

IMT 573: Problem Set 7

Regression

Jenny Skytta

Due: May 15, 2022

Collaborators:

Instructions: Before beginning this assignment, please ensure you have access to R and RStudio; this can be on your own personal computer or on the IMT 573 R Studio Cloud.

1. Download the `07_ps_regression.Rmd` file from Canvas or save a copy to your local directory on RStudio Cloud. Supply your solutions to the assignment by editing `07_ps_regression.Rmd`.
2. Replace the “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. Be sure that each visualization adds value to your written explanation; avoid redundancy – you do not need four different visualizations of the same pattern.
4. Collaboration on problem sets is fun and useful, and we encourage it, but each student must turn in an individual write-up in their own words as well as code/work that is their own. Regardless of whether you work with others, what you turn in must be your own work; this includes code and interpretation of results. The names of all collaborators must be listed on each assignment. Do not copy-and-paste from other students’ responses or code.
5. All materials and resources that you use (with the exception of lecture slides) must be appropriately referenced within your assignment.
6. Remember partial credit will be awarded for each question for which a serious attempt at finding an answer has been shown. Students are *strongly* encouraged to attempt each question and to document their reasoning process even if they cannot find the correct answer. If you would like to include R code to show this process, but it does not run without errors you can do so with the `eval=FALSE` option as follows:

```
a + b # these object don't exist
# if you run this on its own it will give an error
```

7. When you have completed the assignment and have **checked** that your code both runs in the Console and knits correctly when you click Knit, download and rename the knitted PDF file to `ps7_YourLastName_YourFirstName.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(knitr) # this will keep code on the page!
opts_chunk$set(tidy.opts=list(width.cutoff=60),tidy=TRUE)
```

Problem 1: 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. This data is much used in data science and statistics to demonstrate regression problems; and while it has a lot of advantages it comes with concerns. Load this data and use it to answer the following questions.

```
data("Boston")
?Boston
```

(a) Briefly describe where these data come from and why they were collected. Be sure to mention any concerns you have about these data. These data comprise variables that could affect the housing values in suburbs of Boston. The biggest outlier within these data is that there is a variable called, “black” which denotes the proportion of racially Black individuals as a variable affecting price. This is an incredibly overtly racist metric to use to evaluate housing and harkens the problematic practice of Redlining wherein specific geographic areas disallow or discourage racially diverse families from residing within those bounds to uphold white supremacy. The data were actually evaluations of willingness to pay for better air quality to demonstrate value in cleaner air.

```
colnames(Boston)
```

(b) Describe the data and variables that are part of the Boston dataset. Tidy data as necessary.

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

