# INFX 573: Problem Set 6 - Regression

*Shuyang Wu*

*Due: Tuesday, November 15, 2016*

**Collaborators:**

**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 chucks, 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
#look at dataset
str(Boston)
```

```
## 'data.frame':    506 obs. of  14 variables:
##  $ crim   : num  0.00632 0.02731 0.02729 0.03237 0.06905 ...
##  $ zn     : num  18 0 0 0 0 0 12.5 12.5 12.5 12.5 ...
##  $ indus  : num  2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 7.87 ...
##  $ chas   : int  0 0 0 0 0 0 0 0 0 0 ...
##  $ nox    : num  0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 0.524 0.524 ...
##  $ rm     : num  6.58 6.42 7.18 7 7.15 ...
##  $ age    : num  65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9 ...
##  $ dis    : num  4.09 4.97 4.97 6.06 6.06 ...
##  $ rad    : int  1 2 2 3 3 3 5 5 5 5 ...
##  $ tax    : num  296 242 242 222 222 222 311 311 311 311 ...
##  $ ptratio: num  15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
##  $ black  : num  397 397 393 395 397 ...
##  $ lstat  : num  4.98 9.14 4.03 2.94 5.33 ...
##  $ medv   : num  24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...
```

```r
names(Boston)
```

```
## [1] "crim"    "zn"      "indus"   "chas"    "nox"     "rm"      "age"
## [8] "dis"     "rad"     "tax"     "ptratio" "black"   "lstat"   "medv"
```

```r
summary(Boston)
```

```
##       crim                zn              indus            chas
##  Min.   : 0.00632   Min.   :  0.00   Min.   : 0.46   Min.   :0.00000
##  1st Qu.: 0.08204   1st Qu.:  0.00   1st Qu.: 5.19   1st Qu.:0.00000
##  Median : 0.25651   Median :  0.00   Median : 9.69   Median :0.00000
##  Mean   : 3.61352   Mean   : 11.36   Mean   :11.14   Mean   :0.06917
##  3rd Qu.: 3.67708   3rd Qu.: 12.50   3rd Qu.:18.10   3rd Qu.:0.00000
##  Max.   :88.97620   Max.   :100.00   Max.   :27.74   Max.   :1.00000
##       nox               rm             age              dis
##  Min.   :0.3850   Min.   :3.561   Min.   :  2.90   Min.   : 1.130
##  1st Qu.:0.4490   1st Qu.:5.886   1st Qu.: 45.02   1st Qu.: 2.100
##  Median :0.5380   Median :6.208   Median : 77.50   Median : 3.207
##  Mean   :0.5547   Mean   :6.285   Mean   : 68.57   Mean   : 3.795
##  3rd Qu.:0.6240   3rd Qu.:6.623   3rd Qu.: 94.08   3rd Qu.: 5.188
##  Max.   :0.8710   Max.   :8.780   Max.   :100.00   Max.   :12.127
##       rad              tax            ptratio          black
##  Min.   : 1.000   Min.   :187.0   Min.   :12.60   Min.   :  0.32
##  1st Qu.: 4.000   1st Qu.:279.0   1st Qu.:17.40   1st Qu.:375.38
##  Median : 5.000   Median :330.0   Median :19.05   Median :391.44
##  Mean   : 9.549   Mean   :408.2   Mean   :18.46   Mean   :356.67
##  3rd Qu.:24.000   3rd Qu.:666.0   3rd Qu.:20.20   3rd Qu.:396.23
##  Max.   :24.000   Max.   :711.0   Max.   :22.00   Max.   :396.90
##      lstat            medv
##  Min.   : 1.73   Min.   : 5.00
##  1st Qu.: 6.95   1st Qu.:17.02
##  Median :11.36   Median :21.20
##  Mean   :12.65   Mean   :22.53
##  3rd Qu.:16.95   3rd Qu.:25.00
##  Max.   :37.97   Max.   :50.00
```

```r
#Tidy the dataset
Boston$chas <- as.factor(Boston$chas)
Boston$rad <- as.integer(Boston$rad)
```

**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 dataset contains per capita crime rate by town, proportion of residential land zoned for lots over 25,000 sq.ft., proportion of non-retail business acres per town, Charles River dummy variable (= 1 if tract bounds river; 0 otherwise), nitrogen oxides concentration (parts per 10 million), average number of rooms per dwelling, proportion of owner-occupied units built prior to 1940, weighted mean of distances to five Boston employment centres, index of accessibility to radial highways, full-value property-tax rate per ten thousand, pupil-teacher ratio by town, proportion of blacks by town, lower status of the population (percent), 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.

Response variable is the median value of owner-occupied homes. Intuitively, I think number of rooms, distances to five employment centres, index of accessibility to radial highways could be associated with this response.

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.

