

INFX 573: Problem Set 6 - Regression

Amol Surve

Due: Tuesday, November 15, 2016

Collaborators: Abhishek Gupta

Instructions:

Before beginning this assignment, please ensure you have access to R and RStudio.

1. Download the `problemset6.Rmd` file from Canvas. Open `problemset6.Rmd` in RStudio and supply your solutions to the assignment by editing `problemset6.Rmd`.
2. Replace the “Insert Your Name Here” text in the `author:` field with your own full name. Any collaborators must be listed on the top of your assignment.
3. Be sure to include well-documented (e.g. commented) code chunks, figures and clearly written text chunk explanations as necessary. Any figures should be clearly labeled and appropriately referenced within the text.
4. Collaboration on problem sets is acceptable, and even encouraged, but each student must turn in an individual write-up in his or her own words and his or her own work. The names of all collaborators must be listed on each assignment. Do not copy-and-paste from other students’ responses or code.
5. When you have completed the assignment and have **checked** that your code both runs in the Console and knits correctly when you click Knit PDF, rename the R Markdown file to `YourLastName_YourFirstName_ps6.Rmd`, knit a PDF and submit the PDF file on Canvas.

Setup:

In this problem set you will need, at minimum, the following R packages.

```
# Load standard libraries
library(tidyverse)
library(MASS)
# Modern applied statistics functions
library(leaps)
?Boston
boston<-Boston
```

Housing Values in Suburbs of Boston

In this problem we will use the Boston dataset that is available in the `MASS` package. This dataset contains information about median house value for 506 neighborhoods in Boston, MA. Load this data and use it to answer the following questions.

1. Describe the data and variables that are part of the `Boston` dataset. Tidy data as necessary. The Boston data frame has 506 rows and following 14 columns.
 - crim per capita crime rate by town.
 - zn proportion of residential land zoned for lots over 25,000 sq.ft.
 - indus proportion of non-retail business acres per town.

chas Charles River dummy variable (= 1 if tract bounds river; 0 otherwise).
 nox nitrogen oxides concentration (parts per 10 million).
 rm average number of rooms per dwelling.
 age proportion of owner-occupied units built prior to 1940.
 dis weighted mean of distances to five Boston employment centres.
 rad index of accessibility to radial highways.
 tax full-value property-tax rate per \$10,000.
 ptratio pupil-teacher ratio by town.
 black $1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town.
 lstat lower status of the population (percent).
 medv median value of owner-occupied homes in \$1000s.

2. Consider this data in context, what is the response variable of interest? Discuss how you think some of the possible predictor variables might be associated with this response. Here, I am considering “crim” as the response variable and all other variables as predictors. Crime rate will be interesting to analyze with respect to the income, zone, taxes and type of population living in the area.
3. For each predictor, fit a simple linear regression model to predict the response. In which of the models is there a statistically significant association between the predictor and the response? Create some plots to back up your assertions.

```
summary(lm(crim ~ zn, data = boston))

##
## Call:
## lm(formula = crim ~ zn, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.429  -4.222  -2.620   1.250  84.523
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.45369    0.41722  10.675 < 2e-16 ***
## zn          -0.07393    0.01609  -4.594 5.51e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.435 on 504 degrees of freedom
## Multiple R-squared:  0.04019,    Adjusted R-squared:  0.03828
## F-statistic: 21.1 on 1 and 504 DF,  p-value: 5.506e-06
```

#fitting simple linear regression model

```
summary(lm(crim ~ indus, data = boston))

##
## Call:
## lm(formula = crim ~ indus, data = boston)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11.972  -2.698  -0.736   0.712  81.813
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.06374    0.66723  -3.093  0.00209 **
## indus        0.50978    0.05102   9.991 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.866 on 504 degrees of freedom
## Multiple R-squared:  0.1653, Adjusted R-squared:  0.1637
## F-statistic: 99.82 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ chas, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ chas, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.738 -3.661 -3.435   0.018  85.232
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.7444    0.3961   9.453 <2e-16 ***
## chas         -1.8928    1.5061  -1.257   0.209
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.597 on 504 degrees of freedom
## Multiple R-squared:  0.003124, Adjusted R-squared:  0.001146
## F-statistic: 1.579 on 1 and 504 DF,  p-value: 0.2094
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ nox, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ nox, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.371  -2.738  -0.974   0.559  81.728
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept) -13.720      1.699 -8.073 5.08e-15 ***
## nox          31.249      2.999 10.419 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.81 on 504 degrees of freedom
## Multiple R-squared:  0.1772, Adjusted R-squared:  0.1756
## F-statistic: 108.6 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ rm, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ rm, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.604 -3.952 -2.654  0.989 87.197
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   20.482      3.365   6.088 2.27e-09 ***
## rm            -2.684      0.532  -5.045 6.35e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.401 on 504 degrees of freedom
## Multiple R-squared:  0.04807,    Adjusted R-squared:  0.04618
## F-statistic: 25.45 on 1 and 504 DF,  p-value: 6.347e-07
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ age, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ age, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.789 -4.257 -1.230  1.527 82.849
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.77791    0.94398  -4.002 7.22e-05 ***
## age          0.10779    0.01274   8.463 2.85e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.057 on 504 degrees of freedom
## Multiple R-squared:  0.1244, Adjusted R-squared:  0.1227
## F-statistic: 71.62 on 1 and 504 DF,  p-value: 2.855e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ dis, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ dis, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.708 -4.134 -1.527  1.516  81.674
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.4993     0.7304  13.006  <2e-16 ***
## dis          -1.5509     0.1683   -9.213  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.965 on 504 degrees of freedom
## Multiple R-squared:  0.1441, Adjusted R-squared:  0.1425
## F-statistic: 84.89 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ rad, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ rad, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.164  -1.381  -0.141    0.660   76.433
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.28716     0.44348  -5.157 3.61e-07 ***
## rad          0.61791     0.03433  17.998 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.718 on 504 degrees of freedom
## Multiple R-squared:  0.3913, Adjusted R-squared:  0.39
## F-statistic: 323.9 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ tax, data = boston))
```

```
##
```

```
## Call:
## lm(formula = crim ~ tax, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.513  -2.738  -0.194   1.065  77.696
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -8.528369   0.815809  -10.45  <2e-16 ***
## tax          0.029742   0.001847   16.10  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.997 on 504 degrees of freedom
## Multiple R-squared:  0.3396, Adjusted R-squared:  0.3383
## F-statistic: 259.2 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ ptratio, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ ptratio, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
##  -7.654  -3.985  -1.912   1.825  83.353
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -17.6469     3.1473  -5.607 3.40e-08 ***
## ptratio       1.1520     0.1694   6.801 2.94e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.24 on 504 degrees of freedom
## Multiple R-squared:  0.08407, Adjusted R-squared:  0.08225
## F-statistic: 46.26 on 1 and 504 DF,  p-value: 2.943e-11
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ black, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ black, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.756  -2.299  -2.095  -1.296   86.822
```

```
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 16.553529   1.425903  11.609   <2e-16 ***
## black       -0.036280   0.003873  -9.367   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.946 on 504 degrees of freedom
## Multiple R-squared:  0.1483, Adjusted R-squared:  0.1466
## F-statistic: 87.74 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ lstat, data = boston))
```

```
##
## Call:
## lm(formula = crim ~ lstat, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.925  -2.822  -0.664   1.079   82.862
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.33054    0.69376  -4.801 2.09e-06 ***
## lstat        0.54880    0.04776  11.491 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.664 on 504 degrees of freedom
## Multiple R-squared:  0.2076, Adjusted R-squared:  0.206
## F-statistic: 132 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
#fitting simple linear regression model
```

```
summary(lm(crim ~ medv, data = boston))
```

