payout if there is a crash
CAR_AGE	Vehicle Age	Unknown effect on probability of collision, but probably effect the payout if there is a crash
CAR_TYPE	Type of Car	Unknown effect on probability of collision, but probably effect the payout if there is a crash
CAR_USE	Vehicle Use	Commercial vehicles are driven more, so might increase probability of collision
CLM_FREQ	# Claims (Past 5 Years)	The more claims you filed in the past, the more you are likely to file in the future
EDUCATION	Max Education Level	Unknown effect, but in theory more educated people tend to drive more safely
HOMEKIDS	# Children at Home	Unknown effect
HOME_VAL	Home Value	In theory, home owners tend to drive more responsibly
INCOME	Income	In theory, rich people tend to get into fewer crashes
JOB	Job Category	In theory, white collar jobs tend to be safer
KIDSDRIV	# Driving Children	When teenagers drive your car, you are more likely to get into crashes
MSTATUS	Martial Status	In theory, married people drive more safely
MVR_PTS	Motor Vehicle Record Points	If you get lots of traffic tickets, you tend to get into more crashes
OLDCLAIM	Total Claims (Past 5 Years)	If your total payout over the past five years was high, this suggests future payouts will be high
PARENT1	Single Parent	Unknown effect
RED_CAR	A Red Car	Urban legend says that red cars (especially red sports cars) are more risky. Is that true?
REVOKE	License Revoked (Past 7 Years)	If your license was revoked in the past 7 years, you probably are a more risky driver.

VARIABLE NAME:	DEFINITION:	THEORETICAL EFFECT:
SEX	Gender	Urban legend says that women have less crashes than men. Is that true?
TIF	Time in Force	People who have been customers for a long time are usually more safe.
TRAVTIME	Distance to Work	Long drives to work usually suggest greater risk
URBANICITY	Home / Work Area	Unknown
YOJ	Years on Job	People who stay at a job for a long time are usually more safe

Load Libraries:

These are the libraries used to explore, prepare, analyze and build our models

```
library(tidyverse)
library(caret)
library(pROC)
library(corrplot)
library(GGally)
library(psych)
library(car)
library(kableExtra)
library(gridExtra)
library(performance)
library(faraway)
library(jtools)
library(DataExplorer)
library(hrbrthemes)
library(MASS)
```

Load Data set:

We have included the original data sets in our GitHub account and read from this location. Our training data set includes 8,161 records and 26 variables.

```
## Rows: 8,161
## Columns: 26
## $ INDEX      <int> 1, 2, 4, 5, 6, 7, 8, 11, 12, 13, 14, 15, 16, 17, 19, 20, 2~
## $ TARGET_FLAG <int> 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1~
## $ TARGET_AMT  <dbl> 0.000, 0.000, 0.000, 0.000, 2946.000, 0.000, 4021.0~
## $ KIDSDRV    <int> 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ AGE         <int> 60, 43, 35, 51, 50, 34, 54, 37, 34, 50, 53, 43, 55, 53, 45~
## $ HOMEKIDS   <int> 0, 0, 1, 0, 0, 1, 0, 2, 0, 0, 0, 0, 0, 0, 3, 0, 3, 2, 1~
## $ YOJ        <int> 11, 11, 10, 14, NA, 12, NA, 10, 7, 14, 5, 11, 11, 0, 1~
## $ INCOME     <chr> "$67,349", "$91,449", "$16,039", "", "$114,986", "$125,301~
## $ PARENT1    <chr> "No", "No", "No", "No", "No", "Yes", "No", "No", "No", "No~
## $ HOME_VAL   <chr> "$0", "$257,252", "$124,191", "$306,251", "$243,925", "$0"~
## $ MSTATUS    <chr> "z_No", "z_No", "Yes", "Yes", "Yes", "z_No", "Yes", "Yes", ~
```

```

## $ SEX <chr> "M", "M", "z_F", "M", "z_F", "z_F", "M", "z_F", "M"~
## $ EDUCATION <chr> "PhD", "z_High School", "z_High School", "<High School", "~"
## $ JOB <chr> "Professional", "z_Blue Collar", "Clerical", "z_Blue Colla~<int> 14, 22, 5, 32, 36, 46, 33, 44, 34, 48, 15, 36, 25, 64, 48, ~
## $ CAR_USE <chr> "Private", "Commercial", "Private", "Private", "Private", ~<chr> "$14,230", "$14,940", "$4,010", "$15,440", "$18,000", "$17~<int> 11, 1, 4, 7, 1, 1, 1, 1, 7, 1, 7, 7, 6, 1, 6, 6, 7, 4, ~
## $ BLUEBOOK <chr> "Minivan", "Minivan", "z_SUV", "Minivan", "z_SUV", "Sports~<chr> "yes", "yes", "no", "yes", "no", "no", "no", "yes", "no", ~<chr> "$4,461", "$0", "$38,690", "$0", "$19,217", "$0", "$0", "$~<int> 2, 0, 2, 0, 2, 0, 0, 1, 0, 0, 0, 0, 2, 0, 0, 0, 0, 0, 0, 0, 2~<chr> "No", "No", "No", "Yes", "No", "Yes", "Yes", "No", "N~<int> 3, 0, 3, 0, 3, 0, 0, 10, 0, 1, 0, 0, 3, 3, 3, 0, 0, 0, 0, 0, ~<int> 18, 1, 10, 6, 17, 7, 1, 7, 1, 17, 11, 1, 9, 10, 5, 13, 16, ~<chr> "Highly Urban/ Urban", "Highly Urban/ Urban", "Highly Urb~<chr> "Highly Urban/ Urban", "Highly Urban/ Urban", "Highly Urb~
```

Data Exploration:

For insight on the data we use the `summary()` function on the train dataset:

```

##      INDEX      TARGET_FLAG      TARGET_AMT      KIDSDRV
##  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
##                   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

```

Let's look at the statistical summary for the `dfeval` data as well:

```

##          INDEX        TARGET_FLAG        TARGET_AMT        KIDS DRIV        AGE
## Min.  : 3    Mode:logical    Mode:logical    Min.  :0.0000  Min.  :17.00
## 1st Qu.: 2632 NA's:2141     NA's:2141     1st Qu.:0.0000  1st Qu.:39.00
## Median : 5224
## Mean   : 5150
## 3rd Qu.: 7669
## Max.   :10300
## NA's   :1
##
##          HOMEKIDS        YOJ           INCOME        PARENT1
## Min.  :0.0000  Min.  : 0.00  Length:2141      Length:2141
## 1st Qu.:0.0000 1st Qu.: 9.00  Class :character  Class :character
## Median :0.0000  Median :11.00  Mode  :character  Mode  :character
## Mean   :0.7174  Mean   :10.38
## 3rd Qu.:1.0000 3rd Qu.:13.00
## Max.   :5.0000  Max.   :19.00
## NA's   :94
##
##          HOME_VAL        MSTATUS          SEX          EDUCATION
## Length:2141      Length:2141  Length:2141      Length:2141
## Class :character  Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##          #
##          #
##          #
##          #

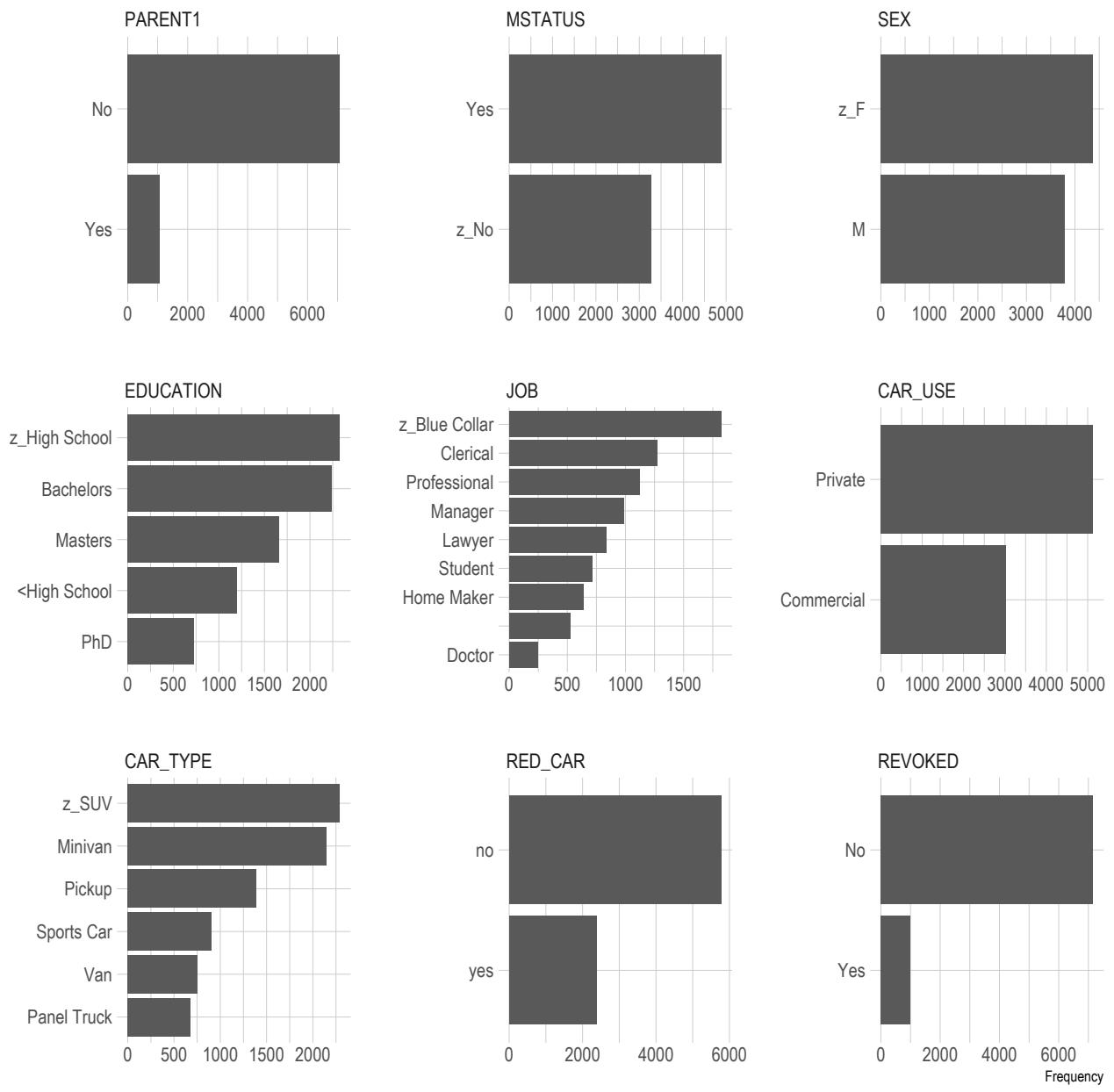
```

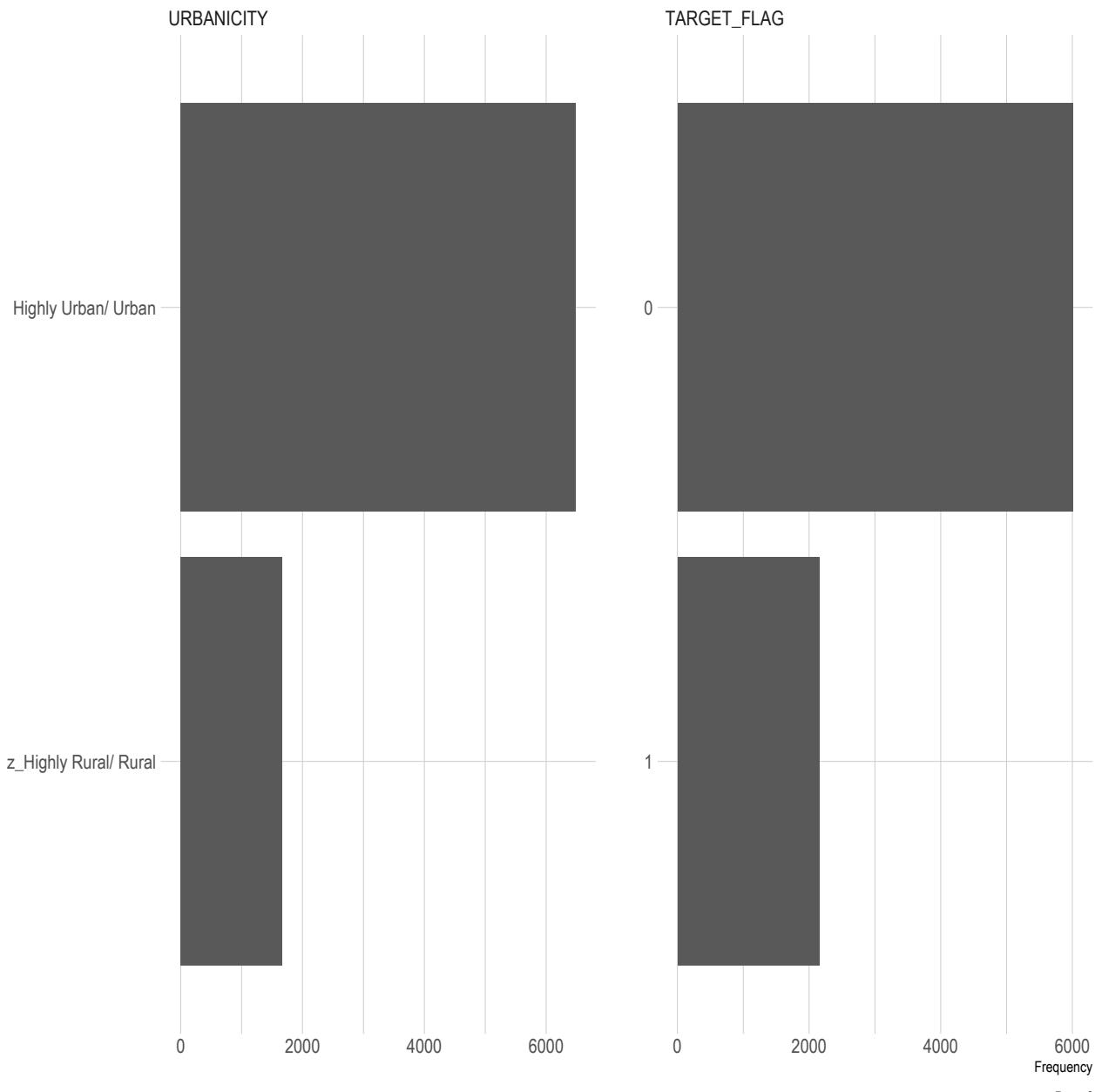
```

##          JOB           TRAVTIME        CAR_USE        BLUEBOOK
## Length:2141      Min. : 5.00  Length:2141  Length:2141
## Class :character  1st Qu.: 22.00  Class :character  Class :character
## Mode  :character   Median : 33.00  Mode  :character  Mode  :character
##                               Mean  : 33.15
##                               3rd Qu.: 43.00
##                               Max.  :105.00
##
##          TIF           CAR_TYPE       RED_CAR        OLDCLAIM
## Min.  : 1.000  Length:2141  Length:2141  Length:2141
## 1st Qu.: 1.000  Class :character  Class :character  Class :character
## Median : 4.000  Mode  :character  Mode  :character  Mode  :character
## Mean   : 5.245
## 3rd Qu.: 7.000
## Max.   :25.000
##
##          CLM_FREQ        REVOKED        MVR_PTS        CAR_AGE
## Min.  :0.000  Length:2141  Min.  : 0.000  Min.  : 0.000
## 1st Qu.:0.000  Class :character  1st Qu.: 0.000  1st Qu.: 1.000
## Median :0.000  Mode  :character  Median : 1.000  Median : 8.000
## Mean   :0.809
## 3rd Qu.:2.000
## Max.   :5.000
##
##          URBANICITY
## Length:2141
## Class :character
## Mode  :character
##
##
```