```
#simple linear regressions
fit.crim <- lm(medv ~ crim, data = Boston)
summary(fit.crim)
```

```
##
## Call:
## lm(formula = medv ~ crim, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -16.957  -5.449  -2.007   2.512  29.800
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 24.03311    0.40914   58.74   <2e-16 ***
## crim        -0.41519    0.04389   -9.46   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.484 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
```

```
fit.zn <- lm(medv ~ zn, data = Boston)
summary(fit.zn) #significant
```

```
##
## Call:
## lm(formula = medv ~ zn, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -15.918  -5.518  -1.006   2.757  29.082
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 20.91758    0.42474  49.248   <2e-16 ***
## zn           0.14214    0.01638   8.675   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.587 on 504 degrees of freedom
## Multiple R-squared:  0.1299, Adjusted R-squared:  0.1282
## F-statistic: 75.26 on 1 and 504 DF,  p-value: < 2.2e-16
```

```
fit.indus <- lm(medv ~ indus, data = Boston)
summary(fit.indus) #significant
```

3

```
## 
## Call:
## lm(formula = medv ~ indus, data = Boston)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max
## -13.017  -4.917  -1.457   3.180  32.943
## 
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 29.75490    0.68345   43.54   <2e-16 ***
## indus       -0.64849    0.05226  -12.41   <2e-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:  0.234,  Adjusted R-squared:  0.2325
## F-statistic:   154 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
fit.nox <- lm(medv ~ nox, data = Boston)
summary(fit.nox) #significant
```

```
## 
## Call:
## lm(formula = medv ~ nox, data = Boston)
## 
## Residuals:
##      Min      1Q  Median      3Q     Max
## -13.691  -5.121  -2.161   2.959  31.310
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   41.346      1.811   22.83   <2e-16 ***
## nox          -33.916      3.196  -10.61   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 8.323 on 504 degrees of freedom
## Multiple R-squared:  0.1826, Adjusted R-squared:  0.181
## F-statistic: 112.6 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
fit.age <- lm(medv ~ age, data = Boston)
summary(fit.age) #significant
```

```
## 
## Call:
## lm(formula = medv ~ age, data = Boston)
## 
## Residuals:
##      Min      1Q  Median      3Q     Max
## -15.097  -5.138  -1.958   2.397  31.338
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 30.97868    0.99911  31.006   <2e-16 ***
```

```
## age          -0.12316    0.01348  -9.137   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.527 on 504 degrees of freedom
## Multiple R-squared:  0.1421, Adjusted R-squared:  0.1404
## F-statistic: 83.48 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
fit.tax <- lm(medv ~ tax, data = Boston)
summary(fit.tax) #significant
```

```
##
## Call:
## lm(formula = medv ~ tax, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -14.091  -5.173  -2.085   3.158  34.058
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 32.970654   0.948296   34.77   <2e-16 ***
## tax         -0.025568   0.002147  -11.91   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.133 on 504 degrees of freedom
## Multiple R-squared:  0.2195, Adjusted R-squared:  0.218
## F-statistic: 141.8 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
fit.ptratio <- lm(medv ~ ptratio, data = Boston)
summary(fit.ptratio) #significant
```

```
##
## Call:
## lm(formula = medv ~ ptratio, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.8342  -4.8262  -0.6426   3.1571  31.2303
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   62.345      3.029   20.58   <2e-16 ***
## ptratio       -2.157      0.163  -13.23   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.931 on 504 degrees of freedom
## Multiple R-squared:  0.2578, Adjusted R-squared:  0.2564
## F-statistic: 175.1 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
fit.black <- lm(medv ~ black, data = Boston)
summary(fit.black) #significant
```