```
glimpse(Boston)
```

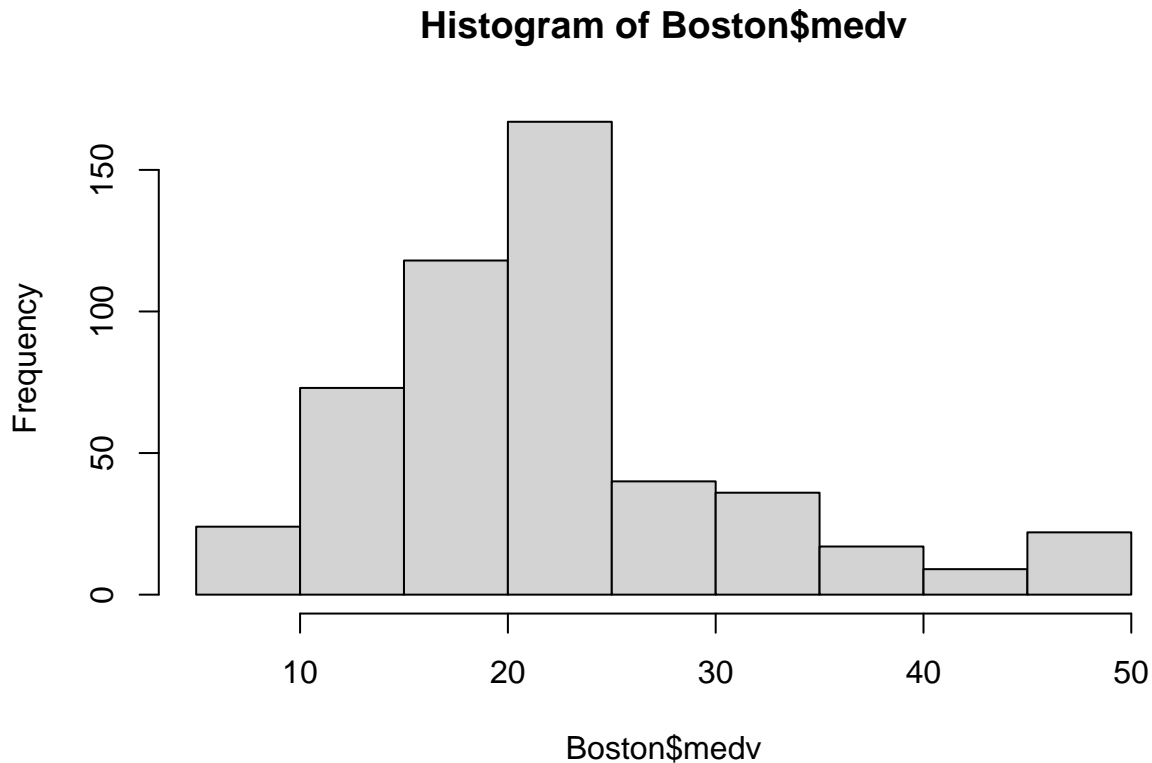
```
## Rows: 506
## Columns: 14
## $ crim    <dbl> 0.00632, 0.02731, 0.02729, 0.03237, 0.06905, 0.02985, 0.08829, ~
## $ zn      <dbl> 18.0, 0.0, 0.0, 0.0, 0.0, 0.0, 12.5, 12.5, 12.5, 12.5, 12.5, 1~
## $ indus   <dbl> 2.31, 7.07, 7.07, 2.18, 2.18, 2.18, 7.87, 7.87, 7.87, 7.87, 7.~
## $ chas    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ nox     <dbl> 0.538, 0.469, 0.469, 0.458, 0.458, 0.458, 0.524, 0.524, 0.524, ~
## $ rm      <dbl> 6.575, 6.421, 7.185, 6.998, 7.147, 6.430, 6.012, 6.172, 5.631, ~
## $ age     <dbl> 65.2, 78.9, 61.1, 45.8, 54.2, 58.7, 66.6, 96.1, 100.0, 85.9, 9~
## $ dis     <dbl> 4.0900, 4.9671, 4.9671, 6.0622, 6.0622, 6.0622, 5.5605, 5.9505~
## $ rad     <int> 1, 2, 2, 3, 3, 3, 5, 5, 5, 5, 5, 5, 5, 4, 4, 4, 4, 4, 4, ~
## $ tax     <dbl> 296, 242, 242, 222, 222, 222, 311, 311, 311, 311, 311, 311, 31~
## $ ptratio <dbl> 15.3, 17.8, 17.8, 18.7, 18.7, 18.7, 15.2, 15.2, 15.2, 15.2, 15~
## $ black   <dbl> 396.90, 396.90, 392.83, 394.63, 396.90, 394.12, 395.60, 396.90~
## $ lstat   <dbl> 4.98, 9.14, 4.03, 2.94, 5.33, 5.21, 12.43, 19.15, 29.93, 17.10~
## $ medv    <dbl> 24.0, 21.6, 34.7, 33.4, 36.2, 28.7, 22.9, 27.1, 16.5, 18.9, 15~
```

The variables are crim: per capita crime rate by town, zn: proportion of residential land zoned for lots over 25,000 square feet, indus: proportion of non-retail business acres per town, chas: Charles River dummy variable that equals 1 if tract bounds the river, otherwise its 0, nox: nitrogen oxides concentration in 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 centers, rad: index of accessibility to radial highways, tax: full-value property tax rate per \$10,000, ptratio: pupil-teacher ratio by town, black: proportion of Black individuals by town, lstat: lower status of the population by percent, medv: median value of owner-occupied homes in 1000 dollars.

(d) Consider this data in context, what is the response variable of interest? The response variable is median value: medv.

```
hist(Boston$medv)
```



```
Mod_nox <- lm(medv ~ nox, data = Boston)
Mod_crim <- lm(medv ~ crim, data = Boston)
Mod_indus <- lm(medv ~ indus, data = Boston)
Mod_chas <- lm(medv ~ chas, data = Boston)
Mod_rm <- lm(medv ~ rm, data = Boston)
Mod_age <- lm(medv ~ age, data = Boston)
Mod_dis <- lm(medv ~ dis, data = Boston)
Mod_rad <- lm(medv ~ rad, data = Boston)
Mod_tax <- lm(medv ~ tax, data = Boston)
Mod_ptratio <- lm(medv ~ ptratio, data = Boston)
Mod_black <- lm(medv ~ black, data = Boston)
Mod_lstat <- lm(medv ~ lstat, data = Boston)
Mod_zn <- lm(medv ~ zn, data = Boston)

summary(Mod_nox)
```

(e) 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.

