

GROUP 2 HW4: Insurance - Data 621 Assignment 4

GROUP 2 MEMBERS: Banu Boopalan, Gregg Maloy, Alexander Moyse, Umais Siddiqui

10/26/2024

Contents

Overview	2
Crash Data Insights	2
Data Exploration	3
Categorical variables	6
Numeric Variables	7
Assessment of Incomplete Data	9
Handling Missing Values And Correlation Analysis	10
Data Preparation for Multiple Linear Regression	14
Removing TARGET_FLAG	14
Handling Missing Data - Multiple Linear Regression	15
Transformations - Multiple Linear Regression	16
Build Models	33
Multiple Linear Regression	33
Binary Logistic Regression	37
Select Models & Prediction	37
Multiple Linear Regression Selection	37
Binary Logistic Regression Selection	37
Prediction	37
Code Appendix	37

Overview

In this assignment, you'll dive into a rich dataset of approximately 8,000 customer records from an auto insurance company. Each record represents a customer and includes two key response variables:

TARGET_FLAG - A binary indicator where a "1" signifies the customer was involved in a car crash, while a "0" means they were not. **TARGET_AMT** - This variable represents the cost incurred in the event of a crash. If there was no crash, this value is zero. If a crash occurred, this variable holds the associated monetary cost, which is greater than zero. Your goal is to develop predictive models that provide insights on two fronts:

The likelihood of a customer being involved in a car crash (using binary logistic regression). The potential cost of a crash, if it occurs (using multiple linear regression). For this task, you'll leverage the variables in the dataset—and any additional variables you derive from them—to create, train, and evaluate your models on a training dataset.

Dataset Variables Overview:

Below, you'll find a brief description of each variable in the dataset to help guide your exploratory analysis and feature engineering efforts.

Crash Data Insights

Target Variables

Attribute	Description	Expected Impact
TARGET_FLAG	Indicates if the customer was involved in a crash (1 = Yes, 0 = No)	None at this stage
TARGET_AMT	Cost incurred in the event of a crash (0 if no crash)	None at this stage

Predictor Variables

Attribute	Description	Theoretical Influence
AGE	Driver's age	Young and very old drivers may have higher risks
BLUEBOOK	Vehicle market value	May affect payout size if a crash occurs
CAR_AGE	Vehicle's age	Possibly influences payout but unclear on crash likelihood
CAR_TYPE	Vehicle type	Potential influence on payout if a crash occurs
CAR_USE	Vehicle's primary use	Commercial usage may increase crash probability
CLM_FREQ	Claims made in past 5 years	More past claims may predict higher future claims
EDUCATION	Highest education level attained	Higher education might correlate with safer driving
HOMEKIDS	Number of children at home	Impact unknown

Attribute	Description	Theoretical Influence
HOME_VAL	Value of home	Homeownership could correlate with responsible driving
INCOME	Annual income	Wealthier individuals may experience fewer crashes
JOB	Employment category	White-collar jobs might suggest safer driving
KIDSDRIV	Number of young drivers in household	Teen drivers could increase crash risk
MSTATUS	Marital status	Married individuals may drive more cautiously
MVR_PTS	Points on motor vehicle record	Higher points suggest increased crash likelihood
OLDCLAIM	Cumulative claims in past 5 years	High past payouts may predict future claims
PARENT1	Single-parent household indicator	Impact unknown
RED_CAR	Indicator for a red car	Potential correlation with risky driving (myth)
REVOKED	Past license revocation (in last 7 years)	Suggests increased risk
SEX	Driver's gender	Myth suggests women may experience fewer crashes
TIF	Policy duration (years)	Long-term policyholders may have safer driving patterns
TRAVTIME	Commute duration	Longer commutes may indicate higher risk
URBANICITY	Urban or rural setting	Impact unknown
YOJ	Years in current job	Stable employment may suggest safer driving habits

Data Exploration

```
## 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, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1~
## $ TARGET_AMT  <dbl> 0.000, 0.000, 0.000, 0.000, 0.000, 2946.000, 0.000, 4021.0~
## $ KIDSDRIV    <int> 0, 0, 0, 0, 0, 0, 0, 0, 1, 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, 0, 3, 0, 3, 2, 1~
## $ YOJ         <int> 11, 11, 10, 14, NA, 12, NA, 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", "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~
## $ TRAVTIME    <int> 14, 22, 5, 32, 36, 46, 33, 44, 34, 48, 15, 36, 25, 64, 48,~
## $ CAR_USE     <chr> "Private", "Commercial", "Private", "Private", "Private", ~
## $ BLUEBOOK    <chr> "$14,230", "$14,940", "$4,010", "$15,440", "$18,000", "$17~
## $ TIF         <int> 11, 1, 4, 7, 1, 1, 1, 1, 1, 7, 1, 7, 7, 6, 1, 6, 6, 7, 4, ~
## $ CAR_TYPE    <chr> "Minivan", "Minivan", "z_SUV", "Minivan", "z_SUV", "Sports~
## $ RED_CAR     <chr> "yes", "yes", "no", "yes", "no", "no", "no", "yes", "no", ~
## $ OLDCLAIM    <chr> "$4,461", "$0", "$38,690", "$0", "$19,217", "$0", "$0", "$~
## $ CLM_FREQ    <int> 2, 0, 2, 0, 2, 0, 0, 1, 0, 0, 0, 0, 2, 0, 0, 0, 0, 0, 0, 2~
```

```
## $ REVOKED      <chr> "No", "No", "No", "No", "Yes", "No", "No", "Yes", "No", "N~
## $ MVR_PTS      <int> 3, 0, 3, 0, 3, 0, 0, 10, 0, 1, 0, 0, 3, 3, 3, 0, 0, 0, ~
## $ CAR_AGE      <int> 18, 1, 10, 6, 17, 7, 1, 7, 1, 17, 11, 1, 9, 10, 5, 13, 16,~
## $ URBANICITY   <chr> "Highly Urban/ Urban", "Highly Urban/ Urban", "Highly Urba~
```

The dataset includes 8,161 records with 23 feature variables and 2 target variables, providing detailed information on customers and their insurance claims history.

```
## INDEX TARGET_FLAG TARGET_AMT KIDSDRIV AGE HOMEKIDS YOJ INCOME PARENT1
## 1 1 0 0 0 60 0 11 $67,349 No
## 2 2 0 0 0 43 0 11 $91,449 No
## 3 4 0 0 0 35 1 10 $16,039 No
## 4 5 0 0 0 51 0 14 No
## 5 6 0 0 0 50 0 NA $114,986 No
## 6 7 1 2946 0 34 1 12 $125,301 Yes
## HOME_VAL MSTATUS SEX EDUCATION JOB TRAVTIME CAR_USE BLUEBOOK
## 1 $0 z_No M PhD Professional 14 Private $14,230
## 2 $257,252 z_No M z_High School z_Blue Collar 22 Commercial $14,940
## 3 $124,191 Yes z_F z_High School Clerical 5 Private $4,010
## 4 $306,251 Yes M <High School z_Blue Collar 32 Private $15,440
## 5 $243,925 Yes z_F PhD Doctor 36 Private $18,000
## 6 $0 z_No z_F Bachelors z_Blue Collar 46 Commercial $17,430
## TIF CAR_TYPE RED_CAR OLDCLAIM CLM_FREQ REVOKED MVR_PTS CAR_AGE
## 1 11 Minivan yes $4,461 2 No 3 18
## 2 1 Minivan yes $0 0 No 0 1
## 3 4 z_SUV no $38,690 2 No 3 10
## 4 7 Minivan yes $0 0 No 0 6
## 5 1 z_SUV no $19,217 2 Yes 3 17
## 6 1 Sports Car no $0 0 No 0 7
## URBANICITY
## 1 Highly Urban/ Urban
## 2 Highly Urban/ Urban
## 3 Highly Urban/ Urban
## 4 Highly Urban/ Urban
## 5 Highly Urban/ Urban
## 6 Highly Urban/ Urban
```

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

```

##      PARENT1          HOME_VAL          MSTATUS          SEX
## Length:8161      Length:8161      Length:8161      Length:8161
## Class :character      Class :character      Class :character      Class :character
## Mode  :character      Mode  :character      Mode  :character      Mode  :character
##
##
##
##      EDUCATION          JOB          TRAVTIME          CAR_USE
## Length:8161      Length:8161      Min.   : 5.00      Length:8161
## Class :character      Class :character      1st Qu.: 22.00      Class :character
## Mode  :character      Mode  :character      Median : 33.00      Mode  :character
##                                     Mean   : 33.49
##                                     3rd Qu.: 44.00
##                                     Max.   :142.00
##
##      BLUEBOOK          TIF          CAR_TYPE          RED_CAR
## Length:8161      Min.   : 1.000      Length:8161      Length:8161
## Class :character      1st Qu.: 1.000      Class :character      Class :character
## Mode  :character      Median : 4.000      Mode  :character      Mode  :character
##                                     Mean   : 5.351
##                                     3rd Qu.: 7.000
##                                     Max.   :25.000
##
##      OLDCLAIM          CLM_FREQ          REVOKED          MVR_PTS
## Length:8161      Min.   :0.0000      Length:8161      Min.   : 0.000
## Class :character      1st Qu.:0.0000      Class :character      1st Qu.: 0.000
## Mode  :character      Median :0.0000      Mode  :character      Median : 1.000
##                                     Mean   :0.7986      Mean   : 1.696
##                                     3rd Qu.:2.0000      3rd Qu.: 3.000
##                                     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

```

On preliminary inspection, we note that several columns contain issues such as incompatible punctuation in financial values, and categorical variables require conversion to factors with clearer labels.

```

##      TARGET_FLAG          TARGET_AMT          KIDSDRIV          AGE
## Min.   :0.0000      Min.   : 0      Min.   :0.0000      Min.   :16.00
## 1st Qu.:0.0000      1st Qu.: 0      1st Qu.:0.0000      1st Qu.:39.00
## Median :0.0000      Median : 0      Median :0.0000      Median :45.00
## Mean   :0.2638      Mean   :1504      Mean   :0.1711      Mean   :44.79
## 3rd Qu.:1.0000      3rd Qu.:1036      3rd Qu.:0.0000      3rd Qu.:51.00
## Max.   :1.0000      Max.   :107586      Max.   :4.0000      Max.   :81.00
##                                     NA's   :6
##      HOMEKIDS          YOJ          INCOME          PARENT1          HOME_VAL
## Min.   :0.0000      Min.   : 0.0      Min.   : 0      No :7084      Min.   : 0

```

```

## 1st Qu.:0.0000 1st Qu.: 9.0 1st Qu.: 28097 Yes:1077 1st Qu.: 0
## Median :0.0000 Median :11.0 Median : 54028 Median :161160
## Mean :0.7212 Mean :10.5 Mean : 61898 Mean :154867
## 3rd Qu.:1.0000 3rd Qu.:13.0 3rd Qu.: 85986 3rd Qu.:238724
## Max. :5.0000 Max. :23.0 Max. :367030 Max. :885282
## NA's :454 NA's :445 NA's :464
## MSTATUS SEX EDUCATION JOB
## No :3267 F:4375 Bachelors :2242 Blue Collar :1825
## Yes:4894 M:3786 High School :2330 Clerical :1271
## Less than High School:1203 Professional:1117
## Masters :1658 Manager : 988
## PhD : 728 Lawyer : 835
## Student : 712
## (Other) :1413
## TRAVTIME CAR_USE BLUEBOOK TIF
## Min. : 5.00 Commercial:3029 Min. : 1500 Min. : 1.000
## 1st Qu.: 22.00 Private :5132 1st Qu.: 9280 1st Qu.: 1.000
## Median : 33.00 Median :14440 Median : 4.000
## Mean : 33.49 Mean :15710 Mean : 5.351
## 3rd Qu.: 44.00 3rd Qu.:20850 3rd Qu.: 7.000
## Max. :142.00 Max. :69740 Max. :25.000
##
## CAR_TYPE RED_CAR OLDCLAIM CLM_FREQ REVOKED
## Minivan :2145 no :5783 Min. : 0 Min. :0.0000 No :7161
## Panel Truck: 676 yes:2378 1st Qu.: 0 1st Qu.:0.0000 Yes:1000
## Pickup :1389 Median : 0 Median :0.0000
## Sports Car : 907 Mean : 4037 Mean :0.7986
## SUV :2294 3rd Qu.: 4636 3rd Qu.:2.0000
## Van : 750 Max. :57037 Max. :5.0000
##
## MVR_PTS CAR_AGE URBANICITY
## Min. : 0.000 Min. : -3.000 Highly Rural/ Rural:1669
## 1st Qu.: 0.000 1st Qu.: 1.000 Highly Urban/ Urban:6492
## Median : 1.000 Median : 8.000
## Mean : 1.696 Mean : 8.328
## 3rd Qu.: 3.000 3rd Qu.:12.000
## Max. :13.000 Max. :28.000
## NA's :510

```

The updated data frame now comprises only numeric and factor columns. It is observed that the car age variable contains values less than 1, including negative values. These will be replaced with a mode value of 1 to ensure data integrity.

Categorical variables

```
## Exploring Categorical Features:
```

```

## Feature: PARENT1
## Levels: No, Yes
##
## Feature: MSTATUS
## Levels: No, Yes
##

```

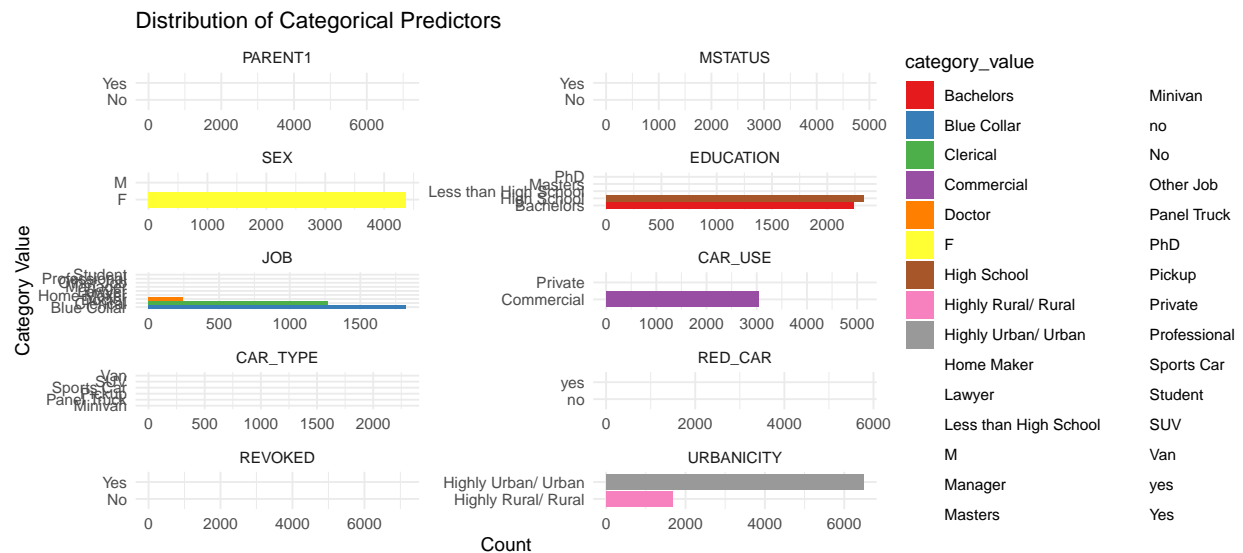
```

## Feature: SEX
## Levels: F, M
##
## Feature: EDUCATION
## Levels: Bachelors, High School, Less than High School, Masters, PhD
##
## Feature: JOB
## Levels: Blue Collar, Clerical, Doctor, Home Maker, Lawyer, Manager, Other Job, Professional, Student
##
## Feature: CAR_USE
## Levels: Commercial, Private
##
## Feature: CAR_TYPE
## Levels: Minivan, Panel Truck, Pickup, Sports Car, SUV, Van
##
## Feature: RED_CAR
## Levels: no, yes
##
## Feature: REVOKED
## Levels: No, Yes
##
## Feature: URBANICITY
## Levels: Highly Rural/ Rural, Highly Urban/ Urban

