

# Data 621 Homework 3: Insurance

Tommy Jenkins, Violeta Stoyanova, Todd Weigel, Peter Kowalchuk, Eleanor R-Secoquian,  
Anthony Pagan

November 6, 2019

## OVERVIEW

In this homework assignment, you will explore, analyze and model a data set containing approximately 8000 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 representing the cost of the crash.

### 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.

## DATA EXPLORATION

### Data Summary

The dataset consists of two data files: training and evaluation. The training dataset contains 26 columns, while the evaluation dataset contains 24. The evaluation dataset is missing columns TARGET\_FLAG, TARGET\_AMT which represent our response variables, respectively whether the person was in a car crash and the cost of the car crash if the person was in an accident. We will start by exploring the training data set since it will be the one used to generate the models.

The columns in the data set are:

| VARIABLE NAME | DEFINITION                               | THEORETICAL EFFECT  |
|---------------|--|---|
| INDEX         | Identification Variable (do not use)     | None  |
| TARGET_FLAG   | Was Car in a crash? 1=YES 0=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       | Marital 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?        |
| REVOKED       | License Revoked (Past 7 Years)           | If your license was revoked in the past 7 years, you probably are a more risky driver.            |
| SEX           | Gender                                   | Urban legend says that women have less crashes then 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                                    |

## Missing Data

An important aspect of any dataset is to determine how much, if any, data is missing. We look at all the variables to see which if any have missing data. We look at the basic descriptive statistics as well as the missing data and percentages.

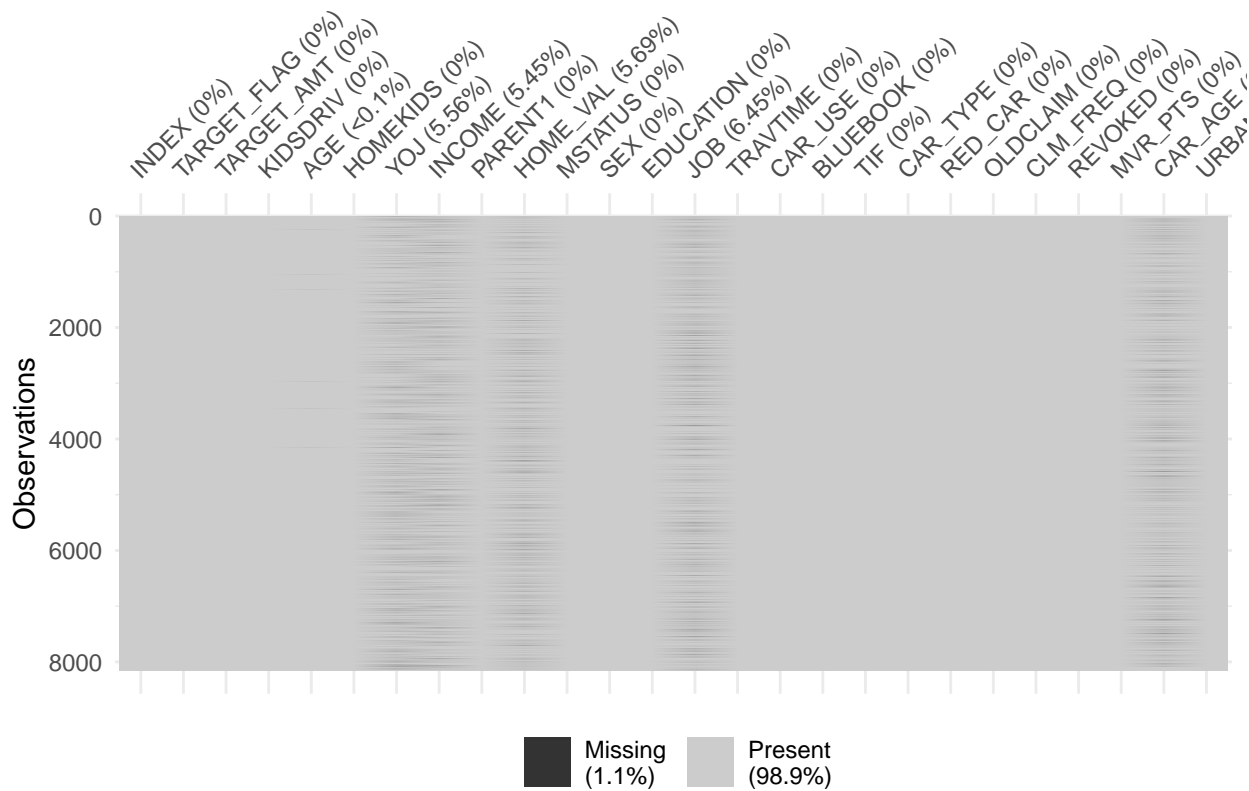
We start by looking at the dataset as a whole and determine how many complete rows, that is rows with data for all predictors, do we have.

```
##      Mode   FALSE    TRUE
## logical   2116    6045
```

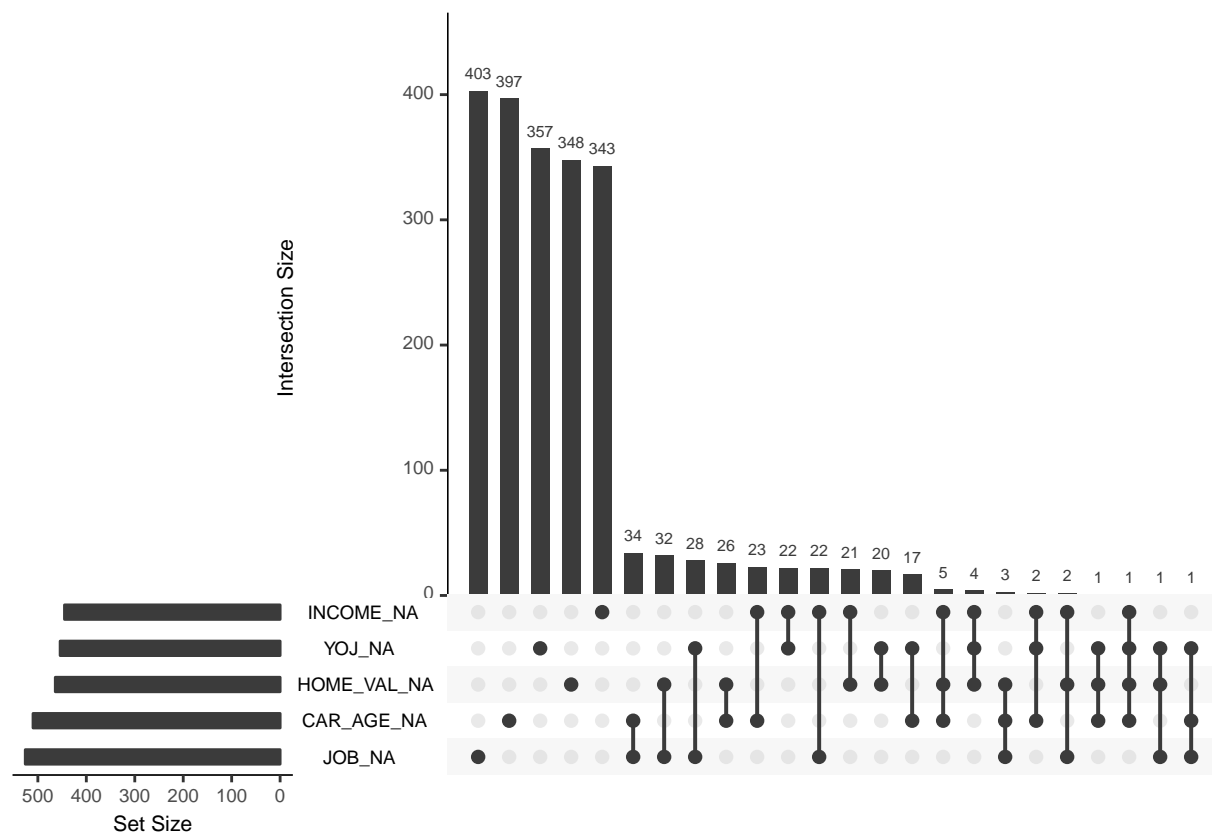
With these results, if we remove all rows with incomplete rows, there will be a total of 6045 rows out of 8161. If we eliminate all non-complete rows and keep only rows with data for all the predictors in the dataset, our new dataset will result in 74% of the total dataset. We create a subset of data with complete cases only to use later in our analysis.

```
## Observations: 6,045
## Variables: 26
## $ INDEX      <int> 1, 2, 4, 7, 12, 13, 14, 15, 16, 19, 20, 22, 23, 24...
## $ TARGET_FLAG <int> 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0,...
## $ TARGET_AMT  <dbl> 0.000, 0.000, 0.000, 2946.000, 2501.000, 0.000, 60...
## $ KIDSDRIV    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
## $ AGE         <int> 60, 43, 35, 34, 34, 50, 53, 43, 55, 45, 39, 42, 34...
## $ HOMEKIDS    <int> 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 3, 0, 3, 2, 1, 0,...
## $ YOJ         <int> 11, 11, 10, 12, 10, 7, 14, 5, 11, 0, 12, 11, 13, 1...
## $ INCOME      <fct> "$67,349", "$91,449", "$16,039", "$125,301", "$62,...
## $ PARENT1     <fct> No, No, No, Yes, No, No, No, No, No, No, Yes, No, ...
## $ HOME_VAL    <fct> "$0", "$257,252", "$124,191", "$0", "$0", "$0", "$...
## $ MSTATUS     <fct> z_No, z_No, Yes, z_No, z_No, z_No, z_No, Yes, Yes,...
## $ SEX         <fct> M, M, z_F, z_F, z_F, M, z_F, z_F, M, z_F, z_F, M, ...
## $ EDUCATION   <fct> PhD, z_High School, z_High School, Bachelors, Bach...
## $ JOB         <fct> Professional, z_Blue Collar, Clerical, z_Blue Coll...
## $ TRAVTIME    <int> 14, 22, 5, 46, 34, 48, 15, 36, 25, 48, 43, 42, 27,...
## $ CAR_USE     <fct> Private, Commercial, Private, Commercial, Private,...
## $ BLUEBOOK    <fct> "$14,230", "$14,940", "$4,010", "$17,430", "$11,20...
## $ TIF         <int> 11, 1, 4, 1, 1, 7, 1, 7, 7, 1, 6, 6, 7, 4, 6, 6, 1...
## $ CAR_TYPE    <fct> Minivan, Minivan, z_SUV, Sports Car, z_SUV, Van, S...
## $ RED_CAR     <fct> yes, yes, no, no, no, no, no, no, yes, no, no, no,...
## $ OLDCLAIM    <fct> "$4,461", "$0", "$38,690", "$0", "$0", "$0", "$0",...
## $ CLM_FREQ    <int> 2, 0, 2, 0, 0, 0, 0, 0, 2, 0, 0, 0, 0, 0, 2, 1, 0,...
## $ REVOKED     <fct> No, No, No, No, No, No, No, No, No, Yes, No, No, No, N...
## $ MVR_PTS     <int> 3, 0, 3, 0, 0, 1, 0, 0, 3, 3, 0, 0, 0, 0, 0, 5, 1,...
## $ CAR_AGE     <int> 18, 1, 10, 7, 1, 17, 11, 1, 9, 5, 13, 16, 20, 7, 1...
## $ URBANICITY  <fct> Highly Urban/ Urban, Highly Urban/ Urban, Highly U...
```

But we can also look at what specific predictors are missing in our dataset. If we do this we can see how there is much more data available, as we find only 5 predictors with missing data. Data missing for these predictors also only accounts for less than 7% of the respective predictors total.



