

Report

Atul Lanka

October 14th, 2016

Abstract

This report will replicate the main results displayed in section 3.2 **Multiple Linear Regression** (chapter 3) of the book *An Introduction to Statistical Learning*.

Introduction

The primary goal of this analysis is to give advice on how to boost sales of the product given the current information on advertising budgets. More specifically, the idea is to determine whether there exists a correlation between advertising and sales, and if so, formulate an accurate model that can be used to predict sales from media budget. For this analysis in particular, a combination of simple linear regression and multiple linear regression.

Data

The Advertising data set comprises of the Sales (in thousands of units) in 200 different markets, along with the advertising budgets (in thousands of dollars) in each market for three different forms of media: TV, Radio, and Newspaper. In this report, the relation between each of them and Sales, and the possible relation between Sales and the three of them combined are observed and studied/analyzed.

We may first look at the table of summary statistics below:

| | TV | Radio | Newspaper | Sales |
|---|----------------|----------------|----------------|---------------|
| 1 | Min. : 0.70 | Min. : 0.000 | Min. : 0.30 | Min. : 1.60 |
| 2 | 1st Qu.: 74.38 | 1st Qu.: 9.975 | 1st Qu.: 12.75 | 1st Qu.:10.38 |
| 3 | Median :149.75 | Median :22.900 | Median : 25.75 | Median :12.90 |
| 4 | Mean :147.04 | Mean :23.264 | Mean : 30.55 | Mean :14.02 |
| 5 | 3rd Qu.:218.82 | 3rd Qu.:36.525 | 3rd Qu.: 45.10 | 3rd Qu.:17.40 |
| 6 | Max. :296.40 | Max. :49.600 | Max. :114.00 | Max. :27.00 |

Table 1: Summary Statistics

Histograms for each variable:

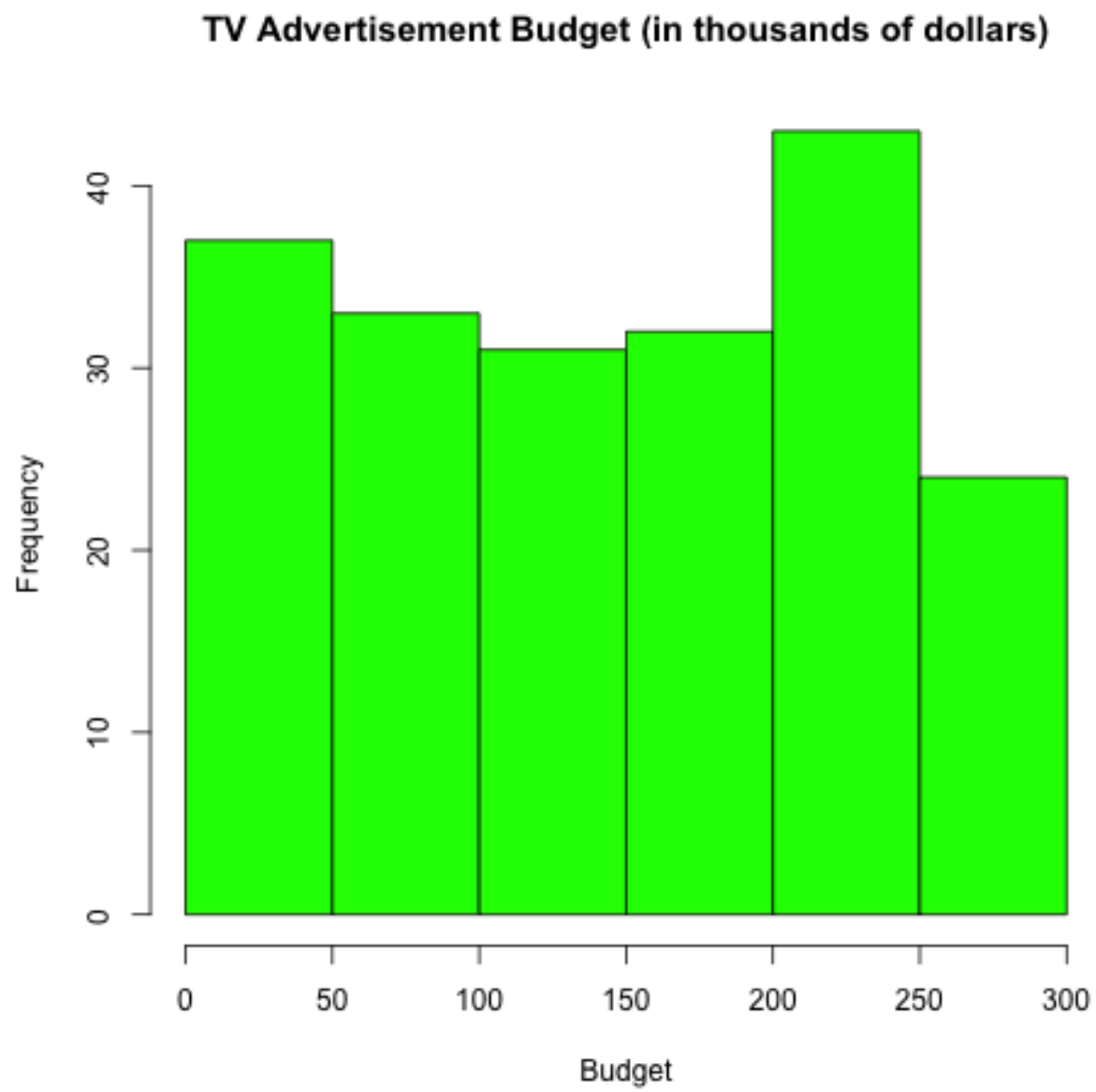


Figure 1: Figure 1: Histogram for TV

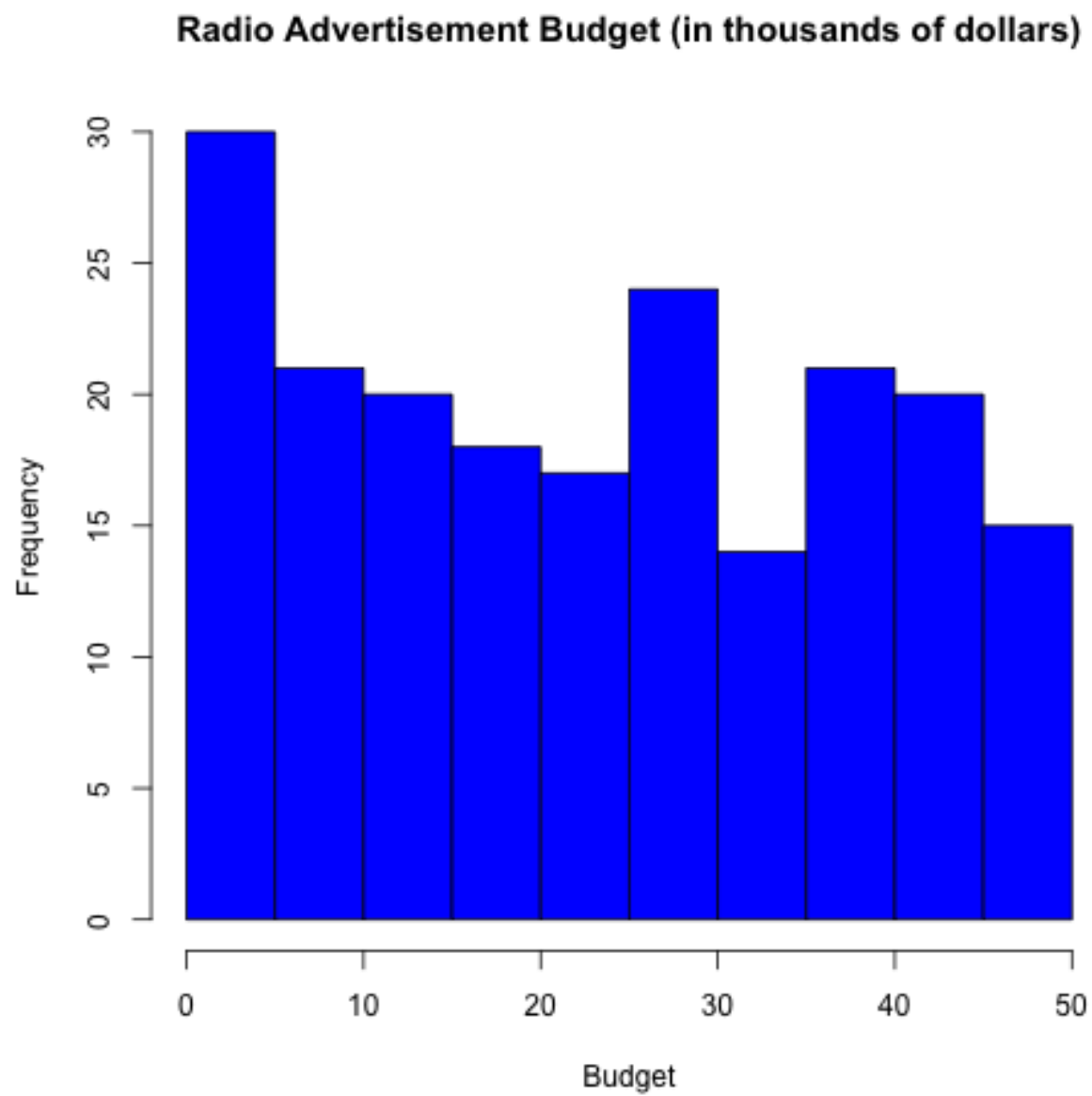


Figure 2: Figure 2: Histogram for Radio

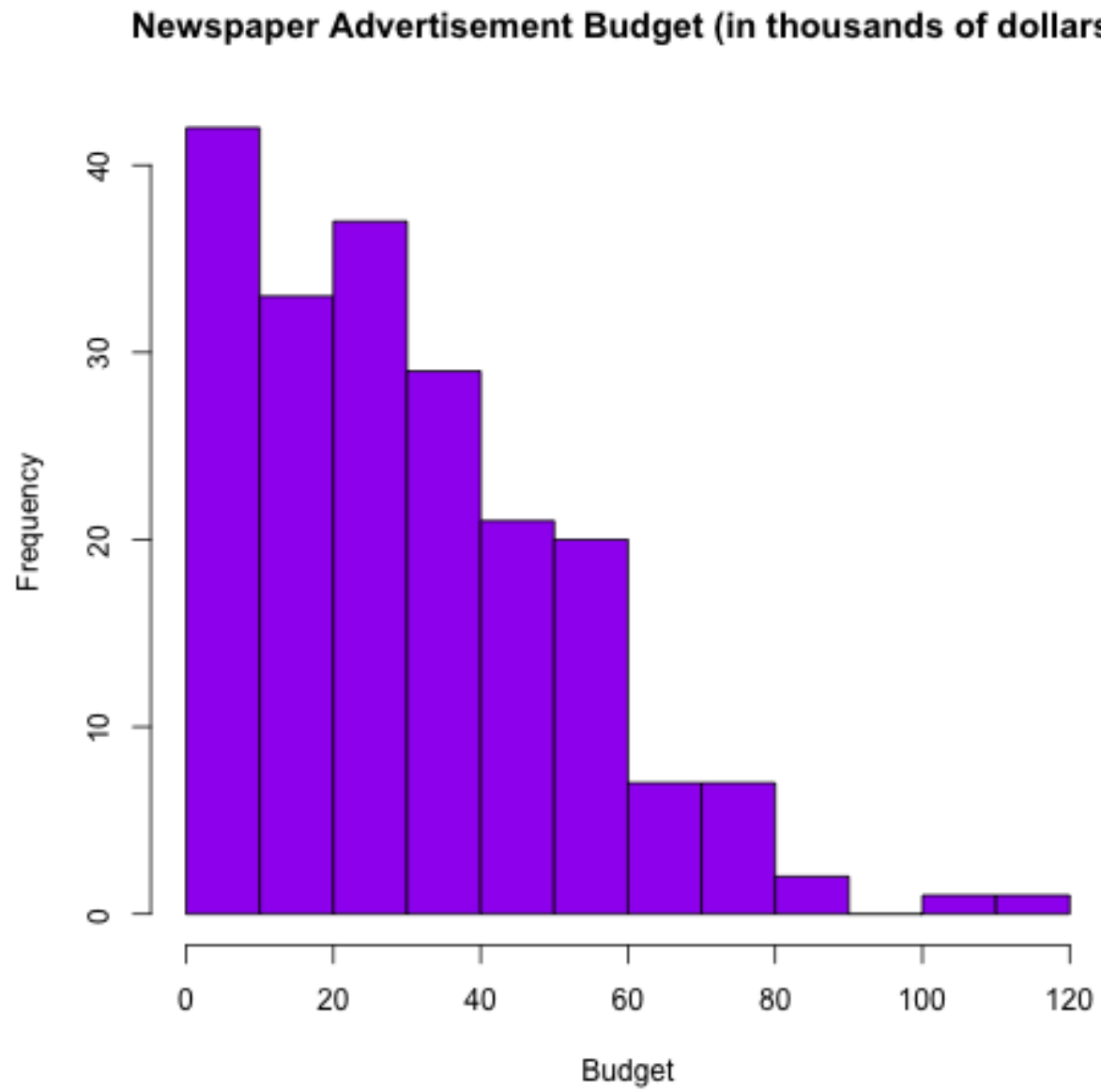


Figure 3: Figure 3: Histogram for Newspaper

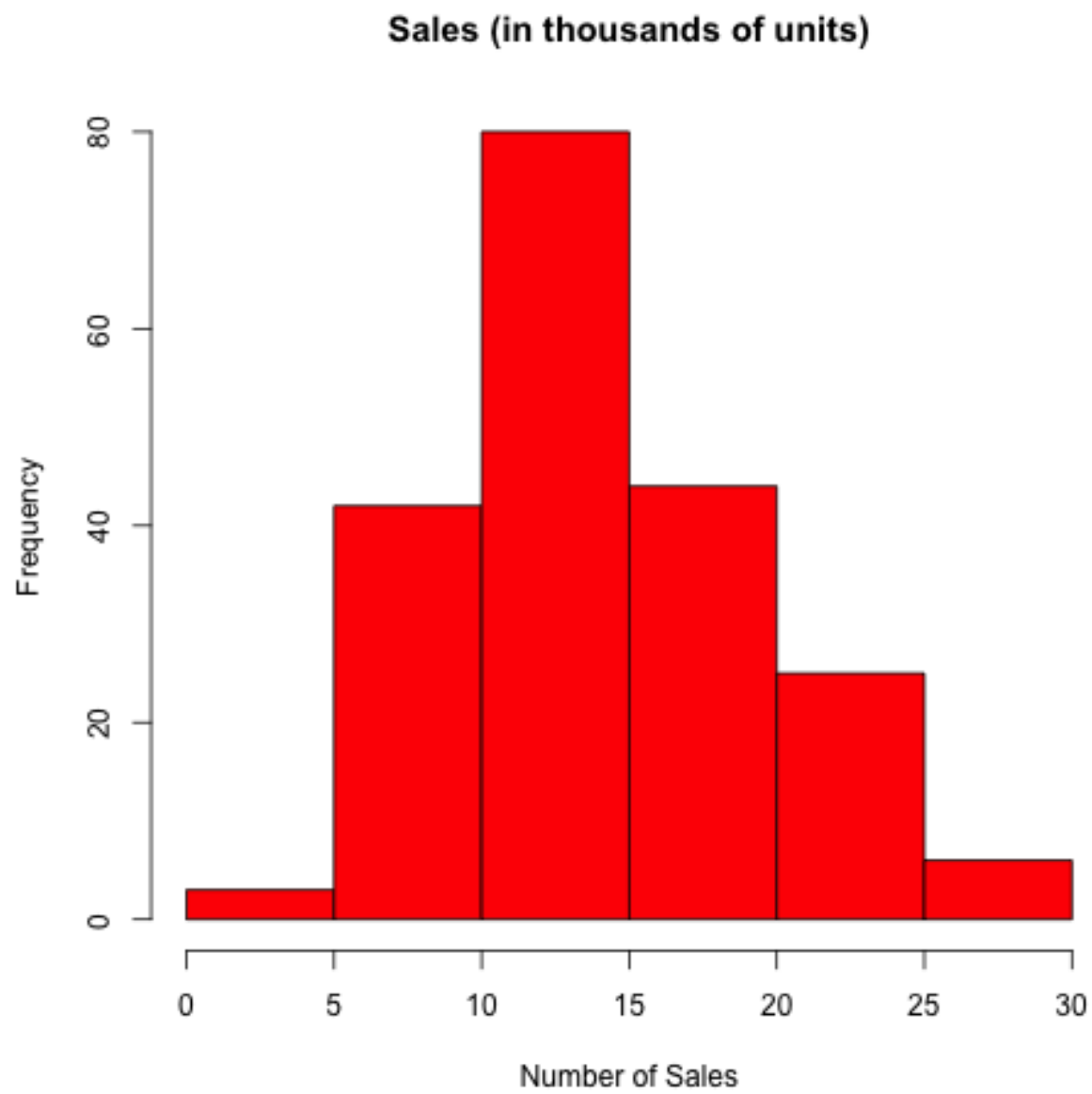


Figure 4: Figure 4: Histogram for Sales

Methodology

Single Linear Regression

We consider each media separately from the data set - TV, Radio and Newspaper - and study its relationship with the dependent variable Sales. The null hypothesis is that each of the independent variables would not have an effect on Sales, and the alternate hypothesis suggests otherwise. Thus a linear model is generated:

$$Sales = \beta_0 + \beta_1(TV/Radio/Newspaper)$$

Multiple Linear Regression

It would be a better approach to expand the model of Sales with multiple predictors rather than single predictors each time, avoiding an excess of linear models. This is done by accomodating each predictor with their respective slope coefficient in a single model. Thus the multiple linear regression model takes the form:

$$Sales = \beta_0 + \beta_1TV + \beta_2Radio + \beta_3Newspaper$$

Results

Single Linear Regression

Using `lm()` to fit the data into a simple linear model, the regression coefficients are as follows:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 7.0326 | 0.4578 | 15.36 | 0.0000 |
| advert\$TV | 0.0475 | 0.0027 | 17.67 | 0.0000 |

Table 2: Regression Coefficients for TV Sales

| | Estimate | Std. Error | t value | Pr(> t) |
|---------------|----------|------------|---------|----------|
| (Intercept) | 9.3116 | 0.5629 | 16.54 | 0.0000 |
| advert\$Radio | 0.2025 | 0.0204 | 9.92 | 0.0000 |