```

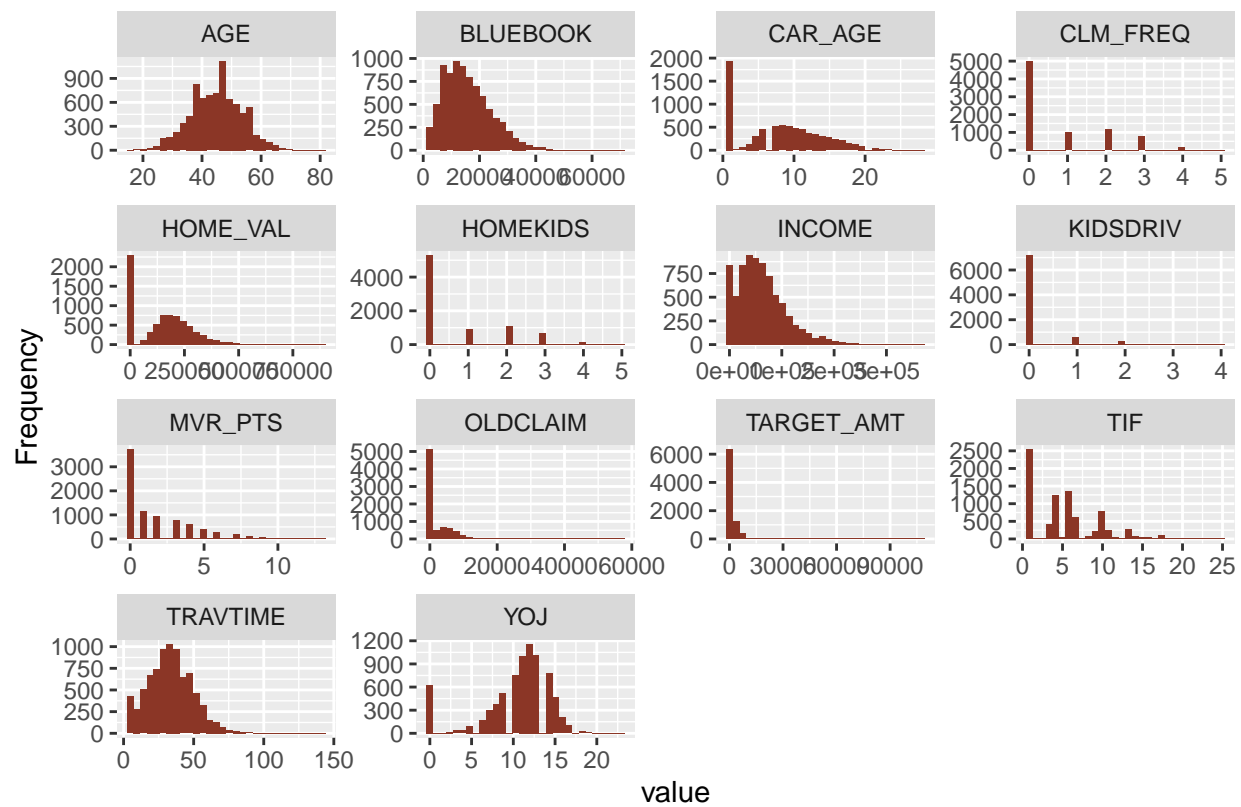
Upon examining the categorical variables, it is observed that the majority of the columns are binary in nature.

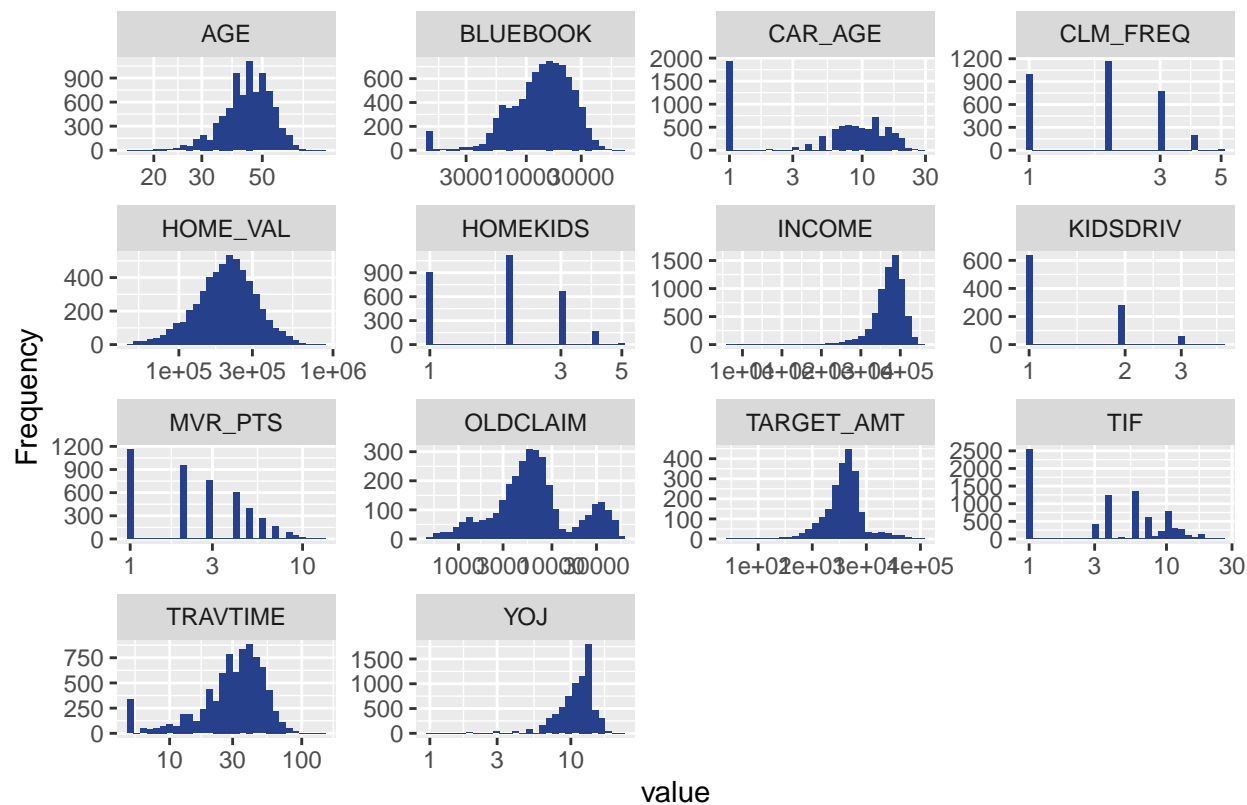
The following graphs illustrate the distribution of all categorical predictors.



Numeric Variables

The following two graphs illustrate the distribution of the numeric variables in our dataset. The first set of histograms, represented in red, displays the distributions on a normal scale, while the second set, depicted in blue, presents the distributions on a log10 scale. Notably, many numeric variables exhibit a mode value of zero, which may warrant further investigation.

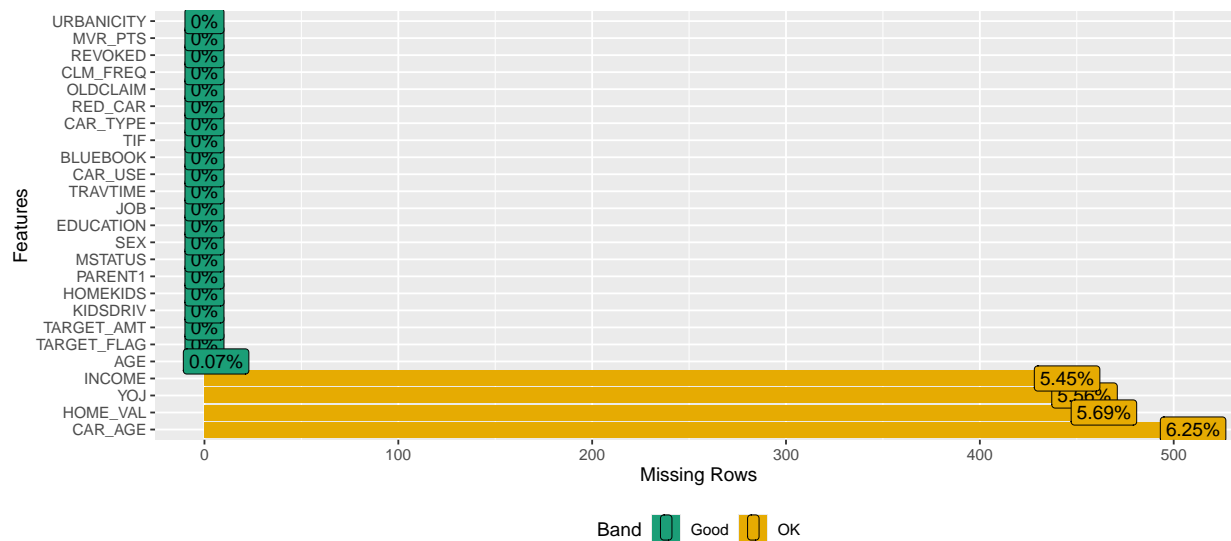




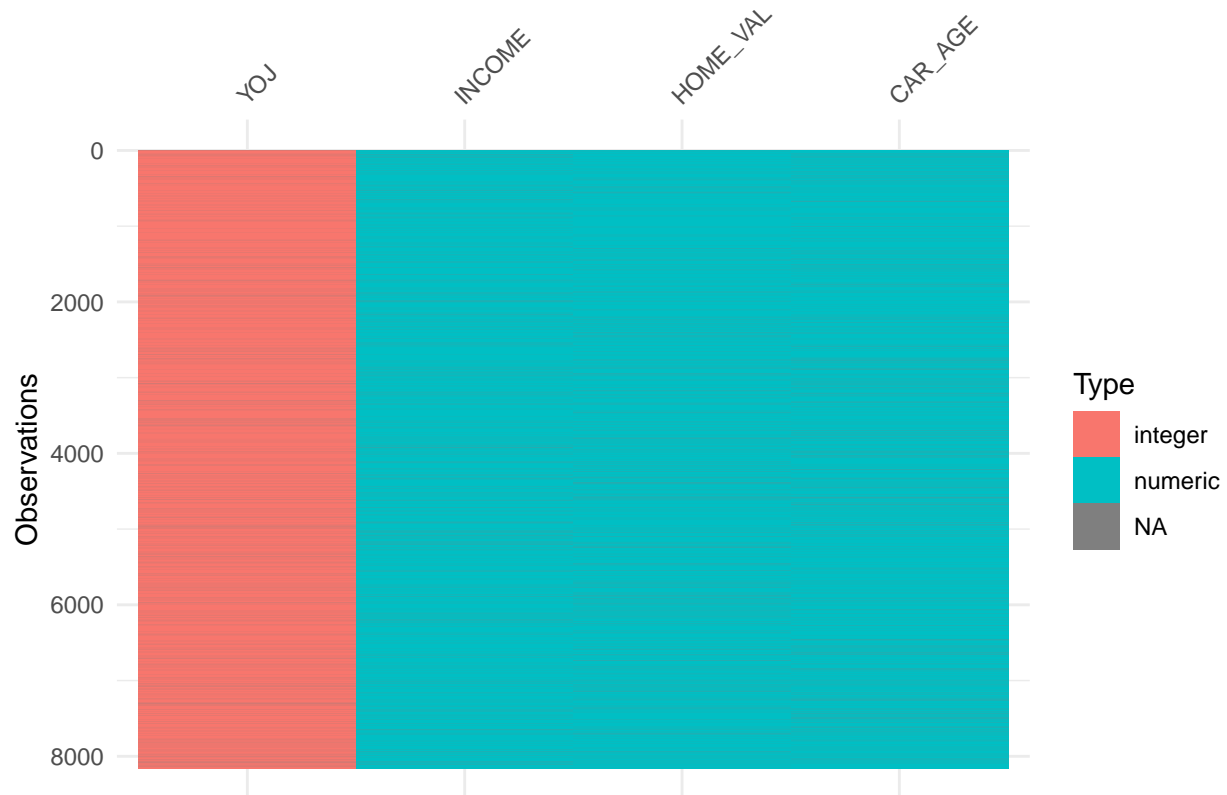
Assessment of Incomplete Data

This section identifies columns within the dataset that contain missing values, denoted as NA:

```
##   AGE YOJ INCOME HOME_VAL CAR_AGE
## 1   6 454   445   464   510
```



```
## TARGET_FLAG TARGET_AMT KIDSDRIV AGE HOMEKIDS YOJ
## 0.000 0.000 0.000 0.001 0.000 0.056
## INCOME PARENT1 HOME_VAL MSTATUS SEX EDUCATION
## 0.055 0.000 0.057 0.000 0.000 0.000
## JOB TRAVTIME CAR_USE BLUEBOOK TIF CAR_TYPE
## 0.000 0.000 0.000 0.000 0.000 0.000
## RED_CAR OLDCLAIM CLM_FREQ REVOKED MVR_PTS CAR_AGE
## 0.000 0.000 0.000 0.000 0.000 0.062
## URBANICITY
## 0.000
```



The analysis reveals that five variables contain missing values. However, there does not appear to be a discernible pattern associated with these missing entries, which suggests they are likely missing at random (MAR). This conclusion allows us to proceed with standard imputation techniques or analyses without significant concern regarding bias introduced by the missing data.

Handling Missing Values And Correlation Analysis

Multiple Imputation by Chained Equations (MICE) is a powerful method for handling missing data, as it generates multiple complete datasets by predicting missing values based on other available data. This method accounts for uncertainty in the imputations and allows for more reliable statistical inference.

```
## [1] "Missing Values Before Imputation:"
```

