

# IMT 573: Problem Set 7 - Regression - Solutions

Vighnesh Misal

Due: Tuesday, November 19, 2019

*Collaborators: Ashish Anand*

## *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 Server.

1. Download the `problemset7.Rmd` file from Canvas or save a copy to your local directory on RStudio Server. Open `problemset7.Rmd` in RStudio and supply your solutions to the assignment by editing `problemset7.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. 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 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. (Note: I am also using the `include=FALSE` option here to not include this code in the PDF, but you need to remove this or change it to `TRUE` if you want to include the code chunk.)
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 PDF, 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(corrplot)
```

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.

Describe the data and variables that are part of the dataset. Tidy data as necessary.

Consider this data in context, what is the response variable of interest?

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.

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  $H_0: \beta_j = 0$ ?

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.

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:

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

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)?

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.

```
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
```

**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 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 ...
```

**nrow**(Boston)

```
## [1] 506
```

**ncol**(Boston)

```
## [1] 14
```

**sum**(**duplicated**(Boston))

```
## [1] 0
```

There are 506 observations in the dataset

There are 14 columns in the dataset

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.

The response variable of interest is **crim** (per capita crime rate by town) because it directly impacts house values for any given area and is a point of interest for me as a data scientist. The residues are also clustered in the upper left and corner and point to the presence of outliers for the presence of skewness in data.

```
attach(Boston)

fit.zn <- lm(crim ~ zn)
summary(fit.zn)

##
## Call:
## lm(formula = crim ~ zn)
##
## 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

fit.indus <- lm(crim ~ indus)
summary(fit.indus)

##
## Call:
## lm(formula = crim ~ indus)
##
## 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

chas <- as.factor(chas)
fit.chas <- lm(crim ~ chas)
summary(fit.chas)

##
## Call:
## lm(formula = crim ~ chas)
##
## 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 ***
## chas1        -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

fit.nox <- lm(crim ~ nox)
summary(fit.nox)

##
## Call:
## lm(formula = crim ~ nox)
##
## 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
```

```

fit.rm <- lm(crim ~ rm)
summary(fit.rm)

##
## Call:
## lm(formula = crim ~ rm)
##
## 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

fit.age <- lm(crim ~ age)
summary(fit.age)

##
## Call:
## lm(formula = crim ~ age)
##
## 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

fit.dis <- lm(crim ~ dis)
summary(fit.dis)

##
## Call:
## lm(formula = crim ~ dis)
##