We look closer at the missing data and look at the intersection of predictors with missing data. We find that the bulk of the missing data is for predictors with no intersection with other missing predictor data.



Having this detail in missing data might be of importance when looking at models. In the next Data Preparation section we will handle these missing cases and build a data set with data for all predictors in all rows.

## Data Exploration

Using TARGET\_FLAG as response variables we confirm when TARGET\_FLAG is 1 TARGET\_AMOUNT >0 and when TARGET\_FLAG is 0 when TARGET\_AMOUNT = 0

```
nrow(subset(InsTrain,TARGET_FLAG == 0))
```

```
## [1] 6008
```

```
nrow(subset(InsTrain,TARGET_AMT == 0))
```

```
## [1] 6008
```

```
nrow(subset(InsTrain,TARGET_FLAG > 0))
```

```
## [1] 2153
```

```
nrow(subset(InsTrain,TARGET_AMT > 0))
```

```
## [1] 2153
```

A glimpse of the data shows that the following columns should be integers and not Factors:

- INCOME
- HOME\_VAL
- BLUEBOOK
- OLDCLAIM
- JOB

We display and view data with all cases and only complete cases

```
## Observations: 8,161
## Variables: 26
## $ INDEX      <int> 1, 2, 4, 5, 6, 7, 8, 11, 12, 13, 14, 15, 16, 17, 1...
## $ TARGET_FLAG <int> 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0,...
## $ TARGET_AMT  <dbl> 0.000, 0.000, 0.000, 0.000, 0.000, 2946.000, 0.000...
## $ KIDSDRIV    <int> 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,...
## $ AGE         <int> 60, 43, 35, 51, 50, 34, 54, 37, 34, 50, 53, 43, 55...
## $ HOMEKIDS    <int> 0, 0, 1, 0, 0, 1, 0, 2, 0, 0, 0, 0, 0, 0, 0, 3, 0,...
## $ YOJ         <int> 11, 11, 10, 14, NA, 12, NA, NA, 10, 7, 14, 5, 11, ...
## $ INCOME      <fct> "$67,349", "$91,449", "$16,039", NA, "$114,986", "...
## $ PARENT1     <fct> No, No, No, No, No, Yes, No, No, No, No, No, No, N...
## $ HOME_VAL    <fct> "$0", "$257,252", "$124,191", "$306,251", "$243,92...
## $ MSTATUS     <fct> z_No, z_No, Yes, Yes, Yes, z_No, Yes, Yes, z_No, z...
## $ SEX         <fct> M, M, z_F, M, z_F, z_F, z_F, M, z_F, M, z_F, z_F, ...
## $ EDUCATION   <fct> PhD, z_High School, z_High School, <High School, P...
## $ JOB         <fct> Professional, z_Blue Collar, Clerical, z_Blue Coll...
## $ TRAVTIME    <int> 14, 22, 5, 32, 36, 46, 33, 44, 34, 48, 15, 36, 25,...
## $ CAR_USE     <fct> Private, Commercial, Private, Private, Private, Co...
## $ BLUEBOOK    <fct> "$14,230", "$14,940", "$4,010", "$15,440", "$18,00...
## $ TIF         <int> 11, 1, 4, 7, 1, 1, 1, 1, 1, 7, 1, 7, 7, 6, 1, 6, 6...
## $ CAR_TYPE    <fct> Minivan, Minivan, z_SUV, Minivan, z_SUV, Sports Ca...
## $ RED_CAR     <fct> yes, yes, no, yes, no, no, no, yes, no, no, no, no...
## $ OLDCLAIM    <fct> "$4,461", "$0", "$38,690", "$0", "$19,217", "$0", ...
## $ CLM_FREQ    <int> 2, 0, 2, 0, 2, 0, 0, 1, 0, 0, 0, 0, 2, 0, 0, 0, 0,...
```

```
## $ REVOKED      <fct> No, No, No, No, Yes, No, No, Yes, No, No, No, No, ...
## $ MVR_PTS      <int> 3, 0, 3, 0, 3, 0, 0, 10, 0, 1, 0, 0, 3, 3, 3, 0, 0...
## $ CAR_AGE      <int> 18, 1, 10, 6, 17, 7, 1, 7, 1, 17, 11, 1, 9, 10, 5,...
## $ URBANICITY   <fct> Highly Urban/ Urban, Highly Urban/ Urban, Highly U...
```

We use `is.na` function to review which columns have NA Values. It display columns and percent of values that are missing.

```
##      INDEX TARGET_FLAG TARGET_AMT KIDSDRIV      AGE  HOMEKIDS
##      0.0          0.0          0.0          0.0      0.1        0.0
##      YOJ          INCOME    PARENT1  HOME_VAL  MSTATUS      SEX
##      5.6          5.5          0.0          5.7      0.0        0.0
## EDUCATION          JOB    TRAVTIME    CAR_USE  BLUEBOOK      TIF
##      0.0          6.4          0.0          0.0      0.0        0.0
## CAR_TYPE    RED_CAR    OLDCLAIM    CLM_FREQ  REVOKED    MVR_PTS
##      0.0          0.0          0.0          0.0      0.0        0.0
## CAR_AGE    URBANICITY
##      6.2          0.0
```

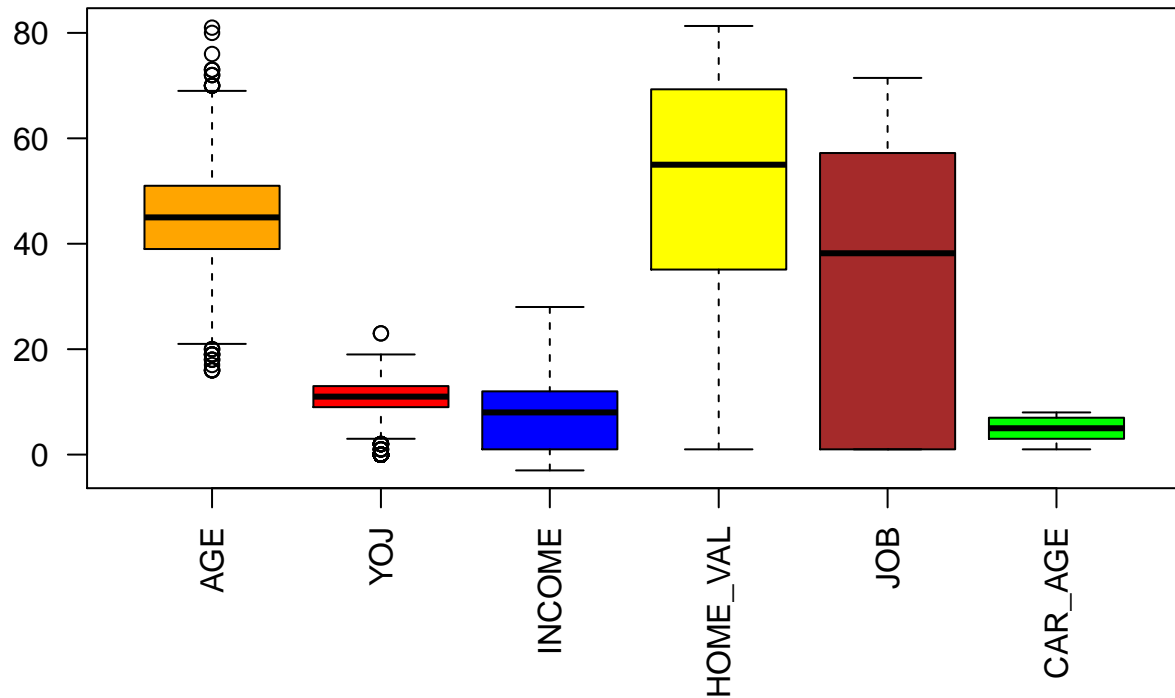
## Data Preperation

As revealed earlier there were a list of columns that we factors that should be integers. We start by converting the columns to numeric.

```
## Observations: 8,161
## Variables: 5
## $ INCOME      <fct> "$67,349", "$91,449", "$16,039", NA, "$114,986", "$12...
## $ HOME_VAL    <fct> "$0", "$257,252", "$124,191", "$306,251", "$243,925",...
## $ JOB         <fct> Professional, z_Blue Collar, Clerical, z_Blue Collar,...
## $ BLUEBOOK    <fct> "$14,230", "$14,940", "$4,010", "$15,440", "$18,000",...
## $ OLDCLAIM    <fct> "$4,461", "$0", "$38,690", "$0", "$19,217", "$0", "$0...

## Observations: 8,161
## Variables: 5
## $ INCOME      <dbl> 5032, 6291, 1249, NA, 508, 745, 1487, 314, 4764, 281,...
## $ HOME_VAL    <dbl> 1, 3258, 347, 3916, 3033, 1, NA, 4166, 1, 1, 1, 2261,...
## $ JOB         <dbl> 6, 8, 1, 8, 2, 8, 8, 8, 1, 6, 4, 6, 5, NA, 3, 1, 4, 4...
## $ BLUEBOOK    <dbl> 434, 503, 2212, 553, 802, 746, 2672, 701, 135, 852, 8...
## $ OLDCLAIM    <dbl> 1449, 1, 1311, 1, 432, 1, 1, 510, 1, 1, 1, 1, 1727, 1...
```

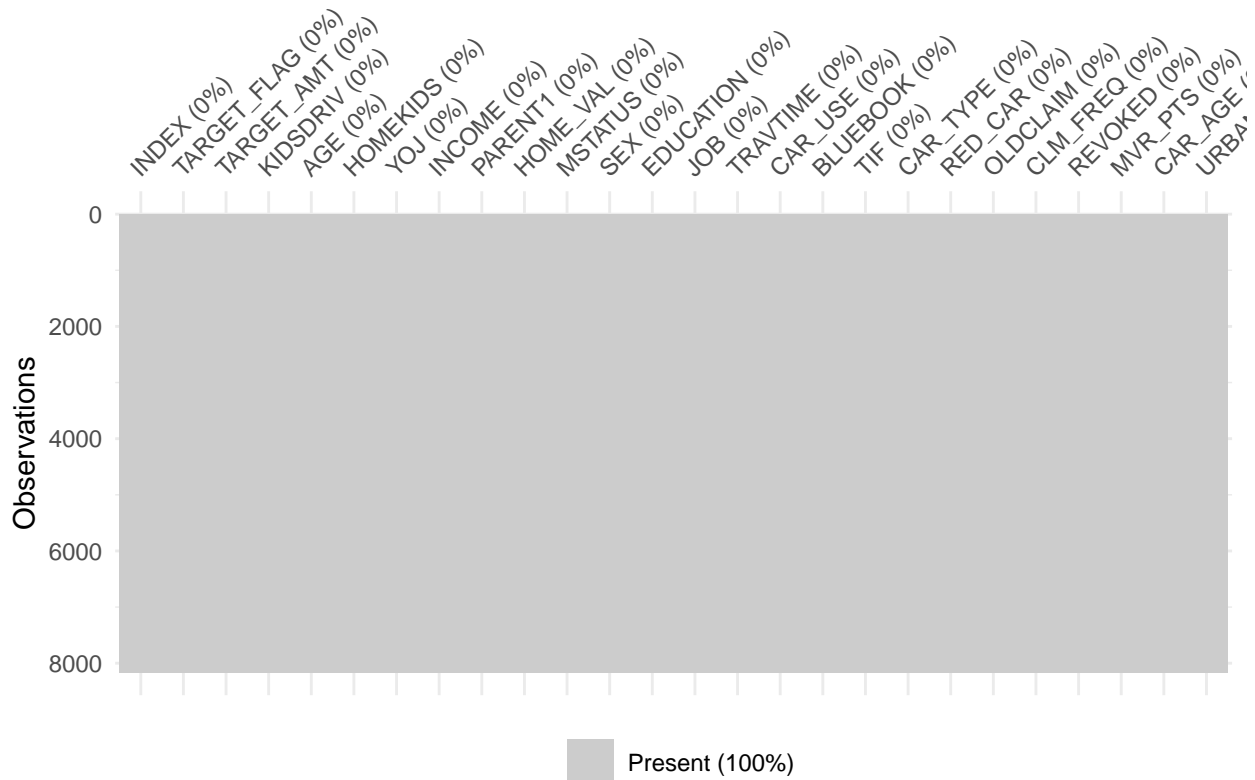
Both boxplot and summary stats with the square root transform of `Home_val` and `Income` to confirm we can use median or mean values to replace NA values if we chose.