```
##
```

```
## Call:
## lm(formula = medv ~ black, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -18.884  -4.862  -1.684   2.932  27.763
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.551034   1.557463   6.775 3.49e-11 ***
## black        0.033593   0.004231   7.941 1.32e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.679 on 504 degrees of freedom
## Multiple R-squared:  0.1112, Adjusted R-squared:  0.1094
## F-statistic: 63.05 on 1 and 504 DF,  p-value: 1.318e-14
```

```r
fit.lstat <- lm(medv ~ lstat, data = Boston)
summary(fit.lstat) #significant
```

```
##
## Call:
## lm(formula = medv ~ lstat, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -15.168  -3.990  -1.318   2.034  24.500
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 34.55384    0.56263   61.41   <2e-16 ***
## lstat       -0.95005    0.03873  -24.53   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.216 on 504 degrees of freedom
## Multiple R-squared:  0.5441, Adjusted R-squared:  0.5432
## F-statistic: 601.6 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
fit.chas <- lm(medv ~ chas, data = Boston)
summary(fit.chas)
```

```
##
## Call:
## lm(formula = medv ~ chas, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -17.094  -5.894  -1.417   2.856  27.906
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 22.0938     0.4176  52.902  < 2e-16 ***
## chas1        6.3462     1.5880   3.996 7.39e-05 ***
```
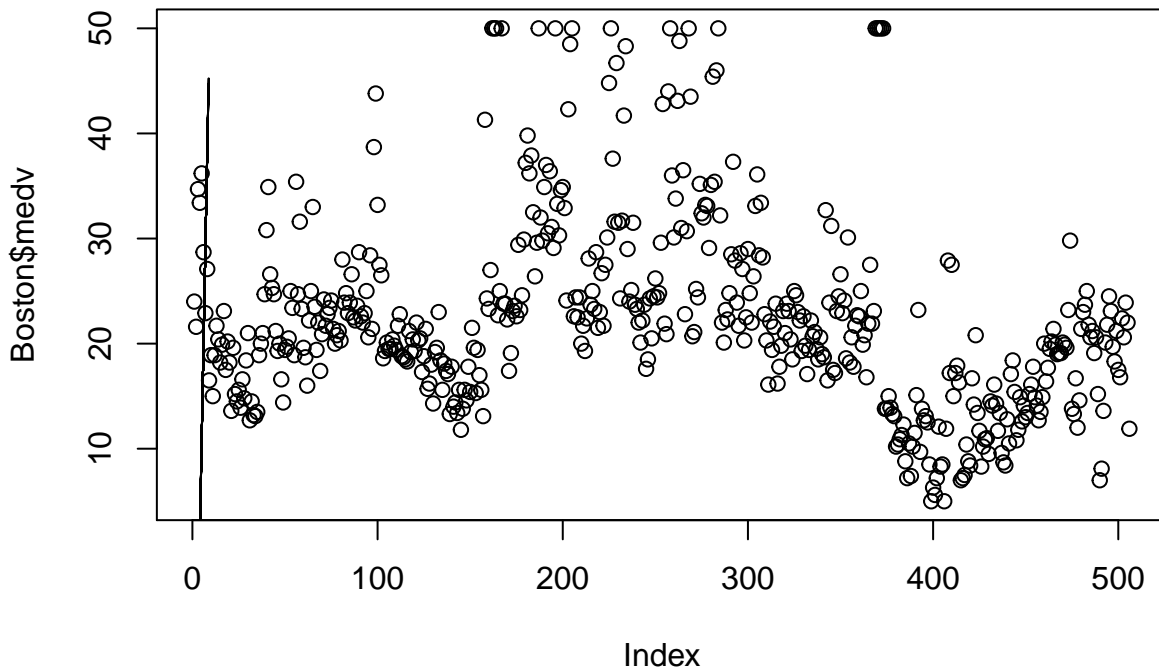
```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9.064 on 504 degrees of freedom
## Multiple R-squared:  0.03072,    Adjusted R-squared:  0.02879
## F-statistic: 15.97 on 1 and 504 DF,  p-value: 7.391e-05
```