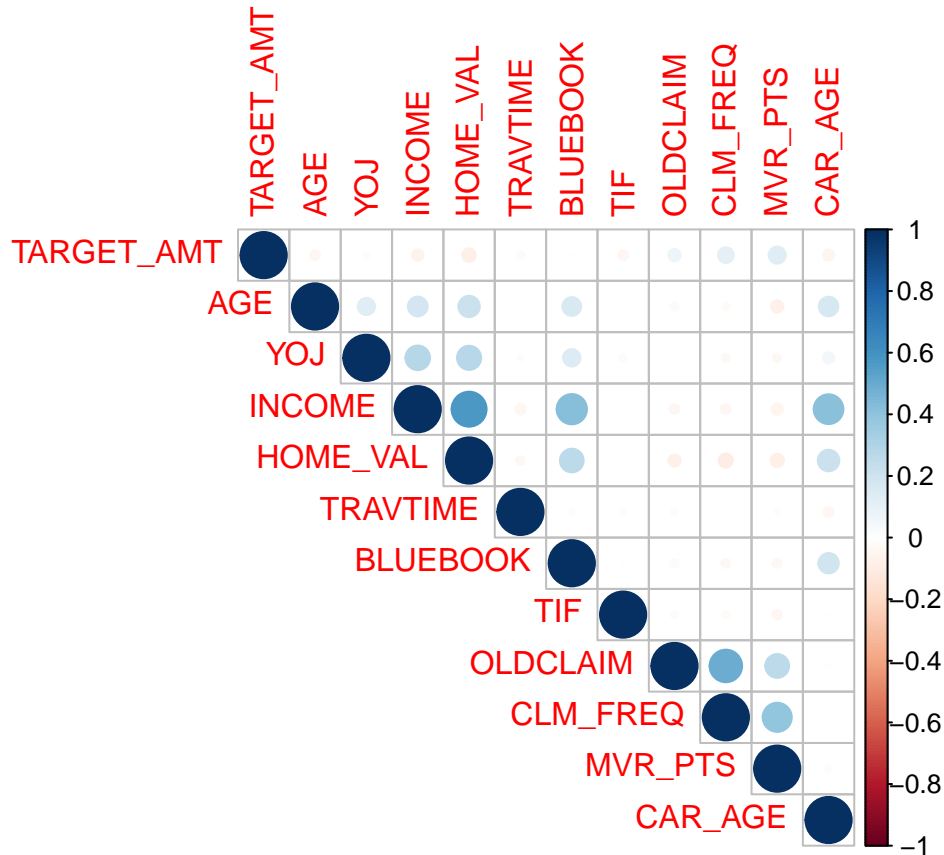
```
## TARGET_AMT AGE YOJ INCOME HOME_VAL TRAVTIME BLUEBOOK
```

```
##          0          6          454          445          464          0          0
##      TIF  OLDCLAIM  CLM_FREQ  MVR_PTS  CAR_AGE
##          0          0          0          0          510
```

```
##
## iter imp variable
## 1 1 AGE YOJ INCOME HOME_VAL CAR_AGE
## 1 2 AGE YOJ INCOME HOME_VAL CAR_AGE
## 1 3 AGE YOJ INCOME HOME_VAL CAR_AGE
## 1 4 AGE YOJ INCOME HOME_VAL CAR_AGE
## 1 5 AGE YOJ INCOME HOME_VAL CAR_AGE
## 2 1 AGE YOJ INCOME HOME_VAL CAR_AGE
## 2 2 AGE YOJ INCOME HOME_VAL CAR_AGE
## 2 3 AGE YOJ INCOME HOME_VAL CAR_AGE
## 2 4 AGE YOJ INCOME HOME_VAL CAR_AGE
## 2 5 AGE YOJ INCOME HOME_VAL CAR_AGE
## 3 1 AGE YOJ INCOME HOME_VAL CAR_AGE
## 3 2 AGE YOJ INCOME HOME_VAL CAR_AGE
## 3 3 AGE YOJ INCOME HOME_VAL CAR_AGE
## 3 4 AGE YOJ INCOME HOME_VAL CAR_AGE
## 3 5 AGE YOJ INCOME HOME_VAL CAR_AGE
## 4 1 AGE YOJ INCOME HOME_VAL CAR_AGE
## 4 2 AGE YOJ INCOME HOME_VAL CAR_AGE
## 4 3 AGE YOJ INCOME HOME_VAL CAR_AGE
## 4 4 AGE YOJ INCOME HOME_VAL CAR_AGE
## 4 5 AGE YOJ INCOME HOME_VAL CAR_AGE
## 5 1 AGE YOJ INCOME HOME_VAL CAR_AGE
## 5 2 AGE YOJ INCOME HOME_VAL CAR_AGE
## 5 3 AGE YOJ INCOME HOME_VAL CAR_AGE
## 5 4 AGE YOJ INCOME HOME_VAL CAR_AGE
## 5 5 AGE YOJ INCOME HOME_VAL CAR_AGE
```

```
## [1] "Missing Values After Imputation:"
```

```
## TARGET_AMT      AGE      YOJ      INCOME  HOME_VAL  TRAVTIME  BLUEBOOK
##          0          0          0          0          0          0          0
##      TIF  OLDCLAIM  CLM_FREQ  MVR_PTS  CAR_AGE
##          0          0          0          0          0
```



```
## [1] "Correlation Matrix for Complete Case Analysis:"
```

```
##          TARGET_AMT      AGE      YOJ      INCOME      HOME_VAL
## TARGET_AMT  1.000000000 -0.052348528 -0.022196571 -0.0562601493 -0.09056112
## AGE         -0.052348528  1.000000000  0.137847876  0.1876862059  0.21598562
## YOJ         -0.022196571  0.137847876  1.000000000  0.2783277152  0.26980907
## INCOME      -0.056260149  0.187686206  0.278327715  1.0000000000  0.57970674
## HOME_VAL    -0.090561124  0.215985625  0.269809074  0.5797067363  1.00000000
## TRAVTIME     0.032287806  0.007807727 -0.015740963 -0.0413200825 -0.03014163
## BLUEBOOK    -0.003183645  0.171170247  0.136335894  0.4332521829  0.26161690
## TIF         -0.041860052  0.000408708  0.030813700  0.0007376252 -0.00460570
## OLDCLAIM     0.080067386 -0.030707066  0.001634368 -0.0377131052 -0.05863833
## CLM_FREQ     0.116939123 -0.027125254 -0.028669411 -0.0451604051 -0.09695212
## MVR_PTS      0.137030840 -0.075556608 -0.035432609 -0.0709892627 -0.09418684
## CAR_AGE     -0.062828101  0.184019005  0.057768248  0.4117386242  0.21531374
##          TRAVTIME      BLUEBOOK      TIF      OLDCLAIM      CLM_FREQ
## TARGET_AMT  0.032287806 -0.003183645 -0.0418600523  0.080067386  0.116939123
## AGE         0.007807727  0.171170247  0.0004087080 -0.030707066 -0.027125254
## YOJ        -0.015740963  0.136335894  0.0308136996  0.001634368 -0.028669411
## INCOME     -0.041320082  0.433252183  0.0007376252 -0.037713105 -0.045160405
## HOME_VAL    -0.030141625  0.261616901 -0.0046056998 -0.058638327 -0.096952119
## TRAVTIME     1.000000000 -0.010979136 -0.0117716399 -0.022707967  0.009510331
## BLUEBOOK    -0.010979136  1.000000000  0.0045237917 -0.032654587 -0.046002944
## TIF         -0.011771640  0.004523792  1.0000000000 -0.018249702 -0.023758956
## OLDCLAIM    -0.022707967 -0.032654587 -0.0182497019  1.000000000  0.494017156
```

```

## CLM_FREQ      0.009510331 -0.046002944 -0.0237589564  0.494017156  1.000000000
## MVR_PTS       0.003815401 -0.061216939 -0.0380976659  0.272706265  0.397847352
## CAR_AGE      -0.030726192  0.185550420  0.0124643954 -0.010610234 -0.006339303
##              MVR_PTS      CAR_AGE
## TARGET_AMT    0.137030840 -0.062828101
## AGE           -0.075556608  0.184019005
## YOJ           -0.035432609  0.057768248
## INCOME        -0.070989263  0.411738624
## HOME_VAL      -0.094186838  0.215313740
## TRAVTIME      0.003815401 -0.030726192
## BLUEBOOK     -0.061216939  0.185550420
## TIF           -0.038097666  0.012464395
## OLDCLAIM      0.272706265 -0.010610234
## CLM_FREQ      0.397847352 -0.006339303
## MVR_PTS       1.000000000 -0.023995843
## CAR_AGE      -0.023995843  1.000000000

```

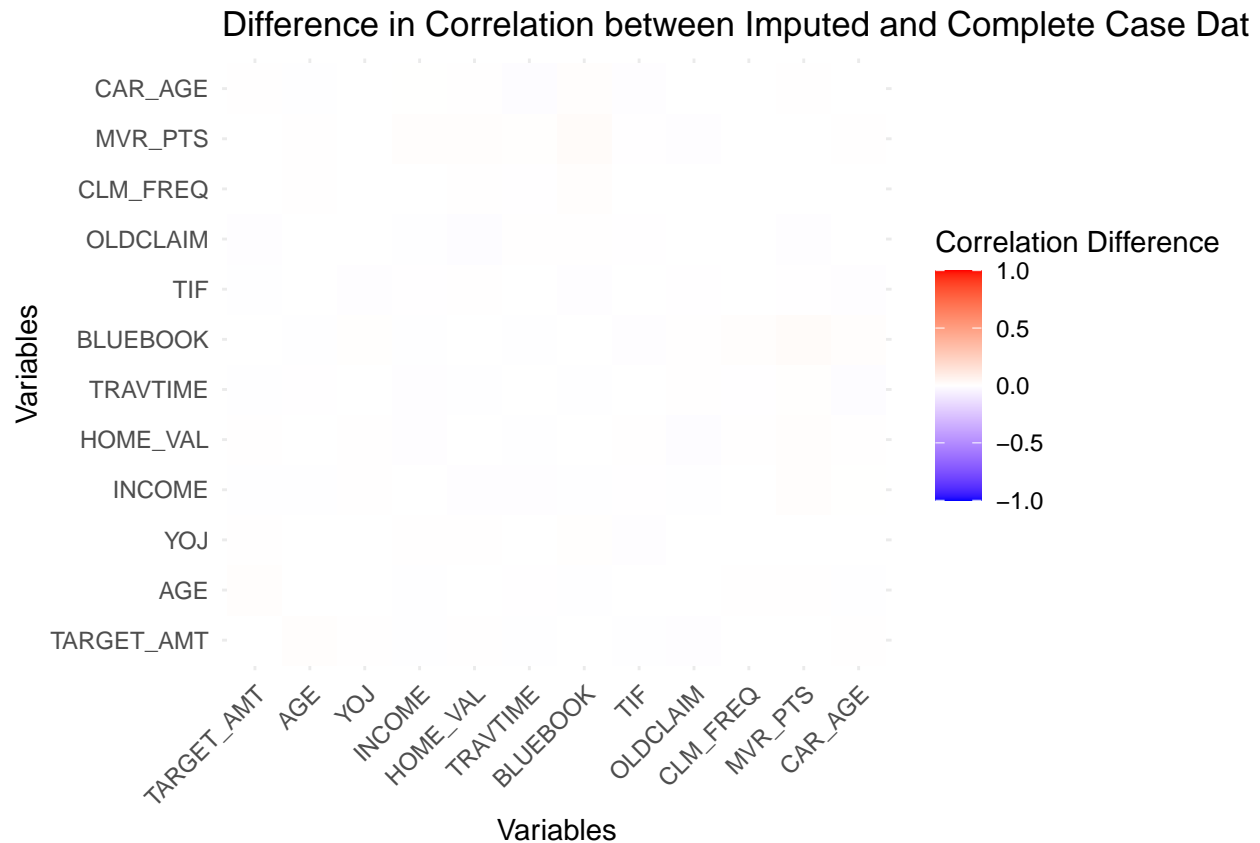
```
## [1] "Correlation Matrix for Imputed Data:"
```

```

##              TARGET_AMT      AGE      YOJ      INCOME      HOME_VAL
## TARGET_AMT  1.000000000 -0.041826419 -0.017860070 -0.060983939 -0.0861878936
## AGE         -0.041826419  1.000000000  0.138761497  0.182928284  0.2143812514
## YOJ         -0.017860070  0.1387614968  1.000000000  0.282659508  0.2733640677
## INCOME      -0.060983939  0.1829282842  0.282659508  1.000000000  0.5723473522
## HOME_VAL    -0.086187894  0.2143812514  0.273364068  0.572347352  1.0000000000
## TRAVTIME    0.027987016  0.0053547772 -0.016038747 -0.048890357 -0.0352608342
## BLUEBOOK   -0.004699523  0.1651777923  0.142660165  0.428970852  0.2630568135
## TIF         -0.046480831 -0.0003363674  0.024330425 -0.002846146  0.0006303218
## OLDCLAIM    0.070953287 -0.0297096301  0.001866237 -0.042264940 -0.0701071116
## CLM_FREQ    0.116419159 -0.0239127328 -0.030361314 -0.044365798 -0.0920016863
## MVR_PTS     0.137865509 -0.0717218955 -0.034559684 -0.058716119 -0.0830885507
## CAR_AGE     -0.058658346  0.1791602948  0.057592905  0.413684204  0.2182023139
##              TRAVTIME    BLUEBOOK      TIF      OLDCLAIM      CLM_FREQ
## TARGET_AMT  0.027987016 -0.004699523 -0.0464808306  0.070953287  0.116419159
## AGE         0.005354777  0.165177792 -0.0003363674 -0.029709630 -0.023912733
## YOJ         -0.016038747  0.142660165  0.0243304249  0.001866237 -0.030361314
## INCOME      -0.048890357  0.428970852 -0.0028461456 -0.042264940 -0.044365798
## HOME_VAL    -0.035260834  0.263056814  0.0006303218 -0.070107112 -0.092001686
## TRAVTIME    1.000000000 -0.017001298 -0.0116046256 -0.019267169  0.006560211
## BLUEBOOK   -0.017001298  1.000000000 -0.0054245723 -0.029517568 -0.036341497
## TIF         -0.011604626 -0.005424572  1.0000000000 -0.021958198 -0.023022955
## OLDCLAIM    -0.019267169 -0.029517568 -0.0219581980  1.000000000  0.495130810
## CLM_FREQ    0.006560211 -0.036341497 -0.0230229550  0.495130810  1.000000000
## MVR_PTS     0.010598511 -0.039130846 -0.0410457340  0.264485025  0.396638373
## CAR_AGE     -0.042936990  0.195606786  0.0058816228 -0.009906046 -0.005909096
##              MVR_PTS      CAR_AGE
## TARGET_AMT  0.13786551 -0.058658346
## AGE         -0.07172190  0.179160295
## YOJ         -0.03455968  0.057592905
## INCOME      -0.05871612  0.413684204
## HOME_VAL    -0.08308855  0.218202314
## TRAVTIME    0.01059851 -0.042936990
## BLUEBOOK   -0.03913085  0.195606786
## TIF         -0.04104573  0.005881623

```

```
## OLDCLAIM    0.26448503 -0.009906046
## CLM_FREQ    0.39663837 -0.005909096
## MVR_PTS     1.00000000 -0.018823869
## CAR_AGE     -0.01882387  1.000000000
```



After completing the data, we have calculated the correlation matrix on the fully imputed dataset. This provides a more accurate representation of the relationships between variables without the bias that could be introduced by simple imputation methods.

It is evident that there are notable positive correlations among the following variables:

Income and Home Value Income and Bluebook Value Income and Car Age Claim Frequency and Old Claims Claim Frequency and MVR Points

The heatmap provides a visual representation of the differences in correlations between the imputed data and complete case data, helping to understand the impact of the missing data handling method.

Data Preparation for Multiple Linear Regression

Removing TARGET_FLAG

Since, for multiple linear regression our objective is to predict the monetary amount of how much it will cost in the event of a crash, we will exclude the TARGET_FLAG variable from our analysis.

Handling Missing Data - Multiple Linear Regression

Before proceeding with imputation, let's assess the missing values in our dataset. We will then handle the missing data using multiple imputation, which is a more robust method than simply replacing missing values with the median.

```
## [1] "Missing Values Before Imputation:"
```

```
## TARGET_AMT  KIDSDRIV      AGE  HOMEKIDS      YOJ      INCOME  PARENT1
##           0           0        5          0      123        110          0
##  HOME_VAL    MSTATUS      SEX  EDUCATION      JOB    TRAVTIME  CAR_USE
##        121          0        0          0        0          0          0
## BLUEBOOK     TIF    CAR_TYPE    RED_CAR  OLDCLAIM  CLM_FREQ  REVOKED
##          0          0        0          0        0          0          0
##  MVR_PTS    CAR_AGE  URBANICITY
##          0        142          0
```

```
##
## iter imp variable
##  1  1 AGE YOJ INCOME HOME_VAL CAR_AGE
##  1  2 AGE YOJ INCOME HOME_VAL CAR_AGE
##  1  3 AGE YOJ INCOME HOME_VAL CAR_AGE
##  1  4 AGE YOJ INCOME HOME_VAL CAR_AGE
##  1  5 AGE YOJ INCOME HOME_VAL CAR_AGE
##  2  1 AGE YOJ INCOME HOME_VAL CAR_AGE
##  2  2 AGE YOJ INCOME HOME_VAL CAR_AGE
##  2  3 AGE YOJ INCOME HOME_VAL CAR_AGE
##  2  4 AGE YOJ INCOME HOME_VAL CAR_AGE
##  2  5 AGE YOJ INCOME HOME_VAL CAR_AGE
##  3  1 AGE YOJ INCOME HOME_VAL CAR_AGE
##  3  2 AGE YOJ INCOME HOME_VAL CAR_AGE
##  3  3 AGE YOJ INCOME HOME_VAL CAR_AGE
##  3  4 AGE YOJ INCOME HOME_VAL CAR_AGE
##  3  5 AGE YOJ INCOME HOME_VAL CAR_AGE
##  4  1 AGE YOJ INCOME HOME_VAL CAR_AGE
##  4  2 AGE YOJ INCOME HOME_VAL CAR_AGE
##  4  3 AGE YOJ INCOME HOME_VAL CAR_AGE
##  4  4 AGE YOJ INCOME HOME_VAL CAR_AGE
##  4  5 AGE YOJ INCOME HOME_VAL CAR_AGE
##  5  1 AGE YOJ INCOME HOME_VAL CAR_AGE
##  5  2 AGE YOJ INCOME HOME_VAL CAR_AGE
##  5  3 AGE YOJ INCOME HOME_VAL CAR_AGE
##  5  4 AGE YOJ INCOME HOME_VAL CAR_AGE
##  5  5 AGE YOJ INCOME HOME_VAL CAR_AGE
```

```
## [1] "Missing Values After Imputation:"
```

```
## TARGET_AMT  KIDSDRIV      AGE  HOMEKIDS      YOJ      INCOME  PARENT1
##           0           0        0          0        0          0          0
##  HOME_VAL    MSTATUS      SEX  EDUCATION      JOB    TRAVTIME  CAR_USE
##          0           0        0          0        0          0          0
## BLUEBOOK     TIF    CAR_TYPE    RED_CAR  OLDCLAIM  CLM_FREQ  REVOKED
##          0           0        0          0        0          0          0
##  MVR_PTS    CAR_AGE  URBANICITY
##          0           0          0
```

Transformations - Multiple Linear Regression

We will be performing transformations and create histograms for several variables, which helps visualize the effect of the transformations on data distribution. Here's a breakdown of how these transformations aid in model building and potential outcomes:

Handling Skewness:

Many of these variables (e.g., INCOME, HOME_VAL, OLDCLAIM) may be right-skewed due to outliers or a large range of values. Transformations like log, square root, and Yeo-Johnson help normalize the distribution, reducing skewness. Normalized distributions (closer to normal) are beneficial for regression-based models, as they assume linear relationships and normally distributed residuals.