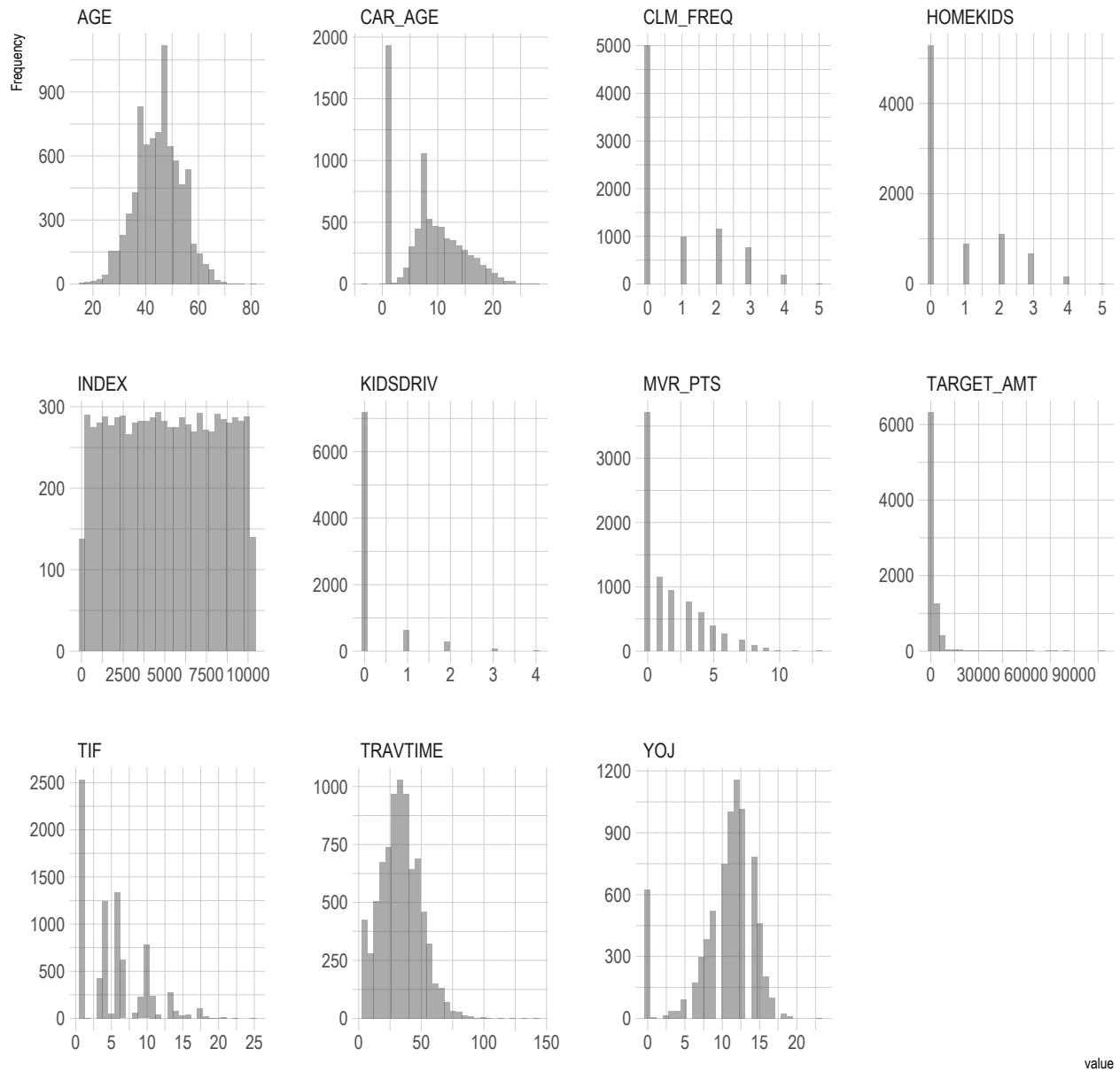
Distributions for training data set:

```
## 4 columns ignored with more than 50 categories.
## INCOME: 6613 categories
## HOME_VAL: 5107 categories
## BLUEBOOK: 2789 categories
## OLDCLAIM: 2857 categories
```



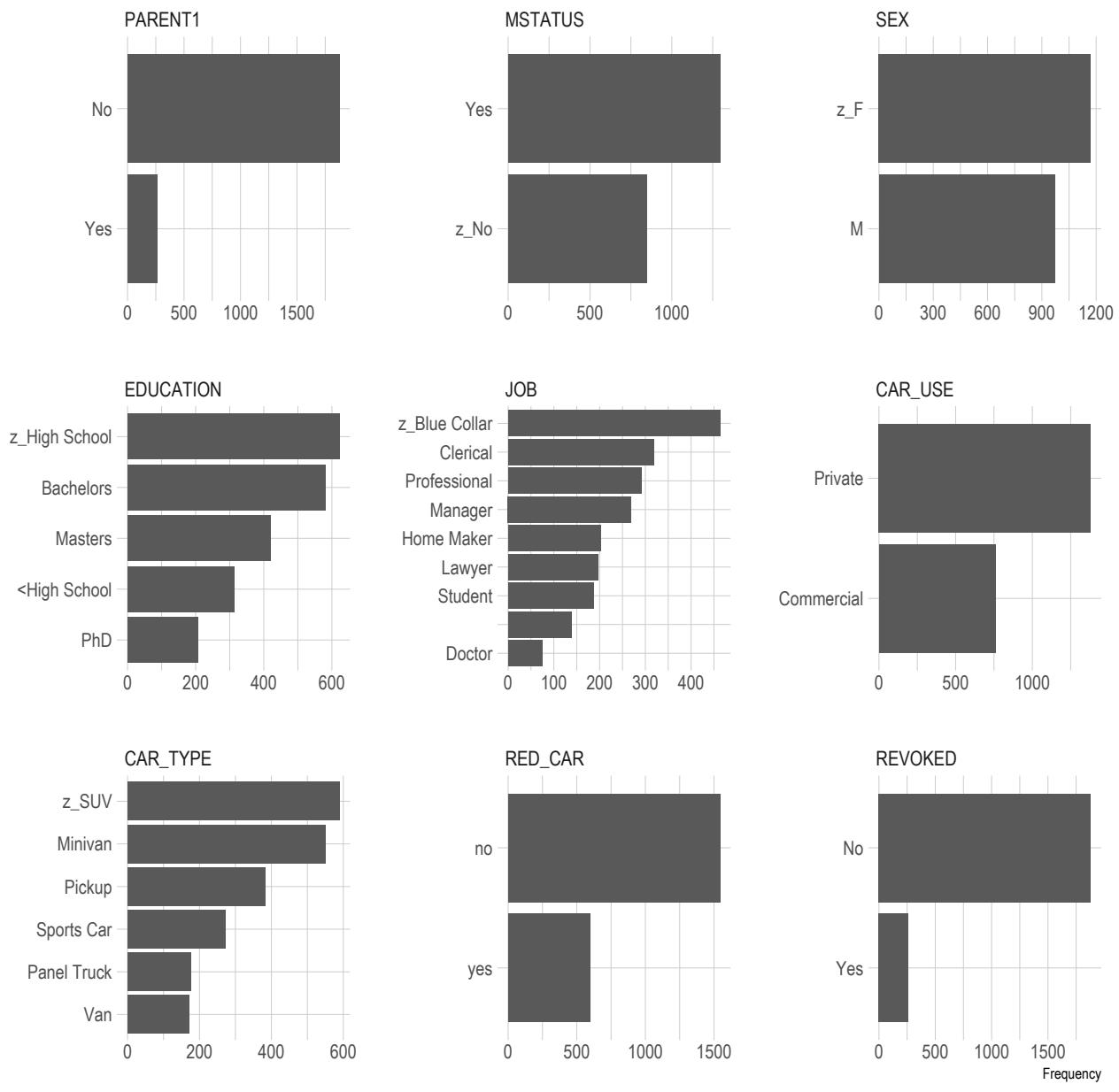


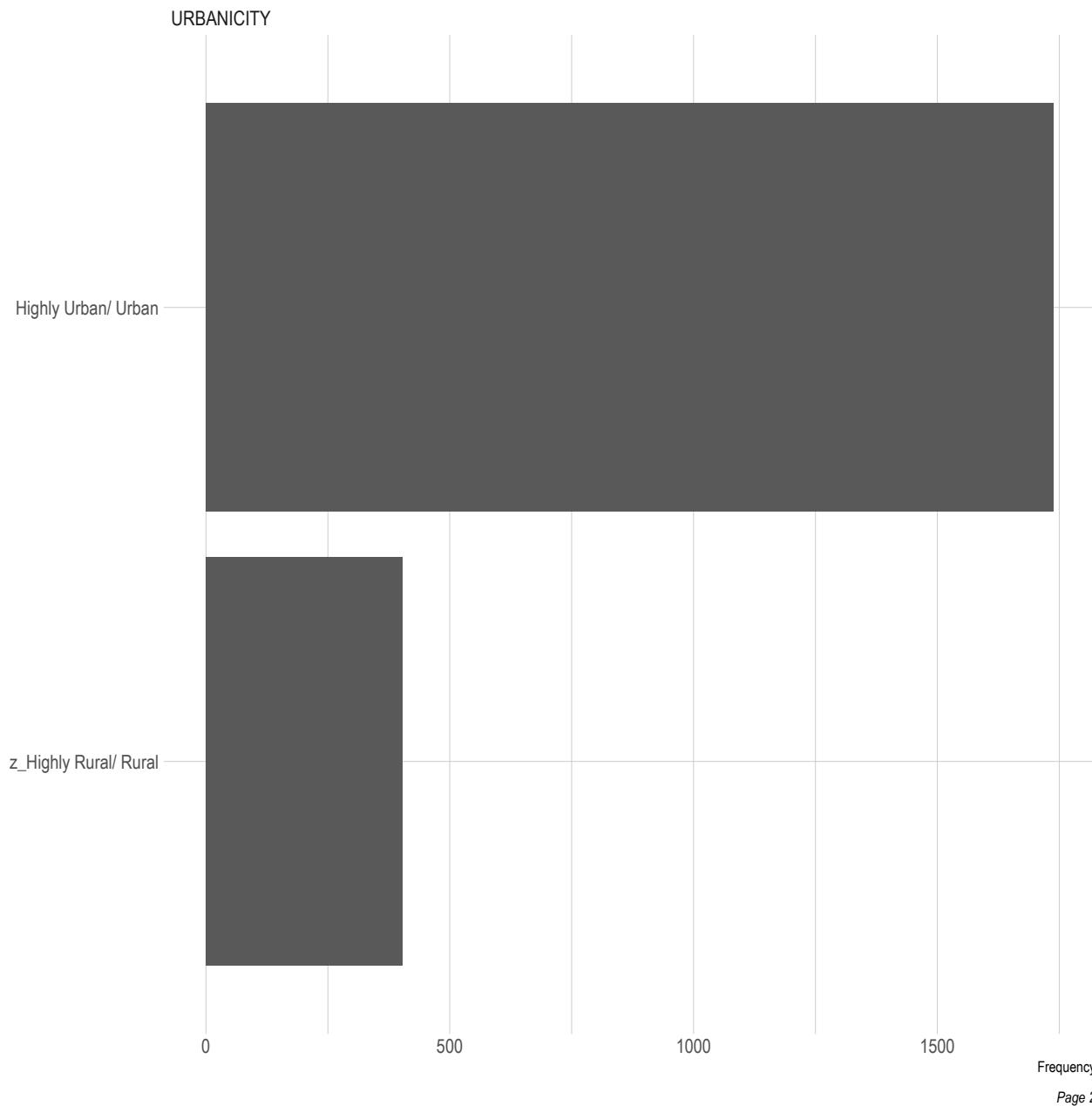
'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



