# 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

```
## 'data.frame':    466 obs. of  14 variables:
## $ zn     : num  0 0 0 30 0 0 0 0 0 80 ...
## $ indus  : num  19.58 19.58 18.1 4.93 2.46 ...
## $ chas   : int  0 1 0 0 0 0 0 0 0 0 ...
## $ nox    : num  0.605 0.871 0.74 0.428 0.488 0.52 0.693 0.693 0.515 0.392 ...
## $ rm     : num  7.93 5.4 6.49 6.39 7.16 ...
## $ age    : num  96.2 100 100 7.8 92.2 71.3 100 100 38.1 19.1 ...
## $ dis    : num  2.05 1.32 1.98 7.04 2.7 ...
## $ rad    : int  5 5 24 6 3 5 24 24 5 1 ...
## $ tax    : int  403 403 666 300 193 384 666 666 224 315 ...
## $ ptratio: num  14.7 14.7 20.2 16.6 17.8 20.9 20.2 20.2 20.2 16.4 ...
## $ black  : num  369 397 387 375 394 ...
## $ lstat  : num  3.7 26.82 18.85 5.19 4.82 ...
## $ medv   : num  50 13.4 15.4 23.7 37.9 26.5 5 7 22.2 20.9 ...
## $ target : int  1 1 1 0 0 0 1 1 0 0 ...


##       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 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 ...
```

