# Home Work Assignment - 03

*Critical Thinking Group 5*

# Contents

# Overview

To attain our objective, we will be following the below best practice steps and guidelines:
1 -Data Exploration
2 -Data Preparation
3 -Build Models
4 -Select Models

```
##       zn              indus              chas              nox
##  Min.   :  0.00   Min.   : 0.460   Min.   :0.00000   Min.   :0.3890
##  1st Qu.:  0.00   1st Qu.: 5.145   1st Qu.:0.00000   1st Qu.:0.4480
##  Median :  0.00   Median : 9.690   Median :0.00000   Median :0.5380
##  Mean   : 11.58   Mean   :11.105   Mean   :0.07082   Mean   :0.5543
##  3rd Qu.: 16.25   3rd Qu.:18.100   3rd Qu.:0.00000   3rd Qu.:0.6240
##  Max.   :100.00   Max.   :27.740   Max.   :1.00000   Max.   :0.8710
##       rm              age              dis              rad
##  Min.   :3.863   Min.   :  2.90   Min.   : 1.130   Min.   : 1.00
##  1st Qu.:5.887   1st Qu.: 43.88   1st Qu.: 2.101   1st Qu.: 4.00
##  Median :6.210   Median : 77.15   Median : 3.191   Median : 5.00
##  Mean   :6.291   Mean   : 68.37   Mean   : 3.796   Mean   : 9.53
##  3rd Qu.:6.630   3rd Qu.: 94.10   3rd Qu.: 5.215   3rd Qu.:24.00
##  Max.   :8.780   Max.   :100.00   Max.   :12.127   Max.   :24.00
##       tax             ptratio           black             lstat
##  Min.   :187.0   Min.   :12.6   Min.   :  0.32   Min.   : 1.730
##  1st Qu.:281.0   1st Qu.:16.9   1st Qu.:375.61   1st Qu.: 7.043
##  Median :334.5   Median :18.9   Median :391.34   Median :11.350
##  Mean   :409.5   Mean   :18.4   Mean   :357.12   Mean   :12.631
##  3rd Qu.:666.0   3rd Qu.:20.2   3rd Qu.:396.24   3rd Qu.:16.930
##  Max.   :711.0   Max.   :22.0   Max.   :396.90   Max.   :37.970
##       medv            target
##  Min.   : 5.00   Min.   :0.0000
##  1st Qu.:17.02   1st Qu.:0.0000
##  Median :21.20   Median :0.0000
##  Mean   :22.59   Mean   :0.4914
##  3rd Qu.:25.00   3rd Qu.:1.0000
##  Max.   :50.00   Max.   :1.0000
```

```
## 'data.frame':    40 obs. of  13 variables:
##  $ zn     : int  0 0 0 0 25 25 0 0 0 ...
##  $ indus  : num  7.07 8.14 8.14 8.14 5.96 5.13 5.13 4.49 4.49 2.89 ...
##  $ chas   : int  0 0 0 0 0 0 0 0 0 0 ...
##  $ nox    : num  0.469 0.538 0.538 0.538 0.499 0.453 0.453 0.449 0.449 0.445 ...
##  $ rm     : num  7.18 6.1 6.5 5.95 5.85 ...
##  $ age    : num  61.1 84.5 94.4 82 41.5 66.2 93.4 56.1 56.8 69.6 ...
##  $ dis    : num  4.97 4.46 4.45 3.99 3.93 ...
##  $ rad    : int  2 4 4 4 5 8 8 3 3 2 ...
##  $ tax    : int  242 307 307 307 279 284 284 247 247 276 ...
##  $ ptratio: num  17.8 21 21 21 19.2 19.7 19.7 18.5 18.5 18 ...
##  $ black  : num  393 380 388 233 397 ...
##  $ lstat  : num  4.03 10.26 12.8 27.71 8.77 ...
##  $ medv   : num  34.7 18.2 18.4 13.2 21 18.7 16 26.6 22.2 21.4 ...
```

Split the full train data set into train and test to validate the model performance
1. Split the data 80% train and 20% for model validation