```
##
## 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
```

```
summary(Mod_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
```

```
summary(Mod_indus)
```

```
##
## 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
```

```
summary(Mod_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 ***
## chas         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
```

```
summary(Mod_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
```

```
summary(Mod_age)
```

```
##
## 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
```

```
summary(Mod_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
```

```
summary(Mod_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
```

```
summary(Mod_tax)
```

```
##
## 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
```

```
summary(Mod_ptratio)
```

```
##
## 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
```

```
summary(Mod_black)
```

```
##
## 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
```

```
summary(Mod_lstat)
```

```
##
## 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
```

```
summary(Mod_zn)
```

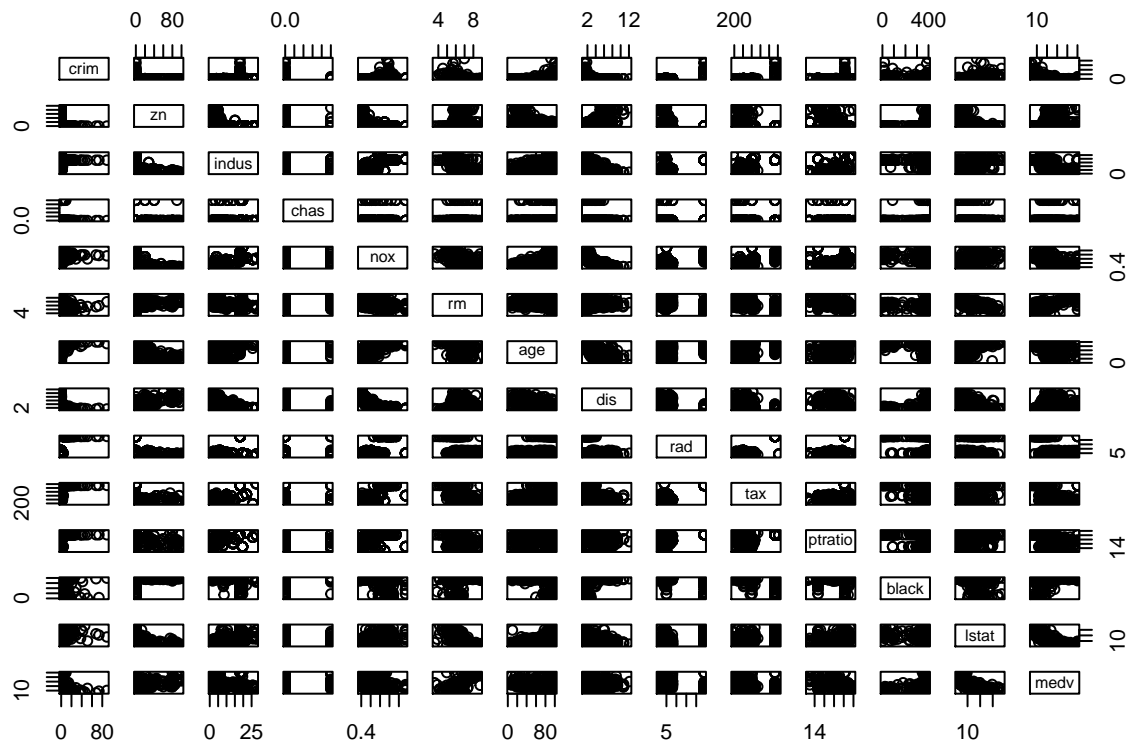
```
##
## Call:
## lm(formula = medv ~ zn, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.918  -5.518  -1.006   2.757  29.082
##
## Coefficients:
```



```
##           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
```

```
# forcing integers to 'as numeric' in Boston dataframe for consistency
Boston$chas <- as.numeric(Boston$chas)
Boston$rad <- as.numeric(Boston$rad)

# plotting the Boston DF to see if there is any visual linear relationships
plot(Boston)
```



There is a statistically significant association between medv and nox, crim, zn, chas, rm, dis, rad, tax, ptratio, black, and lstat.

(f) 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 $H_0 : \beta_j = 0$? The following predictors were statistically significant predictor variables for median value along with their p values: crim 0.001087, zn 0.000778, chas 0.001925, nox 4.25e-06, rm < 2e-16, dis 6.01e-13, rad 5.07e-06, tax 0.001112, ptratio 1.31e-12, black 0.000573, lstat < 2e-16.