```

```
## 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

fit.rad <- lm(crim ~ rad)
summary(fit.rad)

##
## Call:
## lm(formula = crim ~ rad)
##
## 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

fit.tax <- lm(crim ~ tax)
summary(fit.tax)

##
## Call:
## lm(formula = crim ~ tax)
##
## 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

fit.ptratio <- lm(crim ~ ptratio)
summary(fit.ptratio)

##
## Call:
## lm(formula = crim ~ ptratio)
##
## 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

fit.black <- lm(crim ~ black)
summary(fit.black)

##
## Call:
## lm(formula = crim ~ black)
##
## 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
```

```

fit.lstat <- lm(crim ~ lstat)
summary(fit.lstat)

##
## Call:
## lm(formula = crim ~ lstat)
##
## 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

fit.medv <- lm(crim ~ medv)
summary(fit.medv)

##
## Call:
## lm(formula = crim ~ medv)
##
## 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

fit.all <- lm(crim ~ ., data = Boston)
summary(fit.all)

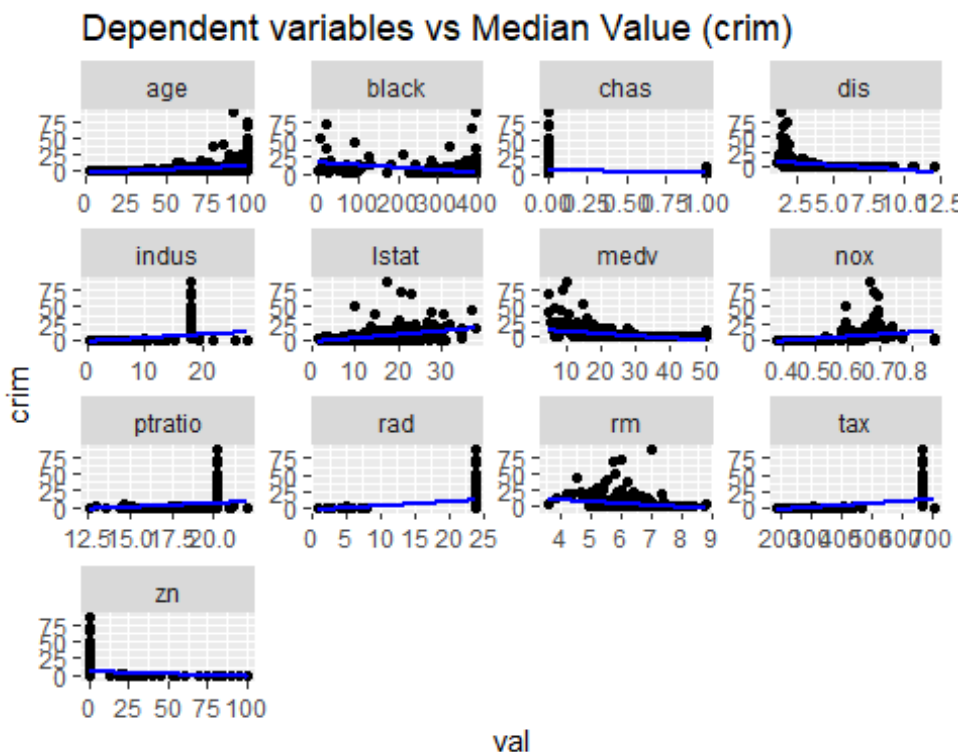
##
## Call:
## lm(formula = crim ~ ., data = Boston)
##

```

```
## 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
```

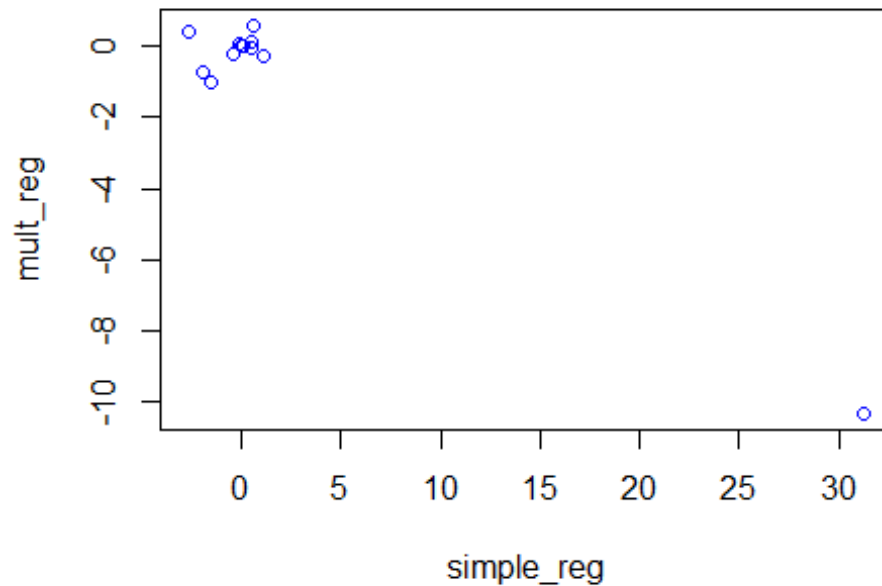
Boston %>%

```
  gather(key, val, -crim) %>%
  ggplot(aes(x = val, y = crim)) +
  geom_point() +
  stat_smooth(method = "lm", se = TRUE, col = "blue") +
  facet_wrap(~key, scales = "free") +
  theme_gray() +
  ggtitle("Dependent variables vs Median Value (crim)")
```