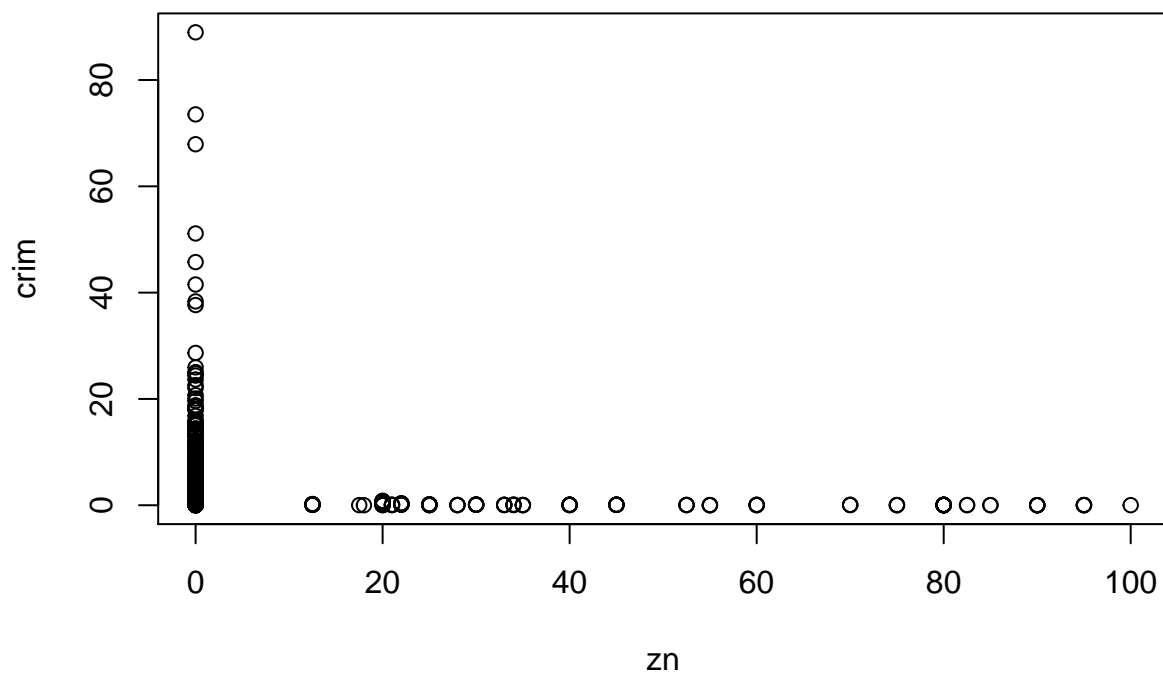
```
##
## Call:
## lm(formula = crim ~ medv, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
##  -9.071  -4.022  -2.343   1.298  80.957
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 11.79654    0.93419   12.63  <2e-16 ***
## medv        -0.36316    0.03839   -9.46  <2e-16 ***
## ---
```