| ## | vars     | n        | mean    | sd      | median | trimmed | mad     | min | max  | range |
|----|----------|----------|---------|---------|--------|---------|---------|-----|------|-------|
| ## | AGE      | 1 8155   | 44.79   | 8.63    | 45.0   | 44.83   | 8.90    | 16  | 81   | 65    |
| ## | YOJ      | 2 7707   | 10.50   | 4.09    | 11.0   | 11.07   | 2.97    | 0   | 23   | 23    |
| ## | INCOME   | 3 7716   | 3040.33 | 2029.52 | 3024.5 | 3018.56 | 2647.18 | 1   | 6612 | 6611  |
| ## | HOME_VAL | 4 7697   | 1785.40 | 1695.15 | 1459.0 | 1639.68 | 2161.63 | 1   | 5106 | 5105  |
| ## | JOB      | 5 7635   | 5.01    | 2.46    | 5.0    | 5.14    | 2.97    | 1   | 8    | 7     |
| ## | CAR_AGE  | 6 7651   | 8.33    | 5.70    | 8.0    | 7.96    | 7.41    | -3  | 28   | 31    |
| ## | skew     | kurtosis | se      |         |        |         |         |     |      |       |
| ## | AGE      | -0.03    | -0.06   | 0.10    |        |         |         |     |      |       |
| ## | YOJ      | -1.20    | 1.18    | 0.05    |        |         |         |     |      |       |
| ## | INCOME   | 0.04     | -1.24   | 23.10   |        |         |         |     |      |       |
| ## | HOME_VAL | 0.43     | -1.24   | 19.32   |        |         |         |     |      |       |
| ## | JOB      | -0.34    | -1.16   | 0.03    |        |         |         |     |      |       |
| ## | CAR_AGE  | 0.28     | -0.75   | 0.07    |        |         |         |     |      |       |

We next replace all NA values with mean values for cases that are missing values and rerun sapply function to confirm there are no longer any missing values.

| ## | INDEX     | TARGET_FLAG | TARGET_AMT | KIDSDRIV | AGE      | HOMEKIDS |
|----|-----------|-------------|------------|----------|----------|----------|
| ## | 0         | 0           | 0          | 0        | 0        | 0        |
| ## | YOJ       | INCOME      | PARENT1    | HOME_VAL | MSTATUS  | SEX      |
| ## | 0         | 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  |            |          |          |          |
| ## | 0         | 0           |            |          |          |          |



| ## | vars     | n      | mean    | sd       | median  | trimmed | mad     | min | max  | range |
|----|----------|--------|---------|----------|---------|---------|---------|-----|------|-------|
| ## | AGE      | 1 8161 | 44.79   | 8.62     | 45.00   | 44.83   | 8.90    | 16  | 81   | 65    |
| ## | YOJ      | 2 8161 | 10.50   | 3.98     | 11.00   | 11.05   | 2.97    | 0   | 23   | 23    |
| ## | INCOME   | 3 8161 | 3040.33 | 1973.41  | 3040.33 | 3020.15 | 2492.74 | 1   | 6612 | 6611  |
| ## | HOME_VAL | 4 8161 | 1785.40 | 1646.25  | 1680.00 | 1642.31 | 2489.29 | 1   | 5106 | 5105  |
| ## | JOB      | 5 8161 | 5.01    | 2.38     | 5.01    | 5.14    | 2.95    | 1   | 8    | 7     |
| ## | CAR_AGE  | 6 8161 | 8.33    | 5.52     | 8.33    | 7.98    | 5.44    | -3  | 28   | 31    |
| ## |          |        | skew    | kurtosis |         |         |         |     |      |       |
| ## | AGE      |        | -0.03   | -0.06    | 0.10    |         |         |     |      |       |
| ## | YOJ      |        | -1.24   | 1.42     | 0.04    |         |         |     |      |       |
| ## | INCOME   |        | 0.05    | -1.14    | 21.84   |         |         |     |      |       |
| ## | HOME_VAL |        | 0.44    | -1.14    | 18.22   |         |         |     |      |       |
| ## | JOB      |        | -0.36   | -1.03    | 0.03    |         |         |     |      |       |
| ## | CAR_AGE  |        | 0.29    | -0.60    | 0.06    |         |         |     |      |       |

We have this way derived a dataset with no missing values. We can use this set of data for our modeling design. We chose to work with this data as opposed to the first “complete” set in which rows with missing data were eliminated.

## Build Model

Modeling design will be divided in two phases. First we will design a model to predict if the person is in a car crash, that is predict TARGET\_FLAG. In a second phase, we will predict TARGET\_AMT, the cost of the crash.

### TARGET\_FLAG Modeling

This response variable being binary, 0 or 1, we will be looking at logistic regression models to find a good fit. We will start with a naive model with all the predictors included as a baseline. First approach will be to

simply the model by reducing the predictors used. We will look at several model metrics such as AIC, BIC. We will also include confusion tables and ROC plot to better understand each model.

### Model 1: all predictors

We start out with a straightforward logit logistical regression with all predictors included. As a note, we need to make sure we do not include the TARGET\_AMT response variable in our model as a predictor.

```
##
## Call:
## glm(formula = TARGET_FLAG ~ . - INDEX - TARGET_AMT, family = binomial(link = "logit"),
##      data = InsTrain)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.5922  -0.7276  -0.4209   0.6419   3.1555
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -4.709e-01  2.687e-01  -1.753  0.079675 .
## KIDSDRIV        3.641e-01  6.053e-02   6.014  1.81e-09 ***
## AGE            -4.635e-03  3.910e-03  -1.185  0.235876
## HOMEKIDS        7.217e-02  3.651e-02   1.977  0.048085 *
## YOJ            -2.674e-02  7.937e-03  -3.369  0.000755 ***
## INCOME         -2.314e-05  1.662e-05  -1.392  0.163862
## PARENT1Yes      3.497e-01  1.082e-01   3.232  0.001229 **
## HOME_VAL       -9.899e-05  2.047e-05  -4.836  1.33e-06 ***
## MSTATUSz_No     4.804e-01  7.717e-02   6.226  4.80e-10 ***
## SEXz_F         -2.579e-01  1.021e-01  -2.526  0.011524 *
## EDUCATIONBachelors -6.887e-01  1.043e-01  -6.603  4.03e-11 ***
## EDUCATIONMasters -7.923e-01  1.291e-01  -6.136  8.47e-10 ***
## EDUCATIONPhD    -1.026e+00  1.545e-01  -6.644  3.06e-11 ***
## EDUCATIONz_High School -1.118e-01  9.212e-02  -1.214  0.224900
## JOB            -2.270e-02  1.406e-02  -1.615  0.106418
## TRAVTIME        1.517e-02  1.868e-03   8.121  4.62e-16 ***
## CAR_USEPrivate  -8.701e-01  8.207e-02 -10.602  < 2e-16 ***
## BLUEBOOK        2.106e-05  3.344e-05   0.630  0.528975
## TIF            -5.453e-02  7.276e-03  -7.495  6.63e-14 ***
## CAR_TYPEPanel Truck 6.109e-02  1.364e-01   0.448  0.654173
## CAR_TYPEPickup    5.281e-01  1.006e-01   5.249  1.53e-07 ***
## CAR_TYPESports Car 1.209e+00  1.217e-01   9.934  < 2e-16 ***
## CAR_TYPEVan       3.726e-01  1.195e-01   3.119  0.001815 **
## CAR_TYPEz_SUV     9.618e-01  1.024e-01   9.392  < 2e-16 ***
## RED_CARyes       1.440e-03  8.558e-02   0.017  0.986572
## OLDCLAIM        8.744e-05  4.214e-05   2.075  0.037989 *
## CLM_FREQ        1.161e-01  3.190e-02   3.639  0.000274 ***
## REVOKEDYes      7.483e-01  7.965e-02   9.395  < 2e-16 ***
## MVRPTS          1.101e-01  1.360e-02   8.092  5.86e-16 ***
## CAR_AGE         -6.893e-04  7.459e-03  -0.092  0.926369
## URBANICITYz_Highly Rural/ Rural -2.281e+00  1.122e-01 -20.331  < 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: 7422.7  on 8130  degrees of freedom
## AIC: 7484.7
##
## Number of Fisher Scoring iterations: 5
```

From the model's summary itself we see that there are several predictors which are not statistically relevant, which suggests a simpler model should be possible. We build a second model without these the non-significant predictors.