# 1 Data Exploration Analysis

In section we will explore and gain some insights into the dataset by pursuing the below high level steps and inquiries:
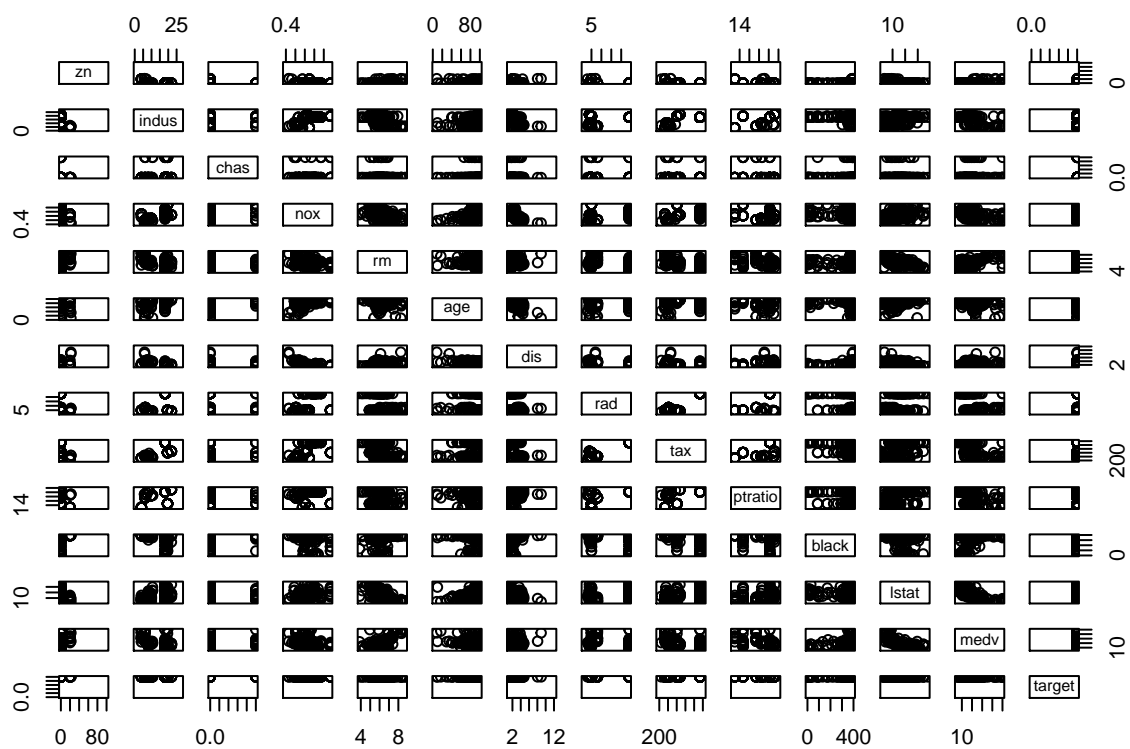-Variable identification
-Variable Relationships
-Data summary analysis
-Outliers and Missing Values Identification