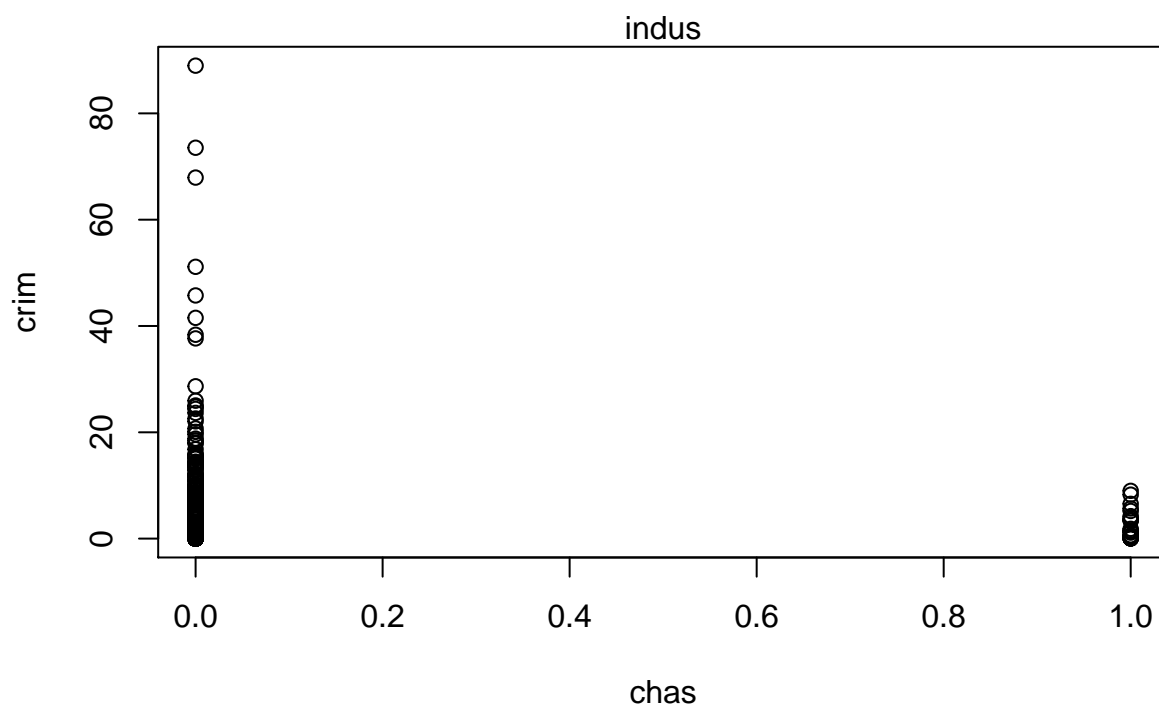
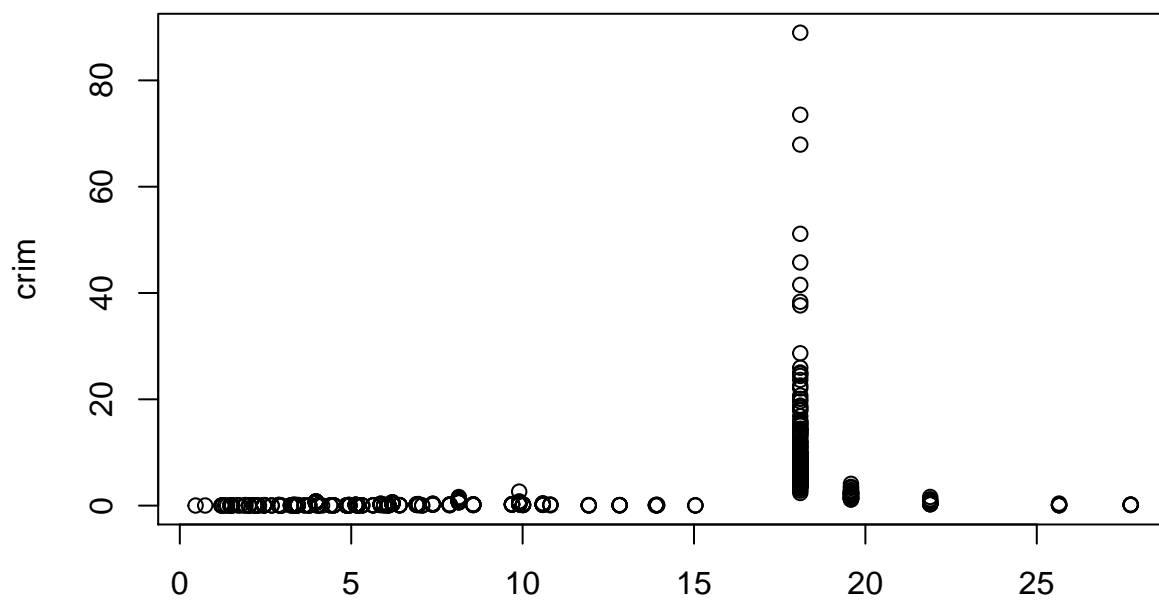
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.934 on 504 degrees of freedom
## Multiple R-squared:  0.1508, Adjusted R-squared:  0.1491
## F-statistic: 89.49 on 1 and 504 DF,  p-value: < 2.2e-16
```

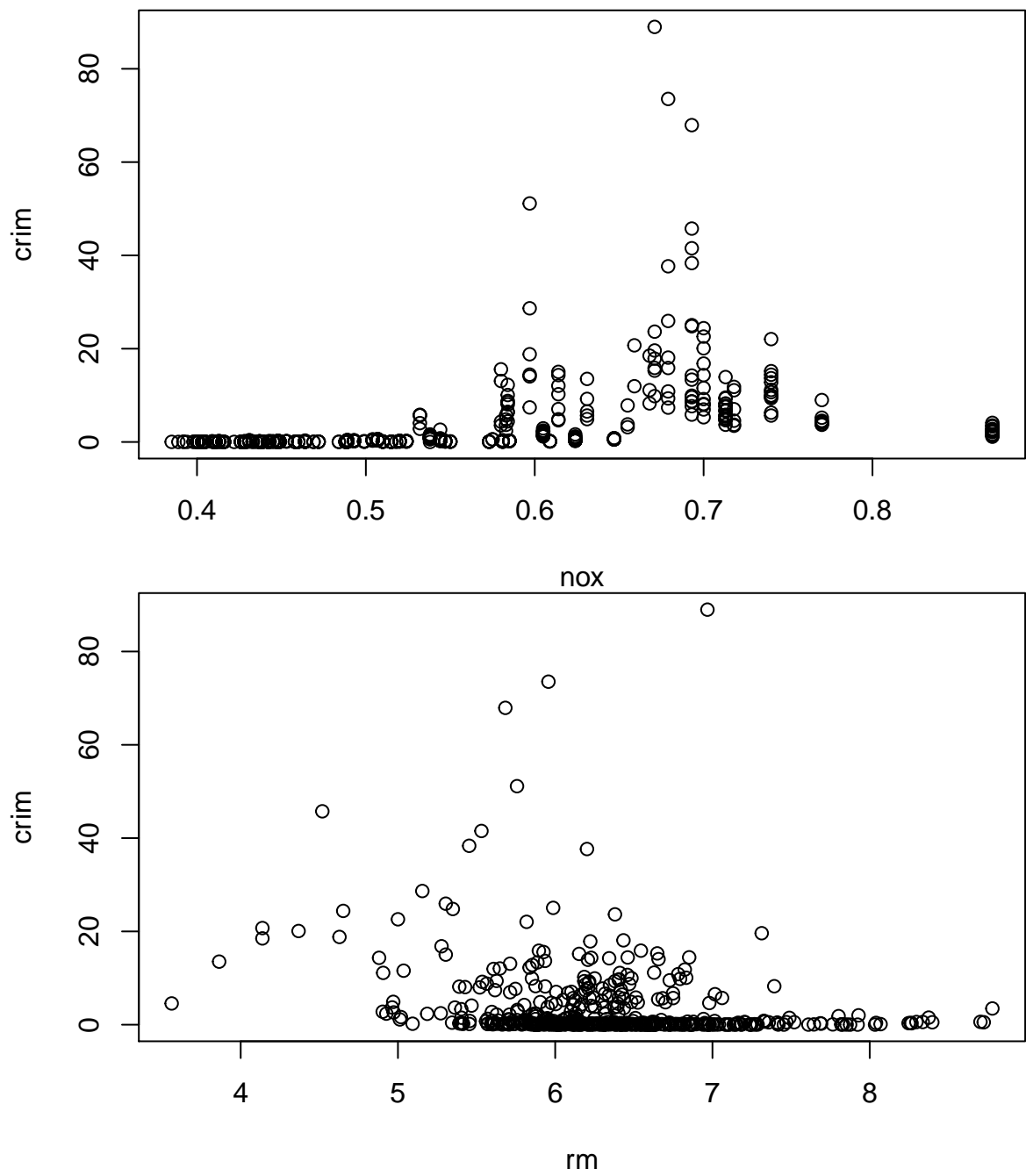
```
#fitting simple linear regression model
```

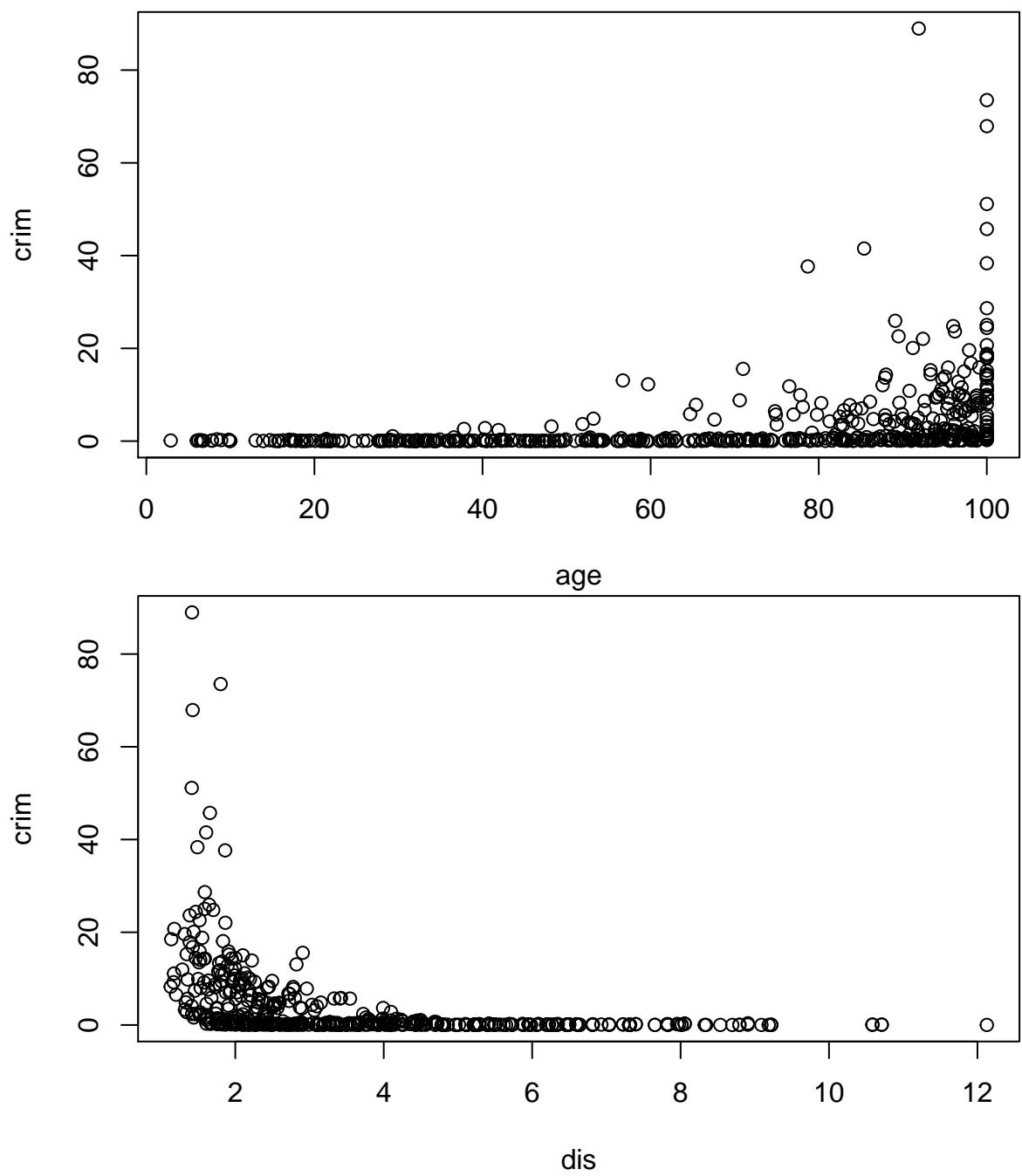
As per the above results, there is statistically significant relationship between