## Model 2: reduced predictors

```
##
## Call:
## glm(formula = TARGET_FLAG ~ . - INDEX - TARGET_AMT - AGE - INCOME -
##      JOB - BLUEBOOK - CAR_AGE - RED_CAR, family = binomial(link = "logit"),
##      data = InsTrain)
##
## Deviance Residuals:
##      Min        1Q    Median        3Q        Max
## -2.5683   -0.7296   -0.4221    0.6443    3.1248
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -8.370e-01  1.689e-01  -4.956 7.21e-07 ***
## KIDSDRIV       3.478e-01  5.947e-02   5.848 4.98e-09 ***
## HOMEKIDS       9.384e-02  3.368e-02   2.786 0.005337 **
## YOJ           -3.225e-02  7.337e-03  -4.395 1.11e-05 ***
## PARENT1Yes     3.612e-01  1.076e-01   3.356 0.000792 ***
## HOME_VAL     -1.010e-04  2.051e-05  -4.922 8.55e-07 ***
## MSTATUSz_No    4.791e-01  7.706e-02   6.217 5.07e-10 ***
## SEXz_F        -2.465e-01  8.784e-02  -2.806 0.005013 **
## EDUCATIONBachelors -7.326e-01  9.501e-02  -7.711 1.25e-14 ***
## EDUCATIONMasters -8.212e-01  1.032e-01  -7.959 1.74e-15 ***
## EDUCATIONPhD    -1.017e+00  1.336e-01  -7.617 2.61e-14 ***
## EDUCATIONz_High School -1.251e-01  9.110e-02  -1.373 0.169855
## TRAVTIME       1.509e-02  1.866e-03   8.088 6.06e-16 ***
## CAR_USEPrivate -8.044e-01  7.261e-02 -11.077 < 2e-16 ***
## TIF           -5.426e-02  7.266e-03  -7.469 8.11e-14 ***
## CAR_TYPEPanel Truck 1.241e-01  1.315e-01   0.943 0.345453
## CAR_TYPEPickup    5.717e-01  9.696e-02   5.897 3.71e-09 ***
## CAR_TYPESports Car 1.213e+00  1.201e-01  10.097 < 2e-16 ***
## CAR_TYPEVan       4.018e-01  1.185e-01   3.390 0.000698 ***
## CAR_TYPEz_SUV     9.653e-01  1.017e-01   9.495 < 2e-16 ***
## OLDCLAIM        8.853e-05  4.209e-05   2.103 0.035443 *
## CLM_FREQ        1.149e-01  3.186e-02   3.606 0.000311 ***
## REVOKEDYes       7.510e-01  7.960e-02   9.435 < 2e-16 ***
## MVR_PTS         1.104e-01  1.358e-02   8.135 4.12e-16 ***
## URBANICITYz_Highly Rural/ Rural -2.268e+00  1.120e-01 -20.259 < 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  on 8160  degrees of freedom
```

```
## Residual deviance: 7430  on 8136  degrees of freedom
## AIC: 7480
##
## Number of Fisher Scoring iterations: 5
```

The new model has a slightly higher AIC which would tell us the first model is slightly less complex.

### AIC Step Method Model 3

Another way of selecting which predictors to use in the model is by calculating the AIC of the model. This metric is similar to the adjusted R-square of a model in that it penalizes models with more predictors over simpler model with few predictors. We use StepAIC function in R to find the lowest AIC with different predictors.

```
## Start:  AIC=7484.66
## TARGET_FLAG ~ (INDEX + TARGET_AMT + KIDSDRIV + AGE + HOMEKIDS +
##      YOJ + INCOME + PARENT1 + HOME_VAL + MSTATUS + SEX + EDUCATION +
##      JOB + TRAVTIME + CAR_USE + BLUEBOOK + TIF + CAR_TYPE + RED_CAR +
##      OLDCLAIM + CLM_FREQ + REVOKED + MVRPTS + CAR_AGE + URBANICITY) -
##      INDEX - TARGET_AMT
##
##      Df Deviance    AIC
## - RED_CAR      1    7422.7 7482.7
## - CAR_AGE      1    7422.7 7482.7
## - BLUEBOOK     1    7423.1 7483.1
## - AGE          1    7424.1 7484.1
## - INCOME       1    7424.6 7484.6
## <none>          1    7422.7 7484.7
## - JOB          1    7425.3 7485.3
## - HOMEKIDS     1    7426.5 7486.5
## - OLDCLAIM     1    7427.0 7487.0
## - SEX          1    7429.1 7489.1
## - PARENT1      1    7433.1 7493.1
## - YOJ          1    7434.0 7494.0
## - CLM_FREQ     1    7435.8 7495.8
## - HOME_VAL     1    7446.3 7506.3
## - KIDSDRIV     1    7458.9 7518.9
## - MSTATUS      1    7461.3 7521.3
## - TIF          1    7480.9 7540.9
## - MVRPTS       1    7488.6 7548.6
## - TRAVTIME     1    7488.9 7548.9
## - EDUCATION    4    7502.1 7556.1
## - REVOKED      1    7509.3 7569.3
## - CAR_USE      1    7537.4 7597.4
## - CAR_TYPE     5    7551.9 7603.9
## - URBANICITY   1    8013.3 8073.3
##
## Step:  AIC=7482.66
## TARGET_FLAG ~ KIDSDRIV + AGE + HOMEKIDS + YOJ + INCOME + PARENT1 +
##      HOME_VAL + MSTATUS + SEX + EDUCATION + JOB + TRAVTIME + CAR_USE +
##      BLUEBOOK + TIF + CAR_TYPE + OLDCLAIM + CLM_FREQ + REVOKED +
##      MVRPTS + CAR_AGE + URBANICITY
##
##      Df Deviance    AIC
## - CAR_AGE      1    7422.7 7480.7
```

```

## - BLUEBOOK      1    7423.1 7481.1
## - AGE            1    7424.1 7482.1
## - INCOME         1    7424.6 7482.6
## <none>           7422.7 7482.7
## - JOB            1    7425.3 7483.3
## - HOMEKIDS       1    7426.5 7484.5
## - OLDCLAIM       1    7427.0 7485.0
## - SEX            1    7431.4 7489.4
## - PARENT1        1    7433.1 7491.1
## - YOJ            1    7434.0 7492.0
## - CLM_FREQ       1    7435.8 7493.8
## - HOME_VAL       1    7446.3 7504.3
## - KIDSDRIV       1    7458.9 7516.9
## - MSTATUS        1    7461.3 7519.3
## - TIF            1    7480.9 7538.9
## - MVR_PTS        1    7488.6 7546.6
## - TRAVTIME       1    7488.9 7546.9
## - EDUCATION      4    7502.1 7554.1
## - REVOKED        1    7509.3 7567.3
## - CAR_USE        1    7537.4 7595.4
## - CAR_TYPE       5    7551.9 7601.9
## - URBANICITY     1    8013.3 8071.3
##
## Step:   AIC=7480.67
## TARGET_FLAG ~ KIDSDRIV + AGE + HOMEKIDS + YOJ + INCOME + PARENT1 +
##               HOME_VAL + MSTATUS + SEX + EDUCATION + JOB + TRAVTIME + CAR_USE +
##               BLUEBOOK + TIF + CAR_TYPE + OLDCLAIM + CLM_FREQ + REVOKED +
##               MVR_PTS + URBANICITY
##
##           Df Deviance    AIC
## - BLUEBOOK      1    7423.1 7479.1
## - AGE            1    7424.1 7480.1
## - INCOME         1    7424.6 7480.6
## <none>           7422.7 7480.7
## - JOB            1    7425.3 7481.3
## - HOMEKIDS       1    7426.6 7482.6
## - OLDCLAIM       1    7427.0 7483.0
## - SEX            1    7431.4 7487.4
## - PARENT1        1    7433.1 7489.1
## - YOJ            1    7434.0 7490.0
## - CLM_FREQ       1    7435.8 7491.8
## - HOME_VAL       1    7446.3 7502.3
## - KIDSDRIV       1    7458.9 7514.9
## - MSTATUS        1    7461.3 7517.3
## - TIF            1    7480.9 7536.9
## - MVR_PTS        1    7488.6 7544.6
## - TRAVTIME       1    7488.9 7544.9
## - REVOKED        1    7509.3 7565.3
## - CAR_USE        1    7537.4 7593.4
## - CAR_TYPE       5    7552.0 7600.0
## - EDUCATION      4    7553.4 7603.4
## - URBANICITY     1    8013.3 8069.3
##
## Step:   AIC=7479.07

```

```

## TARGET_FLAG ~ KIDSDRIV + AGE + HOMEKIDS + YOJ + INCOME + PARENT1 +
##     HOME_VAL + MSTATUS + SEX + EDUCATION + JOB + TRAVTIME + CAR_USE +
##     TIF + CAR_TYPE + OLDCLAIM + CLM_FREQ + REVOKED + MVR_PTS +
##     URBANICITY
##
##           Df Deviance    AIC
## - AGE      1   7424.5 7478.5
## <none>      7423.1 7479.1
## - INCOME    1   7425.2 7479.2
## - JOB       1   7425.7 7479.7
## - HOMEKIDS  1   7427.0 7481.0
## - OLDCLAIM  1   7427.4 7481.4
## - SEX       1   7431.7 7485.7
## - PARENT1   1   7433.4 7487.4
## - YOJ       1   7434.5 7488.5
## - CLM_FREQ  1   7436.1 7490.1
## - HOME_VAL  1   7446.7 7500.7
## - KIDSDRIV  1   7459.2 7513.2
## - MSTATUS   1   7461.9 7515.9
## - TIF       1   7481.3 7535.3
## - MVR_PTS   1   7488.9 7542.9
## - TRAVTIME  1   7489.3 7543.3
## - REVOKED   1   7509.6 7563.6
## - CAR_USE   1   7537.8 7591.8
## - CAR_TYPE  5   7555.5 7601.5
## - EDUCATION 4   7554.7 7602.7
## - URBANICITY 1   8014.0 8068.0
##
## Step:  AIC=7478.53
## TARGET_FLAG ~ KIDSDRIV + HOMEKIDS + YOJ + INCOME + PARENT1 +
##     HOME_VAL + MSTATUS + SEX + EDUCATION + JOB + TRAVTIME + CAR_USE +
##     TIF + CAR_TYPE + OLDCLAIM + CLM_FREQ + REVOKED + MVR_PTS +
##     URBANICITY
##
##           Df Deviance    AIC
## - INCOME    1   7426.5 7478.5
## <none>      7424.5 7478.5
## - JOB       1   7427.3 7479.3
## - OLDCLAIM  1   7428.9 7480.9
## - HOMEKIDS  1   7431.5 7483.5
## - SEX       1   7432.8 7484.8
## - PARENT1   1   7435.7 7487.7
## - CLM_FREQ  1   7437.5 7489.5
## - YOJ       1   7437.7 7489.7
## - HOME_VAL  1   7448.5 7500.5
## - KIDSDRIV  1   7459.2 7511.2
## - MSTATUS   1   7463.8 7515.8
## - TIF       1   7482.5 7534.5
## - TRAVTIME  1   7490.5 7542.5
## - MVR_PTS   1   7491.2 7543.2
## - REVOKED   1   7511.4 7563.4
## - CAR_USE   1   7539.7 7591.7
## - CAR_TYPE  5   7555.8 7599.8
## - EDUCATION 4   7562.4 7608.4