```r
#number of rooms
fit.rm <- lm(medv ~ rm, data = Boston)
summary(fit.rm)
```

```
##
## Call:
## lm(formula = medv ~ rm, data = Boston)
##
## Residuals:
##      Min      1Q  Median      3Q     Max
## -23.346  -2.547   0.090   2.986  39.433
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -34.671      2.650  -13.08   <2e-16 ***
## rm             9.102      0.419   21.72   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.616 on 504 degrees of freedom
## Multiple R-squared:  0.4835, Adjusted R-squared:  0.4825
## F-statistic: 471.8 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
plot(Boston$medv, main = "linear regression plot of rm vs medv")
lines(Boston$rm,predict(fit.rm))
```
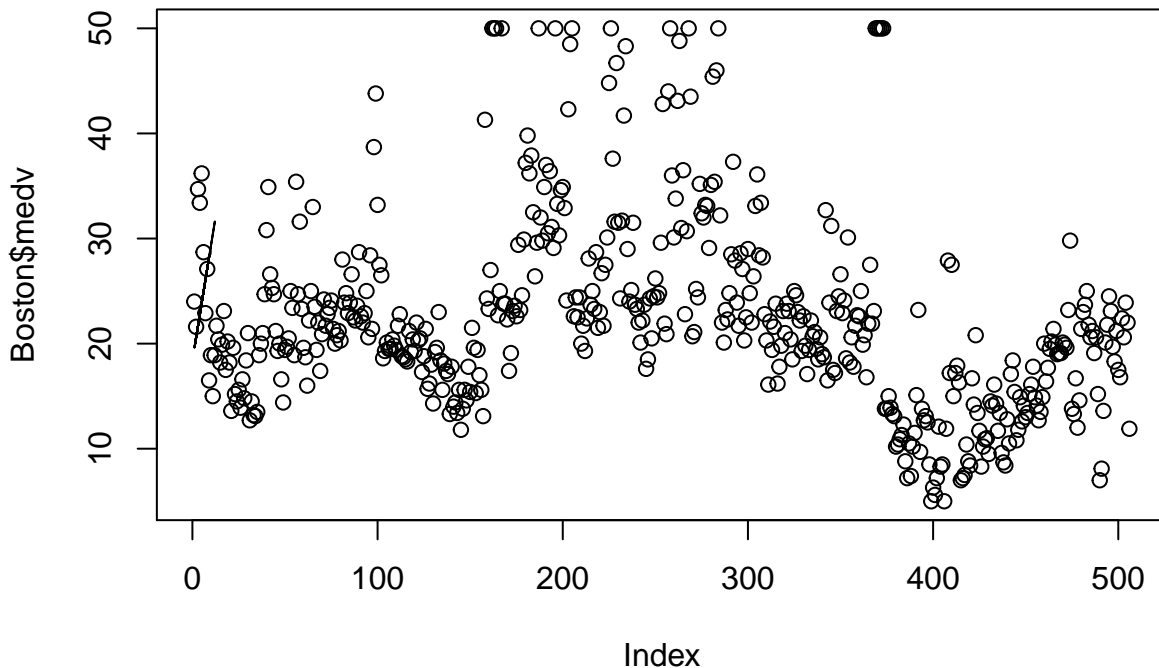
# linear regression plot of rm vs medv



```r
#distances to five employment centres
fit.dis <- lm(medv ~ dis, data = Boston)
summary(fit.dis)
```

```
## 
## Call:
## lm(formula = medv ~ dis, data = Boston)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max
## -15.016  -5.556  -1.865   2.288  30.377
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  18.3901     0.8174  22.499  < 2e-16 ***
## dis           1.0916     0.1884   5.795 1.21e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 8.914 on 504 degrees of freedom
## Multiple R-squared:  0.06246,    Adjusted R-squared:  0.0606
## F-statistic: 33.58 on 1 and 504 DF,  p-value: 1.207e-08
```