every predictor for crim as the response variable except Charles River Dummy. With almost every variable statistically significant and R-squared being low, every predictor describes small amount of variance. Refer to the plots below:

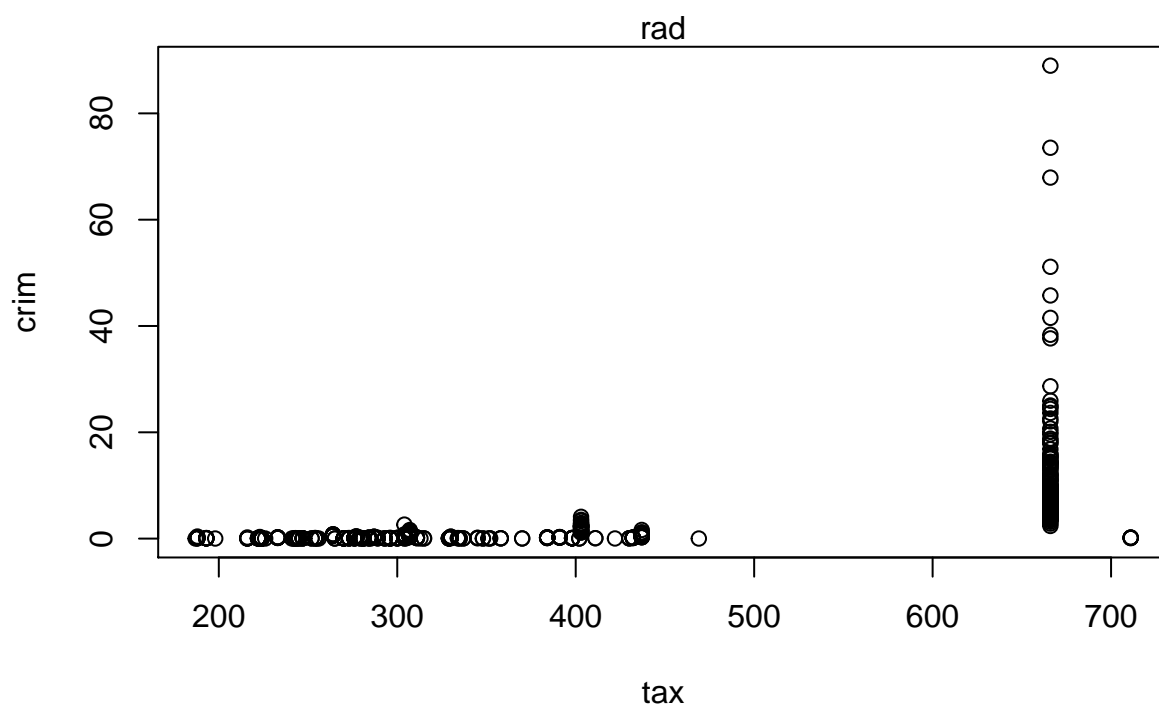
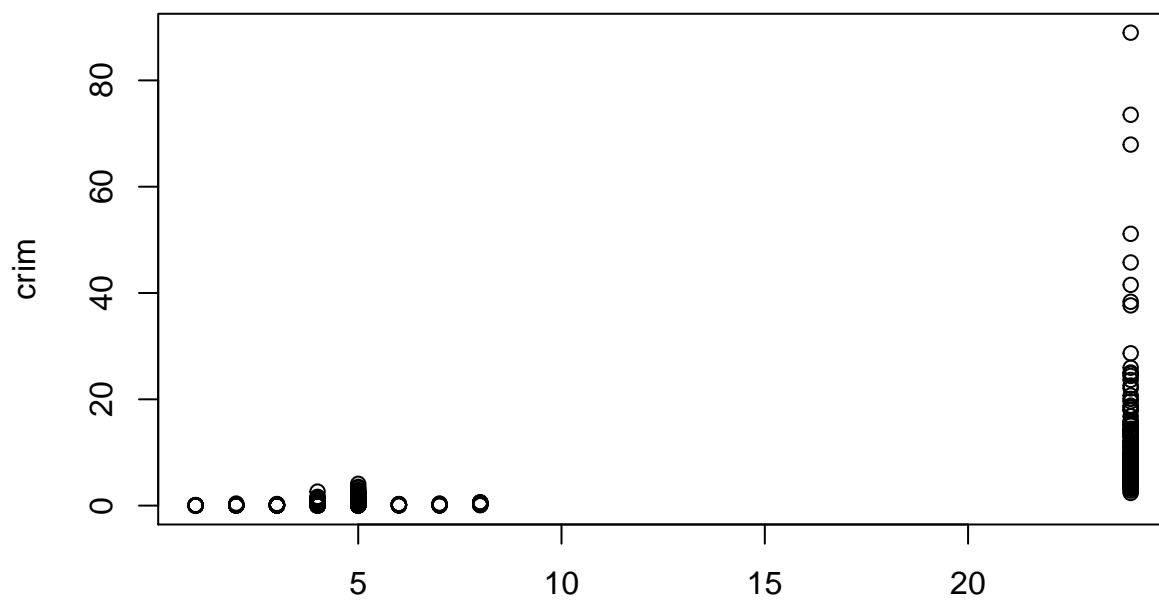
```
plot(crim ~ . - crim, data = boston)
```

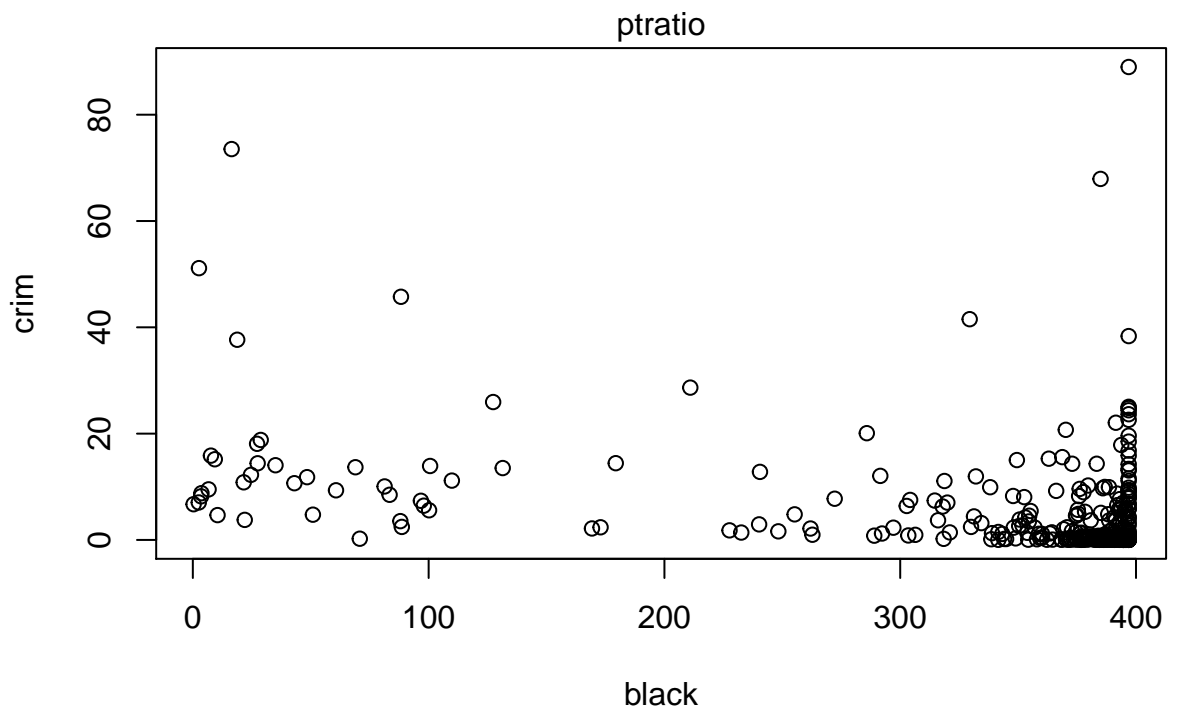
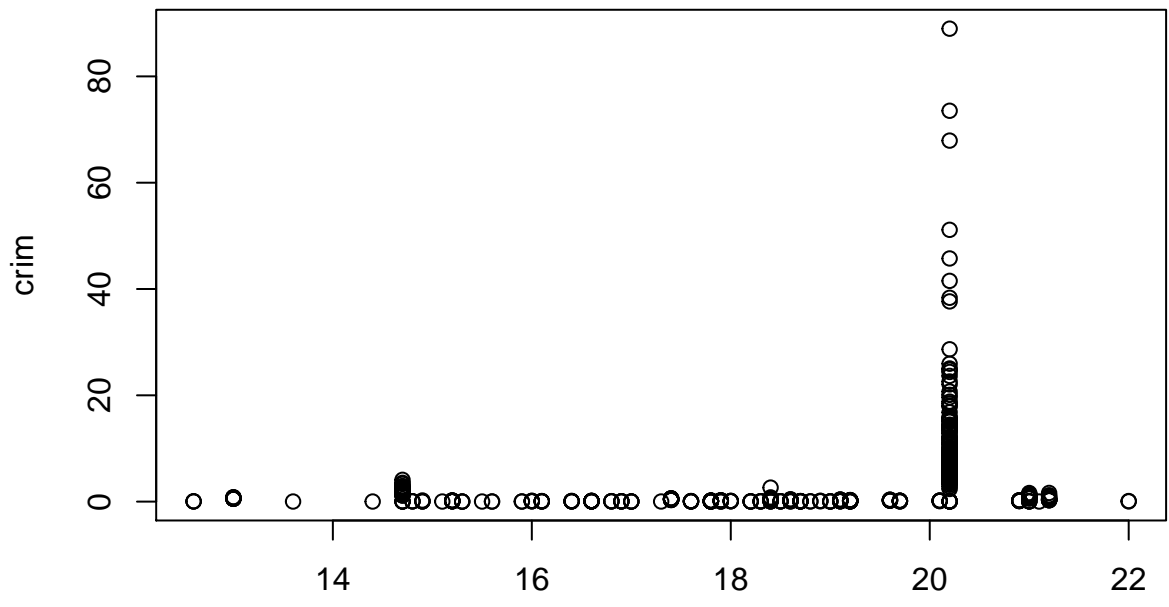


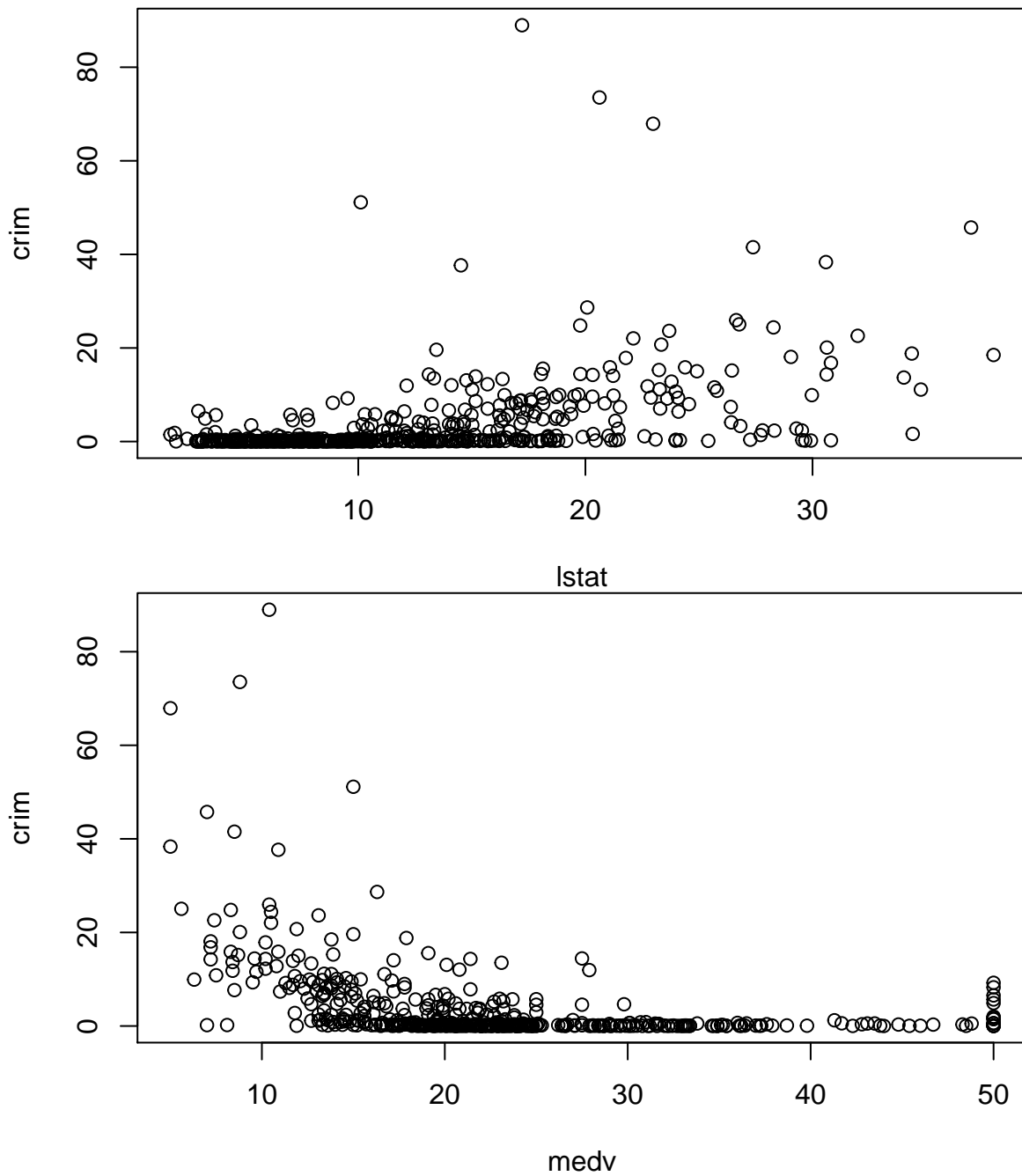