## 1.1 Variable identification

First let's display and examine the data dictionary or the data columns as shown in table 1 and proportion of success and failure cases in target variable.

```
##       zn              indus            chas              nox
##  Min.   :  0.00   Min.   : 0.460   Min.   :0.00000   Min.   :0.3890
##  1st Qu.:  0.00   1st Qu.: 4.945   1st Qu.:0.00000   1st Qu.:0.4480
##  Median :  0.00   Median : 8.560   Median :0.00000   Median :0.5220
##  Mean   : 12.36   Mean   :10.900   Mean   :0.06452   Mean   :0.5512
##  3rd Qu.: 20.00   3rd Qu.:18.100   3rd Qu.:0.00000   3rd Qu.:0.6240
##  Max.   :100.00   Max.   :27.740   Max.   :1.00000   Max.   :0.8710
##       rm              age              dis              rad
##  Min.   :3.863   Min.   :  2.90   Min.   : 1.130   Min.   : 1.000
##  1st Qu.:5.886   1st Qu.: 41.70   1st Qu.: 2.106   1st Qu.: 4.000
##  Median :6.205   Median : 76.50   Median : 3.325   Median : 5.000
##  Mean   :6.295   Mean   : 67.41   Mean   : 3.844   Mean   : 9.204
##  3rd Qu.:6.683   3rd Qu.: 93.85   3rd Qu.: 5.287   3rd Qu.: 8.000
##  Max.   :8.725   Max.   :100.00   Max.   :12.127   Max.   :24.000
##      tax            ptratio           black            lstat
##  Min.   :187.0   Min.   :12.60   Min.   :  0.32   Min.   : 1.730
##  1st Qu.:277.0   1st Qu.:16.60   1st Qu.:376.46   1st Qu.: 6.928
##  Median :330.0   Median :18.60   Median :391.95   Median :10.925
##  Mean   :403.7   Mean   :18.23   Mean   :359.63   Mean   :12.397
##  3rd Qu.:666.0   3rd Qu.:20.20   3rd Qu.:396.21   3rd Qu.:16.672
##  Max.   :711.0   Max.   :22.00   Max.   :396.90   Max.   :37.970
##      medv            target
##  Min.   : 5.00   Min.   :0.0000
##  1st Qu.:17.20   1st Qu.:0.0000
##  Median :21.60   Median :0.0000
##  Mean   :22.85   Mean   :0.4731
##  3rd Qu.:27.02   3rd Qu.:1.0000
##  Max.   :50.00   Max.   :1.0000
```

```
## 
##          0          1
## 0.5268817 0.4731183
```

## 1.2 Data Summary Analysis

In this section, we will create summary data to better understand the initial relationship variables have with our dependent variable using correlation, central tendency, and dispersion As shown in table 2.

```
##            vars   n    mean     sd median trimmed    mad    min    max  range
## zn           1 372   12.36  24.06   0.00    6.04   0.00   0.00 100.00 100.00
## indus        2 372   10.90   6.90   8.56   10.66   7.90   0.46  27.74  27.28
## chas         3 372    0.06   0.25   0.00    0.00   0.00   0.00   1.00   1.00
## nox          4 372    0.55   0.12   0.52    0.54   0.12   0.39   0.87   0.48
## rm           5 372    6.30   0.70   6.21    6.27   0.53   3.86   8.72   4.86
## age          6 372   67.41  28.69  76.50   69.83  30.91   2.90 100.00  97.10
## dis          7 372    3.84   2.13   3.32    3.60   2.05   1.13  12.13  11.00
## rad          8 372    9.20   8.54   5.00    8.28   1.48   1.00  24.00  23.00
## tax          9 372  403.69 167.05 330.00  394.00 108.23 187.00 711.00 524.00
## ptratio     10 372   18.23   2.22  18.60   18.41   2.37  12.60  22.00   9.40
## black       11 372  359.63  88.60 391.96  384.77   7.33   0.32 396.90 396.58
## lstat       12 372   12.40   7.03  10.93   11.62   6.77   1.73  37.97  36.24
## medv        13 372   22.85   9.07  21.60   21.98   6.97   5.00  50.00  45.00
## target      14 372    0.47   0.50   0.00    0.47   0.00   0.00   1.00   1.00
##           skew kurtosis   se
## zn        2.05     3.20 1.25
## indus     0.34    -1.21 0.36
## chas      3.53    10.50 0.01
## nox       0.84     0.09 0.01
## rm        0.39     1.48 0.04
## age      -0.53    -1.09 1.49
## dis       0.96     0.38 0.11
## rad       1.10    -0.67 0.44
## tax       0.72    -1.05 8.66
## ptratio  -0.67    -0.52 0.12
## black    -3.10     8.55 4.59
## lstat     0.95     0.60 0.36
## medv      0.97     1.11 0.47
## target    0.11    -1.99 0.03


##      zn   indus    chas     nox      rm     age     dis     rad     tax
##       0       0       0       0       0       0       0       0       0
## ptratio   black   lstat    medv  target
##       0       0       0       0       0
```