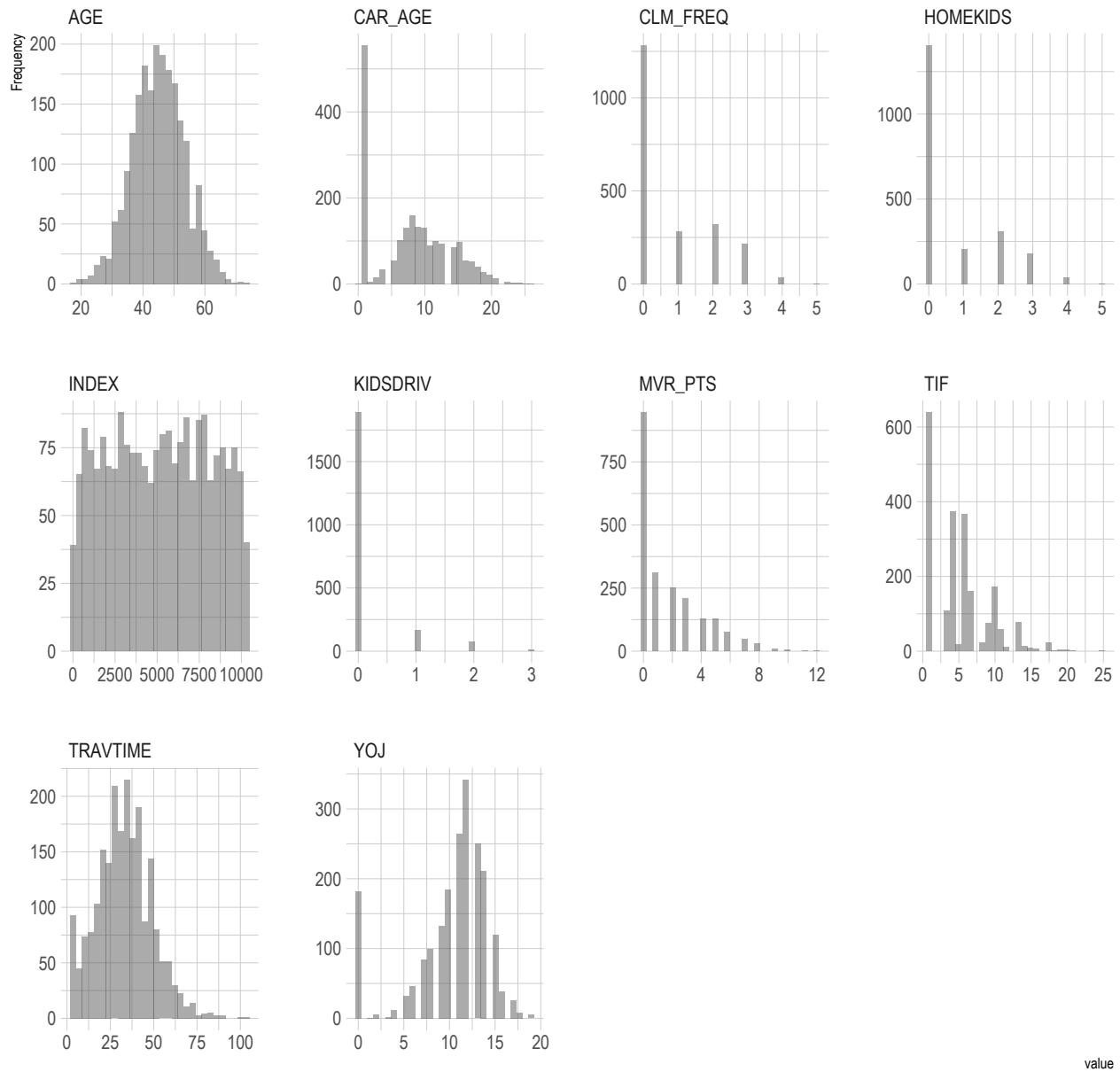
Distributions for evaluation data set:

```
## 4 columns ignored with more than 50 categories.
## INCOME: 1804 categories
## HOME_VAL: 1398 categories
## BLUEBOOK: 1417 categories
## OLDCLAIM: 834 categories
```





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



The following dummy variables are done to both the training and evaluation data set:

PARENT1 The *PARENT1* variable has two values, Yes and No, to indicate if the observation is a single parent. We will construct a dummy variable *SingleParent* = 1 if *PARENT1* = Yes.

```
##   PARENT1     n
## 1      No 7084
## 2      Yes 1077
```

SEX The *SEX* variable has two values, M and z_F. We will create a dummy variable *Male* = 1 if *SEX* = M.

```
##   SEX     n
## 1    M 3786
## 2 z_F 4375
```

MSTATUS The variable *MSTATUS* has two values, Yes and z_No, to indicate the marital status. We will create a dummy variable *Married* = 1 if *MSTATUS* = Yes.

```
##   MSTATUS     n
## 1      Yes 4894
## 2    z_No 3267
```

EDUCATION The *EDUCATION* variable takes on 5 values ranging from less than high school through PHD. We will construct dummy variables: *HighSchool*, *Bachelors*, *Masters*, *PHD*, to indicate the highest level of education completed.

```
##       EDUCATION     n
## 1 <High School 1203
## 2    Bachelors 2242
## 3      Masters 1658
## 4        PhD  728
## 5 z_High School 2330
```

JOB The *JOB* variable takes on 8 values. The *JOB* variable has 526 missing values, so we will construct dummy variables for all 8 values assuming the missing values are not one of the listed professions. The dummy variables we will create are: *Clerical*, *Doctor*, *HomeMaker*, *Lawyer*, *Manager*, *Professional*, *Student*, and *BlueCollar*.

```

##          JOB      n
## 1          526
## 2    Clerical 1271
## 3     Doctor  246
## 4   Home Maker  641
## 5     Lawyer  835
## 6    Manager  988
## 7 Professional 1117
## 8    Student  712
## 9 z_Blue Collar 1825

```

CAR_USE The *CAR_USE* variable has two values, Commercial and Private. We will construct a dummy variable *Commercial* = 1 if Commercial.

```

##          CAR_USE      n
## 1 Commercial 3029
## 2    Private 5132

```

CAR_TYPE The *CAR_TYPE* variable takes on 6 values. We will create dummy variables; *Minivan*, *PanelTruck*, *Pickup*, *SportsCar*, and *Van*.

```

##          CAR_TYPE      n
## 1        Minivan 2145
## 2 Panel Truck  676
## 3     Pickup 1389
## 4 Sports Car  907
## 5       Van  750
## 6     z_SUV 2294

```

RED_CAR The *RED_CAR* variable has two values, yes and no. We will create a dummy variable *RedCar* = 1 if *RED_CAR* = yes.

```

##          RED_CAR      n
## 1        no 5783
## 2      yes 2378

```

REVOKE The *REVOKE* variable has two values, Yes and No. We will create a dummy variable *DLRevoked* = 1 if *REVOKE* = Yes.

```

##          REVOKE      n
## 1        No 7161
## 2      Yes 1000

```

URBANICITY The *URBANICITY* variable has two values, Highly Urban/ Urban and z_Highly Rural/ Rural. We will create a dummy variable *Urban* = 1 if *URBANICITY* = Highly Urban/ Urban.

```
##           URBANICITY      n
## 1  Highly Urban/ Urban 6492
## 2 z_Highly Rural/ Rural 1669
```

Data Preparation:

Performed to both the training and evaluation data sets but will only be displayed for the training data.

Data Cleaning Function Checking for missing values for the training data set. This includes NAs that might be introduced as a result of conversion to numeric.

```
##      INDEX TARGET_FLAG TARGET_AMT KIDSDRV AGE HOMEKIDS
##      0        0        0        0       6        0
##      YOJ     INCOME PARENT1 HOME_VAL MSTATUS SEX
##      454        0        0        0       0        0
## EDUCATION      JOB TRAVTIME CAR_USE BLUEBOOK TIF
##      0        0        0        0       0        0
## CAR_TYPE RED_CAR OLDCLAIM CLM_FREQ REVOKED MVR PTS
##      0        0        0        0       0        0
## CAR AGE URBANICITY
##      510        0

##      BLUEBOOK HOME_VAL INCOME OLDCLAIM
## 1        0      464    445        0
```

- The attributes *BLUEBOOK*, *HOME_VAL*, *INCOME*, and *OLDCLAIM* are dollar amounts stored as characters. Need to convert to int. It is noted that converting blank values to numeric introduces some NAs; therefore our cleaning function will handle these cases.
- Variables with NA: *AGE* (6), *YOJ* (454), *CAR_AGE* (510)
- Consider creating *AGE* groups Under25 and Over65 to account for young and older drivers.
- Consider creating *CAR_AGE* groups to identify new cars. One observation has a *CAR_AGE* = -3, which shouldn't be possible.
- Consider creating *YOJ* (Year on Job) groups to identify job stability; Over5years etc.

```

## INDEX TARGET_FLAG TARGET_AMT KIDSDRV AGE HOMEKIDS YOJ INCOME HOME_VAL
## 1      1          0          0       0  60       0  11  67349       0
## 2      2          0          0       0  43       0  11  91449  257252
## 3      4          0          0       0  35       1  10  16039 124191
## 4      5          0          0       0  51       0  14  54028 306251
## 5      6          0          0       0  50       0  11 114986 243925
## 6      7          1        2946       0  34       1  12 125301       0
## TRAVTIME BLUEBOOK TIF OLDCALL CLM_FREQ MVR PTS CAR AGE SingleParent Male
## 1      14     14230    11     4461       2       3     18       0     1
## 2      22     14940     1       0       0       0       1       0     1
## 3      5      4010     4     38690       2       3     10       0     0
## 4      32     15440     7       0       0       0       6       0     1
## 5      36     18000     1    19217       2       3     17       0     0
## 6      46     17430     1       0       0       0       7       1     0
## Married HighSchool Bachelors Masters PHD Clerical Doctor HomeMaker Lawyer
## 1      0          0          0       0   1       0       0       0     0
## 2      0          1          0       0   0       0       0       0     0
## 3      1          1          0       0   0       1       0       0     0
## 4      1          0          0       0   0       0       0       0     0
## 5      1          0          0       0   1       0       1       0     0
## 6      0          0          1       0   0       0       0       0     0
## Manager Professional Student BlueCollar Commercial Minivan PanelTruck Pickup
## 1      0          1          0       0       0       1       0       0     0
## 2      0          0          0       1       1       1       0       0     0
## 3      0          0          0       0       0       0       0       0     0
## 4      0          0          0       1       0       0       1       0     0
## 5      0          0          0       0       0       0       0       0     0
## 6      0          0          0       1       1       0       0       0     0
## SportsCar Van RedCar DLRevoked Urban
## 1      0    0    1    0    1
## 2      0    0    1    0    1
## 3      0    0    0    0    1
## 4      0    0    1    0    1
## 5      0    0    0    1    1
## 6      1    0    0    0    1

```

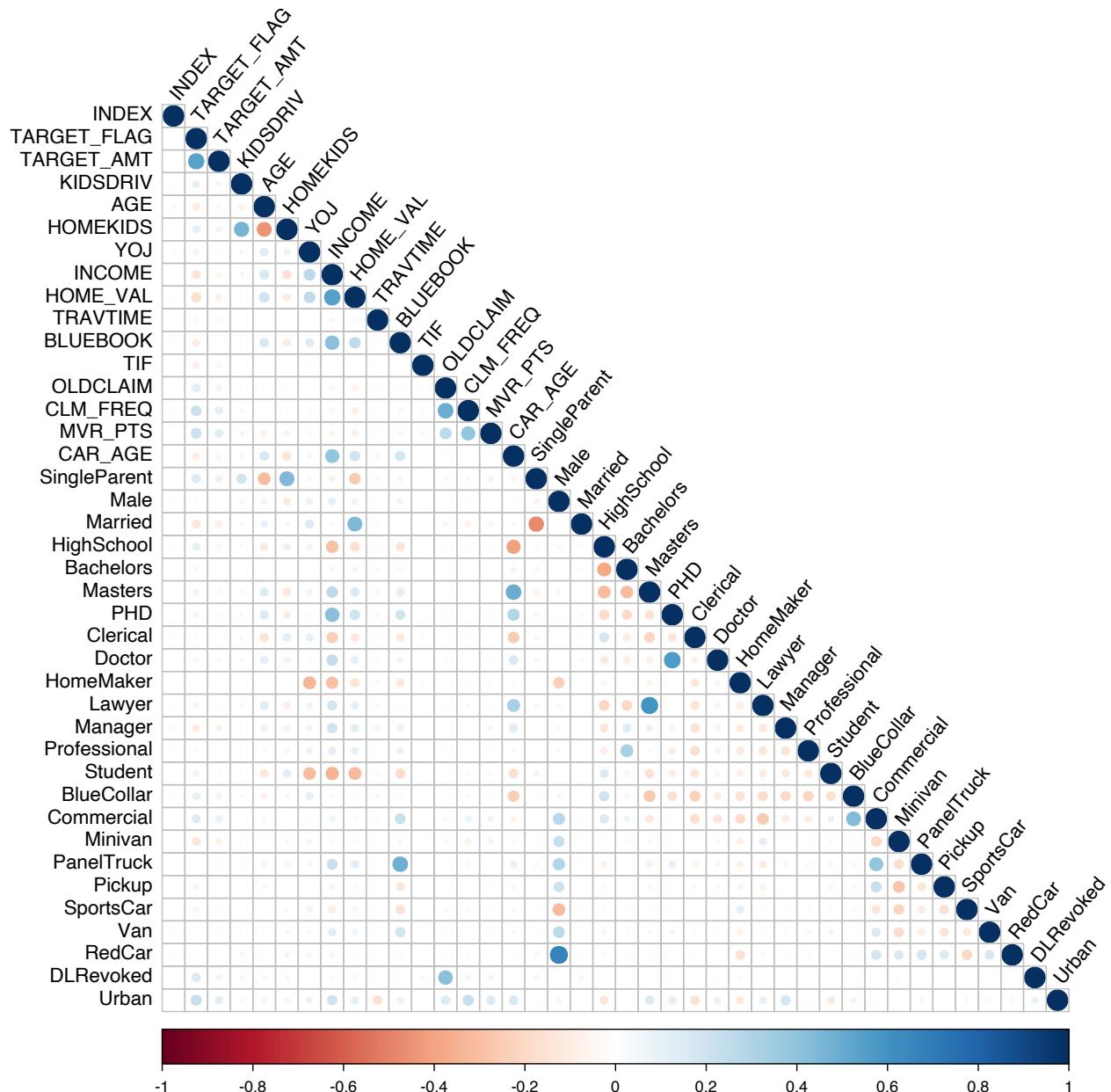
Checking for missing values for the evaluation data set, including NAs that might be introduced as a result of conversion to numeric. The TARGET_FLAG and TARGET_AMT are the response variables so should remain NA for now.

```
##      INDEX TARGET_FLAG TARGET_AMT KIDSDRIV AGE HOMEKIDS
##      0      2141     2141      0      1      0
##      YOJ    INCOME PARENT1 HOME_VAL MSTATUS SEX
##      94      0      0      0      0      0
## EDUCATION      JOB TRAVTIME CAR_USE BLUEBOOK TIF
##      0      0      0      0      0      0
## CAR_TYPE RED_CAR OLDCLAIM CLM_FREQ REVOKED MVR PTS
##      0      0      0      0      0      0
## CAR_AGE URBANICITY
##      129      0

##   BLUEBOOK HOME_VAL INCOME OLDCLAIM
## 1      0      111     125      0
```