```r
plot(Boston$medv, main = "linear regression plot of dis vs medv")
lines(Boston$dis,predict(fit.dis))
```
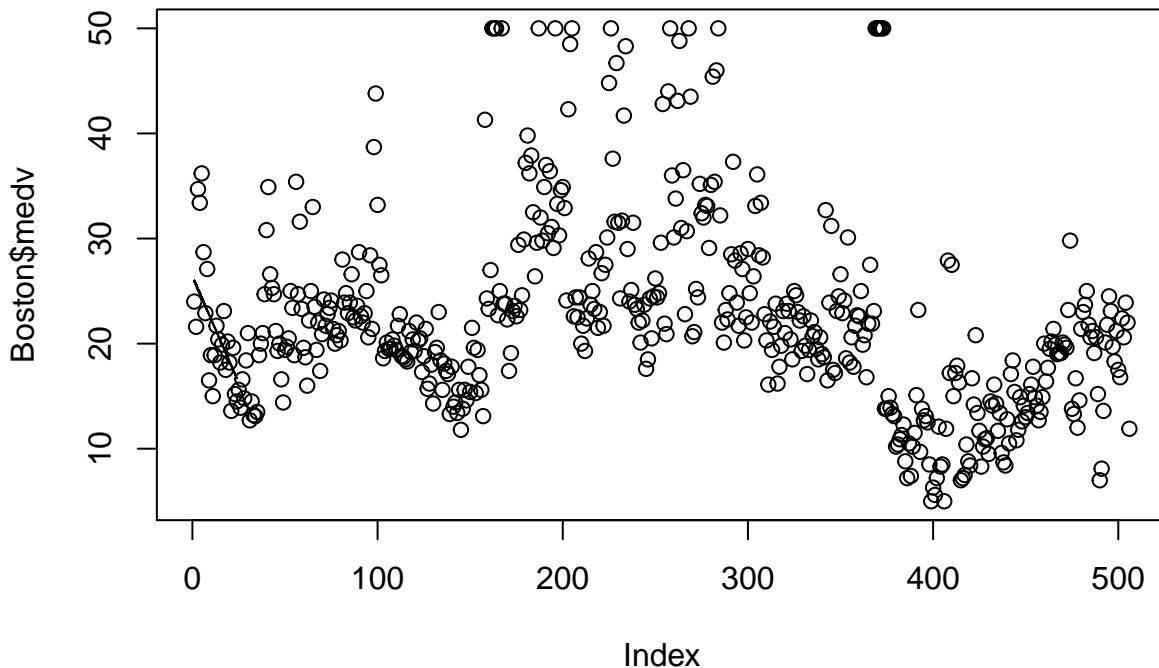
# linear regression plot of dis vs medv



```r
#index of accessibility to radial highways
fit.rad <- lm(medv ~ rad, data = Boston)
summary(fit.rad)
```

```
##
## Call:
## lm(formula = medv ~ rad, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -17.770  -5.199  -1.967   3.321  33.292
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 26.38213    0.56176  46.964   <2e-16 ***
## rad         -0.40310    0.04349  -9.269   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.509 on 504 degrees of freedom
## Multiple R-squared:  0.1456, Adjusted R-squared:  0.1439
## F-statistic: 85.91 on 1 and 504 DF,  p-value: < 2.2e-16
```

```r
plot(Boston$medv, main = "linear regression plot of rad vs medv")
lines(Boston$rad,predict(fit.rad))
```

## linear regression plot of rad vs medv



All of these three models show statistically significant associations between the predictor and the response ($p << 0.05$). But the plots show otherwise.
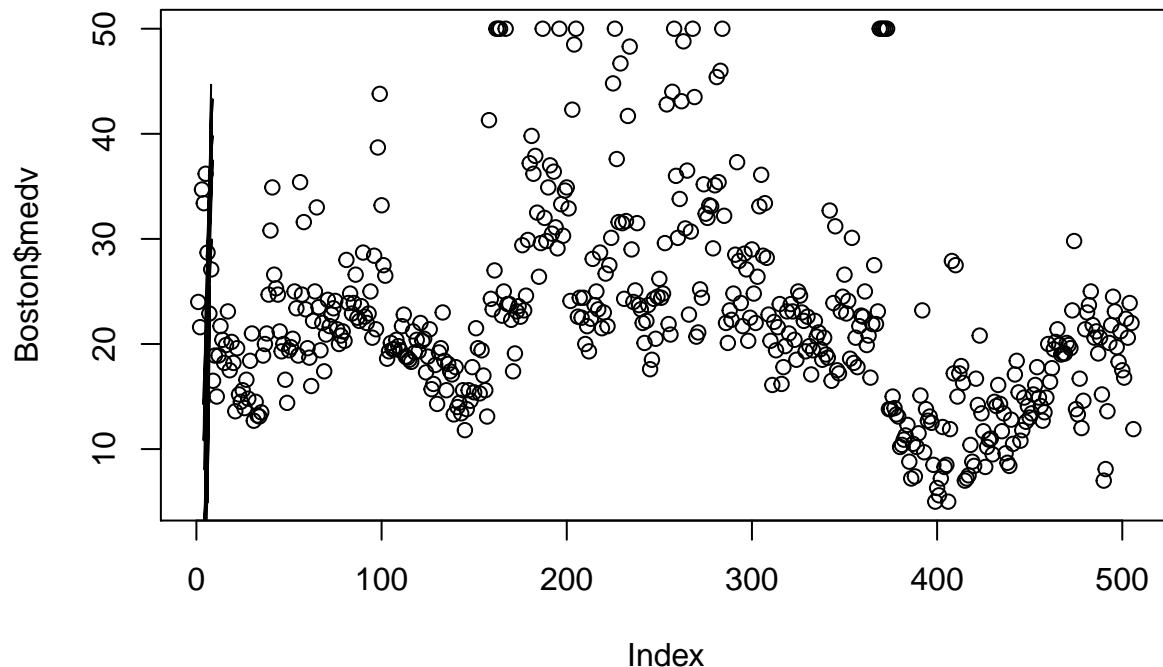
4. Fit a multiple regression model to predict the response using all of the predictors. Describe your results. For which predictors can we reject the null hypothesis $H_0 : \beta_j = 0$?