```

```

## - URBANICITY 1 8016.9 8068.9
##
## Step: AIC=7478.51
## TARGET_FLAG ~ KIDSDRIV + HOMEKIDS + YOJ + PARENT1 + HOME_VAL +
## MSTATUS + SEX + EDUCATION + JOB + TRAVTIME + CAR_USE + TIF +
## CAR_TYPE + OLDCLAIM + CLM_FREQ + REVOKED + MVR_PTS + URBANICITY
##
##           Df Deviance    AIC
## <none>           7426.5 7478.5
## - JOB           1 7430.0 7480.0
## - OLDCLAIM      1 7430.9 7480.9
## - HOMEKIDS      1 7434.0 7484.0
## - SEX           1 7434.8 7484.8
## - PARENT1       1 7437.7 7487.7
## - CLM_FREQ      1 7439.3 7489.3
## - YOJ           1 7445.7 7495.7
## - HOME_VAL      1 7450.6 7500.6
## - KIDSDRIV      1 7461.4 7511.4
## - MSTATUS       1 7465.3 7515.3
## - TIF           1 7485.0 7535.0
## - TRAVTIME      1 7492.5 7542.5
## - MVR_PTS       1 7492.9 7542.9
## - REVOKED       1 7513.7 7563.7
## - CAR_USE       1 7542.8 7592.8
## - CAR_TYPE      5 7558.6 7600.6
## - EDUCATION     4 7567.3 7611.3
## - URBANICITY    1 8017.6 8067.6
##
## Call:
## glm(formula = TARGET_FLAG ~ KIDSDRIV + HOMEKIDS + YOJ + PARENT1 +
## HOME_VAL + MSTATUS + SEX + EDUCATION + JOB + TRAVTIME + CAR_USE +
## TIF + CAR_TYPE + OLDCLAIM + CLM_FREQ + REVOKED + MVR_PTS +
## URBANICITY, family = binomial(link = "logit"), data = InsTrain)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6052  -0.7277  -0.4212   0.6451   3.1351
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -6.416e-01  1.984e-01  -3.234 0.001219 **
## KIDSDRIV        3.513e-01  5.952e-02   5.902 3.59e-09 ***
## HOMEKIDS        9.289e-02  3.368e-02   2.758 0.005818 **
## YOJ            -3.225e-02  7.336e-03  -4.397 1.10e-05 ***
## PARENT1Yes      3.593e-01  1.077e-01   3.337 0.000846 ***
## HOME_VAL       -1.000e-04  2.051e-05  -4.877 1.08e-06 ***
## MSTATUSz_No     4.810e-01  7.708e-02   6.240 4.37e-10 ***
## SEXz_F         -2.517e-01  8.794e-02  -2.862 0.004208 **
## EDUCATIONBachelors -7.253e-01  9.511e-02  -7.626 2.42e-14 ***
## EDUCATIONMasters -8.325e-01  1.034e-01  -8.052 8.17e-16 ***
## EDUCATIONPhD    -1.051e+00  1.347e-01  -7.798 6.28e-15 ***
## EDUCATIONz_High School -1.310e-01  9.115e-02  -1.438 0.150549
## JOB            -2.608e-02  1.392e-02  -1.874 0.060941 .

```

```

## TRAVTIME          1.513e-02  1.867e-03   8.108 5.16e-16 ***
## CAR_USEPrivate    -8.748e-01  8.199e-02 -10.669 < 2e-16 ***
## TIF               -5.459e-02  7.271e-03  -7.508 6.02e-14 ***
## CAR_TYPEPanel Truck  7.944e-02  1.337e-01   0.594 0.552376
## CAR_TYPEPickup     5.445e-01  9.818e-02   5.546 2.92e-08 ***
## CAR_TYPESports Car  1.210e+00  1.201e-01  10.073 < 2e-16 ***
## CAR_TYPEVan        3.766e-01  1.194e-01   3.154 0.001609 **
## CAR_TYPEz_SUV      9.629e-01  1.017e-01   9.469 < 2e-16 ***
## OLDCLAIM          8.789e-05  4.210e-05   2.088 0.036833 *
## CLM_FREQ          1.144e-01  3.186e-02   3.591 0.000329 ***
## REVOKEDYes        7.503e-01  7.963e-02   9.422 < 2e-16 ***
## MVR_PTS           1.102e-01  1.358e-02   8.117 4.77e-16 ***
## URBANICITYz_Highly Rural/ Rural -2.280e+00  1.122e-01 -20.332 < 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: 7426.5  on 8135  degrees of freedom
## AIC: 7478.5
##
## Number of Fisher Scoring iterations: 5

```

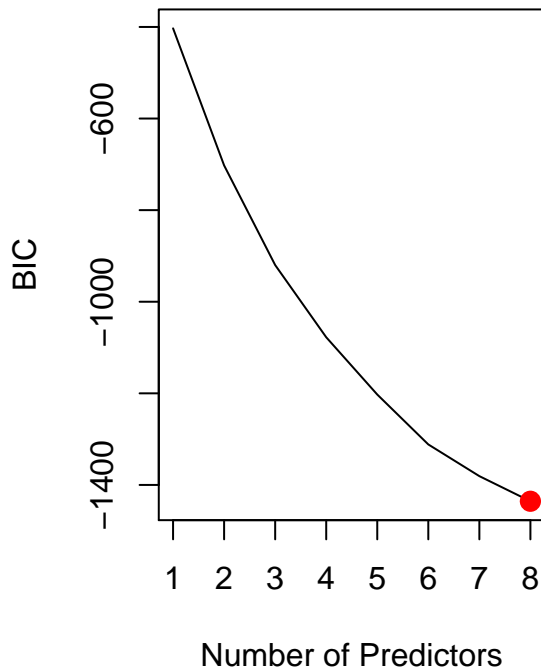
This reduces the predictors used to 25 from 30. The AIC improves from 7482.66 (our original general model) to 7478.5, and we also benefit by having a simpler model less prone to overfitting.

Also, the predictors in the model now are all significant (under 0.05 p level) and all but one under .02 or very significant. Which is much improved over the first model.

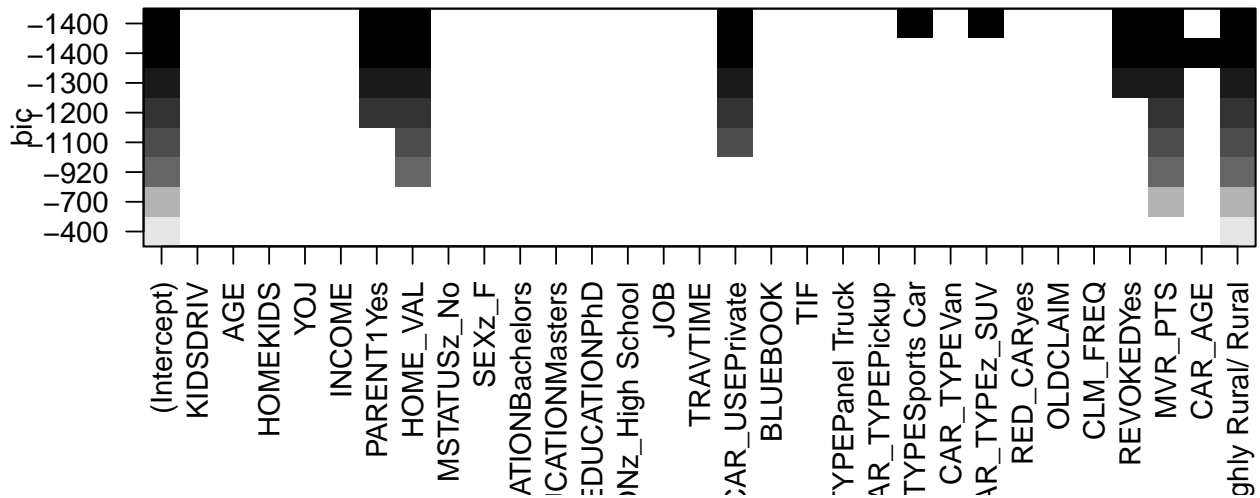
#### BIC Method Model 4

To determine the number of predictors and which predictors to be used we will use the Bayesian Information Criterion (BIC).

## Subset Selection Using BIC



## Predictors vs. BIC



The plot on the right shows that the number of predictors with the lowest BIC are PARENT , HOMEVAL, CAR\_USE, 'CAR\_TYPE', 'REVOKED', 'MVR\_PTS' and 'URBANICITY'. We will use those predictors to build the next model

|                     | Estimate   | Std. Error | z value | Pr(> z )  |
|---------------------|------------|------------|---------|-----------|
| (Intercept)         | -0.9348    | 0.09091    | -10.28  | 8.459e-25 |
| PARENT1Yes          | 0.7962     | 0.07785    | 10.23   | 1.503e-24 |
| HOME_VAL            | -0.0002045 | 1.825e-05  | -11.2   | 3.866e-29 |
| CAR_USEPrivate      | -0.8817    | 0.06719    | -13.12  | 2.409e-39 |
| CAR_TYPEPanel Truck | 0.02916    | 0.1227     | 0.2376  | 0.8122    |

|                                    | Estimate | Std. Error | z value | Pr(> z )  |
|------------------------------------|----------|------------|---------|-----------|
| CAR_TYPEPickup                     | 0.5546   | 0.09339    | 5.938   | 2.881e-09 |
| CAR_TYPESports Car                 | 1.082    | 0.1011     | 10.7    | 1.008e-26 |
| CAR_TYPEVan                        | 0.3126   | 0.1132     | 2.761   | 0.005755  |
| CAR_TYPEz_SUV                      | 0.8356   | 0.08082    | 10.34   | 4.682e-25 |
| REVOKEDYes                         | 0.7567   | 0.07661    | 9.877   | 5.218e-23 |
| MVR_PTS                            | 0.1574   | 0.01221    | 12.89   | 5.303e-38 |
| URBANICITYz_Highly<br>Rural/ Rural | -1.974   | 0.1056     | -18.7   | 5.221e-78 |

(Dispersion parameter for binomial family taken to be 1 )

|                    |                                 |
|--------------------|---------------------------------|
| Null deviance:     | 9418 on 8160 degrees of freedom |
| Residual deviance: | 7853 on 8149 degrees of freedom |

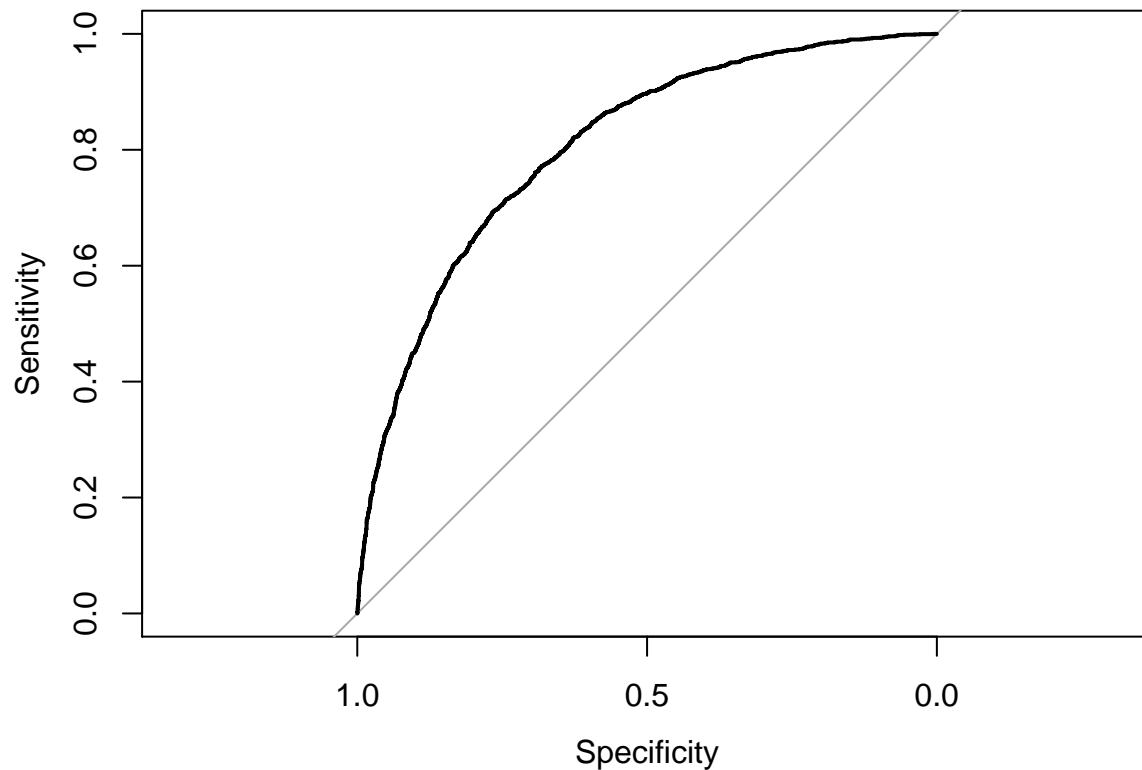
## Select Model

## Compare Model Statistics

## Model 1 - General Model

## ROC Curve

The ROC Curve helps measure true positives and true negative. A high AUC or area under the curve tells us the model is predicting well.