# 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.: 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


##
##         0         1
## 0.5085837 0.4914163
```

## 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            1 466  11.58   23.36   0.00    5.35   0.00   0.00 100.00 100.00
## indus         2 466  11.11    6.85   9.69   10.91   9.34   0.46  27.74  27.28
## chas          3 466   0.07    0.26   0.00    0.00   0.00   0.00   1.00   1.00
## nox           4 466   0.55    0.12   0.54    0.54   0.13   0.39   0.87   0.48
## rm            5 466   6.29    0.70   6.21    6.26   0.52   3.86   8.78   4.92
## age           6 466  68.37   28.32  77.15   70.96  30.02   2.90 100.00  97.10
## dis           7 466   3.80    2.11   3.19    3.54   1.91   1.13  12.13  11.00
## rad           8 466   9.53    8.69   5.00    8.70   1.48   1.00  24.00  23.00
## tax           9 466 409.50  167.90 334.50  401.51 104.52 187.00 711.00 524.00
## ptratio      10 466  18.40    2.20  18.90   18.60   1.93  12.60  22.00   9.40
## black        11 466 357.12   91.32 391.34  383.51   8.24   0.32 396.90 396.58
## lstat        12 466  12.63    7.10  11.35   11.88   7.07   1.73  37.97  36.24
## medv         13 466  22.59    9.24  21.20   21.63   6.00   5.00  50.00  45.00
## target       14 466   0.49    0.50   0.00    0.49   0.00   0.00   1.00   1.00
##            skew kurtosis   se
## zn         2.18     3.81 1.08
## indus      0.29    -1.24 0.32
## chas       3.34     9.15 0.01
## nox        0.75    -0.04 0.01
## rm         0.48     1.54 0.03
## age       -0.58    -1.01 1.31
## dis        1.00     0.47 0.10
## rad        1.01    -0.86 0.40
## tax        0.66    -1.15 7.78
## ptratio   -0.75    -0.40 0.10
## black     -2.92     7.34 4.23
## lstat      0.91     0.50 0.33
## medv       1.08     1.37 0.43
## target     0.03    -2.00 0.02


##      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     | -0.4316818  |
| indus  | 0.6048507   |
| chas   | 0.0800419   |
| nox    | 0.7261062   |
| rm     | -0.1525533  |
| age    | 0.6301062   |
| dis    | -0.6186731  |
| rad    | 0.6281049   |