Improving Model Fit:

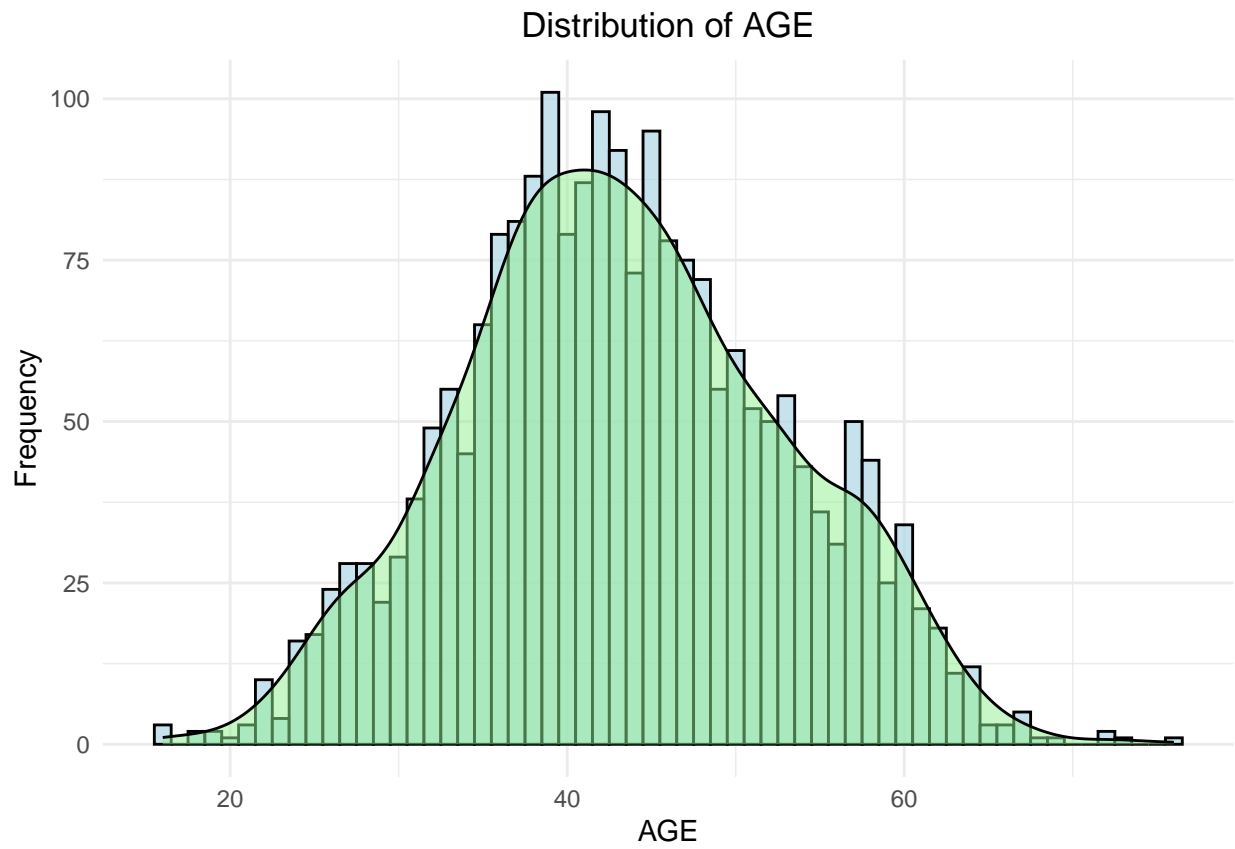
Log and Square Root transformations compress the range of values, which can make the data easier for linear models to handle. For instance, high-income values may dominate the predictive power of INCOME if not transformed. Box-Cox and Yeo-Johnson transformations (which automatically choose an optimal transformation) can help produce more linearly related predictors, which improves linear regression model accuracy.

Comparing the Effect of Transformations:

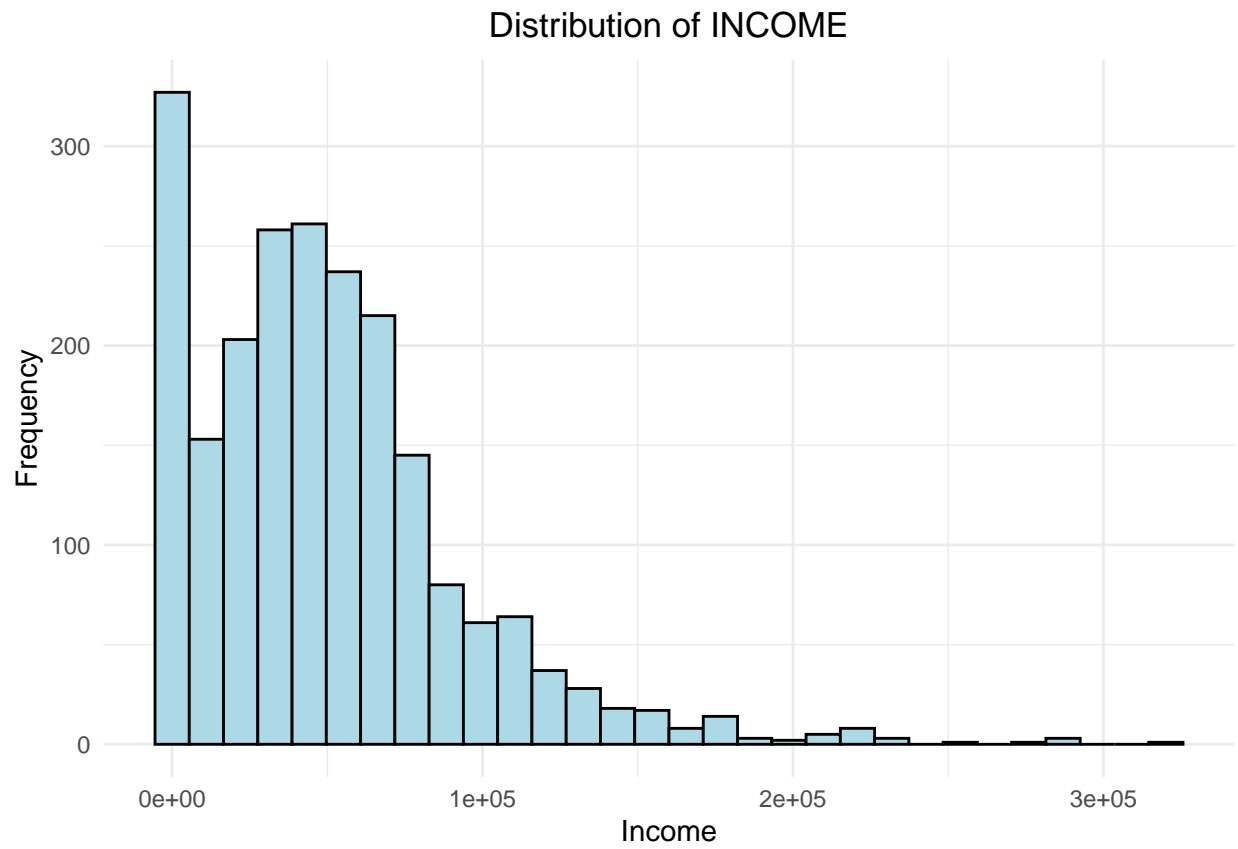
Creating side-by-side histograms allows you to compare the original and transformed distributions. This visual analysis is important for selecting the transformation that brings the distribution closest to normality, which can ultimately improve the performance and interpretability of the model.

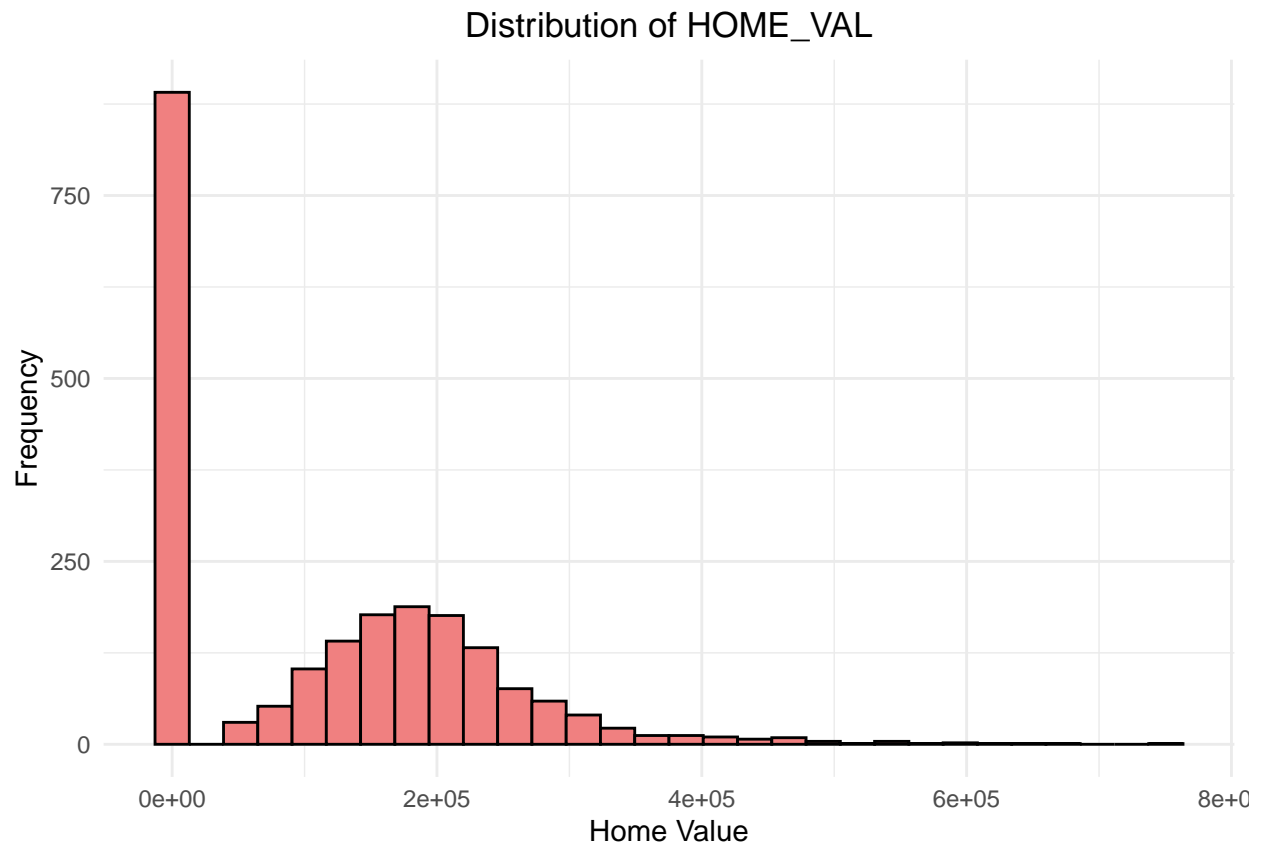
Categorizing Continuous Variables:

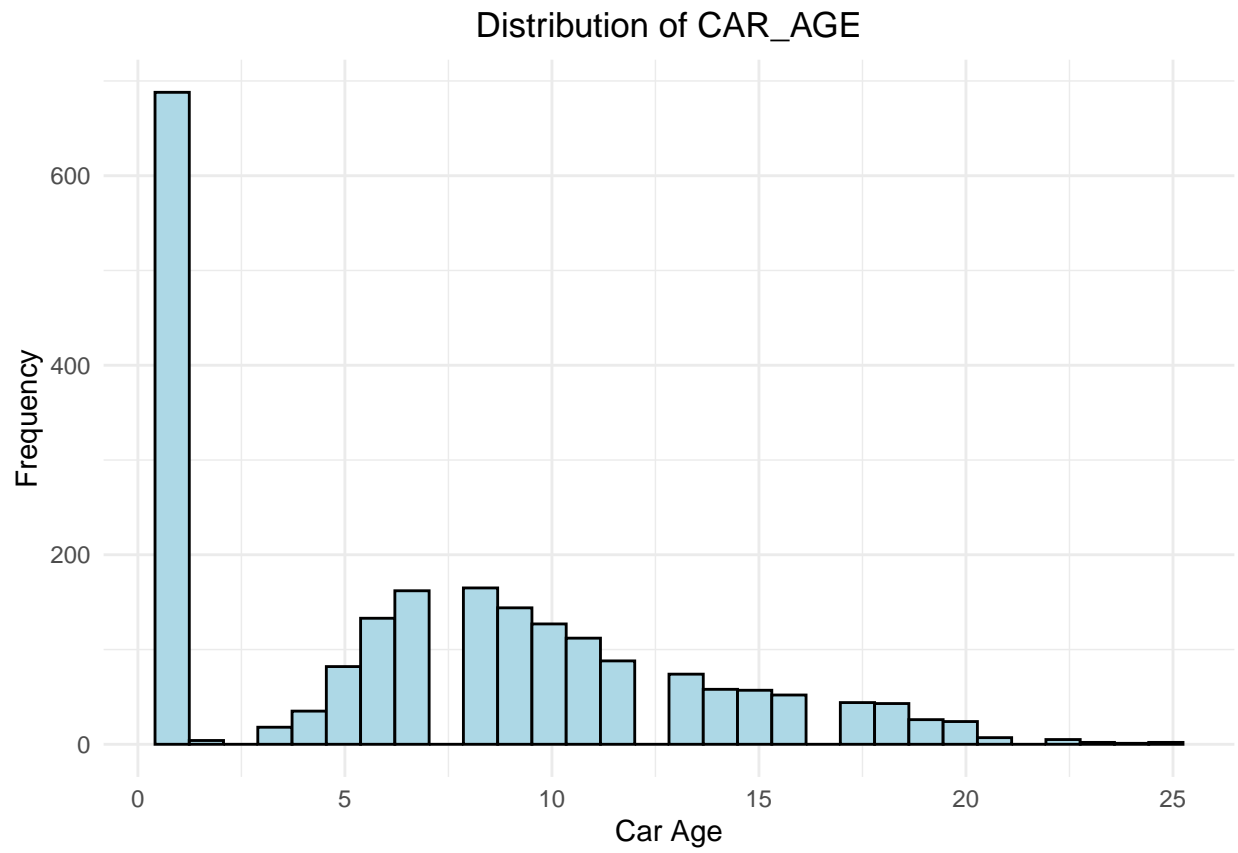
The cut function is used to create binned categories for TIF (Years with Policy) and MVR_PTS (Driving Record Points), which converts continuous variables into categorical bins. This is useful if there are distinct groups within the data that are meaningful (e.g., "Less than 1 year" in TIF). Using Transformed Variables for Modeling After determining the most effective transformation for each variable, we can replace the original variables with the transformed ones in our model. However, it's also useful to keep both versions to allow for comparison in model performance. Here's how to proceed with this:

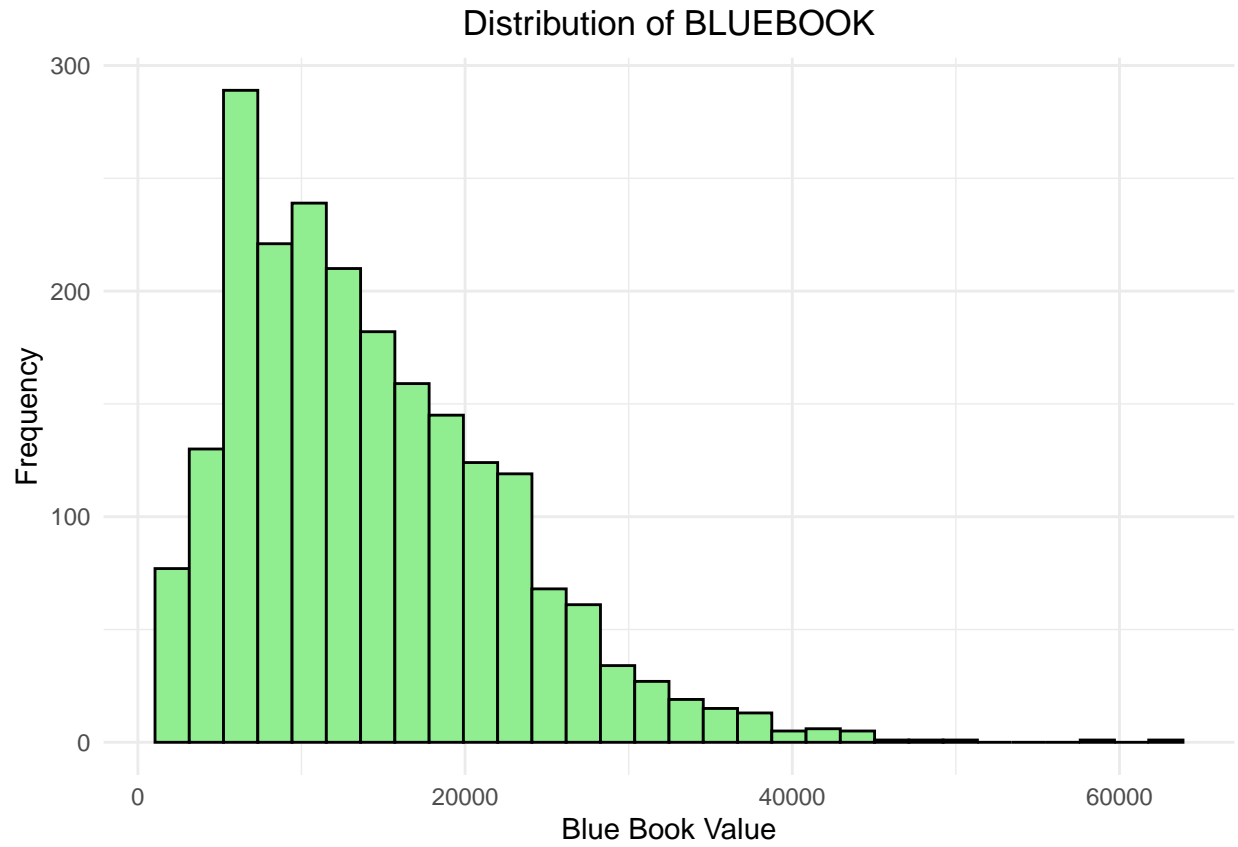


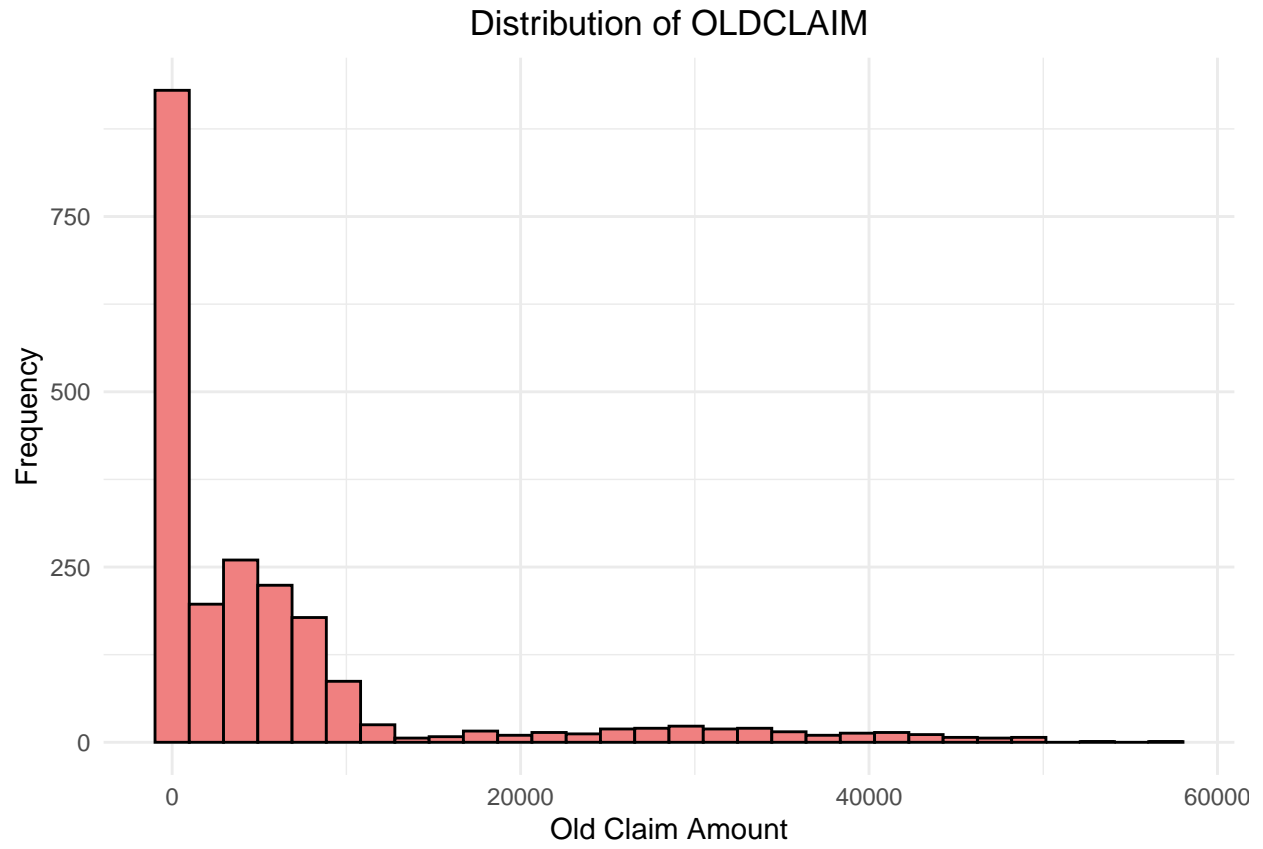
r

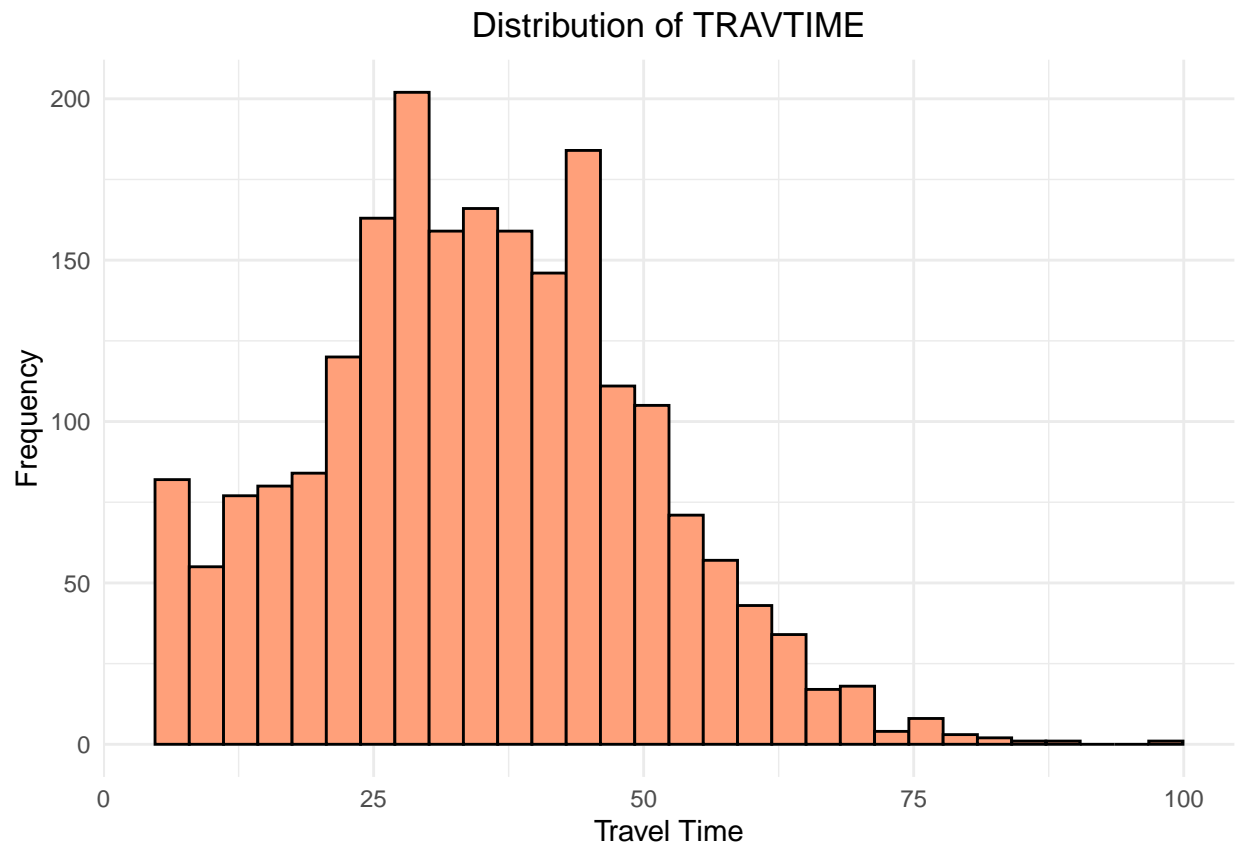


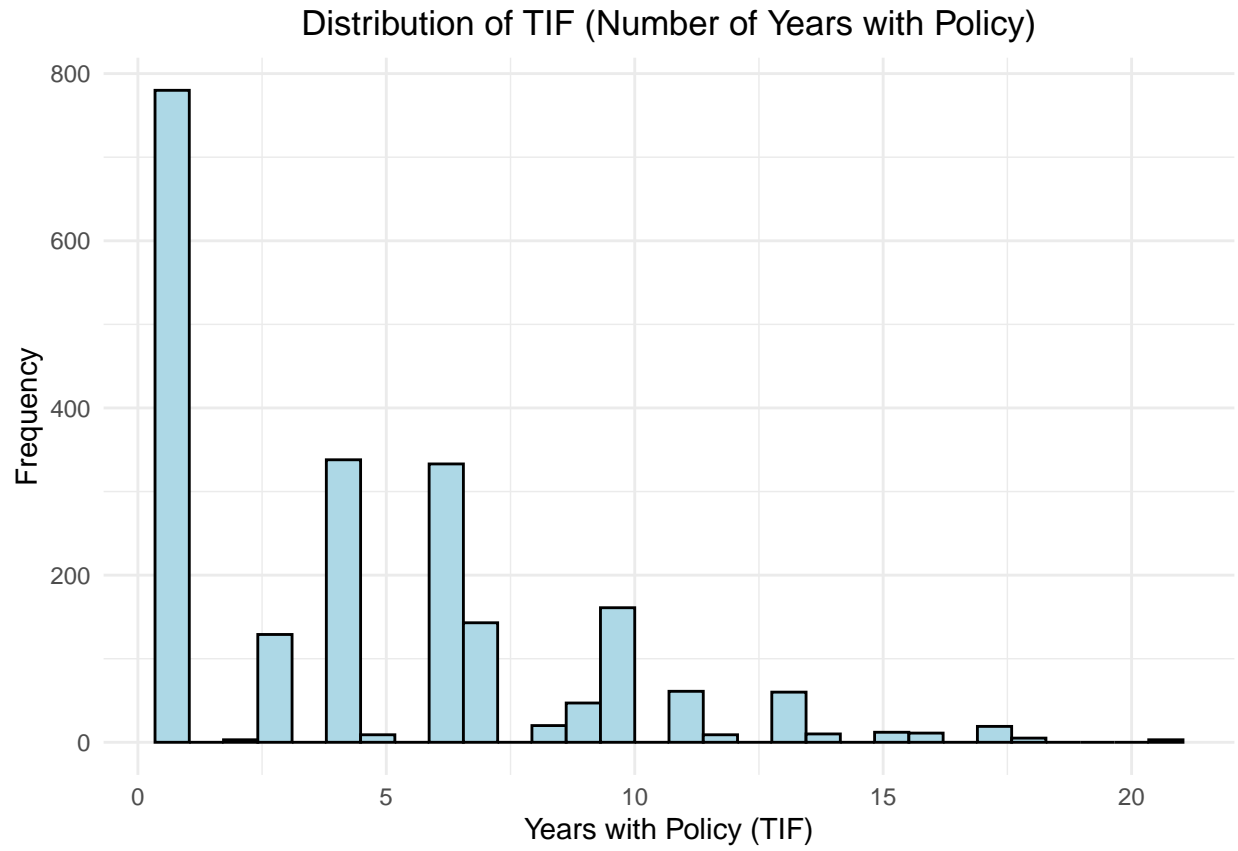


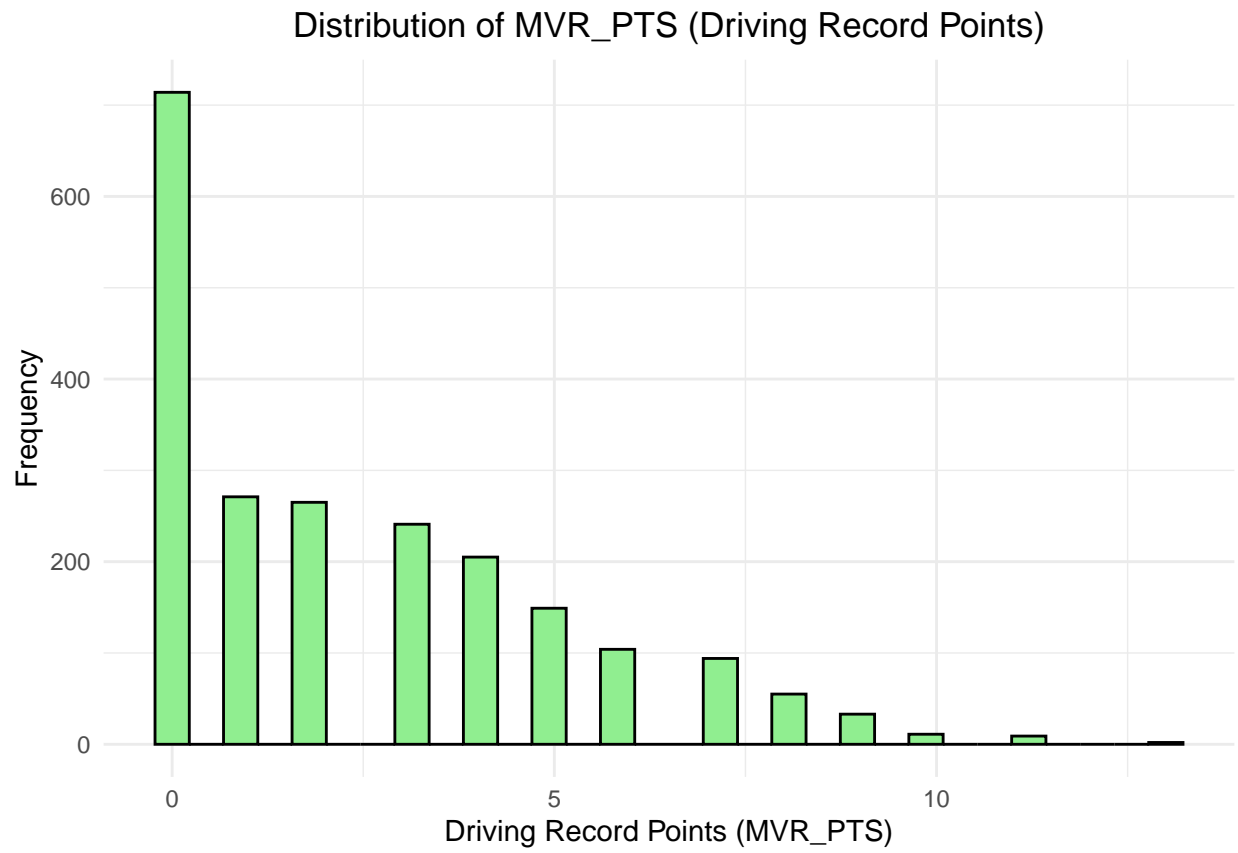


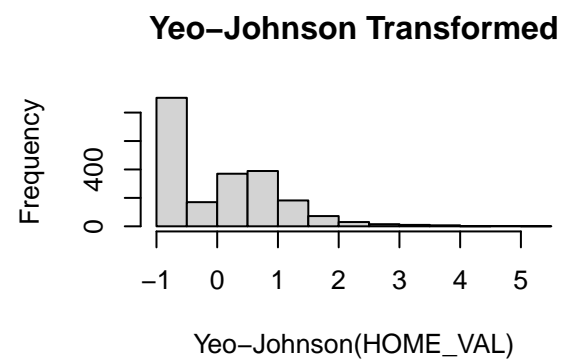
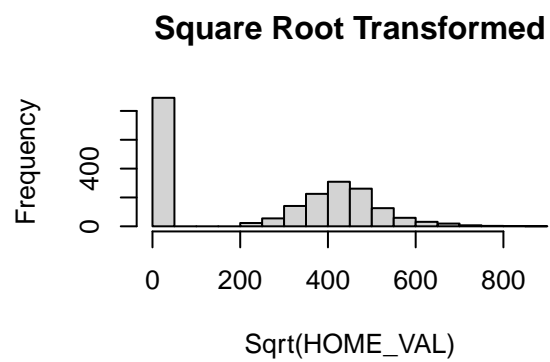
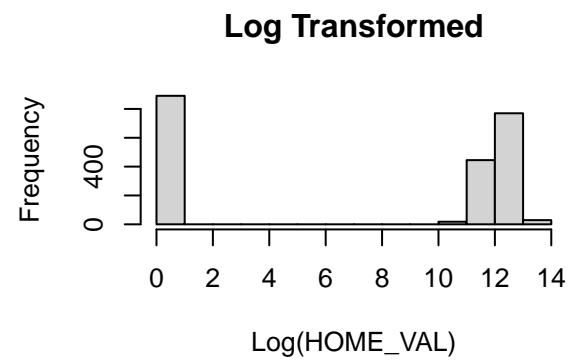
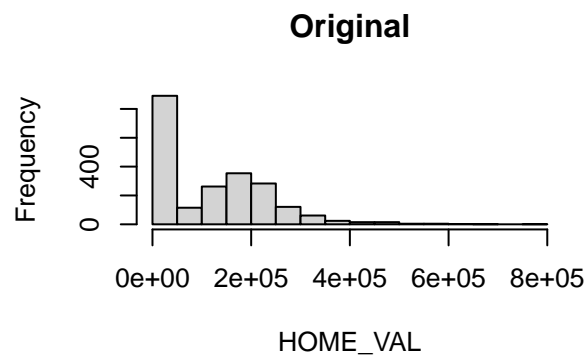