Table 1: Correlation between target and predictor variable

|        | Correlation |
|--------|-------------|
| zn     | -0.4239382  |
| indus  | 0.6034795   |
| chas   | 0.0579716   |
| nox    | 0.7290920   |
| rm     | -0.1605913  |
| age    | 0.6275762   |
| dis    | -0.6167264  |
| rad    | 0.6307187   |
| tax    | 0.6021403   |

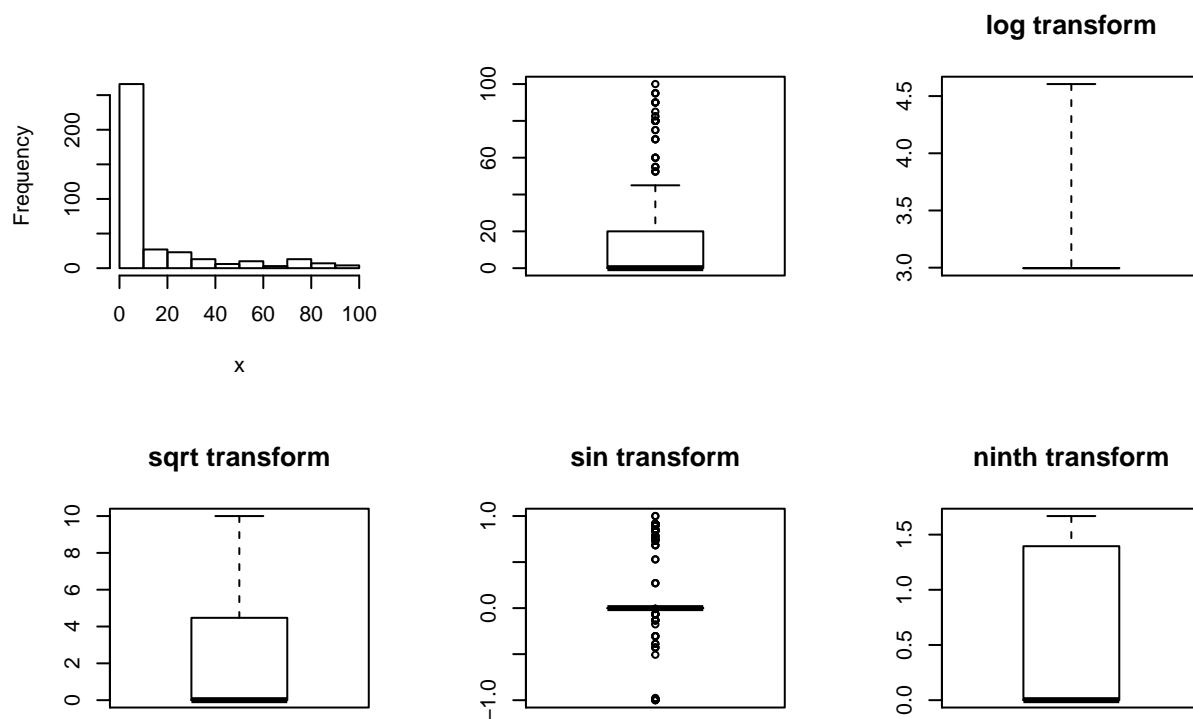| | Correlation |
|---|---|
| ptratio | 0.2198922 |
| black | -0.3463425 |
| lstat | 0.4808888 |
| medv | -0.2724789 |
| target | 1.0000000 |

It is clear from the table that most of the variables are having storng correlation with the target variable.

## 1.3 Outliers and Missing Values Identification

In this section we look at boxplots to determine the outliers in variables and decide on whether to act on the outliers.

Lets do some univariate analysis. We will look at the Histogram and Boxplot for each variable to detect outliers if any and treat it accordingly.

Analysis of variable zn:proportion of residential land zoned for large lots



For zn, we can see that there are large number of values with 0. ninth transformation seem better for this variable..(1)

*

**Please note that we have created similar figures to figure 1 above for each remaining variable. However, we hid the remaining figures for ease of streamlining the report as they have similar shapes. However, we have drawn the below observations from each remaining figure.