```
#multiple linear regressions
fit.all <- lm(medv ~ crim + zn + indus + chas + nox + rm + age + dis + rad + tax + ptratio + black
summary(fit.all)
```

```
##
## Call:
## lm(formula = medv ~ crim + zn + indus + chas + nox + rm + age +
##     dis + rad + tax + ptratio + black + lstat, data = Boston)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -15.595  -2.730  -0.518   1.777  26.199
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.646e+01  5.103e+00   7.144 3.28e-12 ***
## crim        -1.080e-01  3.286e-02  -3.287 0.001087 **
## zn           4.642e-02  1.373e-02   3.382 0.000778 ***
## indus        2.056e-02  6.150e-02   0.334 0.738288
## chas1        2.687e+00  8.616e-01   3.118 0.001925 **
## nox         -1.777e+01  3.820e+00  -4.651 4.25e-06 ***
## rm           3.810e+00  4.179e-01   9.116  < 2e-16 ***
## age          6.922e-04  1.321e-02   0.052 0.958229
## dis         -1.476e+00  1.995e-01  -7.398 6.01e-13 ***
## rad          3.060e-01  6.635e-02   4.613 5.07e-06 ***
```

10

```
## tax           -1.233e-02  3.760e-03  -3.280 0.001112 **
## ptratio       -9.527e-01  1.308e-01  -7.283 1.31e-12 ***
## black          9.312e-03  2.686e-03   3.467 0.000573 ***
## lstat          -5.248e-01  5.072e-02 -10.347  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.745 on 492 degrees of freedom
## Multiple R-squared:  0.7406, Adjusted R-squared:  0.7338
## F-statistic: 108.1 on 13 and 492 DF,  p-value: < 2.2e-16
```