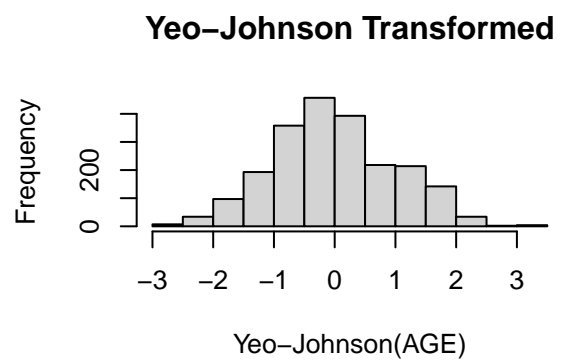
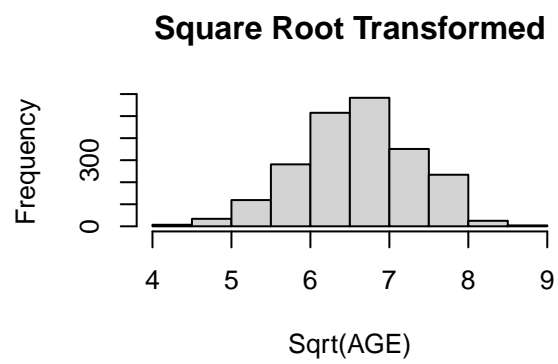
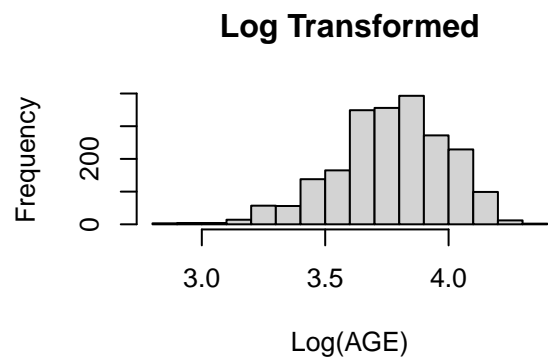
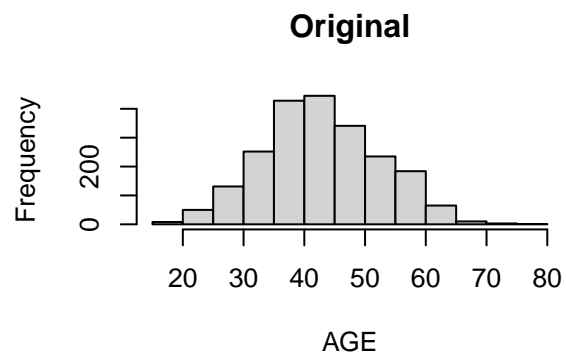




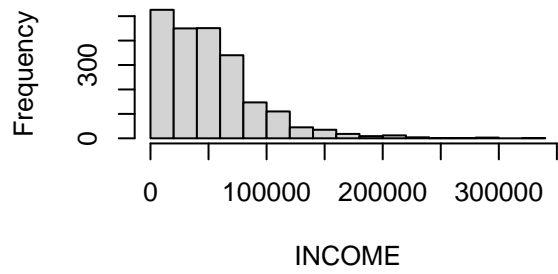




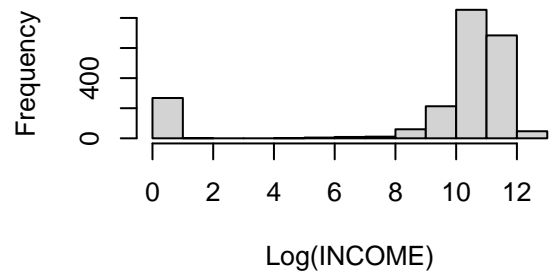




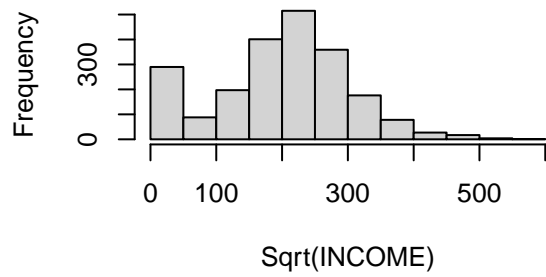
Original



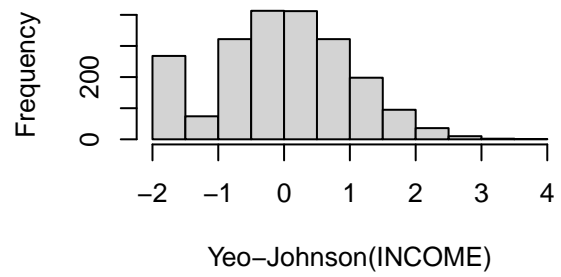
Log Transformed



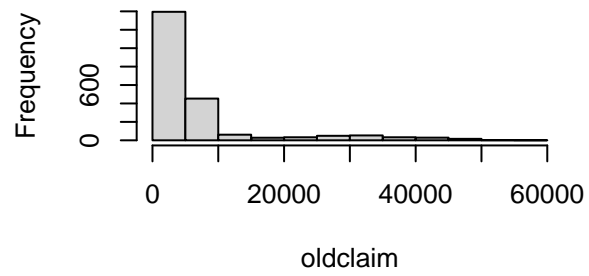
Square Root Transformed



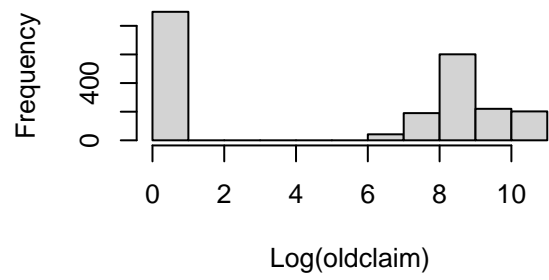
Yeo-Johnson Transformed



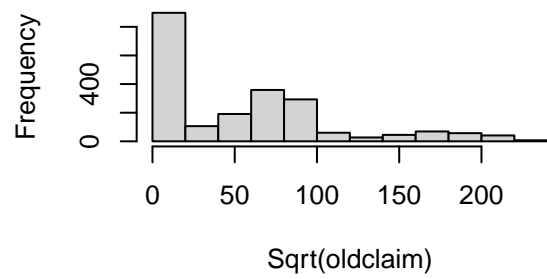
Original



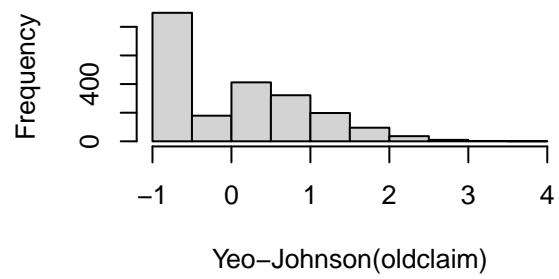
Log Transformed

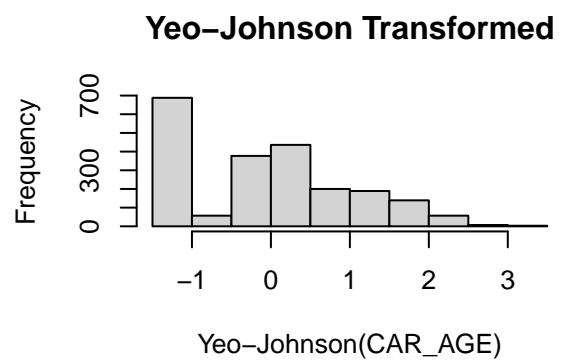
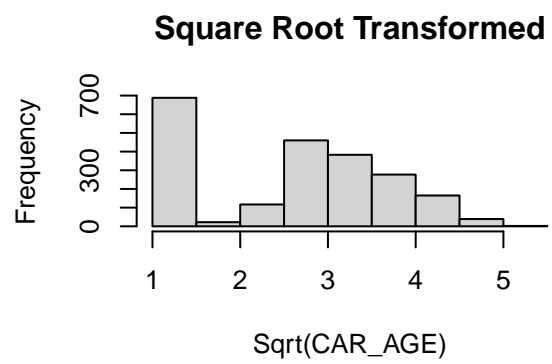
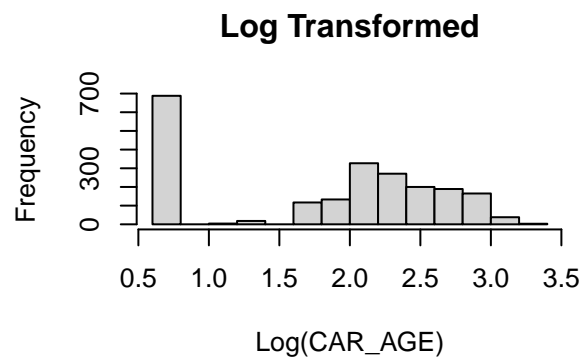
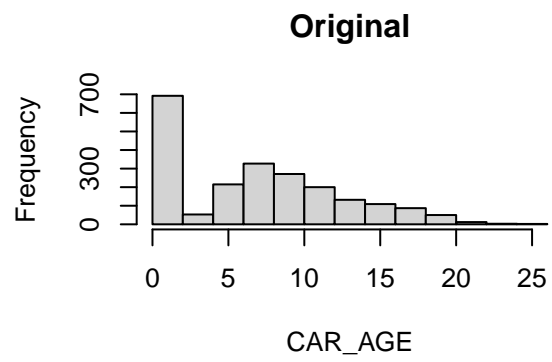


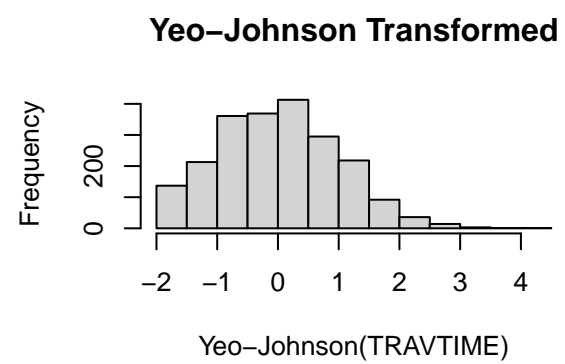
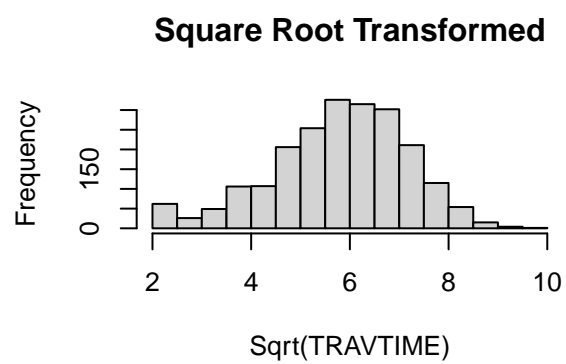
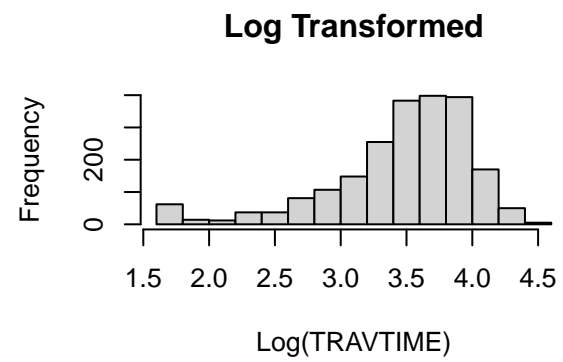
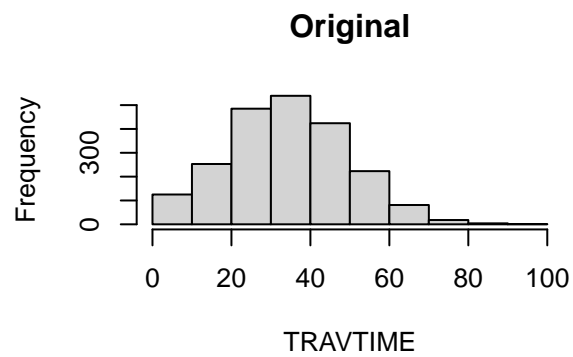
Square Root Transformed