For indus, we can see that there is a spike toward right side of he distribution. Looking at the sqrt transformation it appears that distribution is close to normal and having two peaks after transformation.

For nox, there is a long right tail.

For rm, there are some outliers as we can see from box plot. This variable will need some transformation to handle the outliers.

age of the building variable is skewed heavily towards right side. We will need some transformation for this variable and looks sin transformation is best option for this case

For this variable dis, there are some outliers which needs transformation to handle those outliers. log transformation looks best suited for this scenario.

For rad variable distribution is not uniform as seen from the chart and will need transformation.

For tax variable is not uniformly distributed but there is no outlier for this variable.

For pratio has right aligned peak but no outliers are there in data set.

The variable lstat has long right tail and lef skewed

# 2. Data Preparation

Now that we have completed the preliminary analysis, we will be cleaning and consolidating data into one dataset for use in analysis and modeling. We will be puring the belwo steps as guidlines:
- Outliers treatment
- Missing values treatment
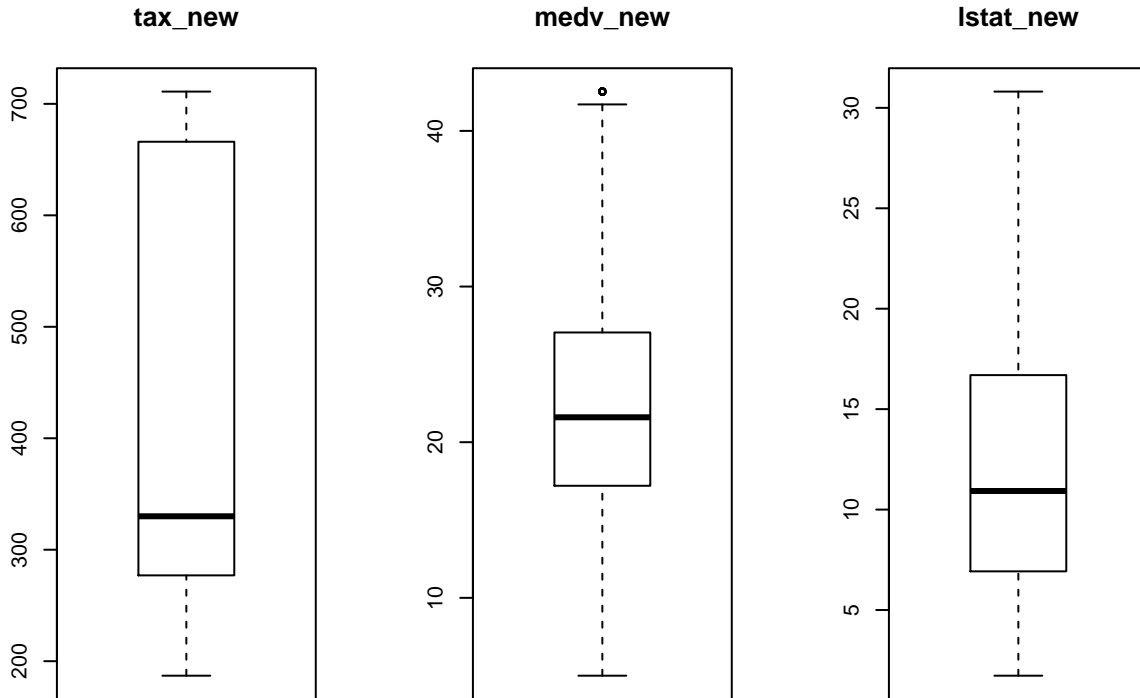- Data transformation

## 2.1 Outliers treatment

For outliers, we will create 2 sets of variables.

The first set uses the capping method. In this method, we will replace all outliers that lie outside the 1.5 times of IQR limits. We will cap it by replacing those observations less than the lower limit with the value of 5th %ile and those that lie above the upper limit with the value of 95th %ile.