```
#plotting crime versus every other variable
```

4. Fit a multiple regression model to predict the response using all of the predictors. Describe your results. For which predictors can we reject the null hypothesis $H_0 : \beta_j = 0$?

```
summary(lm(crim ~ . - crim, data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ . - crim, data = Boston)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.924 -2.120 -0.353  1.019 75.051
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  17.033228   7.234903   2.354 0.018949 *
## zn           0.044855   0.018734   2.394 0.017025 *
## indus        -0.063855   0.083407  -0.766 0.444294
## chas         -0.749134   1.180147  -0.635 0.525867
## nox          -10.313535   5.275536  -1.955 0.051152 .
## rm           0.430131   0.612830   0.702 0.483089
## age          0.001452   0.017925   0.081 0.935488
## dis         -0.987176   0.281817  -3.503 0.000502 ***
## rad          0.588209   0.088049   6.680 6.46e-11 ***
## tax         -0.003780   0.005156  -0.733 0.463793
## ptratio     -0.271081   0.186450  -1.454 0.146611
## black       -0.007538   0.003673  -2.052 0.040702 *
## lstat        0.126211   0.075725   1.667 0.096208 .
## medv        -0.198887   0.060516  -3.287 0.001087 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.439 on 492 degrees of freedom
## Multiple R-squared:  0.454, Adjusted R-squared:  0.4396
## F-statistic: 31.47 on 13 and 492 DF,  p-value: < 2.2e-16
```