| tax    | 0.6111133   |

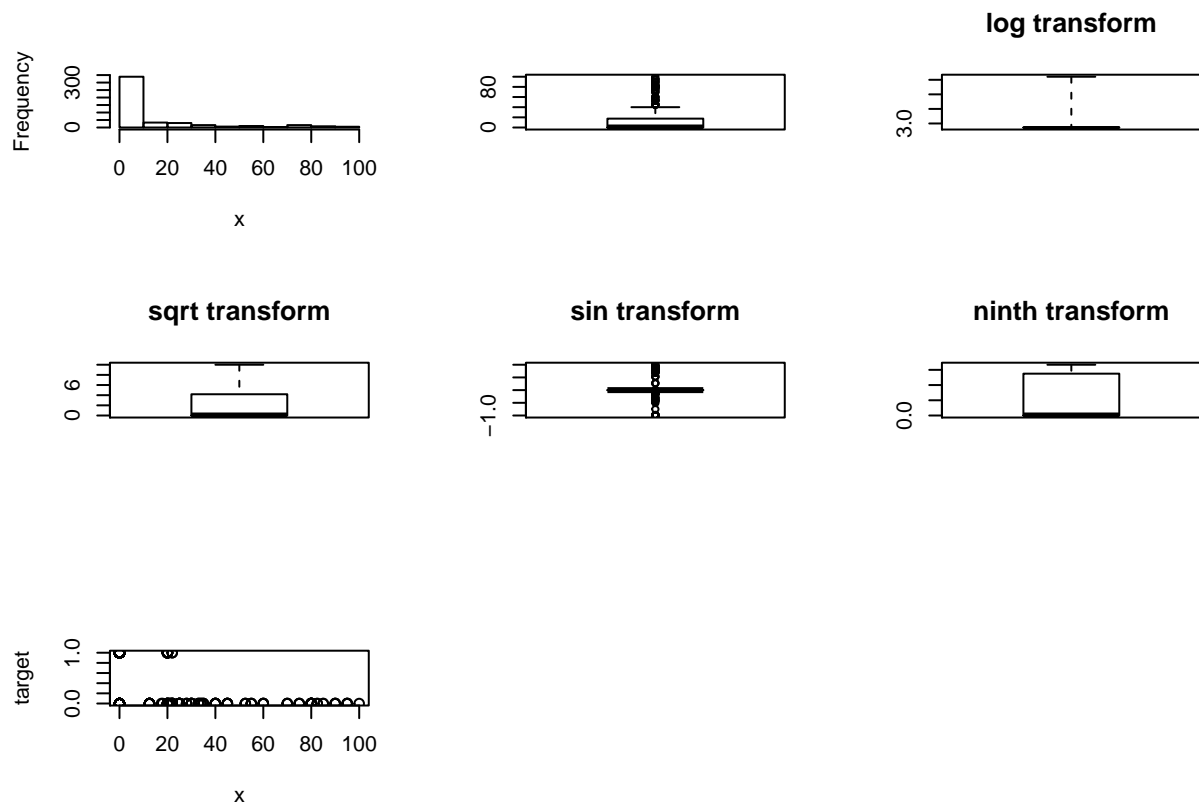|         | Correlation |
|---------|-------------|
| ptratio | 0.2508489   |
| black   | -0.3529568  |
| lstat   | 0.4691270   |
| medv    | -0.2705507  |
| target  | 1.0000000   |

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

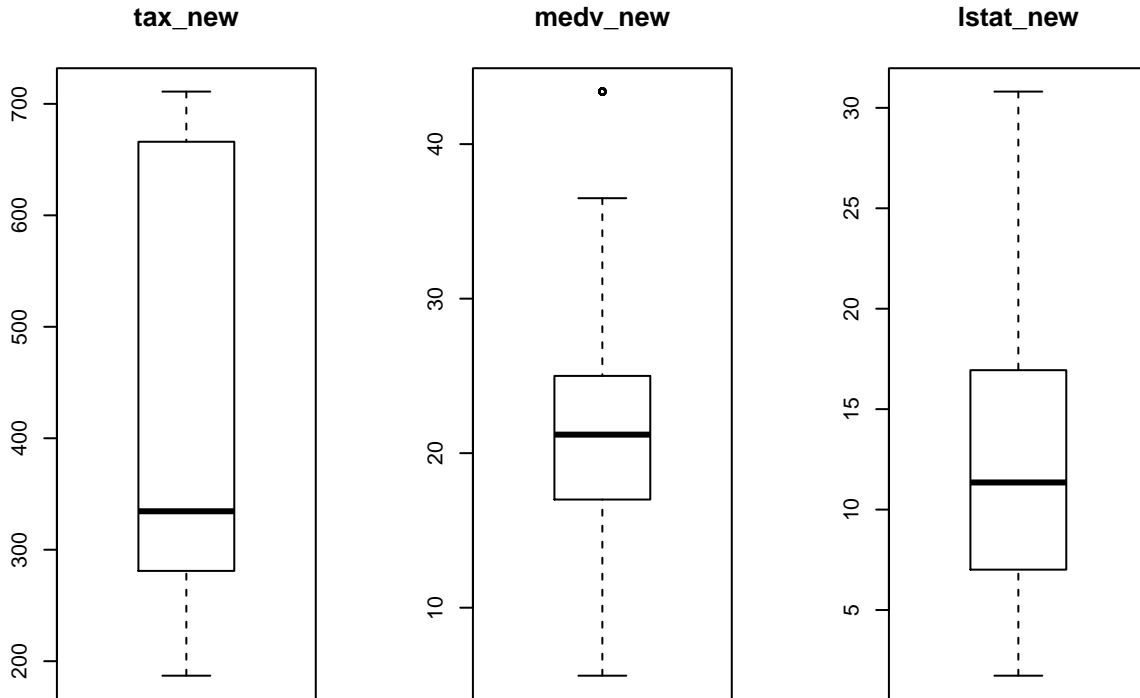For outliers, we will create 2 sets of variables.

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.
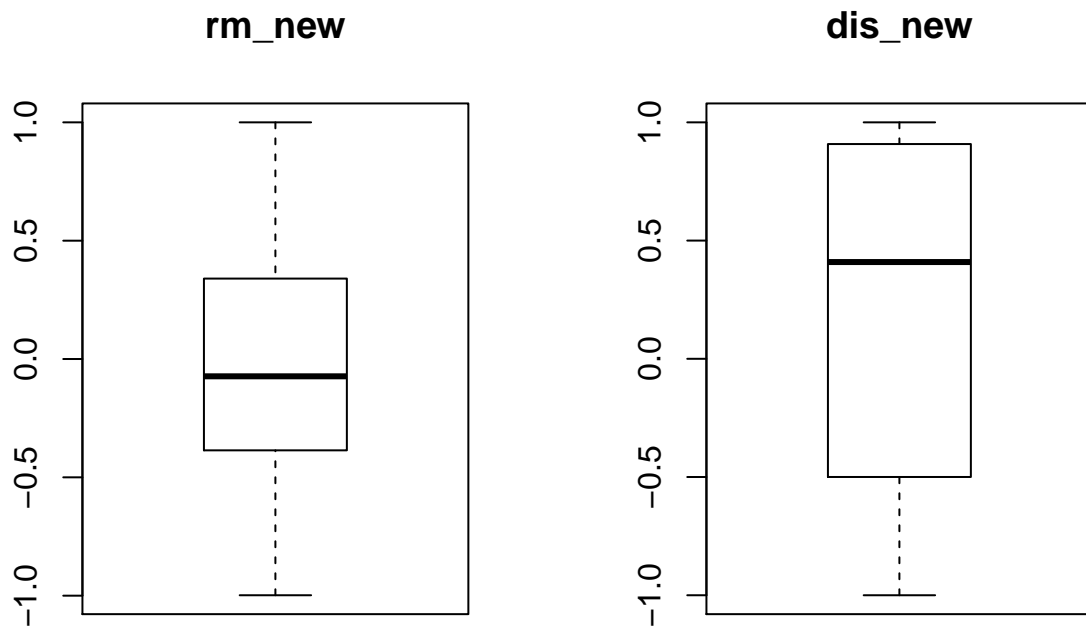
city_crime_train$tax_new city_crime_train$medv_new
city_crime_train$lstat

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$rm_new city_crime_train$dis_new

## rm_new          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 Model One

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 ~ . - zn_new - rm_new - lstat_new - tax_new -
##     medv_new, family = "binomial", data = city_crime_train)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.8469  -0.1389  -0.0017   0.0007   3.3050
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -39.967503   7.480155  -5.343 9.14e-08 ***
## zn           -0.014489   0.028929  -0.501  0.61649
## indus         0.003888   0.055889   0.070  0.94454
## chas1         0.464046   0.743387   0.624  0.53248
## nox          53.269267   8.209555   6.489 8.66e-11 ***
## rm           -0.862730   0.851113  -1.014  0.31075
## age           0.041618   0.015426   2.698  0.00698 **
## dis           0.476938   0.276948   1.722  0.08505 .
## rad           0.800898   0.200223   4.000 6.33e-05 ***
## tax          -0.005040   0.003082  -1.635  0.10201
## ptratio       0.442846   0.142815   3.101  0.00193 **
## black        -0.011963   0.005945  -2.012  0.04421 *
## lstat         0.036461   0.057025   0.639  0.52257
## medv          0.230818   0.078475   2.941  0.00327 **
## dis_new      -1.892538   0.477332  -3.965 7.34e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 165.64  on 451  degrees of freedom
## AIC: 195.64
##
## Number of Fisher Scoring iterations: 9


##
## target FALSE TRUE
##      0   234    3
##      1    35  194
```