- The attributes BLUEBOOK, HOME_VAL, INCOME, and OLDCLAIM are dollar amounts stored as characters. Need to convert to int. It is noted that converting blank values to numeric introduces some NAs; therefore our cleaning function will handle these cases.
- Variables with NA: AGE (1), YOJ (94), CAR_AGE (125)

The correlation plot below is measuring the degree of linear relationship within the cleaned training data. The values in which this is measured falls between -1 and +1, with +1 being a strong positive correlation and -1 a strong negative correlation.



Model Building:

We will be building six different models; three multiple linear regression (MLR) models and three binary logistic regression (BLR) models.

Model 1

MLR: Base model where all variables are tested

```
##  
## Call:  
## lm(formula = TARGET_AMT ~ ., data = cleandf, na.action = na.omit)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -6294    -470     -56     239  101122  
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -6.066e+02 4.746e+02 -1.278 0.201221  
## INDEX        -1.474e-03 1.480e-02 -0.100 0.920691  
## TARGET_FLAG   5.705e+03 1.136e+02 50.240 < 2e-16 ***  
## KIDSDRIV    -3.096e+01 9.910e+01 -0.312 0.754743  
## AGE          5.941e+00 6.175e+00  0.962 0.336051  
## HOMEKIDS    3.970e+01 5.710e+01  0.695 0.486936  
## YOJ          8.275e+00 1.318e+01  0.628 0.530161  
## INCOME       -2.265e-03 1.577e-03 -1.436 0.151171  
## HOME_VAL     4.051e-04 5.164e-04  0.784 0.432786  
## TRAVTIME     5.643e-01 2.825e+00  0.200 0.841655  
## BLUEBOOK    2.933e-02 7.538e-03  3.891 0.000101 ***  
## TIF          -2.946e+00 1.068e+01 -0.276 0.782577  
## OLDCLAIM     3.055e-03 6.502e-03  0.470 0.638492  
## CLM_FREQ     -4.333e+01 4.822e+01 -0.899 0.368906  
## MVR_PTS      5.404e+01 2.277e+01  2.374 0.017641 *  
## CAR_AGE      -2.526e+01 1.118e+01 -2.260 0.023843 *  
## SingleParent 1.401e+02 1.767e+02  0.793 0.427820  
## Male         2.871e+02 1.605e+02  1.788 0.073756 .  
## Married      -1.711e+02 1.268e+02 -1.350 0.177187  
## HighSchool   -1.227e+02 1.502e+02 -0.817 0.414191  
## Bachelors    6.954e+01 1.791e+02  0.388 0.697744  
## Masters      2.263e+02 2.621e+02  0.863 0.387952  
## PHD          4.323e+02 3.110e+02  1.390 0.164631  
## Clerical     -5.782e+00 2.985e+02 -0.019 0.984544  
## Doctor       -2.783e+02 3.571e+02 -0.779 0.435838  
## HomeMaker    -6.523e+01 3.186e+02 -0.205 0.837805  
## Lawyer        7.743e+01 2.583e+02  0.300 0.764368  
## Manager      -1.210e+02 2.521e+02 -0.480 0.631160  
## Professional 1.762e+02 2.698e+02  0.653 0.513667  
## Student      -1.263e+02 3.268e+02 -0.387 0.699045  
## BlueCollar   5.726e+01 2.813e+02  0.204 0.838705  
## Commercial   9.624e+01 1.443e+02  0.667 0.504755  
## Minivan      -1.600e+02 1.571e+02 -1.019 0.308382  
## PanelTruck   -2.123e+02 2.976e+02 -0.713 0.475594
```

```

## Pickup      -1.903e+02  1.740e+02  -1.094  0.274059
## SportsCar   4.487e+01  1.563e+02   0.287  0.774114
## Van        -6.419e+01  2.334e+02  -0.275  0.783290
## RedCar     -2.855e+01  1.302e+02  -0.219  0.826416
## DLRevoked  -3.273e+02  1.526e+02  -2.145  0.031984 *
## Urban      -3.380e+01  1.264e+02  -0.267  0.789103
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3970 on 8121 degrees of freedom
## Multiple R-squared:  0.2913, Adjusted R-squared:  0.2879
## F-statistic: 85.59 on 39 and 8121 DF,  p-value: < 2.2e-16

```

Model 2

MLR: Backward Elimination

```

##
## Call:
## lm(formula = TARGET_AMT ~ KIDSDRV + HOMEKIDS + INCOME + SingleParent +
##      HOME_VAL + Married + TRAVTIME + BLUEBOOK + TIF + CLM_FREQ +
##      MVR PTS + CAR AGE, data = cleandf, na.action = na.omit)
##
## Residuals:
##    Min     1Q Median     3Q    Max
## -5441  -1601   -887    -81 104296
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.377e+03  2.119e+02   6.499 8.54e-11 ***
## KIDSDRV     2.899e+02  1.129e+02   2.568 0.010241 *
## HOMEKIDS    6.561e+01  6.012e+01   1.091 0.275139
## INCOME      -3.171e-03  1.559e-03  -2.034 0.042013 *
## SingleParent 5.607e+02  2.036e+02   2.754 0.005903 **
## HOME_VAL    -4.514e-04  5.841e-04  -0.773 0.439702
## Married     -5.567e+02  1.454e+02  -3.830 0.000129 ***
## TRAVTIME    7.282e+00  3.215e+00   2.265 0.023554 *
## BLUEBOOK    1.632e-02  6.692e-03   2.438 0.014770 *
## TIF         -4.564e+01  1.232e+01  -3.704 0.000214 ***
## CLM_FREQ    2.708e+02  4.814e+01   5.627 1.90e-08 ***
## MVR PTS     2.165e+02  2.600e+01   8.327 < 2e-16 ***
## CAR AGE     -3.422e+01  1.010e+01  -3.388 0.000709 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4611 on 8148 degrees of freedom
## Multiple R-squared:  0.04059, Adjusted R-squared:  0.03918
## F-statistic: 28.73 on 12 and 8148 DF,  p-value: < 2.2e-16

```

Model 3

MLR: BoxCox Transformation

```
##  
## Call:  
## lm(formula = TARGET_AMT ~ ., data = model_bc_transformed)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -6281    -473     -70     240 101112  
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -1.049e+03 5.020e+02 -2.090 0.0366 *  
## INDEX        -2.855e-02 1.799e-01 -0.159 0.8739  
## TARGET_FLAG   5.712e+03 1.136e+02 50.265 < 2e-16 ***  
## KIDSDRV      -3.309e+01 9.907e+01 -0.334 0.7384  
## AGE           5.658e+00 6.170e+00 0.917 0.3591  
## HOMEKIDS     4.026e+01 5.708e+01 0.705 0.4806  
## YOJ           7.698e+00 1.318e+01 0.584 0.5592  
## INCOME        -2.182e-03 1.570e-03 -1.389 0.1648  
## HOME_VAL      4.100e-04 5.163e-04 0.794 0.4271  
## TRAVTIME      1.132e+00 7.904e+00 0.143 0.8861  
## BLUEBOOK      3.914e+00 8.966e-01 4.366 1.28e-05 ***  
## TIF            -1.001e+01 3.648e+01 -0.274 0.7837  
## OLDCLAIM      2.973e-03 6.501e-03 0.457 0.6475  
## CLM_FREQ      -4.300e+01 4.821e+01 -0.892 0.3725  
## MVR PTS       5.407e+01 2.276e+01 2.376 0.0175 *  
## CAR AGE       -2.530e+01 1.118e+01 -2.264 0.0236 *  
## SingleParent  1.404e+02 1.766e+02 0.795 0.4268  
## Male          2.973e+02 1.592e+02 1.867 0.0619 .  
## Married       -1.694e+02 1.268e+02 -1.337 0.1813  
## HighSchool    -1.280e+02 1.502e+02 -0.852 0.3942  
## Bachelors     5.995e+01 1.791e+02 0.335 0.7378  
## Masters        2.143e+02 2.621e+02 0.818 0.4136  
## PHD            4.295e+02 3.110e+02 1.381 0.1673  
## Clerical       -8.202e+00 2.984e+02 -0.027 0.9781  
## Doctor         -2.881e+02 3.570e+02 -0.807 0.4197  
## HomeMaker     -5.052e+01 3.186e+02 -0.159 0.8740  
## Lawyer          7.012e+01 2.582e+02 0.272 0.7860  
## Manager         -1.294e+02 2.521e+02 -0.513 0.6077  
## Professional   1.685e+02 2.697e+02 0.625 0.5321  
## Student         -1.074e+02 3.268e+02 -0.329 0.7425  
## BlueCollar     4.794e+01 2.813e+02 0.170 0.8647  
## Commercial     9.375e+01 1.442e+02 0.650 0.5157  
## Minivan        -1.731e+02 1.558e+02 -1.111 0.2665  
## PanelTruck    -2.116e+02 2.880e+02 -0.735 0.4624  
## Pickup          -1.920e+02 1.730e+02 -1.110 0.2672  
## SportsCar       6.424e+01 1.565e+02 0.411 0.6814  
## Van             -9.857e+01 2.320e+02 -0.425 0.6709  
## RedCar          -2.673e+01 1.302e+02 -0.205 0.8373  
## DLRevoked      -3.258e+02 1.525e+02 -2.136 0.0327 *  
## Urban           -3.716e+01 1.262e+02 -0.294 0.7684
```

```

## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3969 on 8121 degrees of freedom
## Multiple R-squared: 0.2916, Adjusted R-squared: 0.2882
## F-statistic: 85.73 on 39 and 8121 DF, p-value: < 2.2e-16

```

Model 4

BLR: Base model where all variables are tested

```

##
## Call:
## glm(formula = TARGET_FLAG ~ ., family = binomial(link = "logit"),
##      data = reduced_cleandf)
##
## Deviance Residuals:
##    Min      1Q   Median      3Q     Max
## -2.5846 -0.7127 -0.3983  0.6261  3.1524
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.895e+00 3.181e-01 -9.103 < 2e-16 ***
## KIDSDRV      3.862e-01 6.122e-02  6.308 2.82e-10 ***
## AGE         -1.015e-03 4.020e-03 -0.252 0.800672
## HOMEKIDS    4.965e-02 3.713e-02  1.337 0.181119
## YOJ        -1.105e-02 8.582e-03 -1.288 0.197743
## INCOME     -3.423e-06 1.081e-06 -3.165 0.001551 **
## HOME_VAL    -1.306e-06 3.420e-07 -3.819 0.000134 ***
## TRAVTIME    1.457e-02 1.883e-03  7.736 1.03e-14 ***
## BLUEBOOK   -2.084e-05 5.263e-06 -3.959 7.52e-05 ***
## TIF        -5.547e-02 7.344e-03 -7.553 4.26e-14 ***
## OLDCLAIM   -1.389e-05 3.910e-06 -3.554 0.000380 ***
## CLM_FREQ    1.959e-01 2.855e-02  6.864 6.69e-12 ***
## MVR_PTS     1.133e-01 1.361e-02  8.324 < 2e-16 ***
## CAR_AGE     -7.196e-04 7.549e-03 -0.095 0.924053
## SingleParent 3.820e-01 1.096e-01  3.485 0.000492 ***
## Male        8.251e-02 1.120e-01  0.737 0.461416
## Married     -4.938e-01 8.357e-02 -5.909 3.45e-09 ***
## HighSchool   1.764e-02 9.506e-02  0.186 0.852802
## Bachelors   -3.812e-01 1.157e-01 -3.296 0.000981 ***
## Masters     -2.903e-01 1.788e-01 -1.624 0.104397
## PHD        -1.677e-01 2.140e-01 -0.784 0.433295
## Clerical    4.107e-01 1.967e-01  2.088 0.036763 *
## Doctor      -4.458e-01 2.671e-01 -1.669 0.095106 .
## HomeMaker   2.323e-01 2.102e-01  1.106 0.268915
## Lawyer      1.049e-01 1.695e-01  0.619 0.535958
## Manager     -5.572e-01 1.716e-01 -3.248 0.001161 **
## Professional 1.619e-01 1.784e-01  0.907 0.364168
## Student     2.161e-01 2.145e-01  1.007 0.313729
## BlueCollar  3.106e-01 1.856e-01  1.674 0.094158 .