# We can reject the null hypothesis for zn, dis, rad, black, medv as they are much lower than the 5% p-value cutoff threshold.s

```
simple_reg <- vector("numeric",0)
simple_reg <- c(simple_reg, fit.zn$coefficient[2])
simple_reg <- c(simple_reg, fit.indus$coefficient[2])
simple_reg <- c(simple_reg, fit.chas$coefficient[2])
simple_reg <- c(simple_reg, fit.nox$coefficient[2])
simple_reg <- c(simple_reg, fit.rm$coefficient[2])
simple_reg <- c(simple_reg, fit.age$coefficient[2])
simple_reg <- c(simple_reg, fit.dis$coefficient[2])
simple_reg <- c(simple_reg, fit.rad$coefficient[2])
simple_reg <- c(simple_reg, fit.tax$coefficient[2])
simple_reg <- c(simple_reg, fit.ptratio$coefficient[2])
simple_reg <- c(simple_reg, fit.black$coefficient[2])
simple_reg <- c(simple_reg, fit.lstat$coefficient[2])
simple_reg <- c(simple_reg, fit.medv$coefficient[2])
mult_reg <- vector("numeric", 0)
mult_reg <- c(mult_reg, fit.all$coefficients)
mult_reg <- mult_reg[-1]
plot(simple_reg, mult_reg, col = "blue")
```



```
cor(Boston[-c(1, 4)])
```

```
##           zn          indus          nox          rm          age          dis
## zn          1.0000000 -0.5338282 -0.5166037  0.3119906 -0.5695373  0.6644082
## indus      -0.5338282  1.0000000  0.7636514 -0.3916759  0.6447785 -0.7080270
## nox        -0.5166037  0.7636514  1.0000000 -0.3021882  0.7314701 -0.7692301
## rm          0.3119906 -0.3916759 -0.3021882  1.0000000 -0.2402649  0.2052462
## age        -0.5695373  0.6447785  0.7314701 -0.2402649  1.0000000 -0.7478805
## dis         0.6644082 -0.7080270 -0.7692301  0.2052462 -0.7478805  1.0000000
## rad        -0.3119478  0.5951293  0.6114406 -0.2098467  0.4560225 -0.4945879
## tax        -0.3145633  0.7207602  0.6680232 -0.2920478  0.5064556 -0.5344316
## ptratio    -0.3916785  0.3832476  0.1889327 -0.3555015  0.2615150 -0.2324705
## black       0.1755203 -0.3569765 -0.3800506  0.1280686 -0.2735340  0.2915117
## lstat      -0.4129946  0.6037997  0.5908789 -0.6138083  0.6023385 -0.4969958
## medv       0.3604453 -0.4837252 -0.4273208  0.6953599 -0.3769546  0.2499287
##           rad          tax      ptratio          black          lstat          medv
## zn        -0.3119478 -0.3145633 -0.3916785  0.1755203 -0.4129946  0.3604453
## indus      0.5951293  0.7207602  0.3832476 -0.3569765  0.6037997 -0.4837252
## nox        0.6114406  0.6680232  0.1889327 -0.3800506  0.5908789 -0.4273208
## rm        -0.2098467 -0.2920478 -0.3555015  0.1280686 -0.6138083  0.6953599
## age        0.4560225  0.5064556  0.2615150 -0.2735340  0.6023385 -0.3769546
## dis       -0.4945879 -0.5344316 -0.2324705  0.2915117 -0.4969958  0.2499287
## rad        1.0000000  0.9102282  0.4647412 -0.4444128  0.4886763 -0.3816262
## tax        0.9102282  1.0000000  0.4608530 -0.4418080  0.5439934 -0.4685359
## ptratio    0.4647412  0.4608530  1.0000000 -0.1773833  0.3740443 -0.5077867
## black     -0.4444128 -0.4418080 -0.1773833  1.0000000 -0.3660869  0.3334608
```

```
## lstat    0.4886763  0.5439934  0.3740443 -0.3660869  1.0000000 -0.7376627
## medv     -0.3816262 -0.4685359 -0.5077867  0.3334608 -0.7376627  1.0000000
```