```
#fitting multiple regression model
```

Result of fitting Mutiple Regression Model shows that very few variables are statistically significant.

- (a) dis: 0.001 level
- (b) rad: 0.001 level
- (c) medv: 0.01 level
- (d) zn: 0.05 level
- (e) black: 0.05 level

We cannot reject null hypothesis for every other variable & using multiple regression model, R-squared is much higher compared to linear model, indicating that we should explain more of a variance in the outcome.

5. How do your results from (3) compare to your results from (4)? Create a plot displaying the univariate regression coefficients from (3) on the x-axis and the multiple regression coefficients from part (4) on the y-axis. Use this visualization to support your response.

```
#Getting the estimates of the coefficients for all the simple linear regression models
univariateCoefficient <- lm(crim ~ zn, data = boston)$coefficients[2]
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ indus, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ chas, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ nox, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
```

```

lm(crim ~ rm, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ age, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ dis, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ rad, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ tax, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ ptratio, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ black, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ lstat, data = boston)$coefficients[2])
univariateCoefficient <- append(univariateCoefficient,
                                lm(crim ~ medv, data = boston)$coefficients[2])
bostonData <- (lm(crim ~ . - crim, data = boston))
bostonData$coefficients[2:14]

```

```

##          zn          indus          chas          nox          rm
## 0.044855215 -0.063854824 -0.749133611 -10.313534912 0.430130506
##          age          dis          rad          tax          ptratio
## 0.001451643 -0.987175726 0.588208591 -0.003780016 -0.271080558
##          black          lstat          medv
## -0.007537505 0.126211376 -0.198886821

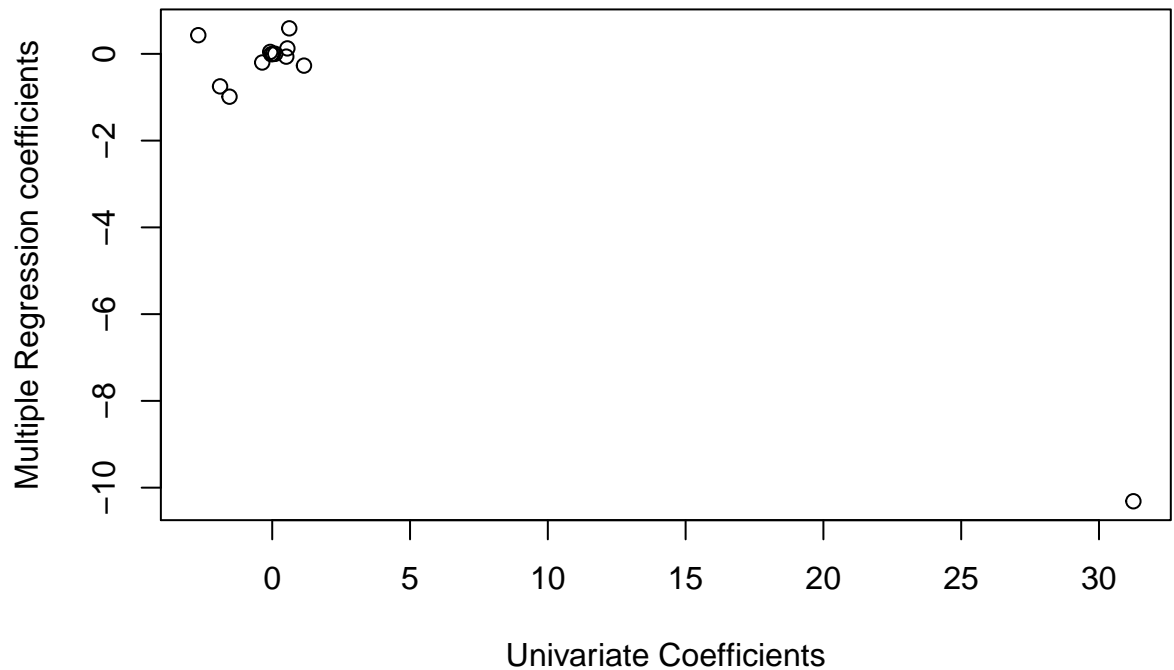
```

```

plot(univariateCoefficient, bostonData$coefficients[2:14], main =
     "Plot of Univariate vs. Multiple Regression Coefficients", xlab =
     "Univariate Coefficients",
     ylab = "Multiple Regression coefficients")

```


Plot of Univariate vs. Multiple Regression Coefficients



6. Is there evidence of a non-linear association between any of the predictors and the response? To answer this question, for each predictor X fit a model of the form:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon$$

```
#fitting the model
summary(lm(crim ~ zn + I(zn^2) + I(zn^3), data = Boston))

##
## Call:
## lm(formula = crim ~ zn + I(zn^2) + I(zn^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.821 -4.614 -1.294   0.473  84.130
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.846e+00  4.330e-01  11.192  < 2e-16 ***
## zn          -3.322e-01  1.098e-01  -3.025  0.00261 **
## I(zn^2)       6.483e-03  3.861e-03   1.679  0.09375 .
## I(zn^3)      -3.776e-05  3.139e-05  -1.203  0.22954
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.372 on 502 degrees of freedom
## Multiple R-squared:  0.05824,    Adjusted R-squared:  0.05261
## F-statistic: 10.35 on 3 and 502 DF,  p-value: 1.281e-06
```

```
#fitting the model
summary(lm(crim ~ indus + I(indus^2) + I(indus^3), data = Boston))

##
## Call:
## lm(formula = crim ~ indus + I(indus^2) + I(indus^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.278 -2.514  0.054  0.764 79.713
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.6625683   1.5739833    2.327  0.0204 *
## indus        -1.9652129   0.4819901   -4.077 5.30e-05 ***
## I(indus^2)    0.2519373   0.0393221    6.407 3.42e-10 ***
## I(indus^3)   -0.0069760   0.0009567   -7.292 1.20e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.423 on 502 degrees of freedom
## Multiple R-squared:  0.2597, Adjusted R-squared:  0.2552
## F-statistic: 58.69 on 3 and 502 DF, p-value: < 2.2e-16
```

```
#fitting the model
summary(lm(crim ~ chas + I(chas^2) + I(chas^3), data = Boston))