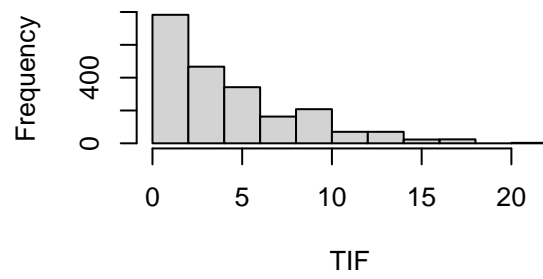
Yeo-Johnson Transformed



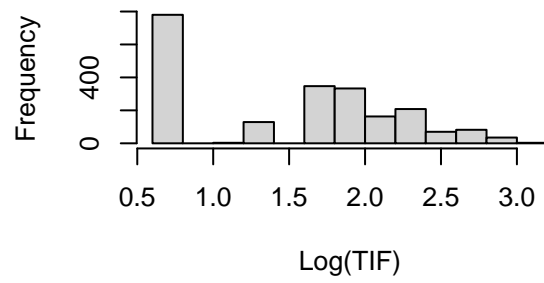




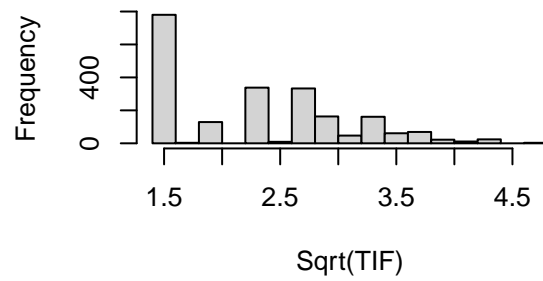
Original



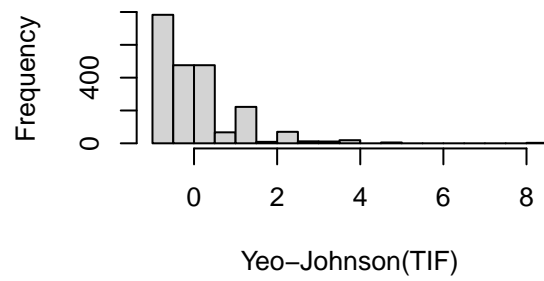
Log Transformed

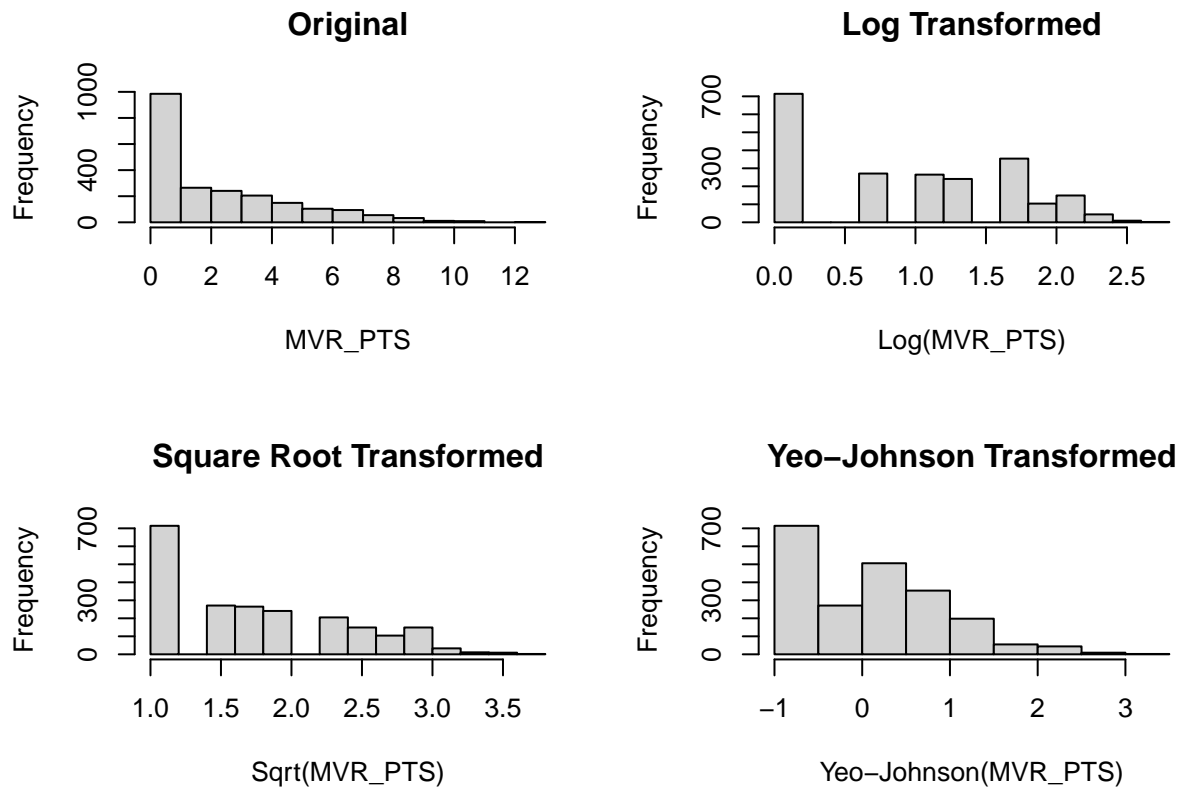


Square Root Transformed



Yeo-Johnson Transformed





Build Models

Multiple Linear Regression

Model 1

Fitting a linear regression model with transformed variables

```
##
## Call:
## lm(formula = TARGET_AMT ~ ., data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10575  -3444  -1603    575   75052
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5.224e+03  2.359e+03   2.214  0.0270 *
## KIDSDRIV     -3.149e+02  4.882e+02  -0.645  0.5191
## AGE           4.000e+00  3.067e+01   0.130  0.8962
## HOMEKIDS      3.510e+02  2.965e+02   1.184  0.2367
## YOJ           5.729e+01  7.352e+01   0.779  0.4360
## INCOME       -1.976e-02  1.017e-02  -1.944  0.0522 .
```

```

## PARENT1Yes          3.034e+01  8.541e+02  0.036  0.9717
## HOME_VAL            1.807e-03  3.024e-03  0.598  0.5501
## MSTATUSYes         -1.442e+03  7.491e+02 -1.925  0.0545 .
## SEXM                2.056e+03  9.381e+02  2.192  0.0286 *
## EDUCATIONHigh School -1.449e+03  7.617e+02 -1.902  0.0574 .
## EDUCATIONLess than High School -8.256e+02  9.352e+02 -0.883  0.3775
## EDUCATIONMasters     5.598e+02  1.316e+03  0.425  0.6707
## EDUCATIONPhD         3.644e+03  1.671e+03  2.181  0.0294 *
## JOBClerical         -9.182e+02  8.482e+02 -1.082  0.2793
## JOBDoctor          -3.980e+03  2.584e+03 -1.540  0.1238
## JOBHome Maker      -6.895e+02  1.295e+03 -0.533  0.5944
## JOBLawyer          -3.381e+02  1.686e+03 -0.200  0.8411
## JOBManager        -1.365e+03  1.319e+03 -1.035  0.3009
## JOBOther Job       2.573e+02  1.680e+03  0.153  0.8783
## JOBProfessional    1.731e+03  1.036e+03  1.672  0.0949 .
## JOBStudent        -5.574e+02  1.101e+03 -0.506  0.6126
## TRAVTIME          -8.015e+00  1.648e+01 -0.486  0.6269
## CAR_USEPrivate     -8.064e+02  7.760e+02 -1.039  0.2989
## BLUEBOOK           1.911e-01  4.333e-02  4.411  1.12e-05 ***
## TIF               -3.095e+01  6.146e+01 -0.504  0.6147
## CAR_TYPEPanel Truck -5.666e+02  1.391e+03 -0.407  0.6838
## CAR_TYPEPickup      1.405e+02  8.757e+02  0.160  0.8725
## CAR_TYPESports Car  1.709e+03  1.061e+03  1.610  0.1076
## CAR_TYPESUV         1.517e+03  9.572e+02  1.585  0.1132
## CAR_TYPEVan        -7.117e+02  1.134e+03 -0.627  0.5306
## RED_CARyes         -9.529e+02  7.307e+02 -1.304  0.1925
## OLDCLAIM           5.793e-02  3.216e-02  1.801  0.0720 .
## CLM_FREQ          -1.726e+02  2.269e+02 -0.761  0.4469
## REVOKEDYes        -1.502e+03  7.594e+02 -1.978  0.0482 *
## MVR_PTS           2.072e+01  9.891e+01  0.209  0.8341
## CAR_AGE           -1.548e+02  6.396e+01 -2.420  0.0157 *
## URBANICITYHighly Urban/ Urban -1.943e+02  1.099e+03 -0.177  0.8597
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8333 on 1151 degrees of freedom
## (318 observations deleted due to missingness)
## Multiple R-squared:  0.05848,    Adjusted R-squared:  0.02821
## F-statistic: 1.932 on 37 and 1151 DF,  p-value: 0.0007587

```

Model Performance on Testing Data:

Mean Absolute Error (MAE): NA

Mean Squared Error (MSE): NA

Root Mean Squared Error (RMSE): NA

Model 2

##

Call:

```
## lm(formula = TARGET_AMT ~ ., data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5070  -1859  -1159    109 103319
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.855e+03  4.207e+02   4.409 1.06e-05 ***
## AGE         -9.890e+00  7.980e+00  -1.239  0.21527
## YOJ          1.731e+01  1.716e+01   1.008  0.31334
## INCOME       -2.030e-03  1.972e-03  -1.029  0.30330
## HOME_VAL     -2.646e-03  6.397e-04  -4.137 3.57e-05 ***
## TRAVTIME      7.296e+00  4.202e+00   1.736  0.08254 .
## BLUEBOOK      2.782e-02  8.789e-03   3.165  0.00156 **
## TIF          -5.059e+01  1.611e+01  -3.141  0.00169 **
## OLDCLAIM      2.099e-03  8.644e-03   0.243  0.80813
## CLM_FREQ      3.030e+02  6.909e+01   4.386 1.18e-05 ***
## MVR_PTS       2.570e+02  3.330e+01   7.718 1.37e-14 ***
## CAR_AGE      -3.948e+01  1.295e+01  -3.049  0.00230 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5218 on 6158 degrees of freedom
## Multiple R-squared:  0.03589,    Adjusted R-squared:  0.03417
## F-statistic: 20.84 on 11 and 6158 DF,  p-value: < 2.2e-16
```

Model Performance on Testing Data:

Mean Absolute Error (MAE): 1721.449

Mean Squared Error (MSE): 3924572

Root Mean Squared Error (RMSE): 1981.053

Model 3

```
##
## Call:
## lm(formula = TARGET_AMT ~ ., data = train_data_scaled)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9595 -0.3505 -0.2187  0.0220 19.4482
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.730e-16  1.251e-02   0.000 1.00000
## AGE         -1.654e-02  1.314e-02  -1.259  0.20821
## YOJ          2.057e-02  1.787e-02   1.151  0.24972
## INCOME       -5.975e-03  3.536e-02  -0.169  0.86583
## HOME_VAL     -6.460e-02  1.555e-02  -4.154 3.31e-05 ***
```

```

## TRAVTIME          2.155e-02  1.255e-02   1.717  0.08597 .
## BLUEBOOK          4.411e-02  1.397e-02   3.158  0.00160 **
## TIF               -3.912e-02  1.255e-02  -3.118  0.00183 **
## OLDCLAIM          3.316e-03  1.447e-02   0.229  0.81870
## CLM_FREQ          6.677e-02  1.525e-02   4.377  1.22e-05 ***
## MVR_PTS           1.064e-01  1.376e-02   7.728  1.26e-14 ***
## CAR_AGE           -5.850e-02  5.203e-02  -1.124  0.26093
## log_income        -1.309e-02  2.130e-02  -0.614  0.53903
## log_car_age        2.030e-02  4.215e-02   0.482  0.63012
## income_car_age_interaction -8.178e-03  4.033e-02  -0.203  0.83934
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.983 on 6155 degrees of freedom
## Multiple R-squared:  0.036, Adjusted R-squared:  0.03381
## F-statistic: 16.42 on 14 and 6155 DF, p-value: < 2.2e-16

##              AGE              YOJ
##      1.102570      2.037874
##      INCOME      HOME_VAL
##      7.982517      1.544118
##      TRAVTIME      BLUEBOOK
##      1.005778      1.245890
##      TIF      OLDCLAIM
##      1.004843      1.336377
##      CLM_FREQ      MVR_PTS
##      1.485765      1.209354
##      CAR_AGE      log_income
##      17.286837      2.897237
##      log_car_age income_car_age_interaction
##      11.341430      10.387322

##
## Call:  glmnet(x = x_train, y = y_train, alpha = 1, lambda = best_lambda)
##
##      Df %Dev   Lambda
## 1 10 3.39 0.001975

## Model Performance on Testing Data:

## Mean Absolute Error (MAE): 0.3245636

## Mean Squared Error (MSE): 0.1369326

## Root Mean Squared Error (RMSE): 0.370044

```

Binary Logistic Regression

Model 1

Model 2

Model 3

Model 4

Model 5

Model 6

Model 7

Select Models & Prediction

Multiple Linear Regression Selection

Third model (with residual standard error: 0.983 and significantly lower MAE, MSE, and RMSE values on the test set). Here's why:

Lower Error Metrics: The first model's error metrics (MAE, MSE, RMSE) are substantially lower than those of the other models, suggesting that its predictions are closer to the actual values on the test data.

Residual Standard Error (RSE): The third model has an RSE of 0.983 compared to the second model's RSE of 5218, indicating tighter residuals, which implies better model fit if both models are evaluated on the same response scale.

F-statistic and R-squared: Both models have similar F-statistics and R-squared values, so there's no distinct advantage for one model over the other in terms of explained variance. However, the significantly lower error metrics of the first model make it the better choice overall for predictive accuracy.

Binary Logistic Regression Selection

Prediction

Code Appendix

```
knitr::opts_chunk$set(echo=FALSE, error=FALSE, warning=FALSE, message=FALSE)

# Libraries

library(stringr)
library(tidyr)
library(DataExplorer)
library(dplyr)
library(visdat)
```

```

library(pROC)
library(mice)
library(corrplot)
library(MASS)
library(caret)
library(e1071)
library(rbin)
library(bestNormalize)
library(GGally)
library(ggplot2)
library(readr)
library(reshape2)
library(purrr)
library(leaps)
# Load necessary package
library(caTools)
library(car) # For VIF
library(glmnet)

# training data
insurance_training_data <- read.csv('https://raw.githubusercontent.com/umais/DATA/refs/heads/main/insurance_training_data.csv')
# test data
insurance_evaluation_data <- read.csv('https://raw.githubusercontent.com/umais/DATA/refs/heads/main/insurance_evaluation_data.csv')
# Check the structure of the data
glimpse(insurance_training_data)

# Display the first few rows and a summary
head(insurance_training_data)
summary(insurance_training_data)
# Remove an index column if present
insurance_training_data_clean <- dplyr::select(insurance_training_data, -INDEX)

# Clean special characters in financial columns
insurance_training_data_clean$HOME_VAL <- substr(insurance_training_data_clean$HOME_VAL, 2, nchar(insurance_training_data_clean$HOME_VAL))
insurance_training_data_clean$HOME_VAL <- as.numeric(str_remove_all(insurance_training_data_clean$HOME_VAL, '[^0-9.]'))

insurance_training_data_clean$BLUEBOOK <- substr(insurance_training_data_clean$BLUEBOOK, 2, nchar(insurance_training_data_clean$BLUEBOOK))
insurance_training_data_clean$BLUEBOOK <- as.numeric(str_remove_all(insurance_training_data_clean$BLUEBOOK, '[^0-9.]'))

insurance_training_data_clean$INCOME <- substr(insurance_training_data_clean$INCOME, 2, nchar(insurance_training_data_clean$INCOME))
insurance_training_data_clean$INCOME <- as.numeric(str_remove_all(insurance_training_data_clean$INCOME, '[^0-9.]'))

insurance_training_data_clean$OLDCLAIM <- substr(insurance_training_data_clean$OLDCLAIM, 2, nchar(insurance_training_data_clean$OLDCLAIM))
insurance_training_data_clean$OLDCLAIM <- as.numeric(str_remove_all(insurance_training_data_clean$OLDCLAIM, '[^0-9.]'))

# Remove 'z_' prefix from marital status and convert to a factor
insurance_training_data_clean$MSTATUS <- as.factor(str_remove(insurance_training_data_clean$MSTATUS, 'z_'))

# Remove 'z_' prefix from parental status and convert to a factor
insurance_training_data_clean$PARENT1 <- as.factor(str_remove(insurance_training_data_clean$PARENT1, 'z_'))

# Replace '<' with 'Less than ' in education level to clarify the meaning

```

```

insurance_training_data_clean$EDUCATION <- str_replace(insurance_training_data_clean$EDUCATION, '<', 'L')

# Remove 'z_' prefix from sex and convert to a factor
insurance_training_data_clean$SEX <- as.factor(str_remove(insurance_training_data_clean$SEX, 'z_'))

# Remove 'z_' prefix from education level and convert to a factor
insurance_training_data_clean$EDUCATION <- as.factor(str_remove(insurance_training_data_clean$EDUCATION, 'z_'))

# Recode empty job entries as 'Other Job' to handle missing data
insurance_training_data_clean$JOB[insurance_training_data_clean$JOB == ""] <- 'Other Job'

# Remove 'z_' prefix from job titles and convert to a factor
insurance_training_data_clean$JOB <- as.factor(str_remove(insurance_training_data_clean$JOB, 'z_'))

# Remove 'z_' prefix from car usage category and convert to a factor
insurance_training_data_clean$CAR_USE <- as.factor(str_remove(insurance_training_data_clean$CAR_USE, 'z_'))

# Remove 'z_' prefix from car type and convert to a factor
insurance_training_data_clean$CAR_TYPE <- as.factor(str_remove(insurance_training_data_clean$CAR_TYPE, 'z_'))

# Remove 'z_' prefix from urbanicity status and convert to a factor
insurance_training_data_clean$URBANICITY <- as.factor(str_remove(insurance_training_data_clean$URBANICITY, 'z_'))

# Remove 'z_' prefix from revoked status and convert to a factor
insurance_training_data_clean$REVOKED <- as.factor(str_remove(insurance_training_data_clean$REVOKED, 'z_'))

# Remove 'z_' prefix from red car indicator and convert to a factor
insurance_training_data_clean$RED_CAR <- as.factor(str_remove(insurance_training_data_clean$RED_CAR, 'z_'))

summary(insurance_training_data_clean)

insurance_training_data_clean$CAR_AGE[insurance_training_data_clean$CAR_AGE < 1] <- 1
# Identify categorical columns and store their names in cat_features
cat_features <- names(insurance_training_data_clean)[map_chr(insurance_training_data_clean, class) == "factor"]

# Display each categorical column and its unique levels
cat("Exploring Categorical Features:\n")
walk(cat_features, ~cat("Feature:", ., "\nLevels:", paste(levels(insurance_training_data_clean[[.]]), collapse = ", ")))

# Select categorical features from the cleaned insurance training data
categorical_data <- insurance_training_data_clean[cat_features]

# Melt the data frame to create a long format suitable for ggplot
melted_data <- melt(categorical_data, measure.vars = cat_features, variable.name = 'category', value.name = 'value')

# Create a bar plot to visualize the distribution of categorical predictors
ggplot(melted_data, aes(x = category_value)) +
  geom_bar(aes(fill = category_value)) +
  scale_fill_brewer(palette = "Set1") +
  facet_wrap(~ category, nrow = 5L, scales = 'free') +