```

```

## Commercial    7.564e-01  9.172e-02   8.247 < 2e-16 ***
## Minivan      -7.682e-01  1.113e-01  -6.904 5.05e-12 ***
## PanelTruck   -2.074e-01  2.002e-01  -1.036 0.300239
## Pickup        -2.142e-01  1.166e-01  -1.838 0.066133 .
## SportsCar     2.570e-01  9.815e-02   2.618 0.008840 **
## Van           -1.496e-01  1.589e-01  -0.942 0.346418
## RedCar        -9.728e-03  8.636e-02  -0.113 0.910313
## DLRevoked     8.874e-01  9.133e-02   9.716 < 2e-16 ***
## Urban         2.390e+00  1.128e-01  21.181 < 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: 7297.6 on 8123 degrees of freedom
## AIC: 7373.6
##
## Number of Fisher Scoring iterations: 5

```

Model 5

BLR: drop term model

The `dropterm()` function from the MASS package tests all models from the current model selected by fitting all models that differ from the current model by dropping a single term, maintaining marginality.

```

##
## Call:
## glm(formula = TARGET_FLAG ~ KIDSDRV + AGE + HOMEKIDS + YOJ +
##       INCOME + HOME_VAL + TRAVTIME + BLUEBOOK + TIF + OLDCLAIM +
##       CLM_FREQ + MVR PTS + CAR_AGE + SingleParent + Male + Married +
##       HighSchool + Bachelors + Masters + PHD + Clerical + Doctor +
##       HomeMaker + Lawyer + Manager + Professional + Student + BlueCollar +
##       Commercial + Minivan + PanelTruck + Pickup + SportsCar +
##       Van + RedCar + DLRevoked + Urban, family = binomial(link = "logit"),
##       data = cleandf)
##
## Deviance Residuals:
##      Min        1Q     Median        3Q       Max
## -2.5846  -0.7127  -0.3983   0.6261   3.1524
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.895e+00 3.181e-01 -9.103 < 2e-16 ***
## KIDSDRV      3.862e-01 6.122e-02  6.308 2.82e-10 ***
## AGE          -1.015e-03 4.020e-03 -0.252 0.800672
## HOMEKIDS     4.965e-02 3.713e-02  1.337 0.181119
## YOJ          -1.105e-02 8.582e-03 -1.288 0.197743
## INCOME       -3.423e-06 1.081e-06 -3.165 0.001551 **
## HOME_VAL     -1.306e-06 3.420e-07 -3.819 0.000134 ***

```

```

## TRAVTIME      1.457e-02  1.883e-03   7.736 1.03e-14 ***
## BLUEBOOK     -2.084e-05  5.263e-06  -3.959 7.52e-05 ***
## TIF          -5.547e-02  7.344e-03  -7.553 4.26e-14 ***
## OLDCLAIM    -1.389e-05  3.910e-06  -3.554 0.000380 ***
## CLM_FREQ     1.959e-01  2.855e-02   6.864 6.69e-12 ***
## MVR_PTS      1.133e-01  1.361e-02   8.324 < 2e-16 ***
## CAR_AGE     -7.196e-04  7.549e-03  -0.095 0.924053
## SingleParent 3.820e-01  1.096e-01   3.485 0.000492 ***
## Male          8.251e-02  1.120e-01   0.737 0.461416
## Married      -4.938e-01  8.357e-02  -5.909 3.45e-09 ***
## HighSchool   1.764e-02  9.506e-02   0.186 0.852802
## Bachelors    -3.812e-01  1.157e-01  -3.296 0.000981 ***
## Masters       -2.903e-01  1.788e-01  -1.624 0.104397
## PHD          -1.677e-01  2.140e-01  -0.784 0.433295
## Clerical      4.107e-01  1.967e-01   2.088 0.036763 *
## Doctor        -4.458e-01  2.671e-01  -1.669 0.095106 .
## HomeMaker    2.323e-01  2.102e-01   1.106 0.268915
## Lawyer         1.049e-01  1.695e-01   0.619 0.535958
## Manager       -5.572e-01  1.716e-01  -3.248 0.001161 **
## Professional  1.619e-01  1.784e-01   0.907 0.364168
## Student       2.161e-01  2.145e-01   1.007 0.313729
## BlueCollar   3.106e-01  1.856e-01   1.674 0.094158 .
## Commercial   7.564e-01  9.172e-02   8.247 < 2e-16 ***
## Minivan       -7.682e-01  1.113e-01  -6.904 5.05e-12 ***
## PanelTruck   -2.074e-01  2.002e-01  -1.036 0.300239
## Pickup         -2.142e-01  1.166e-01  -1.838 0.066133 .
## SportsCar     2.570e-01  9.815e-02   2.618 0.008840 **
## Van           -1.496e-01  1.589e-01  -0.942 0.346418
## RedCar        -9.728e-03  8.636e-02  -0.113 0.910313
## DLRevoked    8.874e-01  9.133e-02   9.716 < 2e-16 ***
## Urban         2.390e+00  1.128e-01  21.181 < 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: 7297.6 on 8123 degrees of freedom
## AIC: 7373.6
##
## Number of Fisher Scoring iterations: 5

```

Model 6

BLR: Backward Elimination

```

##
## Call:
## glm(formula = TARGET_FLAG ~ KIDSDRV + AGE + HOMEKIDS + YOJ +
##      HOME_VAL + TRAVTIME + BLUEBOOK + TIF + OLDCLAIM + CLM_FREQ +
##      CAR_AGE + SingleParent + Male + Married + Bachelors + Clerical +

```

```

## Doctor + HomeMaker + Lawyer + Manager + Student + BlueCollar +
## Commercial + Minivan + PanelTruck + Pickup + SportsCar +
## Van + RedCar + Urban, family = binomial(link = "logit"),
## data = cleandf)
##
## Deviance Residuals:
##      Min     1Q Median     3Q    Max
## -2.6465 -0.7290 -0.4128  0.7133  3.1414
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.644e+00 2.617e-01 -10.104 < 2e-16 ***
## KIDSDRV      4.111e-01 6.003e-02   6.848 7.49e-12 ***
## AGE          -1.873e-03 3.957e-03  -0.473 0.635973
## HOMEKIDS     5.970e-02 3.646e-02   1.638 0.101525
## YOJ          -1.619e-02 8.412e-03  -1.924 0.054317 .
## HOME_VAL     -1.880e-06 3.130e-07  -6.006 1.90e-09 ***
## TRAVTIME     1.461e-02 1.858e-03   7.866 3.66e-15 ***
## BLUEBOOK    -2.337e-05 5.117e-06  -4.568 4.92e-06 ***
## TIF          -5.839e-02 7.235e-03  -8.071 6.97e-16 ***
## OLDCLAIM     5.892e-06 3.315e-06   1.777 0.075547 .
## CLM_FREQ     2.013e-01 2.625e-02   7.670 1.72e-14 ***
## CAR_AGE      -1.453e-02 6.338e-03  -2.293 0.021845 *
## SingleParent 4.113e-01 1.078e-01   3.816 0.000135 ***
## Male          9.082e-02 1.104e-01   0.823 0.410687
## Married      -4.177e-01 8.081e-02  -5.169 2.35e-07 ***
## Bachelors    -2.532e-01 6.944e-02  -3.646 0.000266 ***
## Clerical      4.675e-01 1.083e-01   4.316 1.59e-05 ***
## Doctor        -6.710e-01 2.261e-01  -2.968 0.002996 **
## HomeMaker     2.370e-01 1.404e-01   1.688 0.091338 .
## Lawyer        -7.999e-02 1.289e-01  -0.621 0.534909
## Manager       -7.195e-01 1.180e-01  -6.100 1.06e-09 ***
## Student       2.893e-01 1.370e-01   2.112 0.034714 *
## BlueCollar    3.194e-01 9.985e-02   3.199 0.001380 **
## Commercial   7.415e-01 8.445e-02   8.780 < 2e-16 ***
## Minivan       -8.067e-01 1.096e-01  -7.359 1.85e-13 ***
## PanelTruck   -3.070e-01 1.951e-01  -1.573 0.115635
## Pickup         -2.320e-01 1.138e-01  -2.039 0.041499 *
## SportsCar     2.370e-01 9.677e-02   2.449 0.014306 *
## Van           -1.949e-01 1.558e-01  -1.252 0.210745
## RedCar        -5.288e-03 8.510e-02  -0.062 0.950450
## Urban         2.444e+00 1.118e-01   21.862 < 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: 7479.3 on 8130 degrees of freedom
## AIC: 7541.3
##
## Number of Fisher Scoring iterations: 5

```

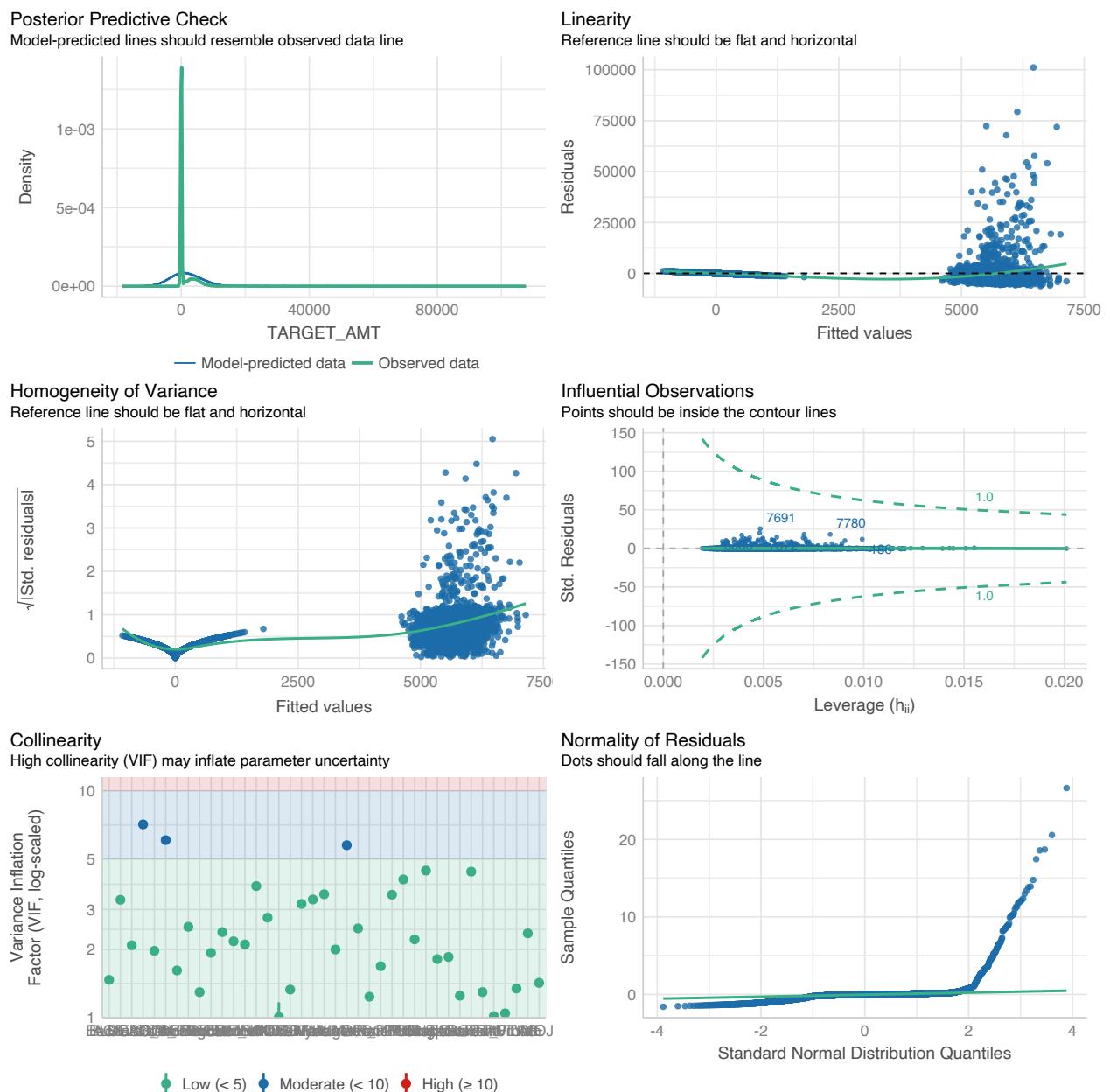
Select Models:

Lets examine our linear models

- **Model 1**

We can see that there may be some moderate collinearity in 3 of our predictors which may lead to instability in our coefficient estimates

At the right tail our residuals are not normally distributed

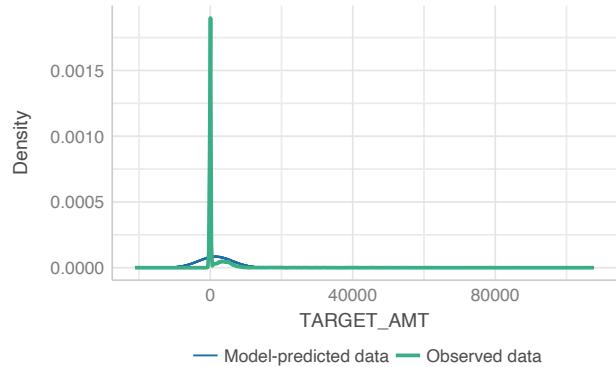


- **Model 2**

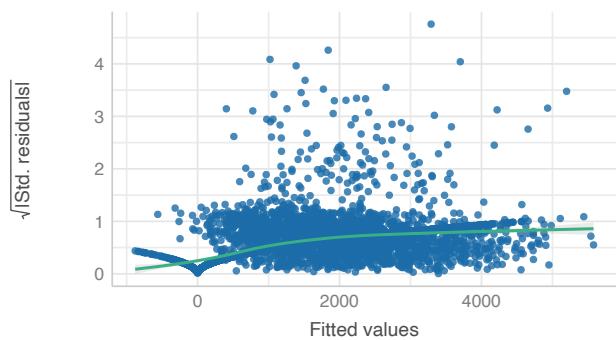
There are no collinearity issues with the second model