##
## Call:
## lm(formula = crim ~ chas + I(chas^2) + I(chas^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.738 -3.661 -3.435  0.018 85.232
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.7444    0.3961    9.453 <2e-16 ***
## chas        -1.8928    1.5061   -1.257  0.209
## I(chas^2)         NA           NA      NA      NA
## I(chas^3)         NA           NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.597 on 504 degrees of freedom
## Multiple R-squared:  0.003124, Adjusted R-squared:  0.001146
## F-statistic: 1.579 on 1 and 504 DF, p-value: 0.2094
```

```
#fitting the model
summary(lm(crim ~ nox + I(nox^2) + I(nox^3), data = Boston))
```

```
##
```

```
## Call:
## lm(formula = crim ~ nox + I(nox^2) + I(nox^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.110 -2.068 -0.255  0.739 78.302
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    233.09      33.64   6.928 1.31e-11 ***
## nox           -1279.37     170.40  -7.508 2.76e-13 ***
## I(nox^2)        2248.54     279.90   8.033 6.81e-15 ***
## I(nox^3)       -1245.70     149.28  -8.345 6.96e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.234 on 502 degrees of freedom
## Multiple R-squared:  0.297, Adjusted R-squared:  0.2928
## F-statistic: 70.69 on 3 and 502 DF, p-value: < 2.2e-16
```

#fitting the model

```
summary(lm(crim ~ rm + I(rm^2) + I(rm^3), data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ rm + I(rm^2) + I(rm^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.485  -3.468  -2.221  -0.015  87.219
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  112.6246    64.5172   1.746  0.0815 .
## rm           -39.1501    31.3115  -1.250  0.2118
## I(rm^2)        4.5509     5.0099   0.908  0.3641
## I(rm^3)       -0.1745     0.2637  -0.662  0.5086
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.33 on 502 degrees of freedom
## Multiple R-squared:  0.06779, Adjusted R-squared:  0.06222
## F-statistic: 12.17 on 3 and 502 DF, p-value: 1.067e-07
```

```
summary(lm(crim ~ age + I(age^2) + I(age^3), data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ age + I(age^2) + I(age^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.762 -2.673 -0.516  0.019 82.842
```

```
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.549e+00  2.769e+00  -0.920  0.35780
## age          2.737e-01  1.864e-01   1.468  0.14266
## I(age^2)     -7.230e-03  3.637e-03  -1.988  0.04738 *
## I(age^3)      5.745e-05  2.109e-05   2.724  0.00668 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.84 on 502 degrees of freedom
## Multiple R-squared:  0.1742, Adjusted R-squared:  0.1693
## F-statistic: 35.31 on 3 and 502 DF,  p-value: < 2.2e-16
```

#fitting the model

```
summary(lm(crim ~ dis + I(dis^2) + I(dis^3), data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ dis + I(dis^2) + I(dis^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.757  -2.588   0.031   1.267  76.378
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  30.0476     2.4459  12.285 < 2e-16 ***
## dis         -15.5543     1.7360  -8.960 < 2e-16 ***
## I(dis^2)       2.4521     0.3464   7.078 4.94e-12 ***
## I(dis^3)      -0.1186     0.0204  -5.814 1.09e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.331 on 502 degrees of freedom
## Multiple R-squared:  0.2778, Adjusted R-squared:  0.2735
## F-statistic: 64.37 on 3 and 502 DF,  p-value: < 2.2e-16
```

#fitting the model

```
summary(lm(crim ~ rad + I(rad^2) + I(rad^3), data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ rad + I(rad^2) + I(rad^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.381  -0.412  -0.269   0.179  76.217
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.605545  2.050108  -0.295  0.768
## rad          0.512736  1.043597   0.491  0.623
```

```
## I(rad^2)    -0.075177    0.148543   -0.506    0.613
## I(rad^3)     0.003209    0.004564    0.703    0.482
##
## Residual standard error: 6.682 on 502 degrees of freedom
## Multiple R-squared:  0.4, Adjusted R-squared:  0.3965
## F-statistic: 111.6 on 3 and 502 DF,  p-value: < 2.2e-16
```

#fitting the model

```
summary(lm(crim ~ tax + I(tax^2) + I(tax^3), data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ tax + I(tax^2) + I(tax^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.273  -1.389   0.046   0.536  76.950
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.918e+01  1.180e+01   1.626   0.105
## tax         -1.533e-01  9.568e-02  -1.602   0.110
## I(tax^2)     3.608e-04  2.425e-04   1.488   0.137
## I(tax^3)    -2.204e-07  1.889e-07  -1.167   0.244
##
## Residual standard error: 6.854 on 502 degrees of freedom
## Multiple R-squared:  0.3689, Adjusted R-squared:  0.3651
## F-statistic:  97.8 on 3 and 502 DF,  p-value: < 2.2e-16
```

#fitting the model

```
summary(lm(crim ~ ptratio + I(ptratio^2) + I(ptratio^3), data = Boston))
```

```
##
## Call:
## lm(formula = crim ~ ptratio + I(ptratio^2) + I(ptratio^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
##  -6.833  -4.146  -1.655   1.408  82.697
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  477.18405  156.79498   3.043  0.00246 **
## ptratio     -82.36054   27.64394  -2.979  0.00303 **
## I(ptratio^2)   4.63535    1.60832   2.882  0.00412 **
## I(ptratio^3)  -0.08476    0.03090  -2.743  0.00630 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.122 on 502 degrees of freedom
## Multiple R-squared:  0.1138, Adjusted R-squared:  0.1085
## F-statistic: 21.48 on 3 and 502 DF,  p-value: 4.171e-13
```

```
#fitting the model
summary(lm(crim ~ black + I(black^2) + I(black^3), data = Boston))

##
## Call:
## lm(formula = crim ~ black + I(black^2) + I(black^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.096  -2.343  -2.128  -1.439   86.790
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.826e+01  2.305e+00   7.924  1.5e-14 ***
## black        -8.356e-02  5.633e-02  -1.483   0.139
## I(black^2)    2.137e-04  2.984e-04   0.716   0.474
## I(black^3)   -2.652e-07  4.364e-07  -0.608   0.544
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.955 on 502 degrees of freedom
## Multiple R-squared:  0.1498, Adjusted R-squared:  0.1448
## F-statistic: 29.49 on 3 and 502 DF, p-value: < 2.2e-16
```

```
#fitting the model
summary(lm(crim ~ lstat + I(lstat^2) + I(lstat^3), data = Boston))

##
## Call:
## lm(formula = crim ~ lstat + I(lstat^2) + I(lstat^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.234  -2.151  -0.486   0.066   83.353
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.2009656  2.0286452   0.592  0.5541
## lstat        -0.4490656  0.4648911  -0.966  0.3345
## I(lstat^2)    0.0557794  0.0301156   1.852  0.0646 .
## I(lstat^3)   -0.0008574  0.0005652  -1.517  0.1299
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.629 on 502 degrees of freedom
## Multiple R-squared:  0.2179, Adjusted R-squared:  0.2133
## F-statistic: 46.63 on 3 and 502 DF, p-value: < 2.2e-16
```

```
#fitting the model
summary(lm(crim ~ medv + I(medv^2) + I(medv^3), data = Boston))
```