```

```

coord_flip() +
  labs(title = "Distribution of Categorical Predictors",
        x = "Category Value",
        y = "Count") +
  theme_minimal()
plot_histogram(insurance_training_data_clean, geom_histogram_args = list("fill" = "tomato4"))

plot_histogram(insurance_training_data_clean, scale_x = "log10", geom_histogram_args = list("fill" = "r

# Summarize the dataset to check for columns with missing values
insurance_training_data_clean %>%
  summarise_all(funs(sum(is.na(.)))) %>%
  select_if(~any(.) > 0)

# Visualize the missing values in the dataset to understand their distribution
plot_missing(insurance_training_data_clean)

# Calculate and display the proportion of missing values for each column
round(colSums(is.na(insurance_training_data_clean)) / nrow(insurance_training_data_clean), 3)

# Visualize specific columns to further investigate missing data patterns
vis_dat(insurance_training_data_clean %>% dplyr::select(YOJ, INCOME, HOME_VAL, CAR_AGE))

# Select numeric columns for correlation analysis
numeric_data <- insurance_training_data_clean[, c('TARGET_AMT', 'AGE', 'YOJ', 'INCOME', 'HOME_VAL', 'TR

# Document missing values before imputation
missing_summary_before <- colSums(is.na(numeric_data))
print("Missing Values Before Imputation:")
print(missing_summary_before)

# Perform multiple imputation
imputed_data <- mice(numeric_data, m = 5, method = 'pmm', seed = 123) # Predictive Mean Matching

# Create a complete dataset by averaging the multiple imputations
completed_data <- complete(imputed_data)

# Document missing values after imputation
missing_summary_after <- colSums(is.na(completed_data))
print("Missing Values After Imputation:")
print(missing_summary_after)

# Generate a correlation matrix and plot it
corrplot(cor(completed_data), type = "upper")

# Sensitivity Analysis

```



```

# Compare correlations from original data (complete case analysis) vs. imputed data

# Complete case analysis (removing rows with NA values)
complete_case_data <- na.omit(numeric_data)
cor_complete_case <- cor(complete_case_data)

# Correlation of imputed data
cor_imputed <- cor(completed_data)

# Print correlation matrices for comparison
print("Correlation Matrix for Complete Case Analysis:")
print(cor_complete_case)

print("Correlation Matrix for Imputed Data:")
print(cor_imputed)

# Visualize the difference in correlations
cor_diff <- cor_imputed - cor_complete_case
ggplot(melt(cor_diff), aes(Var1, Var2, fill = value)) +
  geom_tile() +
  scale_fill_gradient2(low = "blue", high = "red", mid = "white", limit = c(-1, 1), name="Correlation D")
  theme_minimal() +
  labs(title = "Difference in Correlation between Imputed and Complete Case Data", x = "Variables", y = "Variables")
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

crash_data <- subset(filter(insurance_training_data_clean, TARGET_FLAG==1), select = -c(TARGET_FLAG))

# Check for missing values before imputation
missing_summary_before <- colSums(is.na(crash_data))
print("Missing Values Before Imputation:")
print(missing_summary_before)

# Impute missing values
imputed_data <- mice(crash_data, m = 5, method = 'pmm', seed = 123) # Predictive Mean Matching
crash_data_imputed <- complete(imputed_data)

# Check for missing values after imputation
missing_summary_after <- colSums(is.na(crash_data_imputed))
print("Missing Values After Imputation:")
print(missing_summary_after)

# Create a histogram and density plot for the AGE variable
ggplot(crash_data_imputed, aes(x = AGE)) +
  geom_histogram(binwidth = 1, fill = "lightblue", color = "black", alpha = 0.7) +
  geom_density(aes(y = ..count.. * 1), fill = "lightgreen", alpha = 0.5) +
  labs(title = "Distribution of AGE", x = "AGE", y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5))

```

```

# Create a histogram for the INCOME variable
ggplot(data = crash_data_imputed, aes(x = INCOME)) +
  geom_histogram(bins = 30, fill = "lightblue", color = "black") +
  labs(title = "Distribution of INCOME",
       x = "Income",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Create a histogram for the HOME_VAL variable
ggplot(data = crash_data_imputed, aes(x = HOME_VAL)) +
  geom_histogram(bins = 30, fill = "lightcoral", color = "black") +
  labs(title = "Distribution of HOME_VAL",
       x = "Home Value",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Create a histogram for the CAR_AGE variable
ggplot(data = crash_data_imputed, aes(x = CAR_AGE)) +
  geom_histogram(bins = 30, fill = "lightblue", color = "black") +
  labs(title = "Distribution of CAR_AGE",
       x = "Car Age",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Create a histogram for the BLUEBOOK variable
ggplot(data = crash_data_imputed, aes(x = BLUEBOOK)) +
  geom_histogram(bins = 30, fill = "lightgreen", color = "black") +
  labs(title = "Distribution of BLUEBOOK",
       x = "Blue Book Value",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Create a histogram for the OLDCLAIM variable
ggplot(data = crash_data_imputed, aes(x = OLDCLAIM)) +
  geom_histogram(bins = 30, fill = "lightcoral", color = "black") +
  labs(title = "Distribution of OLDCLAIM",
       x = "Old Claim Amount",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Create a histogram for the TRAVTIME variable
ggplot(data = crash_data_imputed, aes(x = TRAVTIME)) +
  geom_histogram(bins = 30, fill = "lightsalmon", color = "black") +
  labs(title = "Distribution of TRAVTIME",
       x = "Travel Time",
       y = "Frequency") +
  theme_minimal() +

```

```

    theme(plot.title = element_text(hjust = 0.5)) # Center the title
# Histogram for TIF (Number of Years with Policy)
ggplot(data = crash_data_imputed, aes(x = TIF)) +
  geom_histogram(bins = 30, fill = "lightblue", color = "black") +
  labs(title = "Distribution of TIF (Number of Years with Policy)",
        x = "Years with Policy (TIF)",
        y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Histogram for MVR_PTS (Driving Record Points)
ggplot(data = crash_data_imputed, aes(x = MVR_PTS)) +
  geom_histogram(bins = 30, fill = "lightgreen", color = "black") +
  labs(title = "Distribution of MVR_PTS (Driving Record Points)",
        x = "Driving Record Points (MVR_PTS)",
        y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5)) # Center the title

# Example variable to transform
home_val_variable <- crash_data_imputed$HOME_VAL # Replace with your actual variable

# 1. Log Transformation
home_val_log_transformed <- log(home_val_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
home_val_sqrt_transformed <- sqrt(home_val_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
home_val_box_cox_transformed <- boxcox(home_val_variable + 1) # Add 1 to handle zeros, need to extract

home_val_yj_transformed <- bestNormalize(home_val_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (home_val_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(home_val_variable, main="Original", xlab="HOME_VAL")
hist(home_val_log_transformed, main="Log Transformed", xlab="Log(HOME_VAL)")
hist(home_val_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(HOME_VAL)")
hist(home_val_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(HOME_VAL)")

# Example variable to transform
age_variable <- crash_data_imputed$AGE # Replace with your actual variable

# 1. Log Transformation
age_log_transformed <- log(age_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
age_sqrt_transformed <- sqrt(age_variable + 1) # Add 1 to handle zeros

```

```

# 3. Box-Cox Transformation
age_box_cox_transformed <- boxcox(age_variable + 1) # Add 1 to handle zeros, need to extract lambda

age_yj_transformed <- bestNormalize(age_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (age_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(age_variable, main="Original", xlab="AGE")
hist(age_log_transformed, main="Log Transformed", xlab="Log(AGE)")
hist(age_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(AGE)")
hist(age_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(AGE)")

# Example variable to transform
income_variable <- crash_data_imputed$INCOME # Replace with your actual variable

# 1. Log Transformation
income_log_transformed <- log(income_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
income_sqrt_transformed <- sqrt(income_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
income_box_cox_transformed <- boxcox(income_variable + 1) # Add 1 to handle zeros, need to extract lam

income_yj_transformed <- bestNormalize(income_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (income_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(income_variable, main="Original", xlab="INCOME")
hist(income_log_transformed, main="Log Transformed", xlab="Log(INCOME)")
hist(income_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(INCOME)")
hist(income_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(INCOME)")

#OldClaim

oldclaim_variable <- crash_data_imputed$OLDCLAIM # Replace with your actual variable

oldclaim_log_transformed <- log(oldclaim_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
oldclaim_sqrt_transformed <- sqrt(oldclaim_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
oldclaim_box_cox_transformed <- boxcox(oldclaim_variable + 1) # Add 1 to handle zeros, need to extract

```

```

oldclaim_yj_transformed <- bestNormalize(oldclaim_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (oldclaim_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(oldclaim_variable, main="Original", xlab="oldclaim")
hist(oldclaim_log_transformed, main="Log Transformed", xlab="Log(oldclaim)")
hist(oldclaim_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(oldclaim)")
hist(oldclaim_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(oldclaim)")

# CAR AGE
car_age_variable <- crash_data_imputed$CAR_AGE # Replace with your actual variable

car_age_log_transformed <- log(car_age_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
car_age_sqrt_transformed <- sqrt(car_age_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
car_age_box_cox_transformed <- boxcox(car_age_variable + 1) # Add 1 to handle zeros, need to extract l

car_age_yj_transformed <- bestNormalize(car_age_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (car_age_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(car_age_variable, main="Original", xlab="CAR_AGE")
hist(car_age_log_transformed, main="Log Transformed", xlab="Log(CAR_AGE)")
hist(car_age_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(CAR_AGE)")
hist(car_age_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(CAR_AGE)")

#TRAVTIME TRANSFORMATIONS

TRAVTIME_variable <- crash_data_imputed$TRAVTIME # Replace with your actual variable

TRAVTIME_log_transformed <- log(TRAVTIME_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
TRAVTIME_sqrt_transformed <- sqrt(TRAVTIME_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
TRAVTIME_box_cox_transformed <- boxcox(TRAVTIME_variable + 1) # Add 1 to handle zeros, need to extract l

TRAVTIME_yj_transformed <- bestNormalize(TRAVTIME_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (TRAVTIME_variable + 1) # Add 1 to handle zeros

```

```

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(TRAVTIME_variable, main="Original", xlab="TRAVTIME")
hist(TRAVTIME_log_transformed, main="Log Transformed", xlab="Log(TRAVTIME)")
hist(TRAVTIME_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(TRAVTIME)")
hist(TRAVTIME_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(TRAVTIME)")

#TIF

TIF_variable <- crash_data_imputed$TIF # Replace with your actual variable

TIF_log_transformed <- log(TIF_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
TIF_sqrt_transformed <- sqrt(TIF_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
TIF_box_cox_transformed <- boxcox(TIF_variable + 1) # Add 1 to handle zeros, need to extract lambda

TIF_yj_transformed <- bestNormalize(TIF_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (TIF_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(TIF_variable, main="Original", xlab="TIF")
hist(TIF_log_transformed, main="Log Transformed", xlab="Log(TIF)")
hist(TIF_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(TIF)")
hist(TIF_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(TIF)")

#MVR_PTS TRANSFORMATIONS

MVR_PTS_variable <- crash_data_imputed$MVR_PTS # Replace with your actual variable

MVR_PTS_log_transformed <- log(MVR_PTS_variable + 1) # Add 1 to handle zeros

# 2. Square Root Transformation
MVR_PTS_sqrt_transformed <- sqrt(MVR_PTS_variable + 1) # Add 1 to handle zeros

# 3. Box-Cox Transformation
MVR_PTS_box_cox_transformed <- boxcox(MVR_PTS_variable + 1) # Add 1 to handle zeros, need to extract lambda

MVR_PTS_yj_transformed <- bestNormalize(MVR_PTS_variable, method = "yeo.johnson")$x.t

# 5. Inverse Transformation
inverse_transformed <- 1 / (MVR_PTS_variable + 1) # Add 1 to handle zeros

# Check the results with histograms
par(mfrow=c(2,2)) # Set up the plotting area
hist(MVR_PTS_variable, main="Original", xlab="MVR_PTS")
hist(MVR_PTS_log_transformed, main="Log Transformed", xlab="Log(MVR_PTS)")
hist(MVR_PTS_sqrt_transformed, main="Square Root Transformed", xlab="Sqrt(MVR_PTS)")

```

```

hist(MVR PTS_yj_transformed, main="Yeo-Johnson Transformed", xlab="Yeo-Johnson(MVR PTS)")

crash_data_imputed_transformed <- crash_data_imputed %>%
  mutate(
    # Log transformation of AGE
    INCOME_transformed = bestNormalize(INCOME, method = "yeo.johnson")$x.t,      # Log transformation of INCOME
    CAR_AGE_transformed = sqrt(CAR_AGE + 1), # Square root transformation of CAR_AGE
    HOME_VAL_transformed = sqrt(HOME_VAL + 1), # Log transformation of HOME_VAL
    OLDCLAIM_transformed=bestNormalize(oldclaim_variable, method = "yeo.johnson")$x.t,
    TRAVTIME_transformed=sqrt(TRAVTIME + 1)

  )

# Set seed for reproducibility
set.seed(123) # You can set any number

# Create a split index
split <- sample.split(crash_data$TARGET_AMT, SplitRatio = 0.7)

# Split data into training and testing sets
train_data <- subset(crash_data, split == TRUE)
test_data <- subset(crash_data, split == FALSE)

# Fit the model on the training data
model <- lm(TARGET_AMT ~ ., data = train_data)
summary(model)

# Predict on the testing data
predictions <- predict(model, newdata = test_data)

# Evaluate model performance
# Calculate Mean Absolute Error (MAE)
MAE <- mean(abs(predictions - test_data$TARGET_AMT))

# Calculate Mean Squared Error (MSE)
MSE <- mean((predictions - test_data$TARGET_AMT)^2)

# Calculate Root Mean Squared Error (RMSE)
RMSE <- sqrt(MSE)

# Print the performance metrics
cat("Model Performance on Testing Data:\n")
cat("Mean Absolute Error (MAE):", MAE, "\n")
cat("Mean Squared Error (MSE):", MSE, "\n")
cat("Root Mean Squared Error (RMSE):", RMSE, "\n")

# Set seed for reproducibility
set.seed(123) # You can set any number

```



```

# Create a split index
split <- sample.split(completed_data$TARGET_AMT, SplitRatio = 0.7)

# Split data into training and testing sets
train_data <- subset(completed_data, split == TRUE)
test_data <- subset(completed_data, split == FALSE)

# Fit the model on the training data
model <- lm(TARGET_AMT ~ ., data = train_data)
summary(model)

# Predict on the testing data
predictions <- predict(model, newdata = test_data)

# Evaluate model performance
# Calculate Mean Absolute Error (MAE)
MAE <- mean(abs(predictions - test_data$TARGET_AMT))

# Calculate Mean Squared Error (MSE)
MSE <- mean((predictions - test_data$TARGET_AMT)^2)

# Calculate Root Mean Squared Error (RMSE)
RMSE <- sqrt(MSE)

# Print the performance metrics
cat("Model Performance on Testing Data:\n")
cat("Mean Absolute Error (MAE):", MAE, "\n")
cat("Mean Squared Error (MSE):", MSE, "\n")
cat("Root Mean Squared Error (RMSE):", RMSE, "\n")

# Feature Engineering: Transformations and Interaction Terms
train_data$log_income <- log(train_data$INCOME + 1)
train_data$log_car_age <- log(train_data$CAR_AGE + 1)
train_data$income_car_age_interaction <- train_data$INCOME * train_data$CAR_AGE

test_data$log_income <- log(test_data$INCOME + 1)
test_data$log_car_age <- log(test_data$CAR_AGE + 1)
test_data$income_car_age_interaction <- test_data$INCOME * test_data$CAR_AGE

# Optional: Scaling if predictors have large variances
preproc <- preProcess(train_data, method = c("center", "scale"))
train_data_scaled <- predict(preproc, train_data)
test_data_scaled <- predict(preproc, test_data)

# Step 1: Fit the initial model on the training data
initial_model <- lm(TARGET_AMT ~ ., data = train_data_scaled)
summary(initial_model)

# Check for multicollinearity
vif_values <- car::vif(initial_model)
print(vif_values)

```



```

# Step 2: Remove high VIF variables if any are >5
high_vif_vars <- names(vif_values[vif_values > 5])
train_data_reduced <- train_data_scaled[, !(names(train_data_scaled) %in% high_vif_vars)]
test_data_reduced <- test_data_scaled[, !(names(test_data_scaled) %in% high_vif_vars)]

# Step 3: Fit a regularized model (Lasso) on the reduced data
x_train <- model.matrix(TARGET_AMT ~ ., data = train_data_reduced)[, -1]
y_train <- train_data_reduced$TARGET_AMT
lasso_cv <- cv.glmnet(x_train, y_train, alpha = 1)
best_lambda <- lasso_cv$lambda.min

final_lasso_model <- glmnet(x_train, y_train, alpha = 1, lambda = best_lambda)
print(final_lasso_model)

# Predictions on the test set
x_test <- model.matrix(TARGET_AMT ~ ., data = test_data_reduced)[, -1]
predictions <- predict(final_lasso_model, newx = x_test)

# Evaluate Model Performance
MAE <- mean(abs(predictions - test_data_reduced$TARGET_AMT))
MSE <- mean((predictions - test_data_reduced$TARGET_AMT)^2)
RMSE <- sqrt(MSE)

# Print performance metrics
cat("Model Performance on Testing Data:\n")
cat("Mean Absolute Error (MAE):", MAE, "\n")
cat("Mean Squared Error (MSE):", MSE, "\n")
cat("Root Mean Squared Error (RMSE):", RMSE, "\n")

```