```
# running a correlation analysis of the Boston DF
cor(Boston[,1:14])
```

```
##           crim           zn           indus           chas           nox
```

```
## crim      1.00000000 -0.20046922  0.40658341 -0.055891582  0.42097171
## zn        -0.20046922  1.00000000 -0.53382819 -0.042696719 -0.51660371
## indus     0.40658341 -0.53382819  1.00000000  0.062938027  0.76365145
## chas      -0.05589158 -0.04269672  0.06293803  1.000000000  0.09120281
## nox       0.42097171 -0.51660371  0.76365145  0.091202807  1.00000000
## rm        -0.21924670  0.31199059 -0.39167585  0.091251225 -0.30218819
## age       0.35273425 -0.56953734  0.64477851  0.086517774  0.73147010
## dis       -0.37967009  0.66440822 -0.70802699 -0.099175780 -0.76923011
## rad       0.62550515 -0.31194783  0.59512927 -0.007368241  0.61144056
## tax       0.58276431 -0.31456332  0.72076018 -0.035586518  0.66802320
## ptratio   0.28994558 -0.39167855  0.38324756 -0.121515174  0.18893268
## black     -0.38506394  0.17552032 -0.35697654  0.048788485 -0.38005064
## lstat     0.45562148 -0.41299457  0.60379972 -0.053929298  0.59087892
## medv      -0.38830461  0.36044534 -0.48372516  0.175260177 -0.42732077
##           rm      age      dis      rad      tax      ptratio
## crim      -0.21924670  0.35273425 -0.37967009  0.625505145  0.58276431  0.28994556
## zn         0.31199059 -0.56953734  0.66440822 -0.311947826 -0.31456332 -0.39167855
## indus      -0.39167585  0.64477851 -0.70802699  0.595129275  0.72076018  0.38324756
## chas       0.09125123  0.08651777 -0.09917578 -0.007368241 -0.03558652 -0.12151512
## nox        -0.30218819  0.73147010 -0.76923011  0.611440563  0.66802320  0.1889327
## rm         1.00000000 -0.24026493  0.20524621 -0.209846668 -0.29204783 -0.3555015
## age        -0.24026493  1.00000000 -0.74788054  0.456022452  0.50645559  0.2615150
## dis        0.20524621 -0.74788054  1.00000000 -0.494587930 -0.53443158 -0.2324705
## rad        -0.20984667  0.45602245 -0.49458793  1.000000000  0.91022819  0.4647412
## tax        -0.29204783  0.50645559 -0.53443158  0.910228189  1.00000000  0.4608530
## ptratio    -0.35550149  0.26151501 -0.23247054  0.464741179  0.46085304  1.0000000
## black      0.12806864 -0.27353398  0.29151167 -0.444412816 -0.44180801 -0.1773833
## lstat      -0.61380827  0.60233853 -0.49699583  0.488676335  0.54399341  0.3740443
## medv       0.69535995 -0.37695457  0.24992873 -0.381626231 -0.46853593 -0.5077867
##           black      lstat      medv
## crim      -0.38506394  0.4556215 -0.3883046
## zn         0.17552032 -0.4129946  0.3604453
## indus      -0.35697654  0.6037997 -0.4837252
## chas       0.04878848 -0.0539293  0.1752602
## nox        -0.38005064  0.5908789 -0.4273208
## rm         0.12806864 -0.6138083  0.6953599
## age        -0.27353398  0.6023385 -0.3769546
## dis        0.29151167 -0.4969958  0.2499287
## rad        -0.44441282  0.4886763 -0.3816262
## tax        -0.44180801  0.5439934 -0.4685359
## ptratio    -0.17738330  0.3740443 -0.5077867
## black      1.00000000 -0.3660869  0.3334608
## lstat      -0.36608690  1.0000000 -0.7376627
## medv       0.33346082 -0.7376627  1.0000000
```

```
# creating a multivariate model
```

```
m1 <- lm(medv ~ crim + zn + indus + chas + nox + rm + age + dis + rad + tax + ptratio + black + lstat, data = medv)
```

```
# summarizing the multivariate model
```

```
summary(m1)
```

```
##
```

```
## 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
## chas         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 ***
## 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
```