```
##
```

```
## Call:
## lm(formula = crim ~ medv + I(medv^2) + I(medv^3), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -24.427  -1.976  -0.437   0.439  73.655
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  53.1655381   3.3563105   15.840 < 2e-16 ***
## medv        -5.0948305   0.4338321  -11.744 < 2e-16 ***
## I(medv^2)     0.1554965   0.0171904    9.046 < 2e-16 ***
## I(medv^3)    -0.0014901   0.0002038   -7.312 1.05e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.569 on 502 degrees of freedom
## Multiple R-squared:  0.4202, Adjusted R-squared:  0.4167
## F-statistic: 121.3 on 3 and 502 DF,  p-value: < 2.2e-16
```

Based on the above results, following are the key observations regarding the squared and cubed terms for the variables below:

- (a) indus- squared term: 3.42e-10 cubed term: 1.20e-12
- (b) nox- squared term: 6.81e-1 cubed term: 6.96e-16
- (c) dis- squared term: 4.94e-12
cubed term: 1.09e-08
- (d) ptracio- squared term: 0.00412
cubed term: 0.00630
- (e) medv- squared term: 2e-16 cubed term: 1.05e-12

The first thing to note is that with the chas variable, we get NA values for the squared and cubed term. This makes sense as chas is a dummy variable, composed of only 0s and 1s, and these values will not change if they are squared or cubed.

With the variables indus, nox, dis, ptracio, and medv, there is evidence of a non-linear relationship, as each of these variables squared and cubed terms is found to be statistically significant (we reject the null hypothesis that the coefficients on these exponentiated variables are zero).

Age also have a non-linear relationship and with squared & cubed-age, linear age becomes insignificant statistically.

For every other variable, we do not find evidence of a non-linear relationship between the predictor and outcome variables.

7. Consider performing a stepwise model selection procedure to determine the best fit model. Discuss your results. How is this model different from the model in (4)?

```
#Forward Stepwise selection
submod<-regsubsets(crim~.,data=boston,method="forward",nvmax = 13)

#Summary of forward stepwise selection
submodsum<-summary(submod)

#Summary of the model
submodsum
```

```
## Subset selection object
## Call: regsubsets.formula(crim ~ ., data = boston, method = "forward",
##      nvmax = 13)
## 13 Variables (and intercept)
##      Forced in Forced out
## zn          FALSE      FALSE
## indus        FALSE      FALSE
## chas          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
## 1 subsets of each size up to 13
## Selection Algorithm: forward
##      zn indus chas 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 ) "*" "*" " " " " " " " " " " " " " " " " " " " " " " " " "
## 12 ( 1 ) "*" "*" "*" " " " " " " " " " " " " " " " " " " " " " " "
## 13 ( 1 ) "*" "*" "*" "*" " " " " " " " " " " " " " " " " " " " " " "
```

Forward stepwise selection is used and summary method in order to get point at which Bayesian Information Criterion (BIC) of the model is the least. We are eliminating chas and age while fitting the next model since these were the last two variables to be fitted.

```
#Fitting the model after the selection
modelForward<-lm(crim~.-chas-age,data=boston)

#Summary of the model
summary(modelForward)
```

```
##
## Call:
## lm(formula = crim ~ . - chas - age, data = boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.906 -2.133 -0.315  1.065 75.055
##
## Coefficients:
```

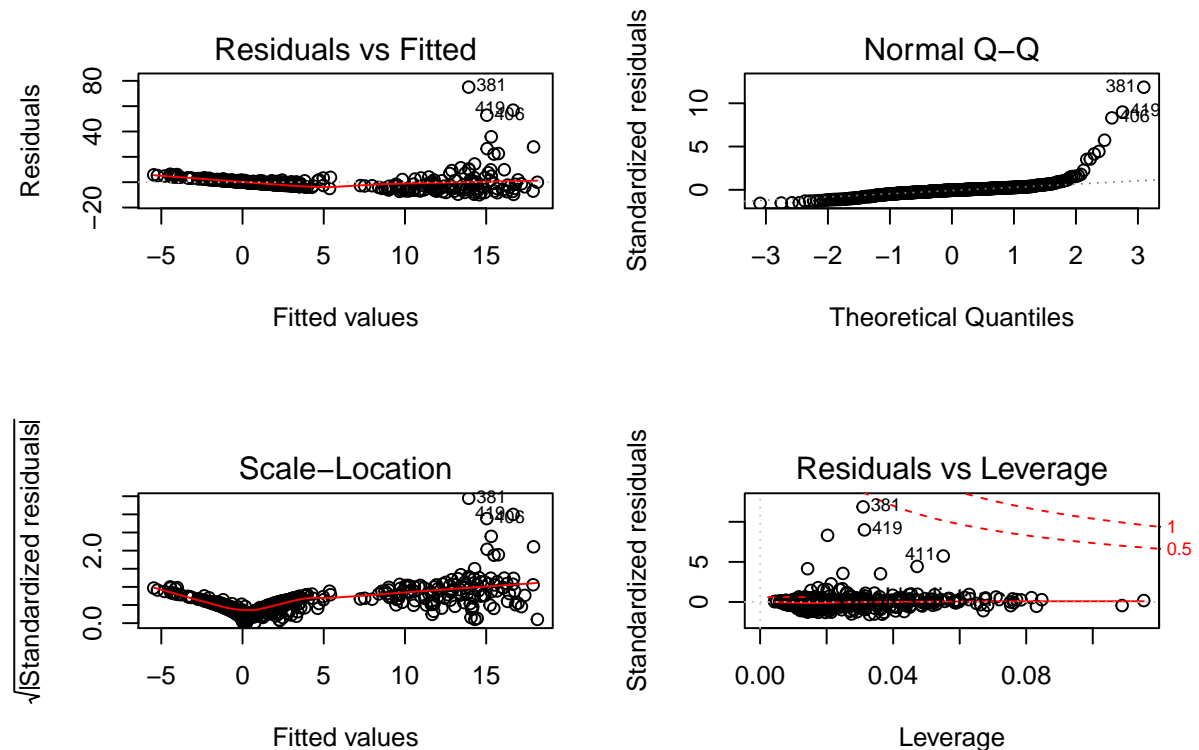


```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 17.096653   7.197290   2.375 0.017909 *
## zn          0.044859   0.018566   2.416 0.016049 *
## indus       -0.069177   0.082845  -0.835 0.404117
## nox        -10.458590   5.068742  -2.063 0.039601 *
## rm          0.445708   0.600550   0.742 0.458339
## dis        -0.997154   0.270082  -3.692 0.000247 ***
## rad         0.583934   0.087449   6.677 6.56e-11 ***
## tax        -0.003455   0.005121  -0.675 0.500266
## ptratio    -0.265328   0.185482  -1.430 0.153212
## black      -0.007599   0.003658  -2.077 0.038277 *
## lstat       0.127215   0.071832   1.771 0.077177 .
## medv       -0.204431   0.059777  -3.420 0.000678 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.429 on 494 degrees of freedom
## Multiple R-squared:  0.4536, Adjusted R-squared:  0.4414
## F-statistic: 37.28 on 11 and 494 DF,  p-value: < 2.2e-16
```

There are two key observations in values of the Residual Standard Error and R-squared values: Residual Standard Error for the model with 11 variables (6.429) is lower than when we fit the model with all the variables (6.439). We also notice the adjusted R-squared value is 0.4414 better for the model in which stepwise selection has occurred than for the model with all the variables where it was 0.4396.

8. Evaluate the statistical assumptions in your regression analysis from (7) by performing a basic analysis of model residuals and any unusual observations. Discuss any concerns you have about your model.

```
par(mfrow=c(2,2))
plot(modelForward)
```



Outliers: Concerning because they are far from the proposed values Residual vs Fitted values: Non-linearity is observed and this is concerning as linear model assumes linearity in the relationship. High leverage points: They have strong effect on estimated regression line.