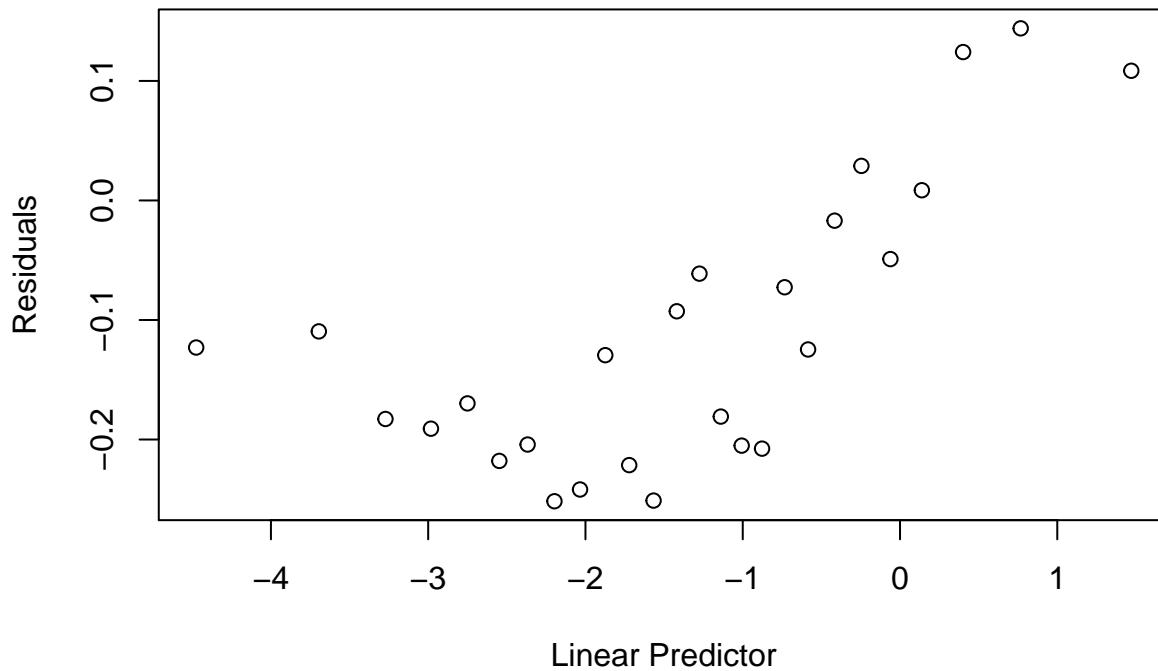
The AUC value of 0.8, tells us this model predicted values are accurate.

## Confusion Matrix

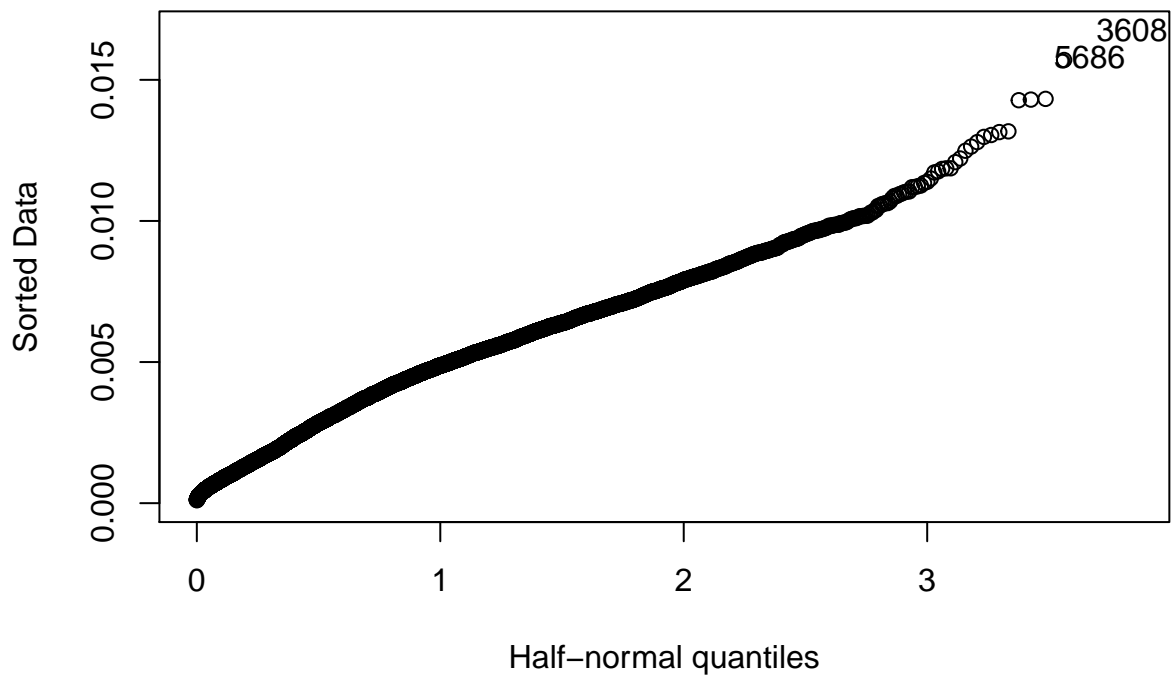


```
##
## targethat    0    1
##           0 5546 1301
##           1  462  852
```

Create a binned diagnostic plot of residuals vs prediction There are definite patterns here, which bear investigating.

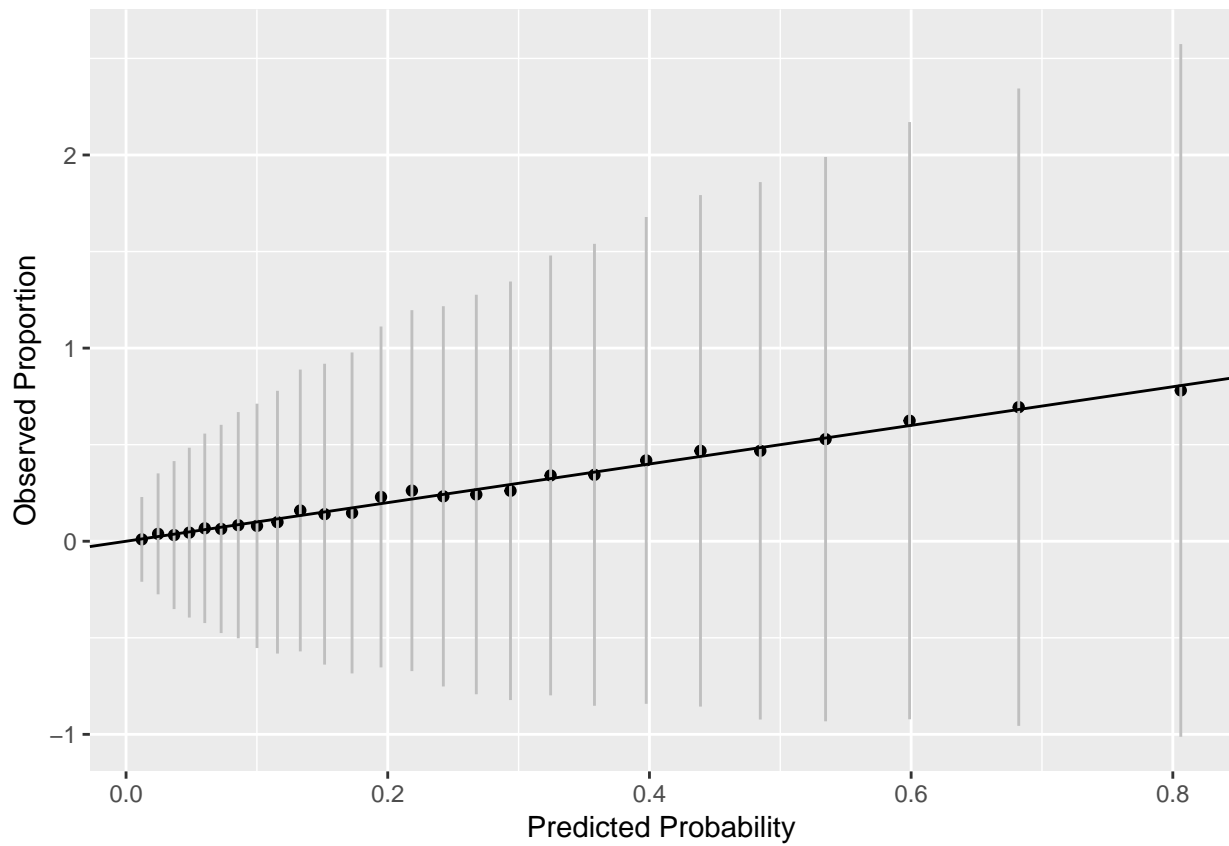


Plot leverages.



We don't see any strong outliers with the leverage plot. The points identified (3608,5686) are essentially in the plot of the line formed, so they are not likely pulling our model in any direction.

### Plot Goodness of fit

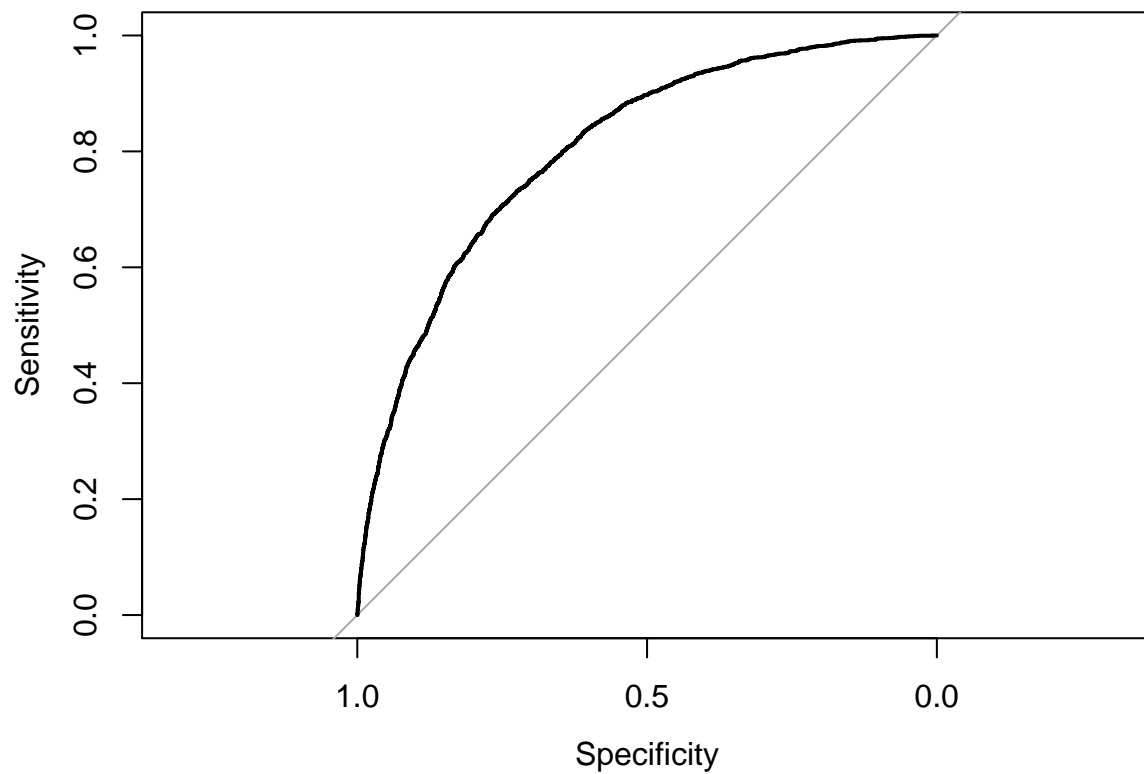


We see that our predictors fall close to the line.