**The linear regression model shows a relationship between each variable and crime rate but when we run the multiple regression we observe that some predictors do not influence the response at all.**

```
fit.zn2 <- lm(crim ~ poly(zn, 3))
summary(fit.zn2)

##
## Call:
## lm(formula = crim ~ poly(zn, 3))
##
## 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)    3.6135     0.3722   9.709  < 2e-16 ***
## poly(zn, 3)1  -38.7498     8.3722  -4.628  4.7e-06 ***
## poly(zn, 3)2   23.9398     8.3722   2.859  0.00442 **
## poly(zn, 3)3  -10.0719     8.3722  -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

fit.indus2 <- lm(crim ~ poly(indus, 3))
summary(fit.indus2)

##
## Call:
## lm(formula = crim ~ poly(indus, 3))
##
## 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)    3.614     0.330  10.950  < 2e-16 ***
## poly(indus, 3)1  78.591     7.423  10.587  < 2e-16 ***
## poly(indus, 3)2  -24.395     7.423  -3.286  0.00109 **
## poly(indus, 3)3  -54.130     7.423  -7.292  1.2e-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

fit.nox2 <- lm(crim ~ poly(nox, 3))
summary(fit.nox2)

##
## Call:
## lm(formula = crim ~ poly(nox, 3))
##
## 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)    3.6135     0.3216  11.237 < 2e-16 ***
## poly(nox, 3)1  81.3720     7.2336  11.249 < 2e-16 ***
## poly(nox, 3)2 -28.8286     7.2336  -3.985 7.74e-05 ***
## poly(nox, 3)3 -60.3619     7.2336  -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

fit.rm2 <- lm(crim ~ poly(rm, 3))
summary(fit.rm2)

##
## Call:
## lm(formula = crim ~ poly(rm, 3))
##
## 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)    3.6135     0.3703   9.758 < 2e-16 ***
## poly(rm, 3)1 -42.3794     8.3297  -5.088 5.13e-07 ***
## poly(rm, 3)2  26.5768     8.3297   3.191 0.00151 **
## poly(rm, 3)3  -5.5103     8.3297  -0.662 0.50858
## ---
## 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

fit.age2 <- lm(crim ~ poly(age, 3))
summary(fit.age2)

##
## Call:
## lm(formula = crim ~ poly(age, 3))
##
## 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)    3.6135     0.3485  10.368 < 2e-16 ***
## poly(age, 3)1  68.1820     7.8397   8.697 < 2e-16 ***
## poly(age, 3)2  37.4845     7.8397   4.781 2.29e-06 ***
## poly(age, 3)3  21.3532     7.8397   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

fit.dis2 <- lm(crim ~ poly(dis, 3))
summary(fit.dis2)

##
## Call:
## lm(formula = crim ~ poly(dis, 3))
##
## 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)    3.6135     0.3259  11.087 < 2e-16 ***
## poly(dis, 3)1 -73.3886     7.3315 -10.010 < 2e-16 ***
## poly(dis, 3)2  56.3730     7.3315   7.689 7.87e-14 ***
## poly(dis, 3)3 -42.6219     7.3315  -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
```



```

fit.rad2 <- lm(crim ~ poly(rad, 3))
summary(fit.rad2)

##
## Call:
## lm(formula = crim ~ poly(rad, 3))
##
## 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)    3.6135     0.2971  12.164 < 2e-16 ***
## poly(rad, 3)1 120.9074     6.6824  18.093 < 2e-16 ***
## poly(rad, 3)2  17.4923     6.6824   2.618  0.00912 **
## poly(rad, 3)3   4.6985     6.6824   0.703  0.48231
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## 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

fit.tax2 <- lm(crim ~ poly(tax, 3))
summary(fit.tax2)

##
## Call:
## lm(formula = crim ~ poly(tax, 3))
##
## 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)    3.6135     0.3047  11.860 < 2e-16 ***
## poly(tax, 3)1 112.6458     6.8537  16.436 < 2e-16 ***
## poly(tax, 3)2  32.0873     6.8537   4.682 3.67e-06 ***
## poly(tax, 3)3  -7.9968     6.8537  -1.167  0.244
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## 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

fit.ptratio2 <- lm(crim ~ poly(ptratio, 3))
summary(fit.ptratio2)