At the right tail our residuals are not normally distributed

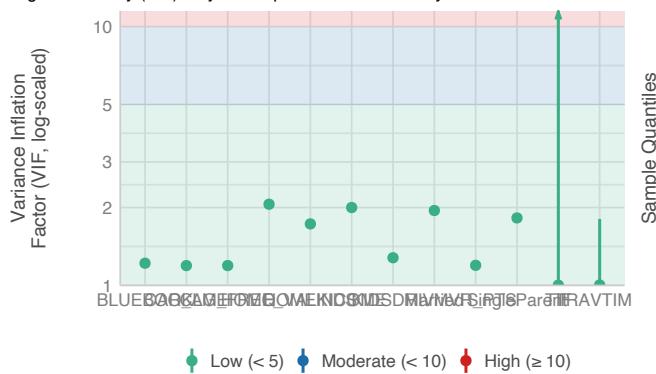
Posterior Predictive Check
Model-predicted lines should resemble observed data line



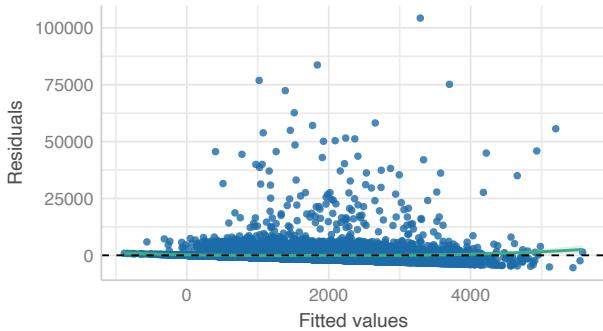
Homogeneity of Variance
Reference line should be flat and horizontal



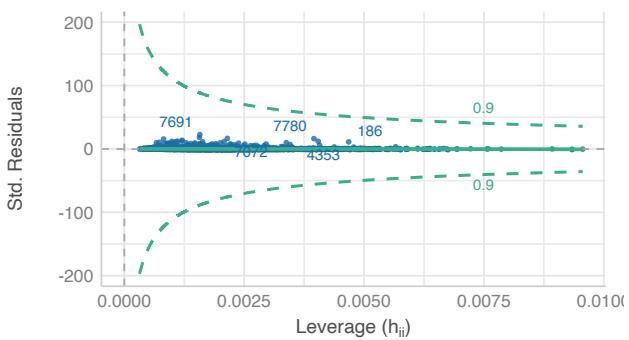
Collinearity
High collinearity (VIF) may inflate parameter uncertainty



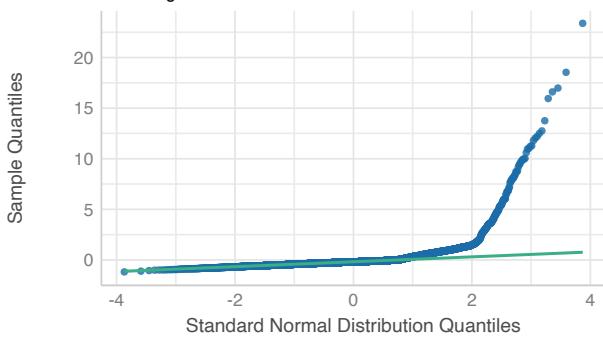
Linearity
Reference line should be flat and horizontal



Influential Observations
Points should be inside the contour lines



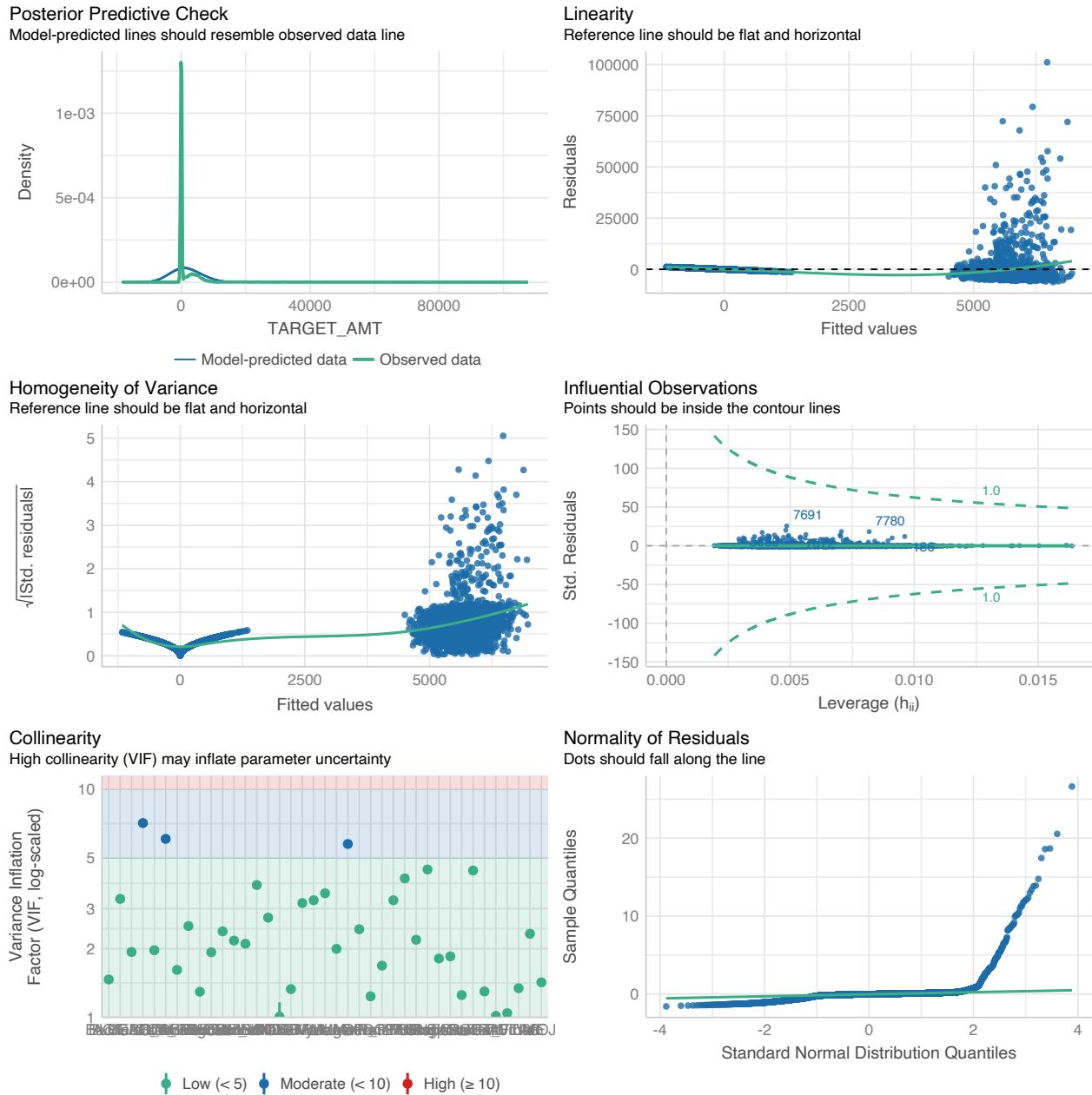
Normality of Residuals
Dots should fall along the line



- **Model 3**

We can see that there may be some moderate collinearity in 3 of our predictors which may lead to instability in our coefficient estimates

At the right tail our residuals are not normally distributed



Let's evaluate the performance of our linear models

Model 1 and Model 3 have identical adjusted R-square. We are going to use Model 3 that performs as well as Model 1 but has a Box Cox Transformation to improve the normality of the distribution of the dependent variable.

- Model 1

```
## # Indices of model performance
##
## AIC | BIC | R2 | R2 (adj.) | RMSE | Sigma
## -----
## 1.585e+05 | 1.587e+05 | 0.291 | 0.288 | 3959.793 | 3969.533
```

- Model 2

```
## # Indices of model performance
##
## AIC | BIC | R2 | R2 (adj.) | RMSE | Sigma
## -----
## 1.609e+05 | 1.610e+05 | 0.041 | 0.039 | 4607.281 | 4610.955
```

- Model 3

```
## # Indices of model performance
##
## AIC | BIC | R2 | R2 (adj.) | RMSE | Sigma
## -----
## 1.584e+05 | 1.587e+05 | 0.292 | 0.288 | 3958.842 | 3968.580
```

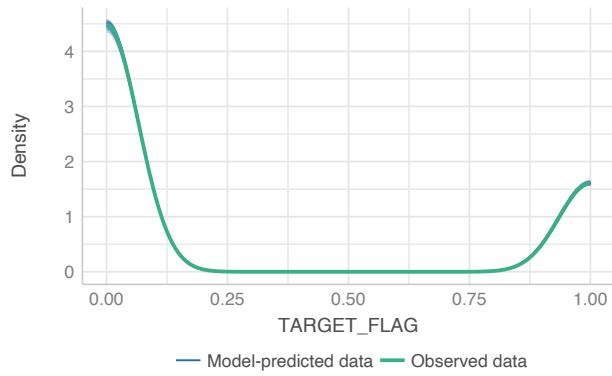
Lets examine our binary logistic models

- **Model 4**

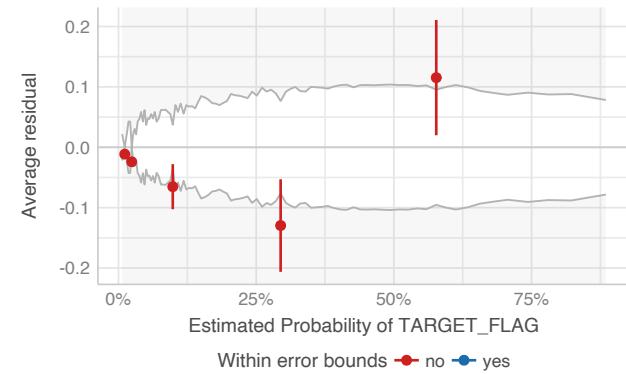
We can see that there may be some moderate collinearity in 3 of our predictors which may lead to instability in our coefficient estimates

At the right tail our residuals are not normally distributed

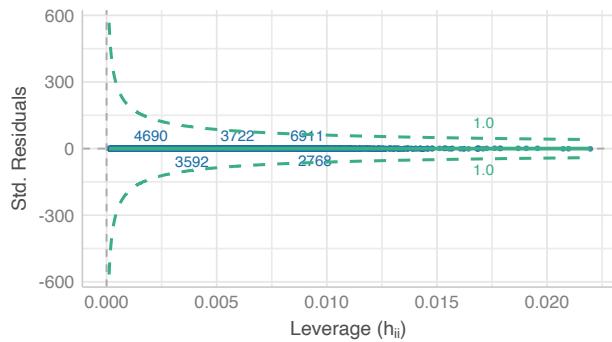
Posterior Predictive Check
Model-predicted lines should resemble observed data line



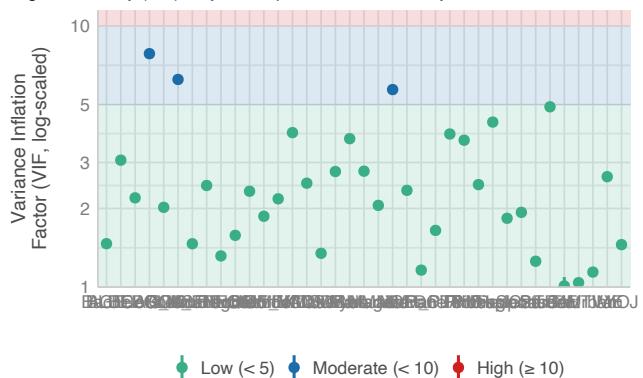
Binned Residuals
Points should be within error bounds



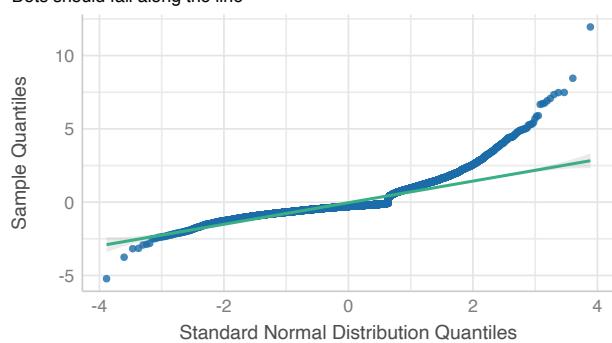
Influential Observations
Points should be inside the contour lines