```
AIC(m1)
```

```
## [1] 3027.609
```

(g) 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. The values are all very closely centered around 0 which denotes significance.

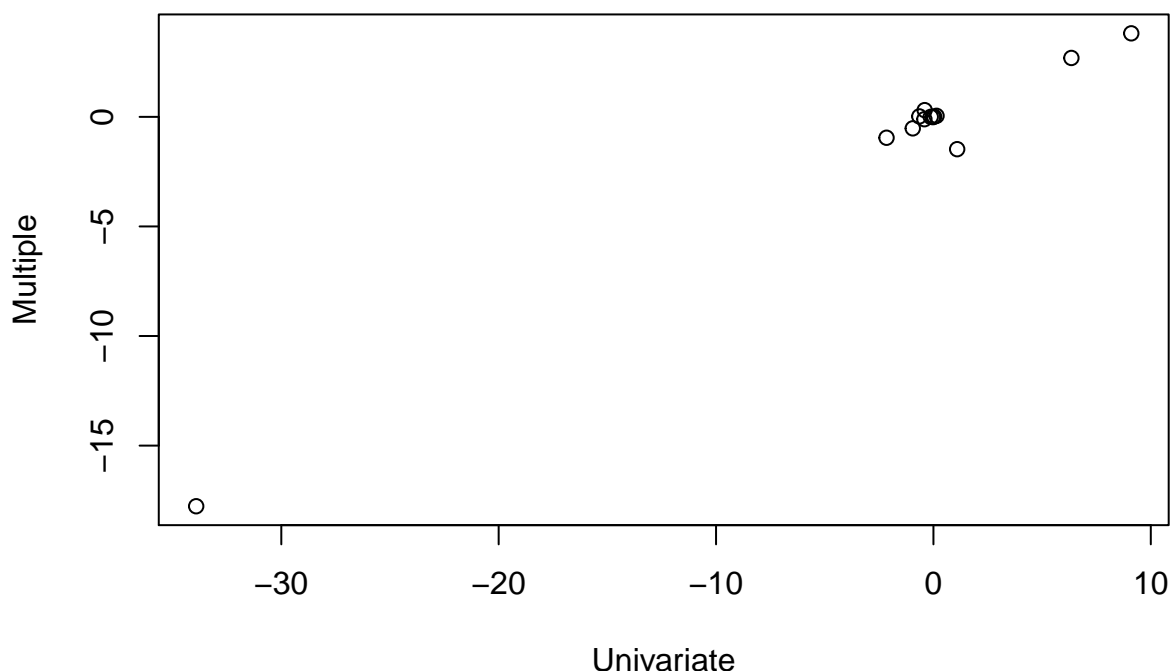
```
# creating a dataframe of univariate coefficient values
Boston_coefficients <- as.tibble(c(Mod_crim$coefficients[2], Mod_zn$coefficients[2], Mod_indus$coefficients[2],
Mod_ptratio$coefficients[2], Mod_black$coefficients[2], Mod_lstat$coefficients[2]))

## Warning: `as.tibble()` was deprecated in tibble 2.0.0.
## Please use `as_tibble()` instead.
## The signature and semantics have changed, see `?as_tibble`.

# adding a column into new dataframe with the multivariate coefficients
Boston_coefficients$mutlicos <- c(m1$coefficients[2], m1$coefficients[3], m1$coefficients[4], m1$coefficients[5],
m1$coefficients[6], m1$coefficients[7], m1$coefficients[8], m1$coefficients[9], m1$coefficients[10], m1$coefficients[11],
m1$coefficients[12], m1$coefficients[13])

#plotting the univariate coefficient values on x and multivariate values
plot(Boston_coefficients, main = "Univariate vs. Multiple Regression Coefficients",
xlab = "Univariate", ylab = "Multiple")
```

Univariate vs. Multiple Regression Coefficients



(h) 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:

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

writing the equation using $X = 1$

```
m1$coefficients[2] + m1$coefficients[3]*1 + m1$coefficients[4]*1^2 + m1$coefficients[5]*1^3 + m1$coeffi
```

```
##      crim
## -13.9604
```

```
stepwise_m1 <- stepAIC(m1, trace = TRUE)
```

(i) 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)?

```
## 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
## - 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
## - 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

```

```
summary(stepwise_m1)
```

```

##
## Call:
## lm(formula = medv ~ crim + zn + chas + nox + rm + dis + rad +
##      tax + ptratio + black + lstat, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max