Accordingly we create the following new variables while retaining the original variables.

city_crime_train$tax $city_crime_train$medv
city_crime_train$lstat

Lets see how the new variables look in boxplots.



In the second set, we will use the sin transformation and create the following variables:

city_crime_train_mod$rm_new $city_crime_train_mod$dis_new

## 2.3 Tranformation for Variables

Following variables will need some transformation:

1. zn: proportion of residential land zoned for large lots (over 25000 square feet) (predictor variable)
2. chas: a dummy var. for whether the suburb borders the Charles River (1) or not (0) (predictor variable)
3. target: whether the crime rate is above the median crime rate (1) or not (0) (response variable)

## 2.6

Lets see how the new variables stack up against wins.

All new variables seem to have a positive correlation with wins. However, some of them do not seem to have a strong correlation. Lets see how they perform while modeling.

# 3 Build Models

Below is a summary table showing models and their respective variables.

### 3.1.1 Model One by using all given variable

In this model, we will be using the original variables. We will create model and we will highlight the variables that being recommended using the AIC value.
First we will produce the summary model as per below:

```
##
## Call:
## glm(formula = target ~ ., family = "binomial", data = city_crime_train)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.8791  -0.1299  -0.0025   0.0011   3.4785
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -41.462153   8.250799  -5.025 5.03e-07 ***
## zn           -0.060580   0.039153  -1.547 0.121799
## indus        -0.063885   0.059335  -1.077 0.281618
## chas          0.789391   0.865818   0.912 0.361912
## nox          53.413503  10.013666   5.334 9.60e-08 ***
## rm           -0.647942   0.904430  -0.716 0.473739
## age           0.028835   0.015680   1.839 0.065915 .
## dis           0.800917   0.268877   2.979 0.002894 **
## rad           0.721751   0.195662   3.689 0.000225 ***
## tax          -0.007065   0.003490  -2.024 0.042948 *
## ptratio       0.440768   0.159366   2.766 0.005679 **
## black        -0.009591   0.006025  -1.592 0.111412
## lstat         0.096941   0.062429   1.553 0.120469
## medv          0.236940   0.091276   2.596 0.009436 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 514.63  on 371  degrees of freedom
## Residual deviance: 140.71  on 358  degrees of freedom
## AIC: 168.71
##
## Number of Fisher Scoring iterations: 9
```

Accuracy=0.9042553

### 3.1.2 Model two- with backward step function with all given variables

```
stepmodel1<- step(model1, direction="backward")
```

```
## Start:  AIC=168.71
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##     ptratio + black + lstat + medv
##
```

```
##            Df Deviance    AIC
## - rm       1    141.22 167.22
## - chas     1    141.55 167.55
## - indus    1    141.93 167.93
## <none>          140.71 168.71
## - lstat    1    143.06 169.06
## - black    1    143.68 169.68
## - zn       1    143.99 169.99
## - age      1    144.45 170.45
## - tax      1    144.93 170.93
## - medv     1    148.67 174.67
## - ptratio  1    149.29 175.29
## - dis      1    150.97 176.97
## - rad      1    171.94 197.94
## - nox      1    195.65 221.65
##
## Step:  AIC=167.22
## target ~ zn + indus + chas + nox + age + dis + rad + tax + ptratio +
##     black + lstat + medv
##
##            Df Deviance    AIC
## - chas     1    142.10 166.10
## - indus    1    142.37 166.37
## <none>          141.22 167.22
## - black    1    144.02 168.02
## - age      1    144.48 168.48
## - zn       1    144.74 168.74
## - lstat    1    145.13 169.13
## - tax      1    145.97 169.97
## - ptratio  1    149.78 173.78
## - dis      1    150.97 174.97
## - medv     1    156.73 180.73
## - rad      1    172.26 196.26
## - nox      1    196.29 220.29
##
## Step:  AIC=166.1
## target ~ zn + indus + nox + age + dis + rad + tax + ptratio +
##     black + lstat + medv
##
##            Df Deviance    AIC
## - indus    1    142.85 164.85
## <none>          142.10 166.10
## - black    1    144.69 166.69
## - age      1    145.65 167.65
## - zn       1    146.09 168.09
## - lstat    1    146.43 168.43
## - tax      1    148.34 170.34
## - ptratio  1    149.90 171.90
## - dis      1    151.42 173.42
## - medv     1    157.16 179.16
## - rad      1    177.68 199.68
## - nox      1    196.44 218.44
##
## Step:  AIC=164.85
```

```
## target ~ zn + nox + age + dis + rad + tax + ptratio + black +
##     lstat + medv
##
##          Df Deviance    AIC
## <none>        142.85 164.85
## - black    1  145.21 165.21
## - age      1  146.69 166.69
## - lstat    1  146.75 166.75
## - zn       1  146.89 166.89
## - ptratio  1  150.46 170.46
## - dis      1  151.87 171.87
## - tax      1  154.08 174.08
## - medv     1  157.59 177.59
## - rad      1  184.71 204.71
## - nox      1  203.12 223.12
```

```
pre_train1_step<-predict(stepmodel1,type="response",newdata=train_test)

table(pre_train1_step>0.5,train_test$target)
```

```
##
##          0  1
##   FALSE 34  5
##   TRUE   7 48
```