### 3.1 Model One with backward step function

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

```
## Start:  AIC=195.64
## 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_new - rm_new - lstat_new -
##     tax_new - medv_new
##
##             Df Deviance    AIC
## - indus      1   165.64 193.64
## - zn         1   165.91 193.91
## - chas       1   166.02 194.02
## - lstat      1   166.04 194.04
## - rm         1   166.67 194.67
## <none>           165.63 195.63
## - tax        1   168.25 196.25
## - dis        1   168.29 196.29
## - black      1   170.84 198.84
## - age        1   173.83 201.83
## - medv       1   175.65 203.65
## - ptratio    1   176.06 204.06
## - dis_new    1   186.15 214.15
## - rad        1   197.68 225.68
## - nox        1   246.22 274.22
##
## Step:  AIC=193.64
## target ~ zn + chas + nox + rm + age + dis + rad + tax + ptratio +
##     black + lstat + medv + dis_new
##
##             Df Deviance    AIC
## - zn         1   165.93 191.93
## - lstat      1   166.05 192.05
## - chas       1   166.08 192.08
## - rm         1   166.68 192.68
## <none>           165.64 193.64
## - dis        1   168.29 194.29
## - tax        1   168.88 194.88
## - black      1   170.88 196.88
## - age        1   173.85 199.85
## - medv       1   175.68 201.68
## - ptratio    1   176.11 202.11
## - dis_new    1   188.52 214.52
## - rad        1   203.74 229.74
## - nox        1   254.38 280.38
##
## Step:  AIC=191.93
## target ~ chas + nox + rm + age + dis + rad + tax + ptratio +
##     black + lstat + medv + dis_new
##
##             Df Deviance    AIC
## - lstat      1   166.24 190.24
## - chas       1   166.44 190.44
## - rm         1   167.27 191.27
## <none>           165.93 191.93
## - dis        1   168.35 192.35
## - tax        1   169.33 193.33
```

```
## - black     1    171.28 195.28
## - age       1    174.85 198.85
## - medv      1    176.18 200.18
## - ptratio   1    178.45 202.45
## - dis_new   1    193.29 217.29
## - rad       1    206.94 230.94
## - nox       1    256.56 280.56
##
## Step:  AIC=190.24
## target ~ chas + nox + rm + age + dis + rad + tax + ptratio +
##      black + medv + dis_new
##
##             Df Deviance    AIC
## - chas       1    166.88 188.88
## <none>            166.24 190.24
## - rm         1    168.56 190.56
## - dis        1    168.93 190.93
## - tax        1    169.45 191.45
## - black      1    171.49 193.49
## - medv       1    176.71 198.71
## - age        1    178.84 200.84
## - ptratio    1    179.38 201.38
## - dis_new    1    193.58 215.58
## - rad        1    207.44 229.44
## - nox        1    258.50 280.50
##
## Step:  AIC=188.88
## target ~ nox + rm + age + dis + rad + tax + ptratio + black +
##      medv + dis_new
##
##             Df Deviance    AIC
## <none>            166.88 188.88
## - dis        1    169.24 189.24
## - rm         1    169.51 189.51
## - tax        1    170.28 190.28
## - black      1    171.96 191.96
## - medv       1    177.75 197.75
## - ptratio    1    179.47 199.47
## - age        1    180.74 200.74
## - dis_new    1    195.77 215.77
## - rad        1    209.89 229.89
## - nox        1    258.55 278.55
```