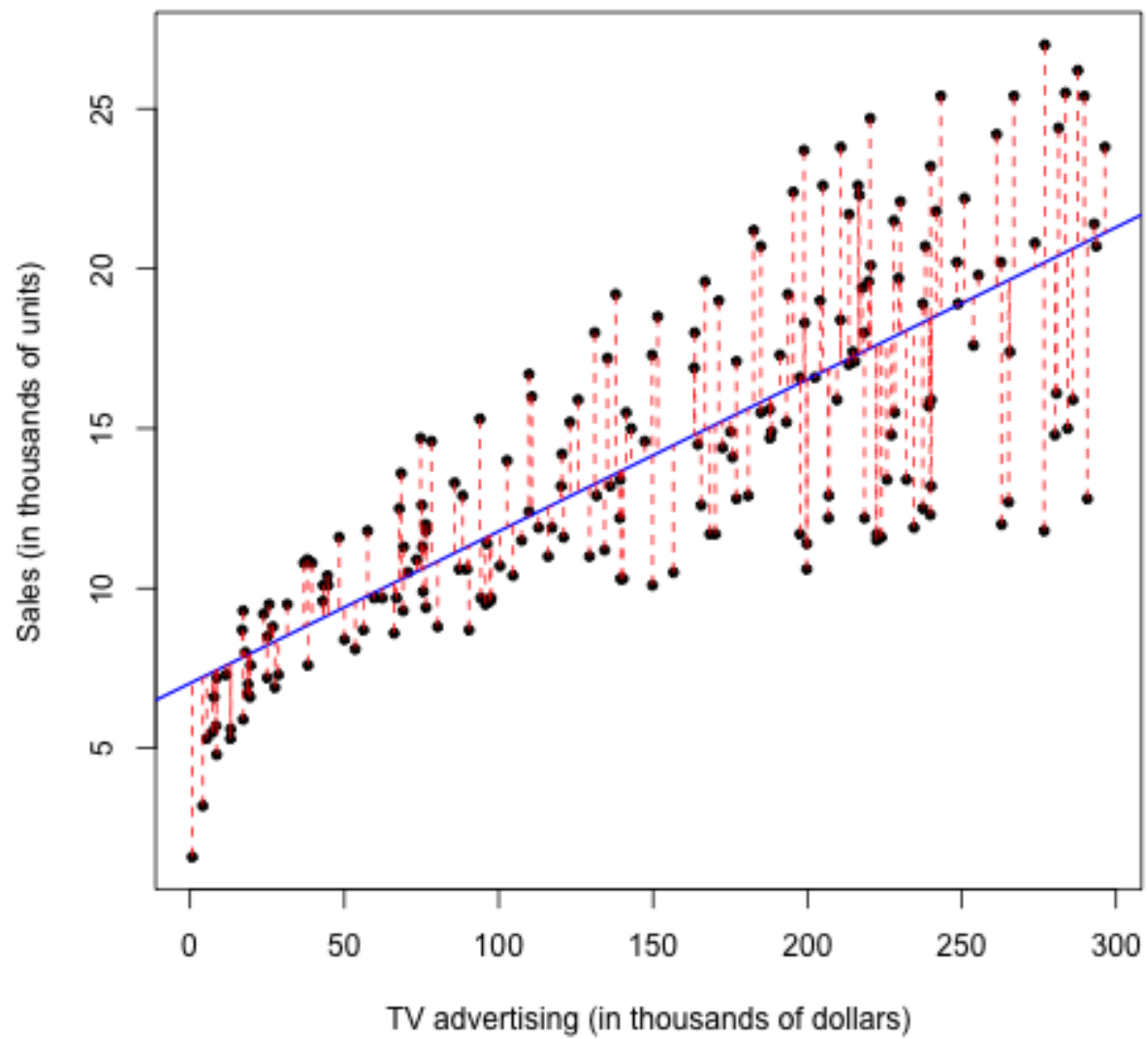
Table 3: Regression Coefficients for TV Radio

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------------|----------|------------|---------|----------|
| (Intercept) | 12.3514 | 0.6214 | 19.88 | 0.0000 |
| advert\$Newspaper | 0.0547 | 0.0166 | 3.30 | 0.0011 |

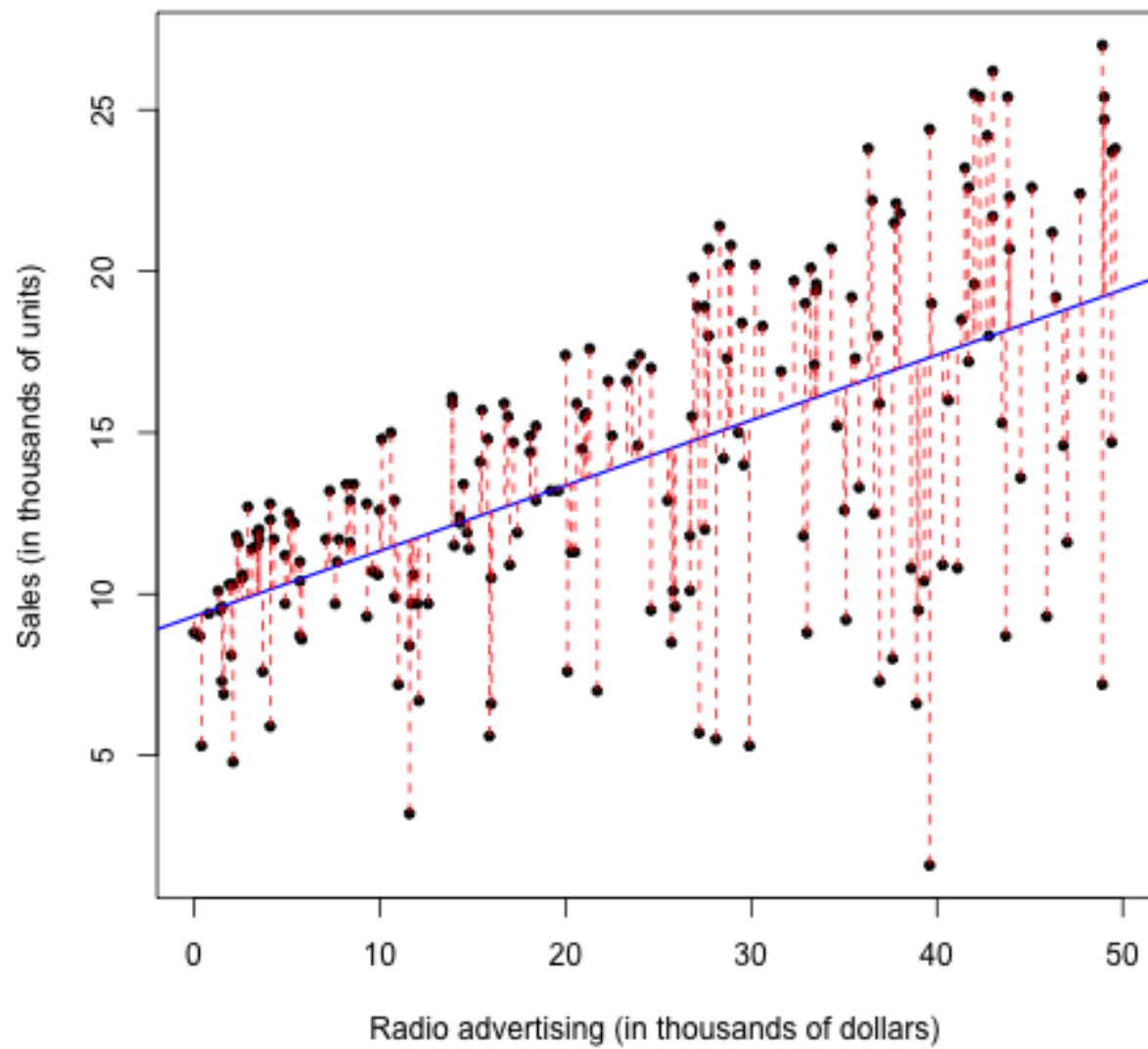
Table 4: Regression Coefficients for TV Newspaper