```

```
##
## Call:
## lm(formula = crim ~ poly(ptratio, 3))
##
## 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)      3.614      0.361  10.008 < 2e-16 ***
## poly(ptratio, 3)1  56.045      8.122   6.901 1.57e-11 ***
## poly(ptratio, 3)2  24.775      8.122   3.050  0.00241 **
## poly(ptratio, 3)3 -22.280      8.122  -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

fit.black2 <- lm(crim ~ poly(black, 3))
summary(fit.black2)

##
## Call:
## lm(formula = crim ~ poly(black, 3))
##
## 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)      3.6135      0.3536  10.218 <2e-16 ***
## poly(black, 3)1 -74.4312      7.9546  -9.357 <2e-16 ***
## poly(black, 3)2   5.9264      7.9546   0.745  0.457
## poly(black, 3)3  -4.8346      7.9546  -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

fit.lstat2 <- lm(crim ~ poly(lstat, 3))
summary(fit.lstat2)

##
## Call:
## lm(formula = crim ~ poly(lstat, 3))
```

```
##
## 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)      3.6135     0.3392  10.654 <2e-16 ***
## poly(lstat, 3)1  88.0697     7.6294  11.543 <2e-16 ***
## poly(lstat, 3)2  15.8882     7.6294   2.082  0.0378 *
## poly(lstat, 3)3 -11.5740     7.6294  -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

fit.medv2 <- lm(crim ~ poly(medv, 3))
summary(fit.medv2)

##
## Call:
## lm(formula = crim ~ poly(medv, 3))
##
## 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)      3.614     0.292  12.374 < 2e-16 ***
## poly(medv, 3)1  -75.058     6.569 -11.426 < 2e-16 ***
## poly(medv, 3)2   88.086     6.569  13.409 < 2e-16 ***
## poly(medv, 3)3  -48.033     6.569  -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
```

The p-value for predictors zn, rm, rad, tax and lstat prove that the cubic coefficient isn't statistically significant.

The p-values for predictors indus, nox, age, dis, ptratio and medv prove the the cubic coefficients are statistically significant.

The p-values for the predictor black prove that quadratic and cubic coefficients aren't statistically significant.

```
nullmodel=lm(medv~1, data=Boston)
fullmodel=lm(medv~., data=Boston)
model.step = step(nullmodel, scope=list(lower=nullmodel, upper=fullmodel),
direction='both')
```

```
## Start: AIC=2246.51
```

```
## medv ~ 1
```

```
##
```

	Df	Sum of Sq	RSS	AIC
## + lstat	1	23243.9	19472	1851.0
## + rm	1	20654.4	22062	1914.2
## + ptratio	1	11014.3	31702	2097.6
## + indus	1	9995.2	32721	2113.6
## + tax	1	9377.3	33339	2123.1
## + nox	1	7800.1	34916	2146.5
## + crim	1	6440.8	36276	2165.8
## + rad	1	6221.1	36495	2168.9
## + age	1	6069.8	36647	2171.0
## + zn	1	5549.7	37167	2178.1
## + black	1	4749.9	37966	2188.9
## + dis	1	2668.2	40048	2215.9
## + chas	1	1312.1	41404	2232.7
## <none>			42716	2246.5

```
##
```

```
## Step: AIC=1851.01
```

```
## medv ~ lstat
```

```
##
```