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

table(target,pre_train1_step >0.75)
```

```
##
## target FALSE TRUE
##      0   234    3
##      1    34  195
```

## 3.2 Model two

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)
##
## Deviance Residuals:
##     Min       1Q    Median       3Q      Max
## -1.7978  -0.1372  -0.0012   0.0006   3.7479
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -41.631145   7.066597  -5.891 3.83e-09 ***
## indus         0.010371   0.055697   0.186 0.852285
## chas1         0.386972   0.716324   0.540 0.589045
## nox          49.978105   7.574874   6.598 4.17e-11 ***
## age           0.039466   0.014272   2.765 0.005689 **
## rad           0.823571   0.203804   4.041 5.32e-05 ***
## ptratio       0.432948   0.150308   2.880 0.003972 **
## black        -0.011718   0.005890  -1.990 0.046641 *
## tax_new      -0.005283   0.003017  -1.751 0.079911 .
## medv_new      0.224594   0.067847   3.310 0.000932 ***
## lstat_new     0.021292   0.063375   0.336 0.736891
## rm_new       -1.395547   0.963844  -1.448 0.147646
## dis_new      -2.328906   0.475160  -4.901 9.52e-07 ***
## zn_new1       0.296890   0.802123   0.370 0.711285
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 167.31  on 452  degrees of freedom
## AIC: 195.31
##
## Number of Fisher Scoring iterations: 9


##
## target FALSE TRUE
##      0   235    2
##      1    34  195
```

## 3.1 Model two with backward step function

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

```
## Start:  AIC=195.31
```

```
## 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
## - indus       1   167.35 193.35
## - lstat_new   1   167.43 193.43
## - zn_new      1   167.45 193.45
## - chas        1   167.60 193.60
## <none>            167.31 195.31
## - rm_new      1   169.44 195.44
## - tax_new     1   170.33 196.33
## - black       1   172.33 198.33
## - age         1   175.88 201.88
## - ptratio     1   176.18 202.18
## - medv_new    1   180.15 206.15
## - dis_new     1   199.36 225.36
## - rad         1   200.70 226.70
## - nox         1   257.91 283.91
##
## Step:  AIC=193.35
## target ~ chas + nox + age + rad + ptratio + black + tax_new +
##     medv_new + lstat_new + rm_new + dis_new + zn_new
##
##               Df Deviance    AIC
## - zn_new      1   167.47 191.47
## - lstat_new   1   167.47 191.47
## - chas        1   167.72 191.72
## <none>            167.35 193.35
## - rm_new      1   169.45 193.45
## - tax_new     1   170.82 194.82
## - black       1   172.40 196.40
## - age         1   175.88 199.88
## - ptratio     1   176.20 200.20
## - medv_new    1   180.20 204.20
## - dis_new     1   201.51 225.51
## - rad         1   207.02 231.02
## - nox         1   269.29 293.29
##
## Step:  AIC=191.46
## target ~ chas + nox + age + rad + ptratio + black + tax_new +
##     medv_new + lstat_new + rm_new + dis_new
##
##               Df Deviance    AIC
## - lstat_new   1   167.72 189.72
## - chas        1   167.76 189.76
## <none>            167.47 191.47
## - rm_new      1   169.49 191.49
## - tax_new     1   170.94 192.94
## - black       1   172.41 194.41
## - age         1   176.02 198.02
## - ptratio     1   178.47 200.47
## - medv_new    1   180.23 202.23
```

```
## - dis_new     1    201.58 223.58
## - rad         1    207.88 229.88
## - nox         1    273.49 295.49
##
## Step:  AIC=189.72
## target ~ chas + nox + age + rad + ptratio + black + tax_new +
##     medv_new + rm_new + dis_new
##
##              Df Deviance    AIC
## - chas        1   168.12 188.12
## <none>            167.72 189.72
## - rm_new      1   171.01 191.01
## - tax_new     1   171.06 191.06
## - black       1   172.58 192.58
## - ptratio     1   178.82 198.82
## - age         1   179.03 199.03
## - medv_new    1   180.24 200.24
## - dis_new     1   201.70 221.70
## - rad         1   208.38 228.38
## - nox         1   273.77 293.77
##
## Step:  AIC=188.11
## target ~ nox + age + rad + ptratio + black + tax_new + medv_new +
##     rm_new + dis_new
##
##              Df Deviance    AIC
## <none>            168.12 188.12
## - tax_new     1   171.58 189.58
## - rm_new      1   171.61 189.61
## - black       1   172.85 190.85
## - ptratio     1   178.87 196.87
## - age         1   180.79 198.79
## - medv_new    1   181.00 199.00
## - dis_new     1   203.04 221.04
## - rad         1   210.44 228.44
## - nox         1   273.82 291.82
```