```r
plot(Boston$medv)
lines(Boston$rm,predict(fit.all))
```



For predictors zn, chas1, nox, rm, dis, rad, ptratio, black and lstat (Proportion of residential land zoned for lots over 25,000 sq.ft, Charles River dummy variable, nitrogen oxides concentration, average number of rooms per dwelling, weighted mean of distances to five Boston employment centres, index of accessibility to radial highways, pupil-teacher ratio, proportion of blacks and lower status of the population(percent) ), we can reject the null hypothesis because their corresponding p values are smaller than 0.05.

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.
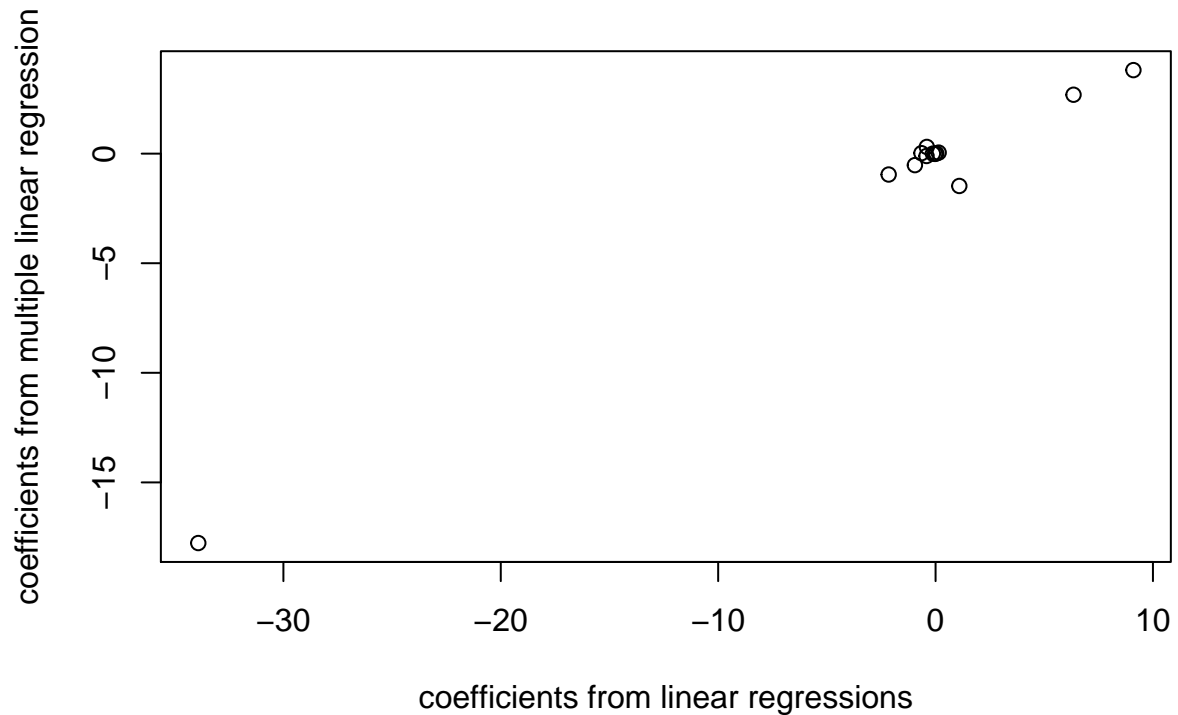
```r
#get simple linear regressions' coefficients

crim <- fit.crim$coefficients
comp <- as.data.frame(crim)
comp$zn <- fit.zn$coefficients
comp$indus <- fit.indus$coefficients
comp$chas <- fit.chas$coefficients
comp$nox <- fit.nox$coefficients
comp$rm <- fit.rm$coefficients
comp$age <- fit.age$coefficients
```

11

```
comp$dis <- fit.dis$coefficients
comp$rad <- fit.rad$coefficients
comp$tax<- fit.tax$coefficients
comp$ptratio <- fit.ptratio$coefficients
comp$black <- fit.black$coefficients
comp$lstat <- fit.lstat$coefficients

#multiple regression
coe <- fit.all$coefficients[c(-1)]
plot(as.numeric(comp[2,]),as.numeric(coe), xlab = "coefficients from linear regressions", ylab = "c
```



Based on the plot, oefficient for all variables from univariate linear regressions is similar to their corresponding coefficient from multiple linear regression, which means [4] supports [3].

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:

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

```
#non-linear regressions' coefficients
predictors <- names(Boston[,-ncol(Boston)])
r2 <- NULL
for(i in predictors){
    tmp <- lm(Boston$medv ~ Boston[,i] + Boston[,i]^2 + Boston[,i]^3)
    r2[i] <- summary(tmp)$r.squared
}
r2

##        crim         zn      indus       chas        nox         rm
## 0.15078047 0.12992084 0.23399003 0.03071613 0.18260304 0.48352546
##         age        dis        rad        tax    ptratio      black
```

```
## 0.14209474 0.06246437 0.14563858 0.21952592 0.25784732 0.11119612
##      lstat
## 0.54414630
```

There is evidence of a cubic polynomial association for rm and lstat because the R-square values are significantly higher, other predictors does not seem to have such association.

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