	Df	Sum of Sq	RSS	AIC
## + rm	1	4033.1	15439	1735.6
## + ptratio	1	2670.1	16802	1778.4
## + chas	1	786.3	18686	1832.2
## + dis	1	772.4	18700	1832.5
## + age	1	304.3	19168	1845.0
## + tax	1	274.4	19198	1845.8
## + black	1	198.3	19274	1847.8
## + zn	1	160.3	19312	1848.8
## + crim	1	146.9	19325	1849.2

```

## + indus      1      98.7 19374 1850.4
## <none>                19472 1851.0
## + rad        1      25.1 19447 1852.4
## + nox         1       4.8 19468 1852.9
## - lstat      1  23243.9 42716 2246.5
##
## Step:  AIC=1735.58
## medv ~ lstat + rm
##
##           Df Sum of Sq  RSS    AIC
## + ptratio  1    1711.3 13728 1678.1
## + chas     1     548.5 14891 1719.3
## + black    1     512.3 14927 1720.5
## + tax      1     425.2 15014 1723.5
## + dis      1     351.2 15088 1725.9
## + crim     1     311.4 15128 1727.3
## + rad      1     180.5 15259 1731.6
## + indus    1      61.1 15378 1735.6
## <none>                15439 1735.6
## + zn       1      56.6 15383 1735.7
## + age      1      20.2 15419 1736.9
## + nox      1      14.9 15424 1737.1
## - rm       1    4033.1 19472 1851.0
## - lstat    1    6622.6 22062 1914.2
##
## Step:  AIC=1678.13
## medv ~ lstat + rm + ptratio
##
##           Df Sum of Sq  RSS    AIC
## + dis      1     499.1 13229 1661.4
## + black    1     389.7 13338 1665.6
## + chas     1     378.0 13350 1666.0
## + crim     1     122.5 13606 1675.6
## + age      1      66.2 13662 1677.7
## <none>                13728 1678.1
## + tax      1      44.4 13684 1678.5
## + nox      1      24.8 13703 1679.2
## + zn       1      15.0 13713 1679.6
## + rad      1       6.1 13722 1679.9
## + indus    1       0.8 13727 1680.1
## - ptratio  1    1711.3 15439 1735.6
## - rm       1    3074.3 16802 1778.4
## - lstat    1    5013.6 18742 1833.7
##
## Step:  AIC=1661.39
## medv ~ lstat + rm + ptratio + dis
##
##           Df Sum of Sq  RSS    AIC
## + nox      1      759.6 12469 1633.5
## + black    1      502.6 12726 1643.8

```

```

## + chas      1      267.4 12962 1653.1
## + indus     1      242.6 12986 1654.0
## + tax       1      240.3 12989 1654.1
## + crim      1      233.5 12995 1654.4
## + zn        1      144.8 13084 1657.8
## + age       1       61.4 13168 1661.0
## <none>             13229 1661.4
## + rad       1       22.4 13206 1662.5
## - dis       1      499.1 13728 1678.1
## - ptratio   1     1859.3 15088 1725.9
## - rm        1     2622.6 15852 1750.9
## - lstat     1     5349.2 18578 1831.2
##
## Step:  AIC=1633.47
## medv ~ lstat + rm + ptratio + dis + nox
##
##           Df Sum of Sq  RSS    AIC
## + chas     1      328.3 12141 1622.0
## + black     1      311.8 12158 1622.7
## + zn        1      151.7 12318 1629.3
## + crim      1      141.4 12328 1629.7
## + rad       1       53.5 12416 1633.3
## <none>             12469 1633.5
## + indus     1       17.1 12452 1634.8
## + tax       1       10.5 12459 1635.0
## + age       1        0.2 12469 1635.5
## - nox       1      759.6 13229 1661.4
## - dis       1     1233.8 13703 1679.2
## - ptratio   1     2116.5 14586 1710.8
## - rm        1     2546.2 15016 1725.5
## - lstat     1     3664.3 16134 1761.8
##
## Step:  AIC=1621.97
## medv ~ lstat + rm + ptratio + dis + nox + chas
##
##           Df Sum of Sq  RSS    AIC
## + black     1      272.8 11868 1612.5
## + zn        1      164.4 11977 1617.1
## + crim      1      116.3 12025 1619.1
## + rad       1       58.6 12082 1621.5
## <none>             12141 1622.0
## + indus     1       26.3 12115 1622.9
## + tax       1        4.2 12137 1623.8
## + age       1        2.3 12139 1623.9
## - chas      1      328.3 12469 1633.5
## - nox       1     820.4 12962 1653.1
## - dis       1    1146.8 13288 1665.6
## - ptratio   1    1924.9 14066 1694.4
## - rm        1    2480.7 14622 1714.0
## - lstat     1    3509.3 15650 1748.5