Accuracy=0.8723404

### 3.1.3 Model three- model with transformed variables

In this model, we will be using the some transformed variables.

First we will produce the summary model as per below:

```
##
## Call:
## glm(formula = target ~ . - zn - rm - dis - tax - lstat - medv,
##     family = "binomial", data = city_crime_train_mod)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.8490  -0.1466  -0.0024   0.0004   3.5826
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -45.541738   8.544894  -5.330 9.84e-08 ***
## indus         0.014531   0.064926   0.224 0.822909
## chas          0.108863   0.811295   0.134 0.893257
## nox          50.472586   9.083435   5.557 2.75e-08 ***
## age           0.036435   0.016117   2.261 0.023780 *
## rad           0.871309   0.241452   3.609 0.000308 ***
## ptratio       0.495086   0.172513   2.870 0.004107 **
## black        -0.010433   0.005881  -1.774 0.076036 .
```

```
## tax_new       -0.005498    0.003495  -1.573 0.115648
## medv_new       0.297542    0.090676   3.281 0.001033 **
## lstat_new      0.053168    0.069612   0.764 0.444998
## rm_new        -1.774497    1.144107  -1.551 0.120904
## dis_new       -2.191201    0.532281  -4.117 3.84e-05 ***
## zn_new         0.465684    0.892871   0.522 0.601978
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 514.63  on 371  degrees of freedom
## Residual deviance: 129.52  on 358  degrees of freedom
## AIC: 157.52
##
## Number of Fisher Scoring iterations: 9


##
##          0  1
##    FALSE 35  3
##    TRUE   6 50
```