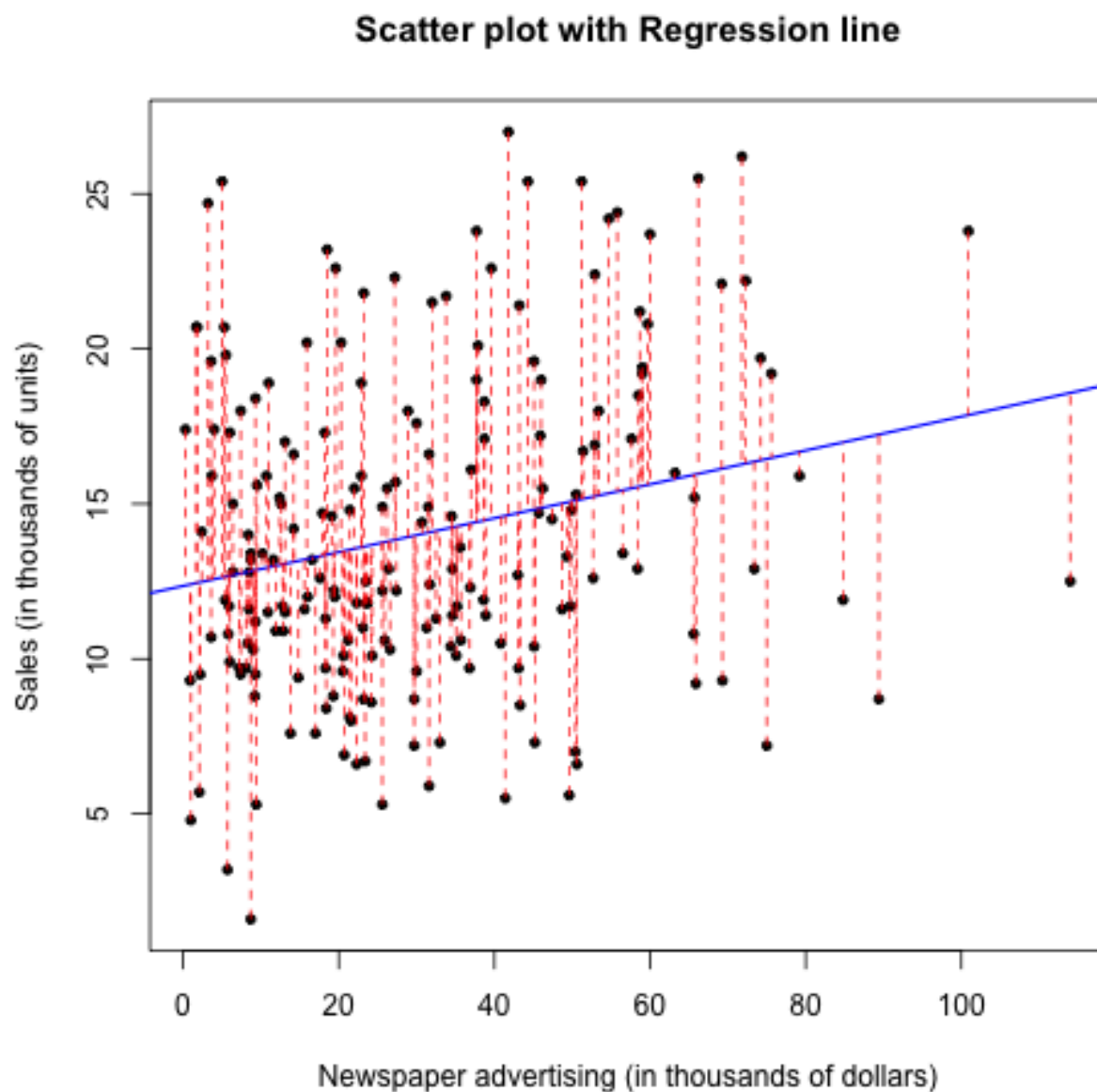
The scatterplots with their respective regression lines are as follows:

Scatter plot with Regression line



Scatter plot with Regression line





Quality Indices calculated with the regression analysis are as follows:

| | Quantity | Value |
|---|-------------|-------------------|
| 1 | RSE | 3.25865636865046 |
| 2 | R^2 | 0.611875050850071 |
| 3 | F-statistic | 312.144994372713 |

Table 5: Quality Indices of Regression of Sales and TV

In the case of TV, the p-value is below 0.05 and the RSE and R^2 values are relatively small. This is reason enough to reject the null hypothesis and establish that TV advertising does impact Sales in a positive, linear correlation (as can be seen in the graph).

Radio advertising budgets too have a similar impact on Sales, although from comparing the two scatter plots, Radio's linear regression is not as accurate as that of TV's.

| | Quantity | Value |
|---|-------------|-------------------|
| 1 | RSE | 4.27494435490106 |
| 2 | R^2 | 0.332032455445295 |
| 3 | F-statistic | 98.4215875667957 |

Table 6: Quality Indices of Regression of Sales and Radio

| | 1 | 2 |
|---|-------------|--------------------|
| 1 | RSE | 5.09248036652019 |
| 2 | R^2 | 0.0521204454443047 |
| 3 | F-statistic | 10.8872990754713 |

Table 7: Quality Indices of Regression of Sales and Newspaper

The same can be said about Newspaper, for the most part. The RSE value is large than that of TV, suggesting that the linear model is not the best fit necessarily. Nevertheless, the low p-value (statistical significance) allows the the null hypothesis to be rejected.

Multiple Regression Analysis

Fitting all 4 variables to a simple linear regression model, the regression coefficients are calculated:

| | Estimate | Std. Error | t value | $\Pr(> t)$ |
|-------------------|----------|------------|---------|-------------|
| (Intercept) | 2.9389 | 0.3119 | 9.42 | 0.0000 |
| advert\$TV | 0.0458 | 0.0014 | 32.81 | 0.0000 |
| advert\$Radio | 0.1885 | 0.0086 | 21.89 | 0.0000 |
| advert\$Newspaper | -0.0010 | 0.0059 | -0.18 | 0.8599 |

Table 8: Simple regression of Sales on Newspaper, TV and Radio

From this table we see that the p-value for Newspaper is not smaller than 0.05 and comparatively larger than TV and Radio. Hence, Newspaper advertising budget, in relation to Sales, is not statistically significant.

The correlation matrix is as follows:

Studying this correlation matrix, keeping in mind that Newspaper advertising is not statistically significant, it is observed that TV has low correlation with Radio and Newspaper, while having a high correlation value with TV. This suggests that TV has a strong impact on Sales while Radio and Newspaper perform differently but equally different. Newspaper has the lowest correlation value with Sales, indicating that Newspaper have the least impact on Sales.

Finally, the matrix of coefficients can be distributed on a scatter plot:

The plots reinforce the conclusions made for the correlation matrix.

The Quality Indices for the Multiple Linear Regression are:

The RSE and the R^2 value are low, significantly lower than the values computed in any of the 3 single linear regressions. This suggests that the multiple regression analysis create a linear model that is a relatively and comparatively good for the data.

| | TV | Radio | Newspaper | Sales |
|-----------|------|-------|-----------|-------|
| TV | 1.00 | 0.05 | 0.06 | 0.78 |
| Radio | 0.05 | 1.00 | 0.35 | 0.58 |
| Newspaper | 0.06 | 0.35 | 1.00 | 0.23 |
| Sales | 0.78 | 0.58 | 0.23 | 1.00 |