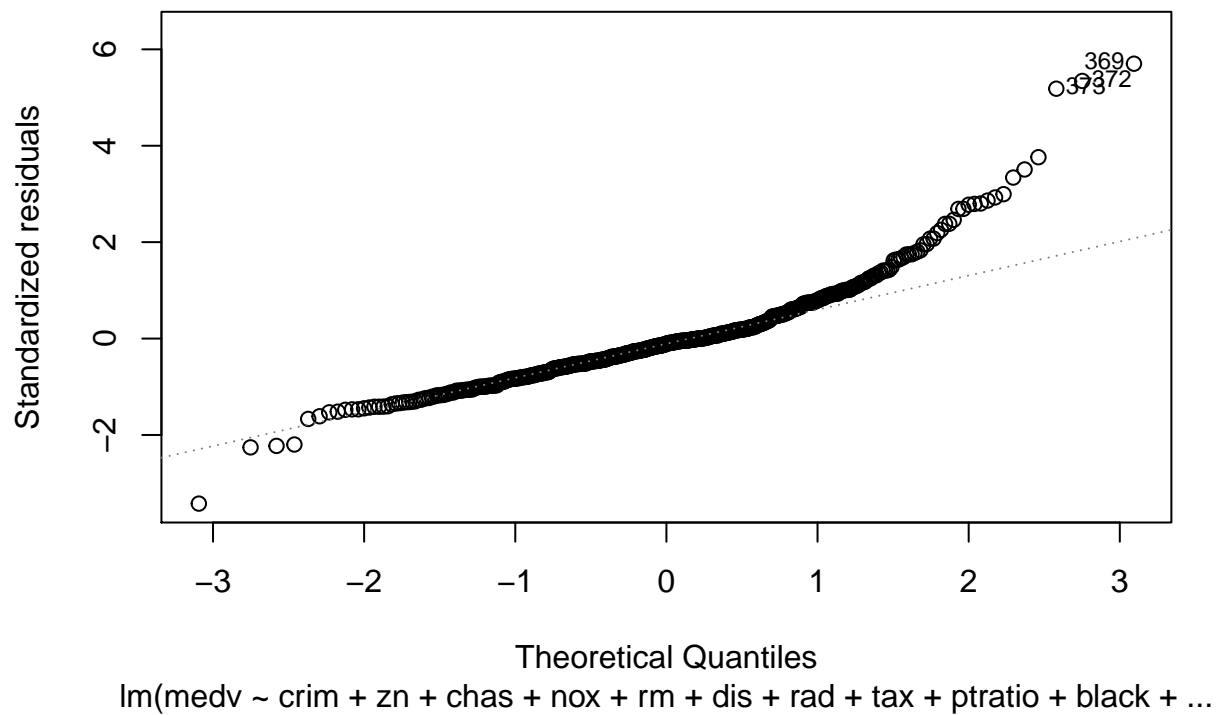
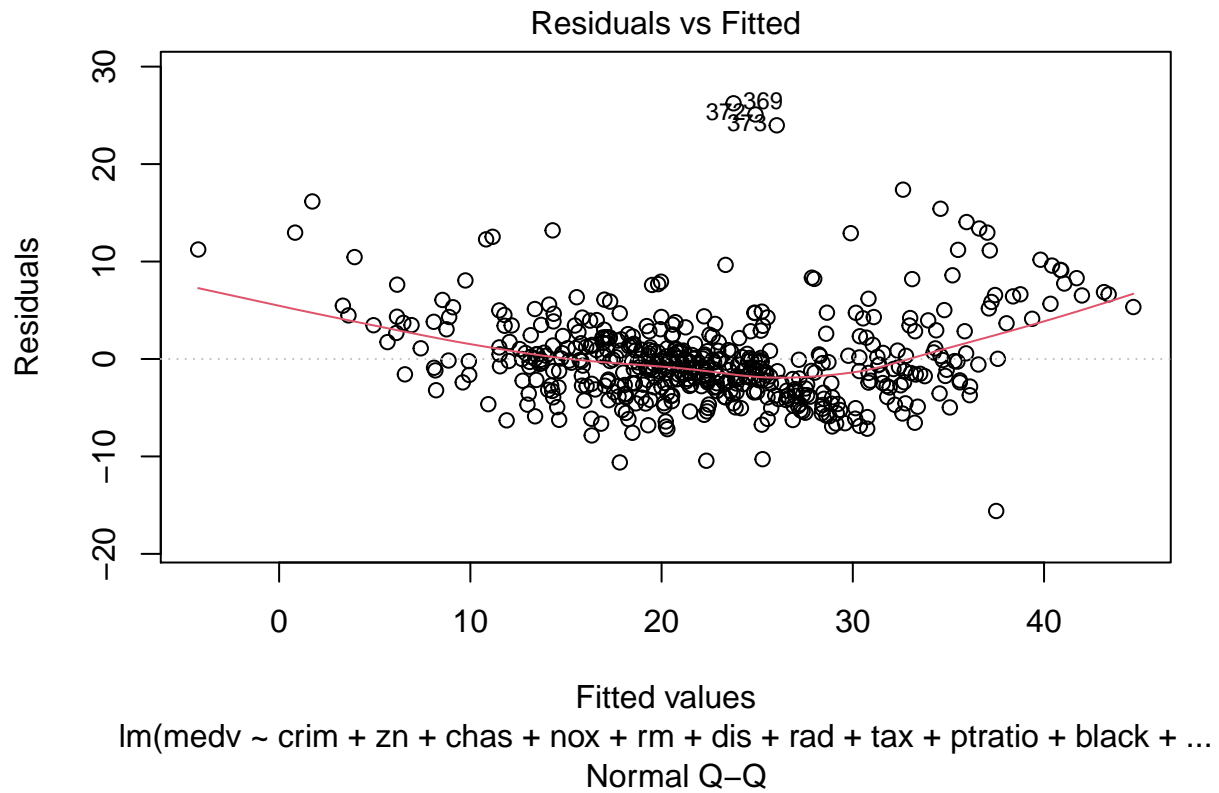
```

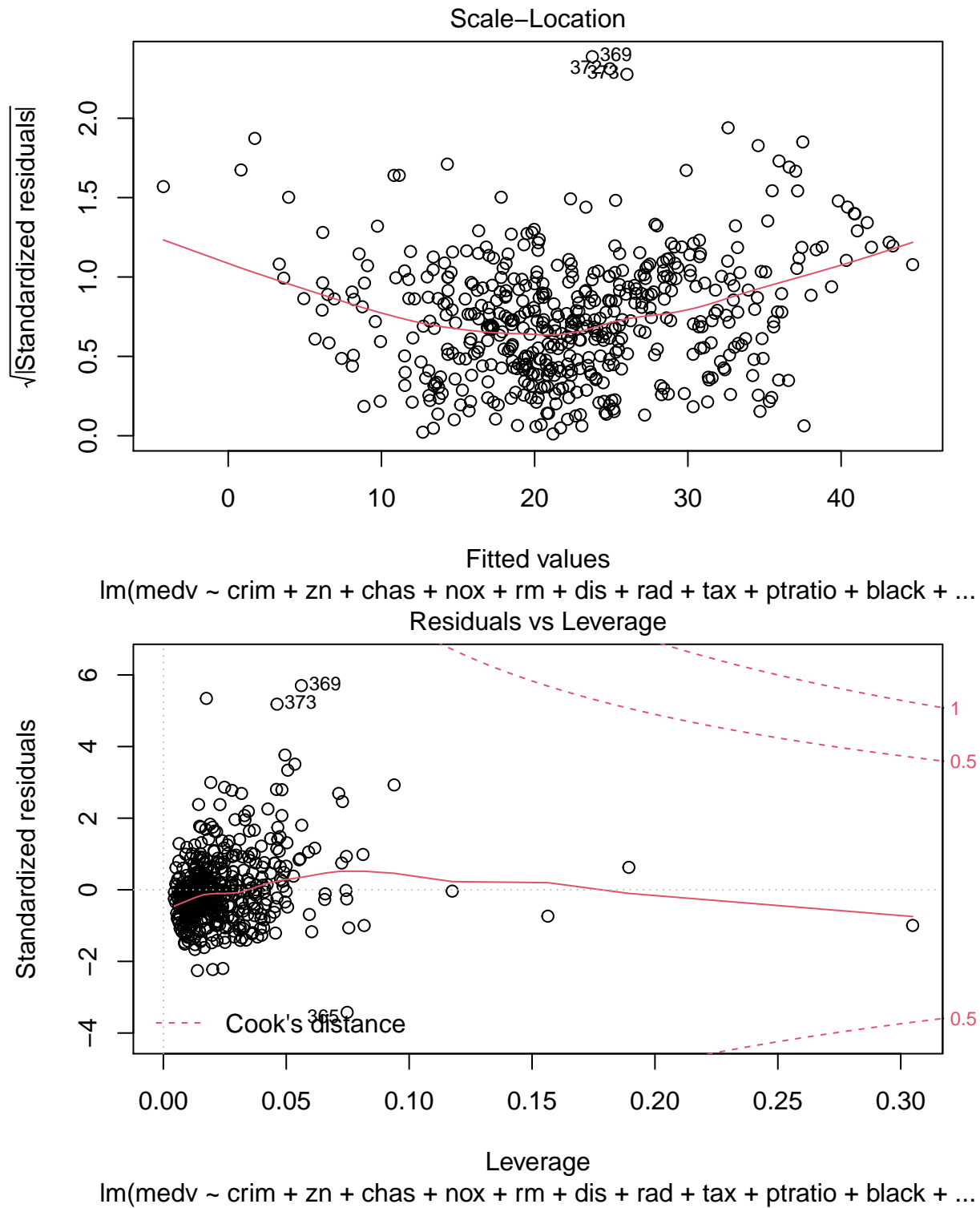
```
## -15.5984 -2.7386 -0.5046 1.7273 26.2373
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  36.341145   5.067492   7.171 2.73e-12 ***
## crim        -0.108413   0.032779  -3.307 0.001010 **
## zn           0.045845   0.013523   3.390 0.000754 ***
## chas         2.718716   0.854240   3.183 0.001551 **
## nox        -17.376023   3.535243  -4.915 1.21e-06 ***
## rm           3.801579   0.406316   9.356 < 2e-16 ***
## dis         -1.492711   0.185731  -8.037 6.84e-15 ***
## rad           0.299608   0.063402   4.726 3.00e-06 ***
## tax         -0.011778   0.003372  -3.493 0.000521 ***
## ptratio     -0.946525   0.129066  -7.334 9.24e-13 ***
## black        0.009291   0.002674   3.475 0.000557 ***
## lstat       -0.522553   0.047424 -11.019 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.736 on 494 degrees of freedom
## Multiple R-squared:  0.7406, Adjusted R-squared:  0.7348
## F-statistic: 128.2 on 11 and 494 DF, p-value: < 2.2e-16
```

In the stepwise model the crim, zn, chas, nox, rm, dis, rad, tax, ptratio, black and lstat variables were kept to minimize the AIC value. The model removed indus and age variables. The earlier model (m1) had an AIC over 3000.

(j) 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. The residuals are not in a cloud like pattern or well dispersed. They are grouped into a clear pattern. There is also a U-shape in the scale. There is also a curvilinear bend at the end of my Q-Q plot meaning my data is skewing. Based on this, I believe the data would need further delving to assess parameters.