Collinearity
High collinearity (VIF) may inflate parameter uncertainty



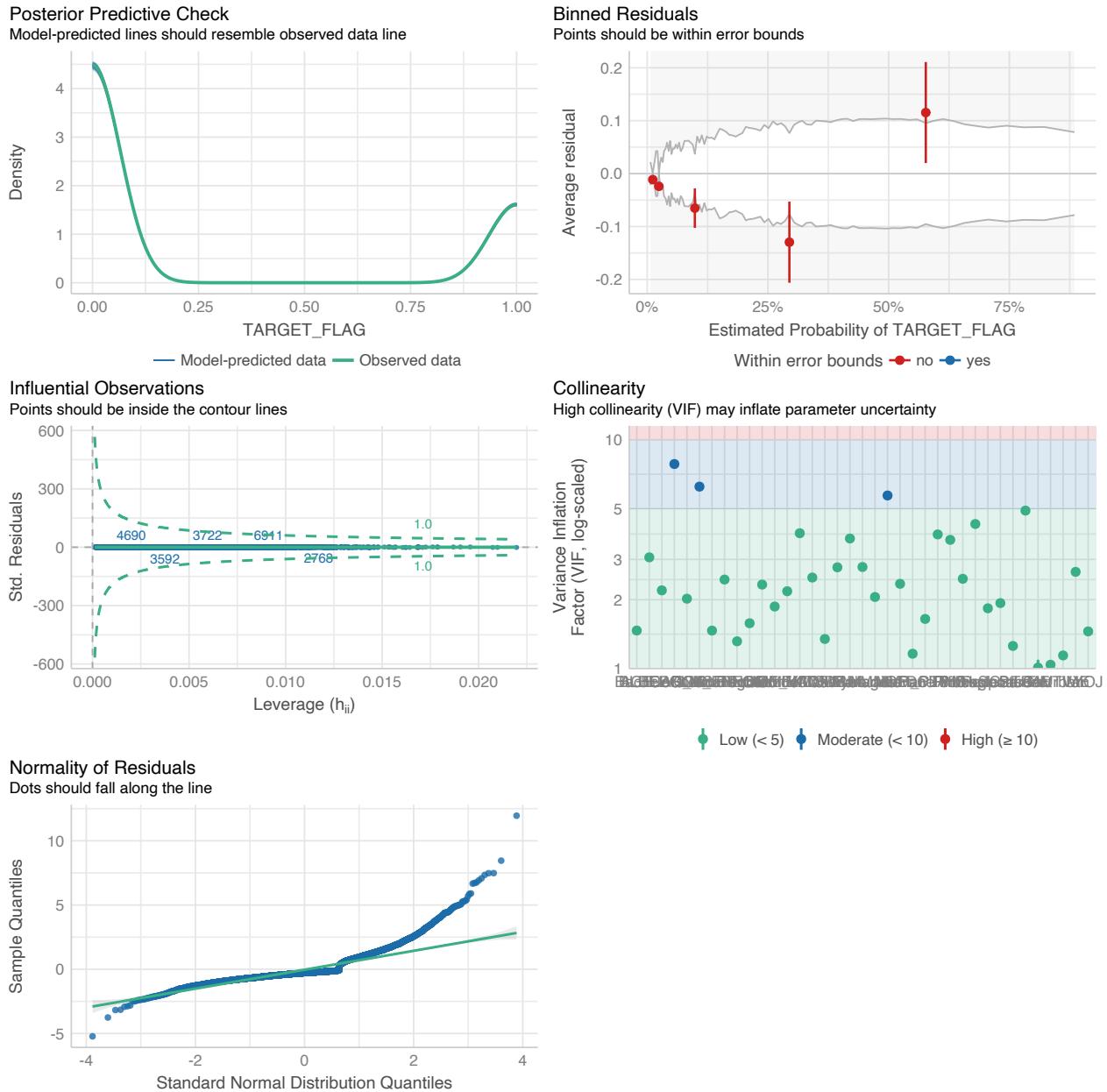
Normality of Residuals
Dots should fall along the line



- **Model 5**

We can see that there may be some moderate collinearity in 3 of our predictors which may lead to instability in our coefficient estimates

At the right tail our residuals are not normally distributed



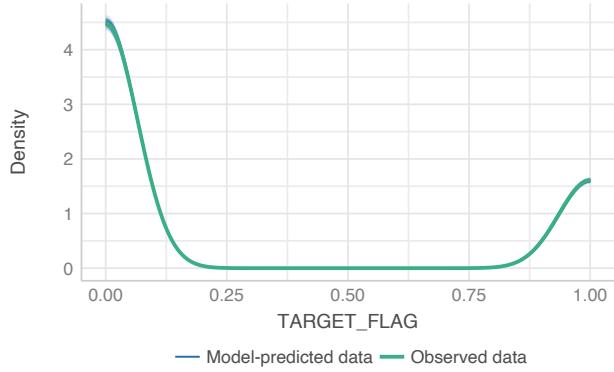
- **Model 6**

There is no collinearity issues in Model 6

Like all our models, at the right tail our residuals are not normally distributed

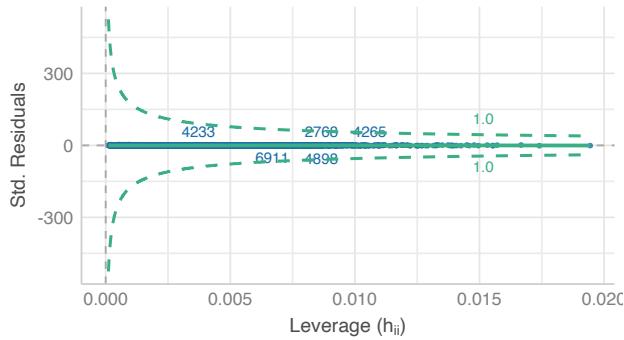
Posterior Predictive Check

Model-predicted lines should resemble observed data line



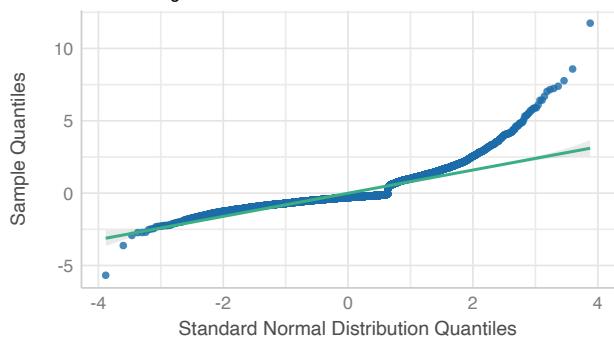
Influential Observations

Points should be inside the contour lines



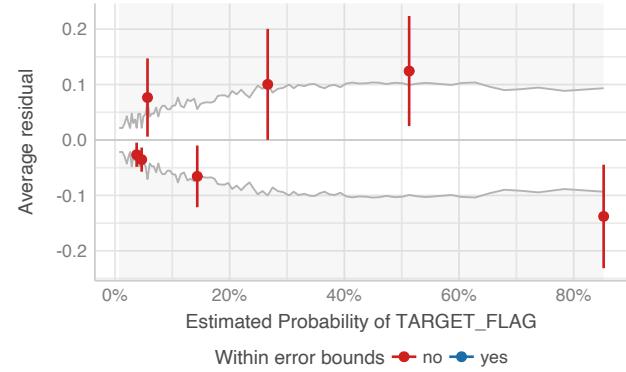
Normality of Residuals

Dots should fall along the line



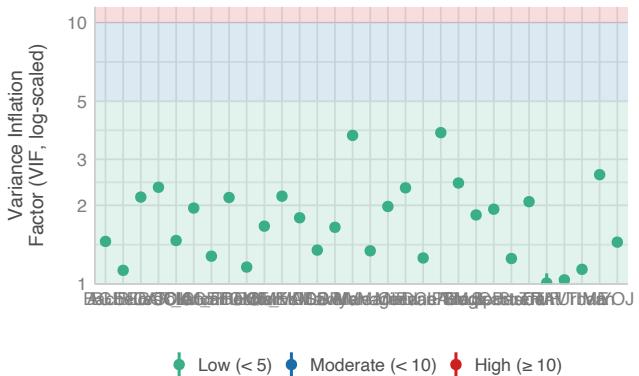
Binned Residuals

Points should be within error bounds



Collinearity

High collinearity (VIF) may inflate parameter uncertainty



Let's evaluate the performance of our binary logistic models

- Model 4

```
## # Indices of model performance
##
## AIC      |     BIC | Tjur's R2 |   RMSE | Sigma | Log_loss | Score_log | Score_spherical |   PCP
## -----
## 7373.640 | 7639.910 |     0.253 | 0.381 | 0.948 |     0.447 |      -Inf | 3.687e-04 | 0.710

## Setting levels: control = 0, case = 1

## Setting direction: controls < cases

## Area under the curve: 0.6752

##      pred
## true 0 1
## 0 5551 457
## 1 1235 918

Accuracy:  $\frac{5551+918}{8161} = 79\%$ 
Classification Error Rate:  $\frac{1235+457}{8161} = 20\%$ 
Precision:  $\frac{918}{918+457} = 67$ 
Sensitivity:  $\frac{918}{918+1235} = 43\%$ 
Specificity:  $\frac{5551}{5551+457} = 92\%$ 
F1 Score:  $\frac{2*918}{2*918+457+1235} = 52\%$ 
```

- Model 5

```
## # Indices of model performance
##
## AIC      |     BIC | Tjur's R2 |   RMSE | Sigma | Log_loss | Score_log | Score_spherical |   PCP
## -----
## 7373.640 | 7639.910 |     0.253 | 0.381 | 0.948 |     0.447 |      -Inf | 3.687e-04 | 0.710

## Setting levels: control = 0, case = 1

## Setting direction: controls < cases

## Area under the curve: 0.6752

##      pred
## true 0 1
## 0 5551 457
## 1 1235 918
```

Accuracy: $\frac{5551+918}{8161} = 79\%$
 Classification Error Rate: $\frac{1235+457}{8161} = 20\%$
 Precision: $\frac{918}{918+457} = 67$
 Sensitivity: $\frac{918}{918+1235} = 43\%$
 Specificity: $\frac{5551}{5551+457} = 92\%$
 F1 Score: $\frac{2*918}{2*918+457+1235} = 52\%$

- Model 6

```

## # Indices of model performance
##
## AIC      |      BIC | Tjur's R2 |   RMSE | Sigma | Log_loss | Score_log | Score_spherical |   PCP
## -----
## 7541.301 | 7758.522 |      0.229 | 0.387 | 0.959 |     0.458 |      -Inf | 3.184e-04 | 0.700

## Setting levels: control = 0, case = 1

## Setting direction: controls < cases

## Area under the curve: 0.6576

##      pred
## true 0 1
## 0 5535 473
## 1 1305 848

Accuracy:  $\frac{5535+848}{8161} = 78\%$   

Classification Error Rate:  $\frac{1305+473}{8161} = 22\%$   

Precision:  $\frac{848}{848+473} = 64$   

Sensitivity:  $\frac{848}{848+1305} = 39\%$   

Specificity:  $\frac{5535}{5535+473} = 92\%$   

F1 Score:  $\frac{2*848}{2*848+473+1305} = 49\%$ 

```

The AUC for all 3 of our models are very similar, we are going to chose Model 6 due to the lack of colinearity issues.

Predictions:

Lets make predictions with our two chosen models (Linear Model 3 and Logistic Model 6):

We are able to apply the coefficients of our model to novel data to make predictions on accident incidence.

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

Now that we have predictions for accidents, we can make predictions on the recovery amounts. Anywhere no accident was predicted, the recovery amount should be set to zero.

```
## [1] 85131.01 73240.93 22207.27 35103.38 59627.88 99700.80 43421.68
## [8] 93365.11 105495.15 132895.94 90836.37 55693.33 41552.42 40966.35
## [15] 56288.83 37558.53 11026.17 89936.38 34697.37 41472.62

## [1] 0.00 0.00 0.00 0.00 0.00 0.00 0.00
## [9] 0.00 0.00 0.00 0.00 41552.42 0.00 0.00
## [17] 11026.17 0.00 0.00
```

Appendix:

Below is all the code used in this homework

```
# load libraries
library(tidyverse)
library(caret)
library(pROC)
library(corrplot)
library(GGally)
library(psych)
library(car)
library(kableExtra)
library(gridExtra)
library(performance)
library(faraway)
library(jtools)
library(DataExplorer)
library(hrbrthemes)
library(MASS)

# load data
dftrain <- read.csv("https://raw.githubusercontent.com/letisalba/Data_621/master/
Homework_4/csv/insurance_training_data.csv")
dfeval <- read.csv("https://raw.githubusercontent.com/letisalba/Data_621/master/
Homework_4/csv/insurance_training_data.csv")
glimpse(dftrain)

# summary of training data
summary(dftrain)

# bar plots for training set
DataExplorer::plot_bar(data = dftrain, order_bar = T, ggtheme = theme_ipsum())

# histograms for training set
DataExplorer::plot_histogram(geom_histogram_args = list(alpha = 0.5),
data = dftrain, ggtheme = theme_ipsum())

# summary of evaluation data
summary(dfeval)

# bar plots for evaluation set
DataExplorer::plot_bar(data = dfeval, order_bar = T, ggtheme = theme_ipsum())

# histograms for evaluation set
DataExplorer::plot_histogram(geom_histogram_args = list(alpha = 0.5),
data = dfeval, ggtheme = theme_ipsum())

# creating dummy variables for some variables in training
# and evaluation data sets
dftrain %>%
  count(PARENT1)
dfeval %>%
  count(PARENT1)
```

```

dftrain %>%
  count(SEX)
dfeval %>%
  count(SEX)

dftrain %>%
  count(MSTATUS)
dfeval %>%
  count(MSTATUS)

dftrain %>%
  count(EDUCATION)
dfeval %>%
  count(EDUCATION)

dftrain %>%
  count(JOB)
dfeval %>%
  count(JOB)

dftrain %>%
  count(CAR_USE)
dfeval %>%
  count(CAR_USE)

dftrain %>%
  count(CAR_TYPE)
dfeval %>%
  count(CAR_TYPE)

dftrain %>%
  count(RED_CAR)
dfeval %>%
  count(RED_CAR)

dftrain %>%
  count(REVOKED)
dfeval %>%
  count(REVOKED)

dftrain %>%
  count(URBANICITY)
dfeval %>%
  count(URBANICITY)

# checking for missing values in training data
colSums(is.na(dftrain))

# Look for blank values in columns that will be converted
# to numeric, which would introduce NAs
dollar_cols <- c("BLUEBOOK", "HOME_VAL", "INCOME", "OLDCLAIM")
ctVals <- c()
for (dollar_col in dollar_cols) {