```

```
##
## Step: AIC=1612.47
## medv ~ lstat + rm + ptratio + dis + nox + chas + black
##
##           Df Sum of Sq  RSS    AIC
## + zn       1    189.94 11678 1606.3
## + rad       1    144.32 11724 1608.3
## + crim      1     55.63 11813 1612.1
## <none>             11868 1612.5
## + indus     1     15.58 11853 1613.8
## + age       1      9.45 11859 1614.1
## + tax       1      2.70 11866 1614.4
## - black     1    272.84 12141 1622.0
## - chas      1    289.27 12158 1622.7
## - nox       1    626.85 12495 1636.5
## - dis       1   1103.33 12972 1655.5
## - ptratio   1   1804.30 13672 1682.1
## - rm        1   2658.21 14526 1712.7
## - lstat     1   2991.55 14860 1724.2
##
## Step: AIC=1606.31
## medv ~ lstat + rm + ptratio + dis + nox + chas + black + zn
##
##           Df Sum of Sq  RSS    AIC
## + crim      1     94.71 11584 1604.2
## + rad       1     93.61 11585 1604.2
## <none>             11678 1606.3
## + indus     1     16.05 11662 1607.6
## + tax       1      3.95 11674 1608.1
## + age       1      1.49 11677 1608.2
## - zn       1    189.94 11868 1612.5
## - black     1    298.37 11977 1617.1
## - chas      1    300.42 11979 1617.2
## - nox       1    627.62 12306 1630.8
## - dis       1   1276.45 12955 1656.8
## - ptratio   1   1364.63 13043 1660.2
## - rm        1   2384.55 14063 1698.3
## - lstat     1   3052.50 14731 1721.8
##
## Step: AIC=1604.19
## medv ~ lstat + rm + ptratio + dis + nox + chas + black + zn +
##      crim
##
##           Df Sum of Sq  RSS    AIC
## + rad       1    228.60 11355 1596.1
## <none>             11584 1604.2
## + indus     1     15.77 11568 1605.5
## + age       1      2.47 11581 1606.1
## + tax       1      1.31 11582 1606.1
## - crim      1     94.71 11678 1606.3
```

```

## - black      1      222.18 11806 1611.8
## - zn         1      229.02 11813 1612.1
## - chas       1      284.34 11868 1614.5
## - nox        1      578.44 12162 1626.8
## - ptratio    1     1192.90 12776 1651.8
## - dis        1     1345.70 12929 1657.8
## - rm         1     2419.57 14003 1698.2
## - lstat      1     2753.42 14337 1710.1
##
## Step:  AIC=1596.1
## medv ~ lstat + rm + ptratio + dis + nox + chas + black + zn +
##      crim + rad
##
##           Df Sum of Sq  RSS    AIC
## + tax      1      273.62 11081 1585.8
## <none>                        11355 1596.1
## + indus    1       33.89 11321 1596.6
## + age      1        0.10 11355 1598.1
## - zn       1     171.14 11526 1601.7
## - rad      1     228.60 11584 1604.2
## - crim     1     229.70 11585 1604.2
## - chas     1     272.67 11628 1606.1
## - black    1     295.78 11651 1607.1
## - nox      1     785.16 12140 1627.9
## - dis      1    1341.37 12696 1650.6
## - ptratio  1    1419.77 12775 1653.7
## - rm       1    2182.57 13538 1683.1
## - lstat    1    2785.28 14140 1705.1
##
## Step:  AIC=1585.76
## medv ~ lstat + rm + ptratio + dis + nox + chas + black + zn +
##      crim + rad + tax
##
##           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