```
plot(stepwise_m1)
```





Problem 2: A Critical Perspective to the Boston Housing Data

(a) When were these data collected? Did you note this in your descriptive above? Did the date surprise you? The data were collected in the late 70s which truthfully didn't surprise me.

(b) Amidst data features like number of rooms and access to highways are features like crime rate, and percentage Black per town. Whether intentional or not, someone looking at this data might infer a link between crime and race just due to the variables present; or even worse might use the data to support harsh policing policies based on race. Suppose for a moment we have a modern version of this dataset; the “Seattle Housing Data.” Discuss, in a few paragraphs, how this hypothetical dataset could be used (1) in a harmful way, and (2) in a beneficial way for society. Utilizing crime as a variable leaves out the etiology of how these data are generated as well as the problematic metric of disproportionate enforcement and abuse of power from law enforcement against predominantly marginalized communities. Police presence and abuse of power happens at a higher rate in Black communities which take at data value would show an increase measure of crime. Additionally, systemic barriers to advancement have historically restricted Black communities from advancement and denied the improvement of community through legacy and home ownership. Many of the communities within Seattle enforced redlining rules up until the 1960s which prevented the advancement of POC individuals within higher SES communities. This forced these POC regardless of SES status to reside in a lower equity zone. Additionally these same policies prevented loans to POC individuals regardless of ability to pay a mortgage. By relegating these POC individuals into urban areas that were allocated as “less than”, they lost on multiple advancement opportunities. Additionally, urban planners allowed for known equity disrupters to be built within or in the bounds of these communities including electrical power plants, freeway offramps and other easements and urban fixtures that led to lower property values.

Using crime statistics to suggest an additional need for enforcement is an absurdity when considering that its stating the $X = X$. The measurement of crime is recorded by policing data which is generated by their direct enforcement. Using their own disproportionate data to justify additional enforcement means they are creating a loop.

One beneficial way in which these data could be used is to create targeted programs to alleviate the burdens that exist in these communities. Noticing a trend in crime, a program to have police ombudsmen within these communities to measure police interactions could be a step towards assessing whether bias exist within the data. Another beneficial way this data could be used is to have afterschool programs or recidivism programs for people in the community who have criminal histories. These programs could help to foster positive outcomes or help career development. Crime is highly correlated with lower SES and a need to meet one’s basic needs. When you can create a scenario wherein people’s needs are met and they are safe, the crime rates will go down. There have been very successful ventures of community gardens that have created sense of community and also fed those who experienced food scarcity.

Citations

Harrison, D. and Rubinfeld, D.L. (1978) Hedonic prices and the demand for clean air. J. Environ. Economics and Management 5, 81–102.

Belsley D.A., Kuh, E. and Welsch, R.E. (1980) Regression Diagnostics. Identifying Influential Data and Sources of Collinearity. New York: Wiley.

Code written above is from the previous course IMT 511 which used the below text to support class scripts. https://www.google.com/books/edition/Programming_Skills_for_Data_Science/BnB6DwAAQBAJ?hl=en&gbpv=1&printsec=frontcover —