```
#model selection
step <- stepAIC(fit.all, direction="both")
```

```
## Start:  AIC=1589.64
## medv ~ crim + zn + indus + chas + nox + rm + age + dis + rad +
##      tax + ptratio + black + lstat
##
##             Df Sum of Sq   RSS    AIC
## - age        1      0.06 11079 1587.7
## - indus      1      2.52 11081 1587.8
## <none>                   11079 1589.6
## - chas       1    218.97 11298 1597.5
## - tax        1    242.26 11321 1598.6
## - crim       1    243.22 11322 1598.6
## - zn         1    257.49 11336 1599.3
## - black      1    270.63 11349 1599.8
## - rad        1    479.15 11558 1609.1
## - nox        1    487.16 11566 1609.4
## - ptratio    1   1194.23 12273 1639.4
## - dis        1   1232.41 12311 1641.0
## - rm         1   1871.32 12950 1666.6
## - lstat      1   2410.84 13490 1687.3
##
## Step:  AIC=1587.65
## medv ~ crim + zn + indus + chas + nox + rm + dis + rad + tax +
##      ptratio + black + lstat
##
##             Df Sum of Sq   RSS    AIC
## - indus      1      2.52 11081 1585.8
## <none>                   11079 1587.7
## + age        1      0.06 11079 1589.6
## - chas       1    219.91 11299 1595.6
## - tax        1    242.24 11321 1596.6
## - crim       1    243.20 11322 1596.6
## - zn         1    260.32 11339 1597.4
## - black      1    272.26 11351 1597.9
## - rad        1    481.09 11560 1607.2
## - nox        1    520.87 11600 1608.9
## - ptratio    1   1200.23 12279 1637.7
## - dis        1   1352.26 12431 1643.9
## - rm         1   1959.55 13038 1668.0
## - lstat      1   2718.88 13798 1696.7
##
## Step:  AIC=1585.76
## medv ~ crim + zn + chas + nox + rm + dis + rad + tax + ptratio +
##      black + lstat
```

```
## 
##           Df Sum of Sq   RSS     AIC
## <none>                   11081 1585.8
## + indus     1     2.52 11079 1587.7
## + age       1     0.06 11081 1587.8
## - chas      1   227.21 11309 1594.0
## - crim      1   245.37 11327 1594.8
## - zn        1   257.82 11339 1595.4
## - black     1   270.82 11352 1596.0
## - tax       1   273.62 11355 1596.1
## - rad       1   500.92 11582 1606.1
## - nox       1   541.91 11623 1607.9
## - ptratio   1  1206.45 12288 1636.0
## - dis       1  1448.94 12530 1645.9
## - rm        1  1963.66 13045 1666.3
## - lstat     1  2723.48 13805 1695.0
```

```
step$anova # display results
```

```
## Stepwise Model Path
## Analysis of Deviance Table
##
## Initial Model:
## medv ~ crim + zn + indus + chas + nox + rm + age + dis + rad +
##     tax + ptratio + black + lstat
##
## Final Model:
## medv ~ crim + zn + chas + nox + rm + dis + rad + tax + ptratio +
##     black + lstat
##
##
##       Step Df   Deviance Resid. Df Resid. Dev      AIC
## 1                              492   11078.78 1589.643
## 2    - age  1 0.06183435       493   11078.85 1587.646
## 3 - indus  1 2.51754013       494   11081.36 1585.761
```

The predictors suggested by stepwise model selection are completely different from the ones suggested by multiple linear regression. (Stepwise model selection suggests: age, indus, whereas multiple linear regression suggests: zn, chas1, nox, rm, dis, rad, ptratio, black and lstat. No overlap.)

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.

Stepwise Linear Regression assumption: Linear relationship, multivariate normality, no or little multicollinearity, no auto-correlation, homoscedasticity. Because the Residual vs Fitted plot shows a non-linear relationship, it means that our model has some non-linearity, unequal error variances, and outliers.

```
#model residuals
plot(step)
```

## Residuals vs Fitted



Fitted values
lm(medv ~ crim + zn + chas + nox + rm + dis + rad + tax + ptratio + black + ...

## Normal Q–Q



Theoretical Quantiles
lm(medv ~ crim + zn + chas + nox + rm + dis + rad + tax + ptratio + black + ...

Scale–Location

√|Standardized residuals|

Fitted values
lm(medv ~ crim + zn + chas + nox + rm + dis + rad + tax + ptratio + black + ...

Residuals vs Leverage

Standardized residuals

Cook's distance

Leverage
lm(medv ~ crim + zn + chas + nox + rm + dis + rad + tax + ptratio + black + ...