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

table(target,pre_train2_step >0.75)
```

```
##
## target FALSE TRUE
##      0   235    2
##      1    34  195
```

## 3.3 Model three with leap package

```r
# install.packages("ISLR")
# install.packages("leaps")
```

17

```r
par(mfrow=c(1,1))
library(ISLR)
library(leaps)

#We will now use the package leaps to evaluate all the best-subset models.
#It gives by default best-subsets up to size 8; lets increase that to 18, i.e. all the variables
regfit <- regsubsets(target~., data = city_crime_train, nvmax = 18)
```

```
## Warning in leaps.setup(x, y, wt = wt, nbest = nbest, nvmax = nvmax,
## force.in = force.in, : 1 linear dependencies found
```

```
## Reordering variables and trying again:
```

```r
summary(regfit)
```

```
## Subset selection object
## Call: regsubsets.formula(target ~ ., data = city_crime_train, nvmax = 18)
## 19 Variables  (and intercept)
##            Forced in Forced out
## zn             FALSE      FALSE
## indus          FALSE      FALSE
## chas1          FALSE      FALSE
## nox            FALSE      FALSE
## rm             FALSE      FALSE
## age            FALSE      FALSE
## dis            FALSE      FALSE
## rad            FALSE      FALSE
## tax            FALSE      FALSE
## ptratio        FALSE      FALSE
## black          FALSE      FALSE
## lstat          FALSE      FALSE
## medv           FALSE      FALSE
## medv_new       FALSE      FALSE
## lstat_new      FALSE      FALSE
## rm_new         FALSE      FALSE
## dis_new        FALSE      FALSE
## zn_new1        FALSE      FALSE
## tax_new        FALSE      FALSE
## 1 subsets of each size up to 18
## Selection Algorithm: exhaustive
##           zn  indus chas1 nox rm  age dis rad tax ptratio black lstat medv
## 1  ( 1 )  " " " "   " "   " " "*" " " " " " " " " " "     " "   " "   " "
## 2  ( 1 )  " " " "   " "   " " "*" " " " " " " "*" " "     " "   " "   " "
## 3  ( 1 )  " " " "   " "   " " "*" " " "*" " " "*" " "     " "   " "   " "
## 4  ( 1 )  " " " "   " "   " " "*" " " "*" " " "*" " "     " "   " "   "*"
## 5  ( 1 )  " " " "   " "   " " "*" " " "*" " " "*" " "     " "   " "   "*"
## 6  ( 1 )  " " " "   " "   " " "*" " " "*" " " "*" " "     "*"   " "   "*"
## 7  ( 1 )  " " " "   " "   " " "*" "*" "*" " " "*" " "     " "   " "   "*"
## 8  ( 1 )  " " " "   " "   " " "*" "*" "*" " " "*" " "     "*"   " "   "*"
## 9  ( 1 )  " " " "   " "   " " "*" "*" "*" " " "*" " "     "*"   "*"   "*"
## 10  ( 1 ) " " " "   " "   " " "*" "*" "*" "*" "*" " "     "*"   " "   "*"
## 11  ( 1 ) " " " "   " "   " " "*" "*" "*" "*" "*" " "     "*"   "*"   "*"
```

18