```

```

    ctVals <- c(ctVals, sum(str_length(dftrain[, dollar_col]) ==
      0))
  }
ctVals <- data.frame(t(ctVals))
colnames(ctVals) <- dollar_cols
ctVals

# data cleaning functions for training data
clean_df <- function(df) {

  df$BLUEBOOK <- as.numeric(gsub("[,$]", "", df$BLUEBOOK))
  df$HOME_VAL <- as.numeric(gsub("[,$]", "", df$HOME_VAL))
  df$INCOME <- as.numeric(gsub("[,$]", "", df$INCOME))
  df$OLDCLAIM <- as.numeric(gsub("[,$]", "", df$OLDCLAIM))

  df <- df %>%
    mutate(SingleParent = ifelse(PARENT1 == "Yes", 1, 0)) %>%
    dplyr::select(-PARENT1) %>%
    mutate(Male = ifelse(SEX == "M", 1, 0)) %>%
    dplyr::select(-SEX) %>%
    mutate(Married = ifelse(MSTATUS == "Yes", 1, 0)) %>%
    dplyr::select(-MSTATUS) %>%
    mutate(HighSchool = ifelse(EDUCATION == "z_High School",
      1, 0)) %>%
    mutate(Bachelors = ifelse(EDUCATION == "Bachelors", 1,
      0)) %>%
    mutate(Masters = ifelse(EDUCATION == "Masters", 1, 0)) %>%
    mutate(PhD = ifelse(EDUCATION == "PhD", 1, 0)) %>%
    dplyr::select(-EDUCATION) %>%
    mutate(Clerical = ifelse(JOB == "Clerical", 1, 0)) %>%
    mutate(Doctor = ifelse(JOB == "Doctor", 1, 0)) %>%
    mutate(HomeMaker = ifelse(JOB == "Home Maker", 1, 0)) %>%
    mutate(Lawyer = ifelse(JOB == "Lawyer", 1, 0)) %>%
    mutate(Manager = ifelse(JOB == "Manager", 1, 0)) %>%
    mutate(Professional = ifelse(JOB == "Professional", 1,
      0)) %%%
    mutate(Student = ifelse(JOB == "Student", 1, 0)) %>%
    mutate(BlueCollar = ifelse(JOB == "z_Blue Collar", 1,
      0)) %>%
    dplyr::select(-JOB) %>%
    mutate(Commercial = ifelse(CAR_USE == "Commercial", 1,
      0)) %>%
    dplyr::select(-CAR_USE) %>%
    mutate(Minivan = ifelse(CAR_TYPE == "Minivan", 1, 0)) %>%
    mutate(PanelTruck = ifelse(CAR_TYPE == "Panel Truck",
      1, 0)) %>%
    mutate(Pickup = ifelse(CAR_TYPE == "Pickup", 1, 0)) %>%
    mutate(SportsCar = ifelse(CAR_TYPE == "Sports Car", 1,
      0)) %>%
    mutate(Van = ifelse(CAR_TYPE == "Van", 1, 0)) %>%
    dplyr::select(-CAR_TYPE) %>%
    mutate(REDCar = ifelse(RED_CAR == "yes", 1, 0)) %>%
    dplyr::select(-RED_CAR) %>%

```

```

    mutate(DLRevoked = ifelse(REVOKED == "Yes", 1, 0)) %>%
  dplyr::select(-REVOKED) %>%
  mutate(Urban = ifelse(URBANICITY == "Highly Urban/ Urban",
    1, 0)) %>%
  dplyr::select(-URBANICITY)

# Change negative values of car age to NA
df$CAR_AGE <- ifelse(df$CAR_AGE < 0, NA, df$CAR_AGE)

# Handle NAs
df$AGE <- ifelse(is.na(df$AGE), median(df$AGE, na.rm = T),
  df$AGE)
df$Y0J <- ifelse(is.na(df$Y0J), median(df$Y0J, na.rm = T),
  df$Y0J)
df$INCOME <- ifelse(is.na(df$INCOME), median(df$INCOME, na.rm = T),
  df$INCOME)
df$HOME_VAL <- ifelse(is.na(df$HOME_VAL), median(df$HOME_VAL,
  na.rm = T), df$HOME_VAL)
df$CAR_AGE <- ifelse(is.na(df$CAR_AGE), median(df$CAR_AGE,
  na.rm = T), df$CAR_AGE)

return(df)
}

cleandf <- clean_df(dftrain)
head(cleandf)

# checking for missing values in evaluation data
colSums(is.na(dfeval))

# Look for blank values in columns that will be converted
# to numeric, which would introduce NAs
dollar_cols <- c("BLUEBOOK", "HOME_VAL", "INCOME", "OLDCLAIM")
ctVals <- c()
for (dollar_col in dollar_cols) {
  ctVals <- c(ctVals, sum(str_length(dfeval[, dollar_col]) ==
    0))
}
ctVals <- data.frame(t(ctVals))
colnames(ctVals) <- dollar_cols
ctVals

# data cleaning functions for evaluation data
eval_clean_df <- function(df) {

  df$BLUEBOOK <- as.numeric(gsub("[,$]", "", df$BLUEBOOK))
  df$HOME_VAL <- as.numeric(gsub("[,$]", "", df$HOME_VAL))
  df$INCOME <- as.numeric(gsub("[,$]", "", df$INCOME))
  df$OLDCLAIM <- as.numeric(gsub("[,$]", "", df$OLDCLAIM))

  df <- df %>%
    mutate(SingleParent = ifelse(PARENT1 == "Yes", 1, 0)) %>%
    dplyr::select(-PARENT1) %>%

```

```

mutate(Male = ifelse(SEX == "M", 1, 0)) %>%
dplyr::select(-SEX) %>%
mutate(Married = ifelse(MSTATUS == "Yes", 1, 0)) %>%
dplyr::select(-MSTATUS) %>%
mutate(HighSchool = ifelse(EDUCATION == "z_High School",
  1, 0)) %>%
mutate(Bachelors = ifelse(EDUCATION == "Bachelors", 1,
  0)) %>%
mutate(Masters = ifelse(EDUCATION == "Masters", 1, 0)) %>%
mutate(PhD = ifelse(EDUCATION == "PhD", 1, 0)) %>%
dplyr::select(-EDUCATION) %>%
mutate(Clerical = ifelse(JOB == "Clerical", 1, 0)) %>%
mutate(Doctor = ifelse(JOB == "Doctor", 1, 0)) %>%
mutate(HomeMaker = ifelse(JOB == "Home Maker", 1, 0)) %>%
mutate(Lawyer = ifelse(JOB == "Lawyer", 1, 0)) %>%
mutate(Manager = ifelse(JOB == "Manager", 1, 0)) %>%
mutate(Professional = ifelse(JOB == "Professional", 1,
  0)) %>%
mutate(Student = ifelse(JOB == "Student", 1, 0)) %>%
mutate(BlueCollar = ifelse(JOB == "z_Blue Collar", 1,
  0)) %>%
dplyr::select(-JOB) %>%
mutate(Commercial = ifelse(CAR_USE == "Commercial", 1,
  0)) %>%
dplyr::select(-CAR_USE) %>%
mutate(Minivan = ifelse(CAR_TYPE == "Minivan", 1, 0)) %>%
mutate(PanelTruck = ifelse(CAR_TYPE == "Panel Truck",
  1, 0)) %>%
mutate(Pickup = ifelse(CAR_TYPE == "Pickup", 1, 0)) %>%
mutate(SportsCar = ifelse(CAR_TYPE == "Sports Car", 1,
  0)) %>%
mutate(Van = ifelse(CAR_TYPE == "Van", 1, 0)) %>%
dplyr::select(-CAR_TYPE) %>%
mutate(RedCar = ifelse(RED_CAR == "yes", 1, 0)) %>%
dplyr::select(-RED_CAR) %>%
mutate(DLRevoked = ifelse(REVOKED == "Yes", 1, 0)) %>%
dplyr::select(-REVOKED) %>%
mutate(Urban = ifelse(URBANICITY == "Highly Urban/ Urban",
  1, 0)) %>%
dplyr::select(-URBANICITY)

# Handle NAs
df$AGE <- ifelse(is.na(df$AGE), median(df$AGE, na.rm = T),
  df$AGE)
df$YOJ <- ifelse(is.na(df$YOJ), median(df$YOJ, na.rm = T),
  df$YOJ)
df$INCOME <- ifelse(is.na(df$INCOME), median(df$INCOME, na.rm = T),
  df$INCOME)
df$HOME_VAL <- ifelse(is.na(df$HOME_VAL), median(df$HOME_VAL,
  na.rm = T), df$HOME_VAL)
df$CAR_AGE <- ifelse(is.na(df$CAR_AGE), median(df$CAR_AGE,
  na.rm = T), df$CAR_AGE)

```

```

    return(df)
}

eval_cleandf <- eval_clean_df(dfeval)
head(eval_cleandf)

# correlation plot for training data
cor_res <- cor(cleandf, use = "na.or.complete")

corrplot(cor_res, type = "lower", order = "original", tl.col = "black",
         tl.srt = 50, tl.cex = 1)
cor_res <- data.frame(cor_res)

# model building for multiple linear regression
model1 <- lm(TARGET_AMT ~ ., cleandf, na.action = na.omit)
summary(model1)
# summ(model1)

model2 <- lm(TARGET_AMT ~ KIDSDRV + HOMEKIDS + INCOME + SingleParent +
              HOME_VAL + Married + TRAVTIME + BLUEBOOK + TIF + CLM_FREQ +
              MVR_PTS + CAR_AGE, cleandf, na.action = na.omit)
summary(model2)
# summ(model2)

model_boxcox <- preProcess(cleandf, c("BoxCox"))
model_bc_transformed <- predict(model_boxcox, cleandf)
model3 <- lm(TARGET_AMT ~ ., model_bc_transformed)
summary(model3)
# summ(model3)

# model building for binary logistic regression
model4 <- glm(TARGET_FLAG ~ ., cleandf, family = binomial(link = "logit"))
summary(model4)

dropterm(model4, test = "F")
model5 <- glm(TARGET_FLAG ~ INDEX + TARGET_AMT + KIDSDRV + AGE +
               HOMEKIDS + YOJ + INCOME + HOME_VAL + TRAVTIME + BLUEBOOK +
               TIF + OLDCLAIM + CLM_FREQ + MVR_PTS + CAR_AGE + SingleParent +
               Male + Married + HighSchool + Bachelors + Masters + PHD +
               Clerical + Doctor + HomeMaker + Lawyer + Manager + Professional +
               Student + BlueCollar + Commercial + Minivan + PanelTruck +
               Pickup + SportsCar + Van + RedCar + DLRevoked + Urban, cleandf,
               family = binomial(link = "logit"))
summary(model5)

model6 <- glm(TARGET_FLAG ~ INDEX + TARGET_AMT + KIDSDRV + AGE +
               HOMEKIDS + YOJ + HOME_VAL + TRAVTIME + BLUEBOOK + TIF + OLDCLAIM +
               CLM_FREQ + CAR_AGE + SingleParent + Male + Married + Bachelors +
               Clerical + Doctor + HomeMaker + Lawyer + Manager + Student +
               BlueCollar + Commercial + Minivan + PanelTruck + Pickup +
               SportsCar + Van + RedCar + Urban, data = cleandf, family = binomial(link = "logit"))
summary(model6)

```

```

# model selection

# checking our model's performance
performance::check_model(model1)
performance::check_model(model2)
performance::check_model(model3)

performance::model_performance(model1)
performance::model_performance(model2)
performance::model_performance(model3)

performance::check_model(model4)
performance::check_model(model5)
performance::check_model(model6)

# performance for model 4 only
performance::model_performance(model4)

pred = round(fitted(model4))
pROC::auc(cleandf$TARGET_FLAG, pred)

table(true = cleandf$TARGET_FLAG, pred = round(fitted(model4)))

# performance for model 5 only
performance::model_performance(model5)

pred = round(fitted(model5))
pROC::auc(cleandf$TARGET_FLAG, pred)

table(true = cleandf$TARGET_FLAG, pred = round(fitted(model5)))

# performance for model 6 only
performance::model_performance(model6)

pred = round(fitted(model6))
pROC::auc(cleandf$TARGET_FLAG, pred)

table(true = cleandf$TARGET_FLAG, pred = round(fitted(model6)))

# start our predictions for our chosen models (3 and 6) and
# start with model 6
clean_evaldf$TARGET_FLAG <- as.numeric(clean_evaldf$TARGET_FLAG)
reduced_clean_evaldf <- clean_evaldf %>%
  dplyr::select(-c("TARGET_AMT"))

prediction <- broom::augment(model6, newdata = reduced_clean_evaldf)
head(as.integer(prediction$.fitted > 0.5), 20)

TARGET_FLAG <- as.data.frame(as.integer(prediction$.fitted >
  0.5))
names(TARGET_FLAG) <- "TARGET_FLAG"

```

```

# predictions on the recovery amounts model 3
clean_evaldf$TARGET_AMT <- as.numeric(clean_evaldf$TARGET_AMT)
clean_evaldf$TARGET_FLAG <- NULL

clean_evaldf <- bind_cols(TARGET_FLAG, clean_evaldf)
prediction <- broom::augment(model3, newdata = clean_evaldf)
head(prediction$.fitted, 20)

recovery <- as.data.frame(prediction$.fitted)
names(recovery) <- "recovery"

clean_evaldf <- bind_cols(recovery, clean_evaldf) %>%
  mutate(recovery = ifelse(TARGET_FLAG == 1, recovery, 0))

# print prediction
head(clean_evaldf$recovery, 20)

```

References:

- <https://rdrr.io/cran/MASS/man/dropterm.html>