```



```

model1 = lm(medv~lstat + rm + ptratio + dis + nox + chas + black + zn + crim
+ rad + tax, data=Boston)
summary(model1)

##
## Call:
## lm(formula = medv ~ lstat + rm + ptratio + dis + nox + chas +
##      black + zn + crim + rad + tax, 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 ***
## lstat       -0.522553   0.047424  -11.019 < 2e-16 ***
## rm          3.801579   0.406316   9.356 < 2e-16 ***
## ptratio     -0.946525   0.129066   -7.334 9.24e-13 ***
## dis         -1.492711   0.185731   -8.037 6.84e-15 ***
## nox        -17.376023   3.535243   -4.915 1.21e-06 ***
## chas         2.718716   0.854240    3.183 0.001551 **
## black        0.009291   0.002674    3.475 0.000557 ***
## zn          0.045845   0.013523    3.390 0.000754 ***
## crim        -0.108413   0.032779   -3.307 0.001010 **
## rad          0.299608   0.063402    4.726 3.00e-06 ***
## tax         -0.011778   0.003372   -3.493 0.000521 ***
## ---
## 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

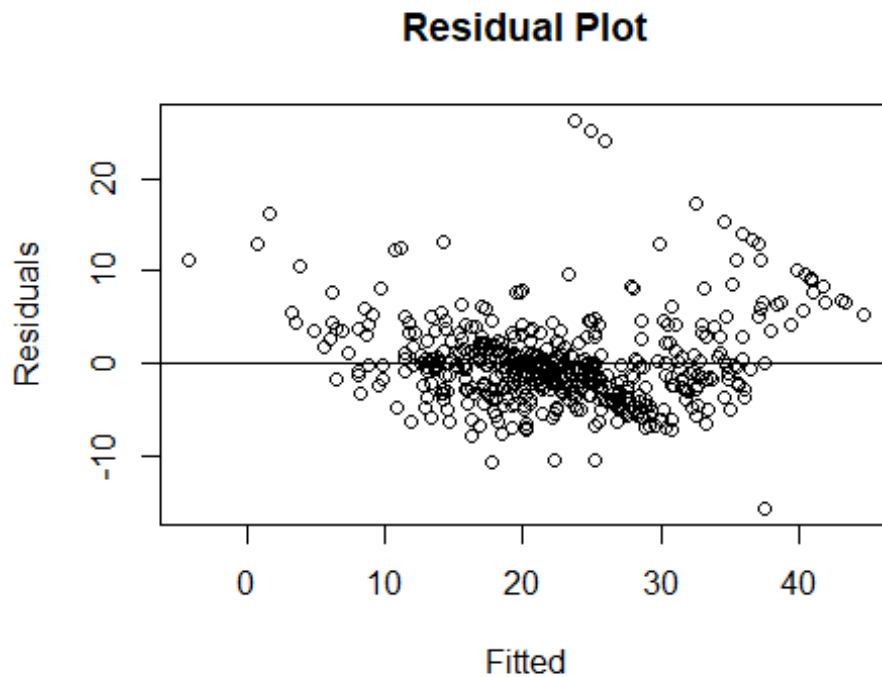
```

The stepwise model selects a model by automatically adding or removing predictors based upon the statistical significance. We usually end up with a single model whereas in multiple regression we tend to compare all possible models for a list of predictors and the model that fits the best might contain one or multiple predictors. This ends up with us observing a number of models and their corresponding summary characteristics.

```

plot(model1$fitted.values, model1$res, xlab="Fitted", ylab="Residuals",
main="Residual Plot")
abline(h=0)

```



The stepwise model gives us more accurate results than the multiple regression model. The plot for residuals v/s fitted for the stepwise model are pretty symmetrically distributed and tend to cluster around the  $y=0$  line and we don't observe any clear patterns. The multiple regression model is highly skewed and its residual plot shows the presence of outliers which cause the model to give an incorrect output due to the presence of outliers