Accuracy=0.9042553

**3.1.4 Model with transformed variable and with with backward step function**

```
stepmodel2<- step(model2, direction="backward")
```

```
## Start:  AIC=157.52
## target ~ (zn + indus + chas + nox + rm + age + dis + rad + tax +
##     ptratio + black + lstat + medv + tax_new + medv_new + lstat_new +
##     rm_new + dis_new + zn_new) - zn - rm - dis - tax - lstat -
##     medv
##
##             Df Deviance    AIC
## - chas       1   129.54 155.54
## - indus      1   129.57 155.57
## - zn_new     1   129.79 155.79
## - lstat_new  1   130.08 156.08
## <none>           129.52 157.52
## - tax_new    1   131.92 157.92
## - rm_new     1   131.97 157.97
## - black      1   132.86 158.86
## - age        1   135.31 161.31
## - ptratio    1   138.64 164.64
## - medv_new   1   142.81 168.81
## - dis_new    1   151.54 177.54
## - rad        1   155.24 181.24
## - nox        1   197.04 223.04
##
## Step:  AIC=155.54
```

```
## target ~ indus + nox + age + rad + ptratio + black + tax_new +
##     medv_new + lstat_new + rm_new + dis_new + zn_new
##
##              Df Deviance    AIC
## - indus       1   129.61 153.61
## - zn_new      1   129.79 153.79
## - lstat_new   1   130.13 154.13
## <none>           129.54 155.54
## - rm_new      1   131.99 155.99
## - tax_new     1   132.13 156.13
## - black       1   132.86 156.86
## - age         1   135.51 159.51
## - ptratio     1   138.79 162.79
## - medv_new    1   142.84 166.84
## - dis_new     1   152.03 176.03
## - rad         1   156.60 180.60
## - nox         1   197.61 221.61
##
## Step:  AIC=153.61
## target ~ nox + age + rad + ptratio + black + tax_new + medv_new +
##     lstat_new + rm_new + dis_new + zn_new
##
##              Df Deviance    AIC
## - zn_new      1   129.82 151.82
## - lstat_new   1   130.28 152.28
## <none>           129.61 153.61
## - rm_new      1   132.04 154.04
## - tax_new     1   132.51 154.51
## - black       1   132.99 154.99
## - age         1   135.51 157.51
## - ptratio     1   138.80 160.80
## - medv_new    1   143.10 165.10
## - dis_new     1   152.60 174.60
## - rad         1   161.77 183.77
## - nox         1   209.86 231.86
##
## Step:  AIC=151.82
## target ~ nox + age + rad + ptratio + black + tax_new + medv_new +
##     lstat_new + rm_new + dis_new
##
##              Df Deviance    AIC
## - lstat_new   1   130.87 150.87
## <none>           129.82 151.82
## - rm_new      1   132.04 152.04
## - tax_new     1   132.69 152.69
## - black       1   133.06 153.06
## - age         1   135.52 155.52
## - ptratio     1   139.74 159.74
## - medv_new    1   143.10 163.10
## - dis_new     1   152.65 172.65
## - rad         1   162.06 182.06
## - nox         1   212.46 232.46
##
## Step:  AIC=150.86
```

```
## target ~ nox + age + rad + ptratio + black + tax_new + medv_new +
##     rm_new + dis_new
##
##           Df Deviance    AIC
## <none>          130.87 150.87
## - tax_new   1  133.34 151.34
## - black     1  133.89 151.89
## - rm_new    1  135.44 153.44
## - age       1  139.74 157.74
## - ptratio   1  141.03 159.03
## - medv_new  1  143.94 161.94
## - dis_new   1  154.34 172.34
## - rad       1  163.53 181.53
## - nox       1  213.91 231.91
```

```
pre_train2_step<-predict(stepmodel2,type="response",newdata=train_test_mod)

table(pre_train2_step>0.5,train_test_mod$target)
```

```
##
##          0  1
##   FALSE 35  4
##   TRUE   6 49
```

Accuracy= 0.893617

**3.1,5 Model three with Linear discrement analysis**

```
##    class  posterior.0 posterior.1        LD1
## 3      1 0.0005609314  0.99943907  2.9179352
## 6      0 0.8593842086  0.14061579 -0.5664873
## 7      1 0.0040700562  0.99592994  2.1737359
## 8      1 0.0014576826  0.99854232  2.5596162
## 23     0 0.9672384727  0.03276153 -1.1568765
```

```
##
##      0  1
##   0 39 14
##   1  2 39
```

Accuracy=0.8297872

**3.1.6 Model with Linear discrement analysis with transformed data**

```
##
##      0  1
##   0 39 14
##   1  2 39
```

Accuracy=0.7978723

17

# 4 Model Selection

In section we will further examine all six models. We will apply a model selection strategy defined below to compare the models.

## 4.1 Model selection strategy:

Following model selection strategy has been used for this assignment:

(1) Compare accuracy of the models & confusion matrix
(2) Compare Precision,Sensitivity,Specificity,F1 score
(3) Compare AUC curve for the models

```
##
##           0  1
##   FALSE  36  4
##   TRUE    5 49
```

ROC Curve w/ AUC=0.954901058444547