### Model 2 - Reduced General Model

#### ROC Curve

The ROC Curve helps measure true positives and true negative. A high AUC or area under the curve tells us the model is predicting well.

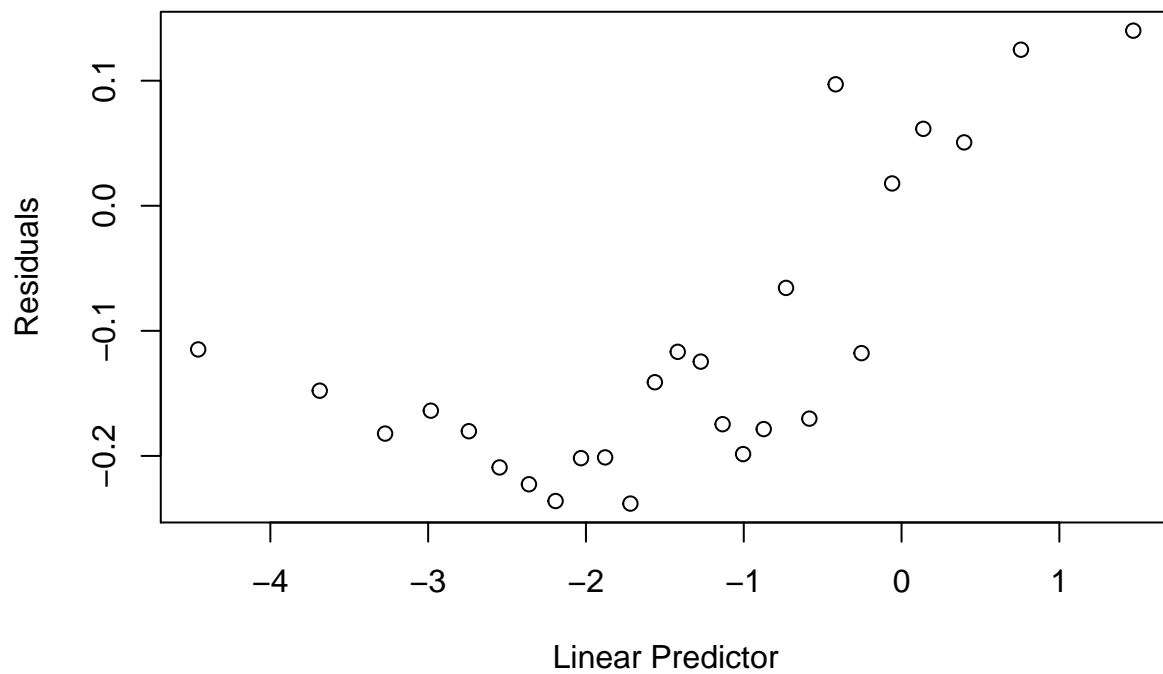


The AUC value of 0.8, tells us this model predicted values are accurate.

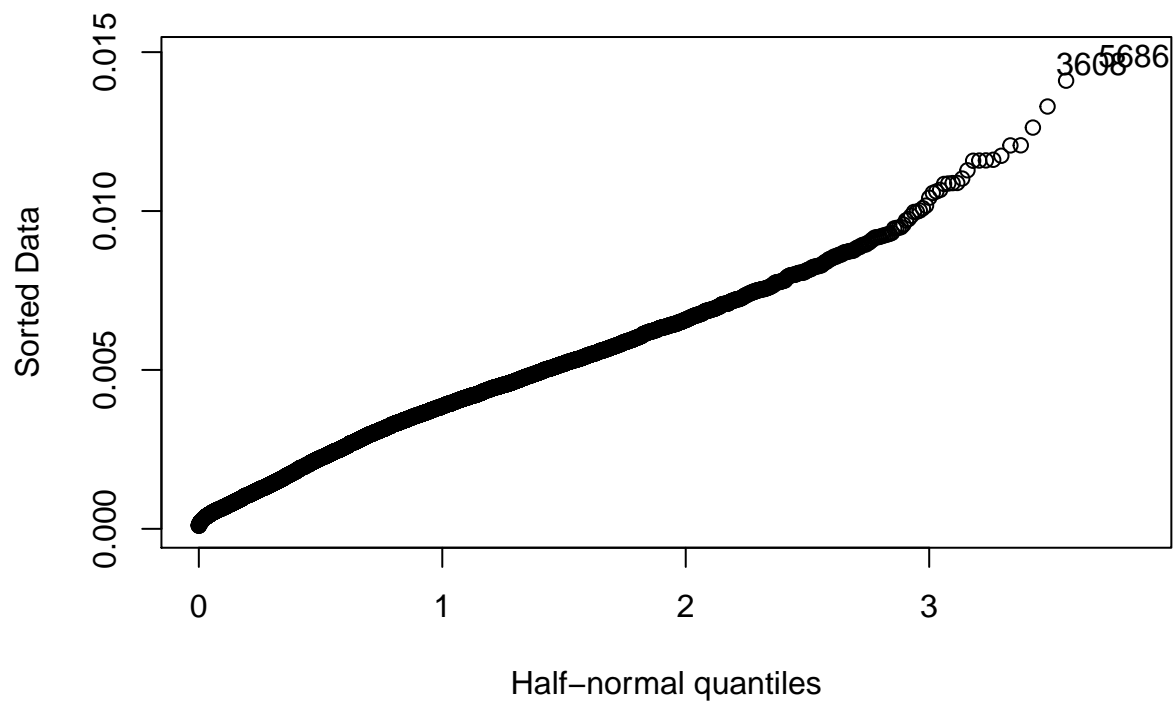
#### Confusion Matrix

```
##
## targethat    0    1
##           0 5551 1296
##           1  457  857
```

**Create a binned diagnostic plot of residuals vs prediction** There are definite patterns here, which bear investigating.

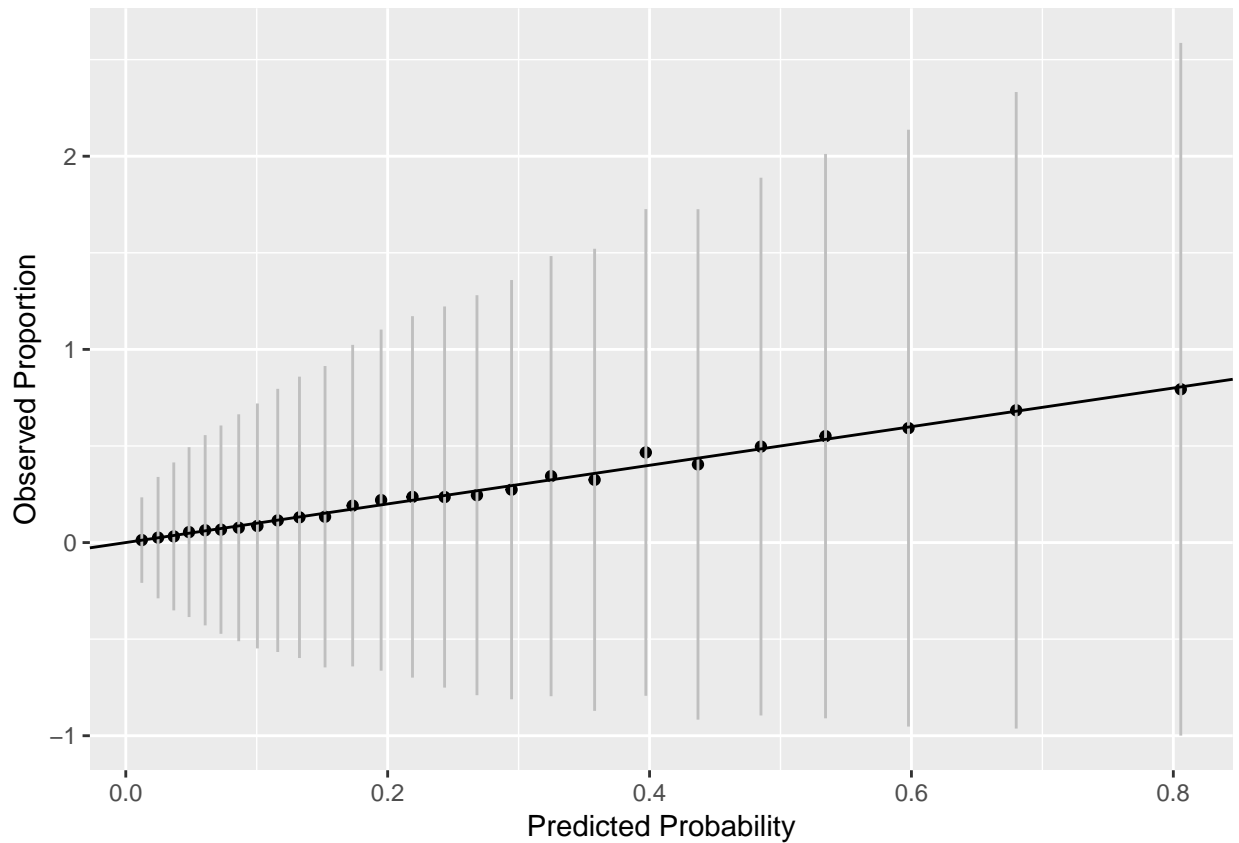


Plot leverages.



We don't see any strong outliers with the leverage plot. The points identified (3608,5686) are essentially in the plot of the line formed, so they are not likely pulling our model in any direction.