```
## 12  ( 1 ) " " " " " "     " "    "*" "*" "*" "*" "*" " " " " " "     "*"    "*"    "*"
## 13  ( 1 ) " " "*"    " "    "*" "*" "*" "*" "*" " " " " " "     "*"    "*"    "*"
## 14  ( 1 ) "*" "*"    " "    "*" "*" "*" "*" "*" " " " " " "     "*"    "*"    "*"
## 15  ( 1 ) "*" "*"    " "    "*" "*" "*" "*" "*" "*" " " " "     "*"    "*"    "*"
## 16  ( 1 ) "*" "*"    "*"    "*" "*" "*" "*" "*" " " " " " "     "*"    "*"    "*"
## 17  ( 1 ) "*" "*"    "*"    "*" "*" "*" "*" "*" "*" "*"         "*"    "*"    "*"
## 18  ( 1 ) "*" "*"    "*"    "*" "*" "*" "*" "*" "*" "*"         "*"    "*"    "*"
##           tax_new medv_new lstat_new rm_new dis_new zn_new1
## 1   ( 1 ) " "     " "      " "       " "    " "     " "
## 2   ( 1 ) " "     " "      " "       " "    " "     " "
## 3   ( 1 ) " "     " "      " "       " "    " "     " "
## 4   ( 1 ) " "     " "      " "       " "    " "     " "
## 5   ( 1 ) " "     " "      " "       " "    "*"     " "
## 6   ( 1 ) " "     " "      " "       " "    "*"     " "
## 7   ( 1 ) " "     " "      " "       "*"    "*"     " "
## 8   ( 1 ) " "     " "      " "       "*"    "*"     " "
## 9   ( 1 ) " "     " "      " "       "*"    "*"     " "
## 10  ( 1 ) " "     " "      " "       "*"    "*"     "*"
## 11  ( 1 ) " "     " "      " "       "*"    "*"     "*"
## 12  ( 1 ) " "     " "      "*"       "*"    "*"     "*"
## 13  ( 1 ) " "     " "      "*"       "*"    "*"     "*"
## 14  ( 1 ) " "     " "      "*"       "*"    "*"     "*"
## 15  ( 1 ) " "     " "      "*"       "*"    "*"     "*"
## 16  ( 1 ) "*"     " "      "*"       "*"    "*"     "*"
## 17  ( 1 ) " "     " "      "*"       "*"    "*"     "*"
## 18  ( 1 ) " "     "*"      "*"       "*"    "*"     "*"
```
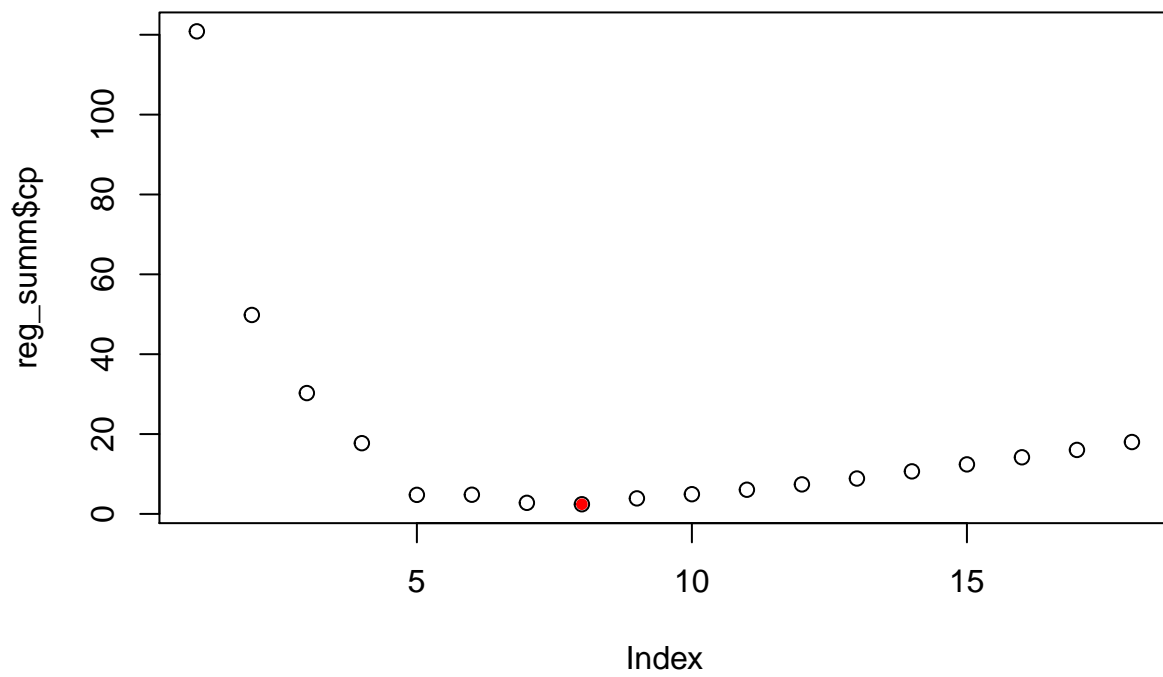
```r
reg_summ <- summary(regfit)
names(reg_summ)
```