**Plot Goodness of fit**

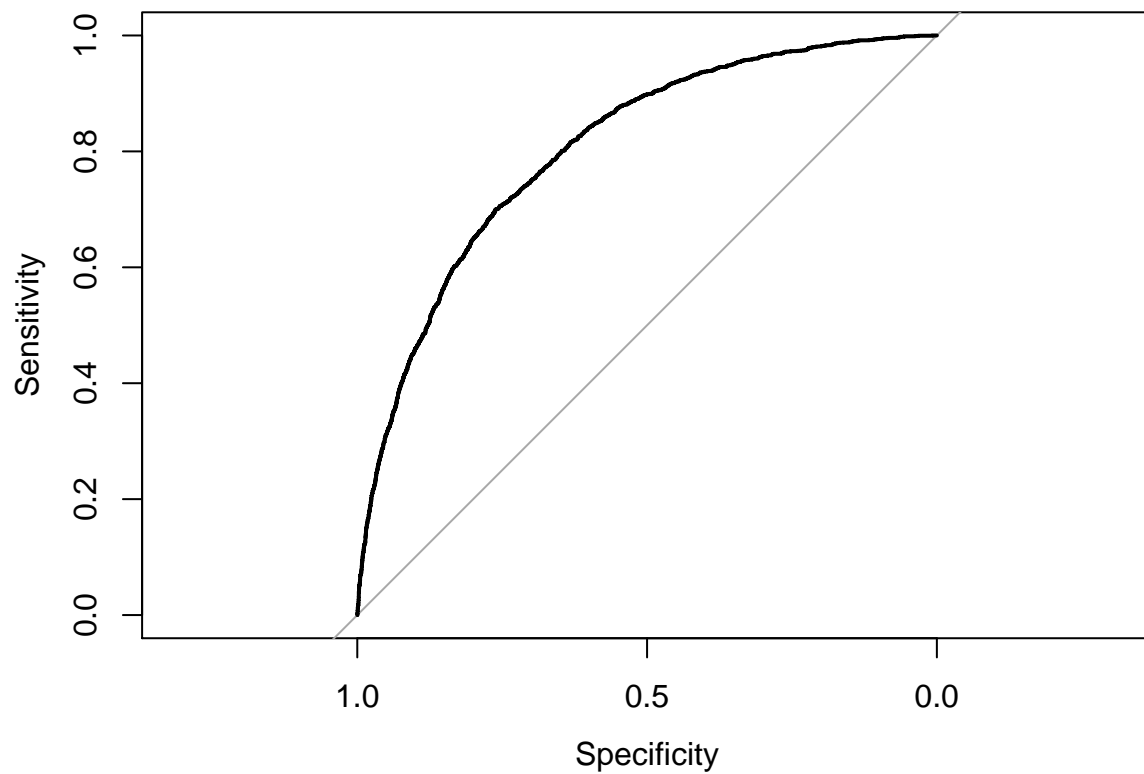


We see that our predictors fall close to the line.

### Model 3 - Srep AIC Model

#### ROC Curve

The ROC Curve helps measure true positives and true negative. A high AUC or area under the curve tells us the model is predicting well.

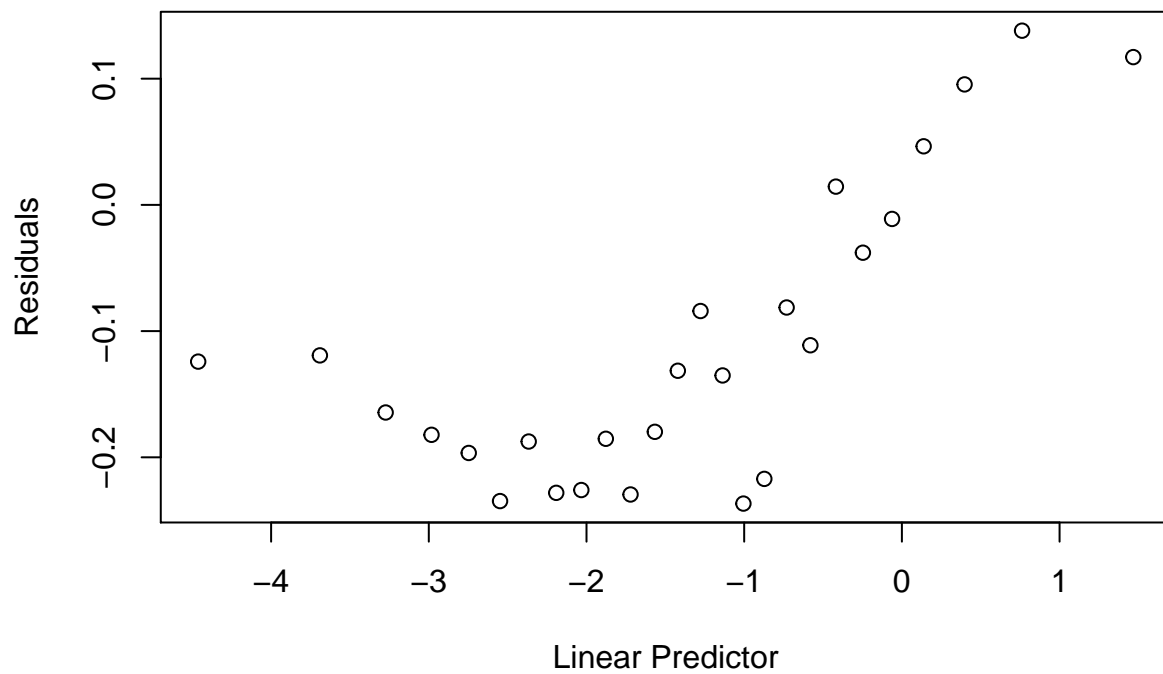


The AUC value of 0.8, tells us this model predicted values are accurate.

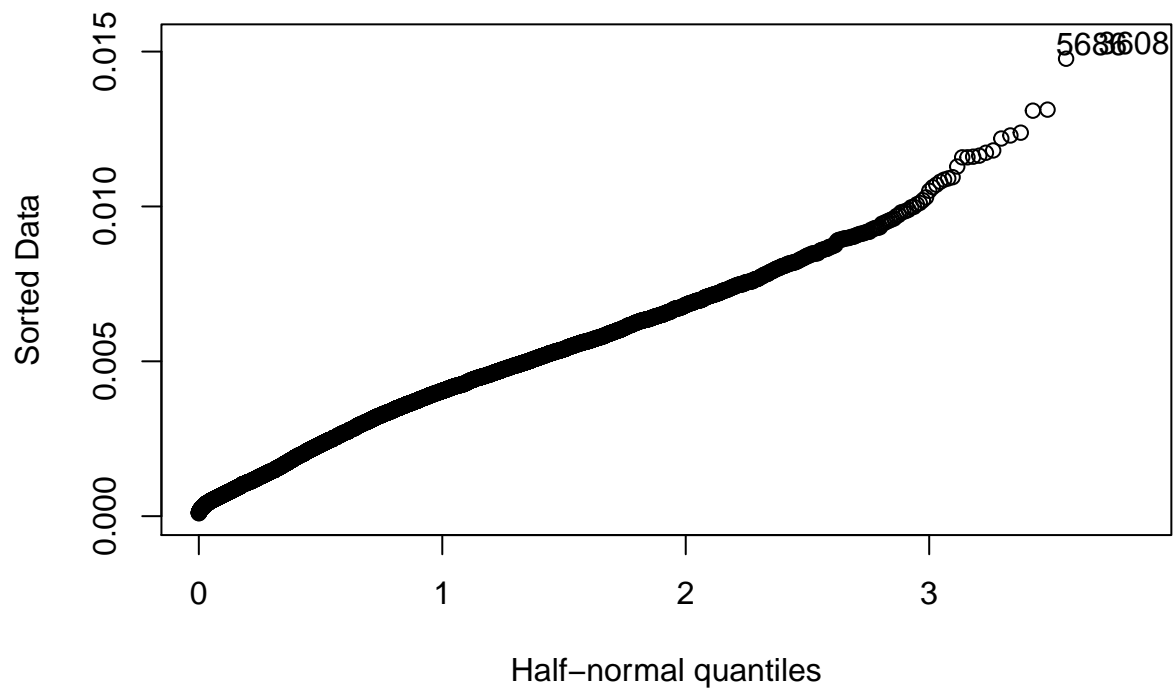
#### Confusion Matrix

```
##
## targethat    0    1
##           0 5556 1298
##           1  452  855
```

**Create a binned diagnostic plot of residuals vs prediction** There are definite patterns here, which bear investigating.

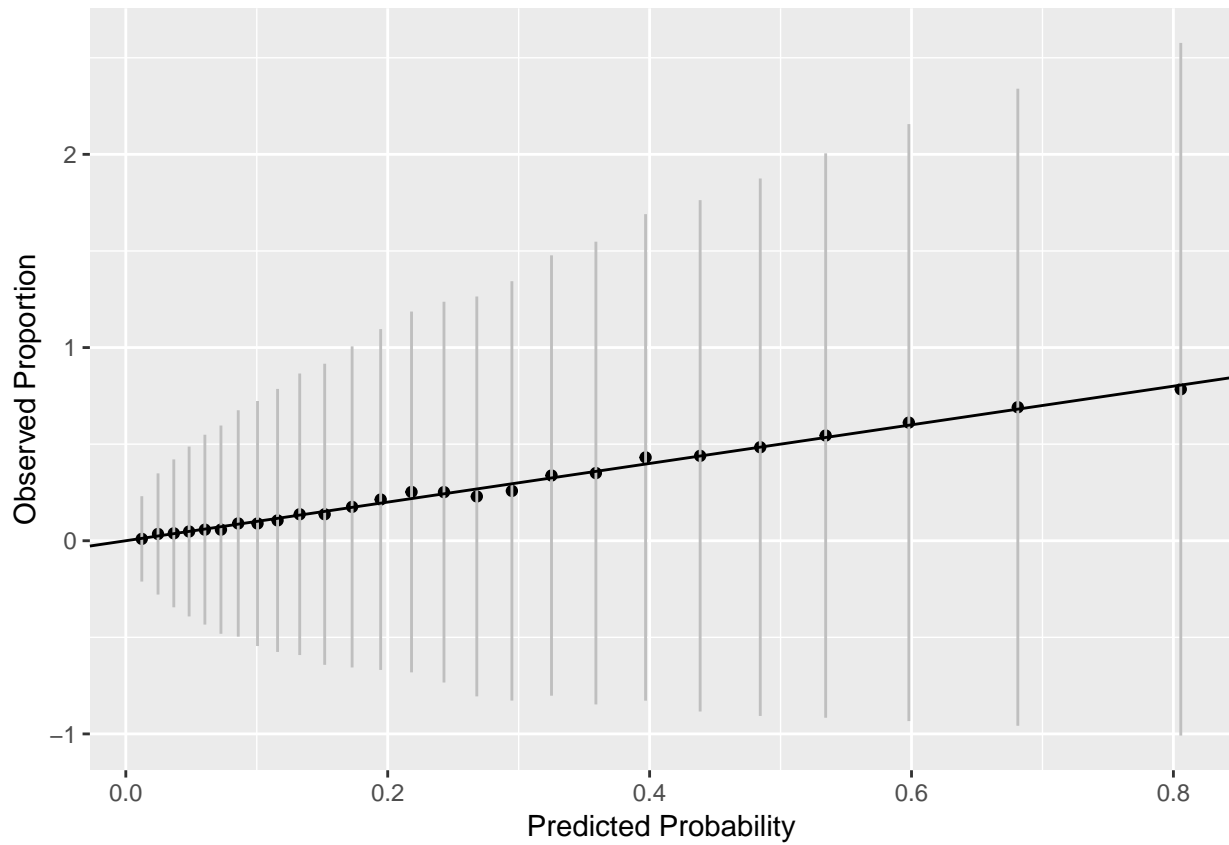


Plot leverages.



We don't see any strong outliers with the leverage plot. The points identified (3608,5686) are essentially in the plot of the line formed, so they are not likely pulling our model in any direction.

**Plot Goodness of fit**



We see that our predictors fall close to the line.

## TARGET\_AMT Modeling

Select Model

Compare Model Statistics

Conclusion

## APPENDIX