```
## [1] "which"  "rsq"    "rss"    "adjr2"  "cp"     "bic"    "outmat" "obj"
```
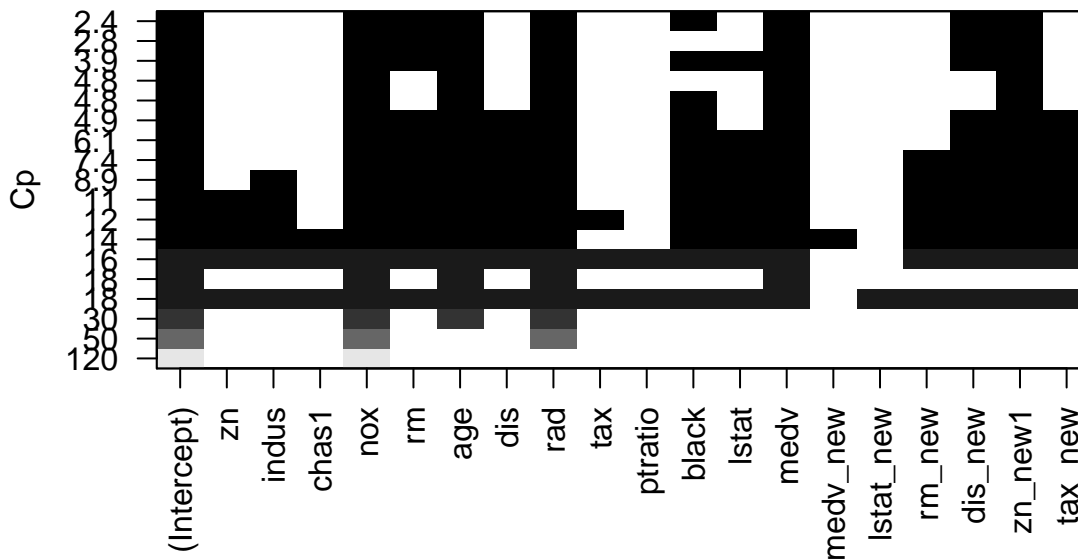
```r
#finding the lowest cp value
#cp or adjr2 or r2 is the value of the chosen model selection statistic for each model
plot(reg_summ$cp)
which.min(reg_summ$cp)
```

```
## [1] 8
```

```r
points(8, reg_summ$cp[8], pch=20,col="red")
```

```
#There is a plot method for the regsubsets object
plot(regfit, scale = "Cp")
```

```r
coef(regfit, 8)
```

```
##   (Intercept)            nox            rm            age            rad
## -0.1591910366   2.0020380540   0.0075020987   0.0039607779   0.0188066131
##         black           medv        dis_new       zn_new1
## -0.0002437521   0.0071614965  -0.0913255905  -0.0398451027
```

```r
model3 <- glm(target ~ nox+rm+age+rad+black+medv+dis_new+zn_new, data = city_crime_train, family = "bin
summary(model3)
```

```
##
## Call:
## glm(formula = target ~ nox + rm + age + rad + black + medv +
##     dis_new + zn_new, family = "binomial", data = city_crime_train)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.0426  -0.2088  -0.0045   0.0028   4.0464
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -26.885976    5.293377  -5.079 3.79e-07 ***
## nox          44.462712    6.611660   6.725 1.76e-11 ***
## rm           -0.304011    0.652704  -0.466 0.641379
```

```
## age            0.032627    0.011786    2.768 0.005635 **
## rad            0.582932    0.151465    3.849 0.000119 ***
## black         -0.009052    0.005210   -1.737 0.082313 .
## medv           0.124413    0.057666    2.157 0.030967 *
## dis_new       -2.529076    0.440512   -5.741 9.40e-09 ***
## zn_new1       -0.839678    0.625105   -1.343 0.179188
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 180.15  on 457  degrees of freedom
## AIC: 198.15
##
## Number of Fisher Scoring iterations: 9
```

```r
pre_train3 <-predict(model3,type="response")

table(target,pre_train3 > 0.5)
```

```
##
## target FALSE TRUE
##      0   222   15
##      1    15  214
```