Table 9: Correlation Matrix for Sales on Newspaper, TV and Radio

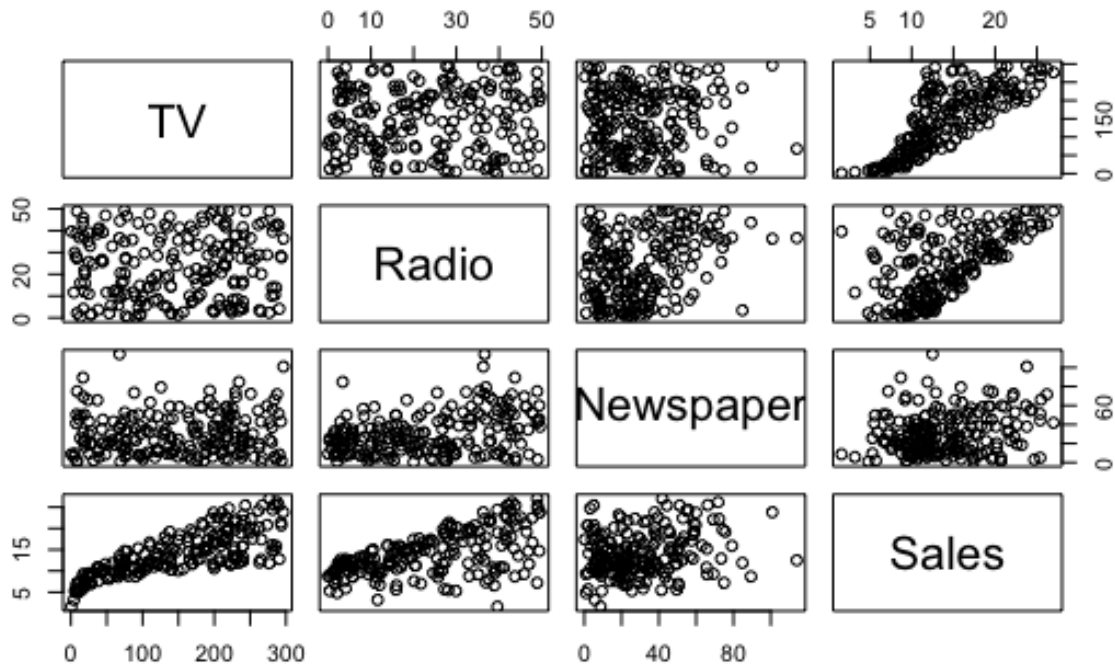
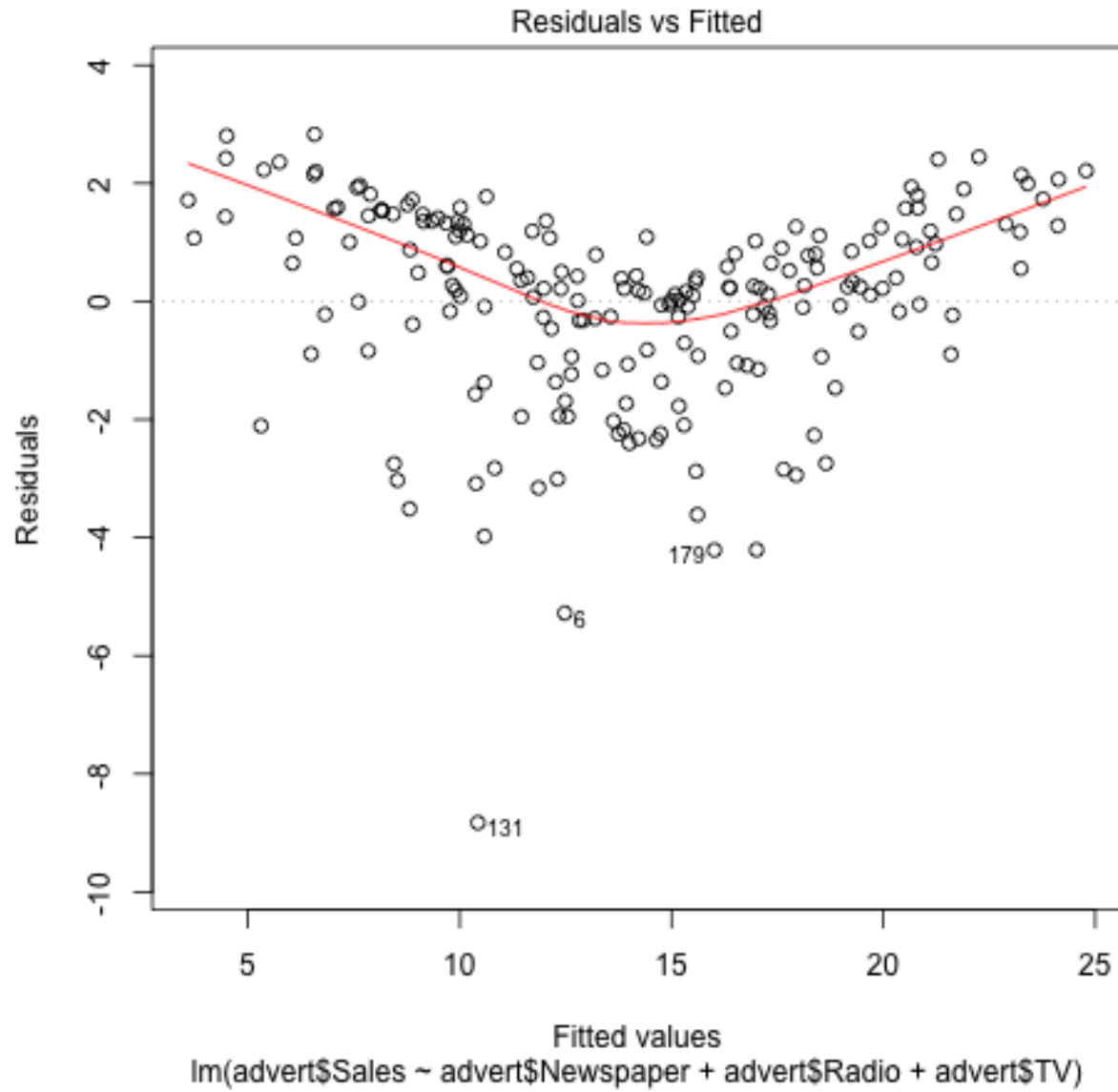


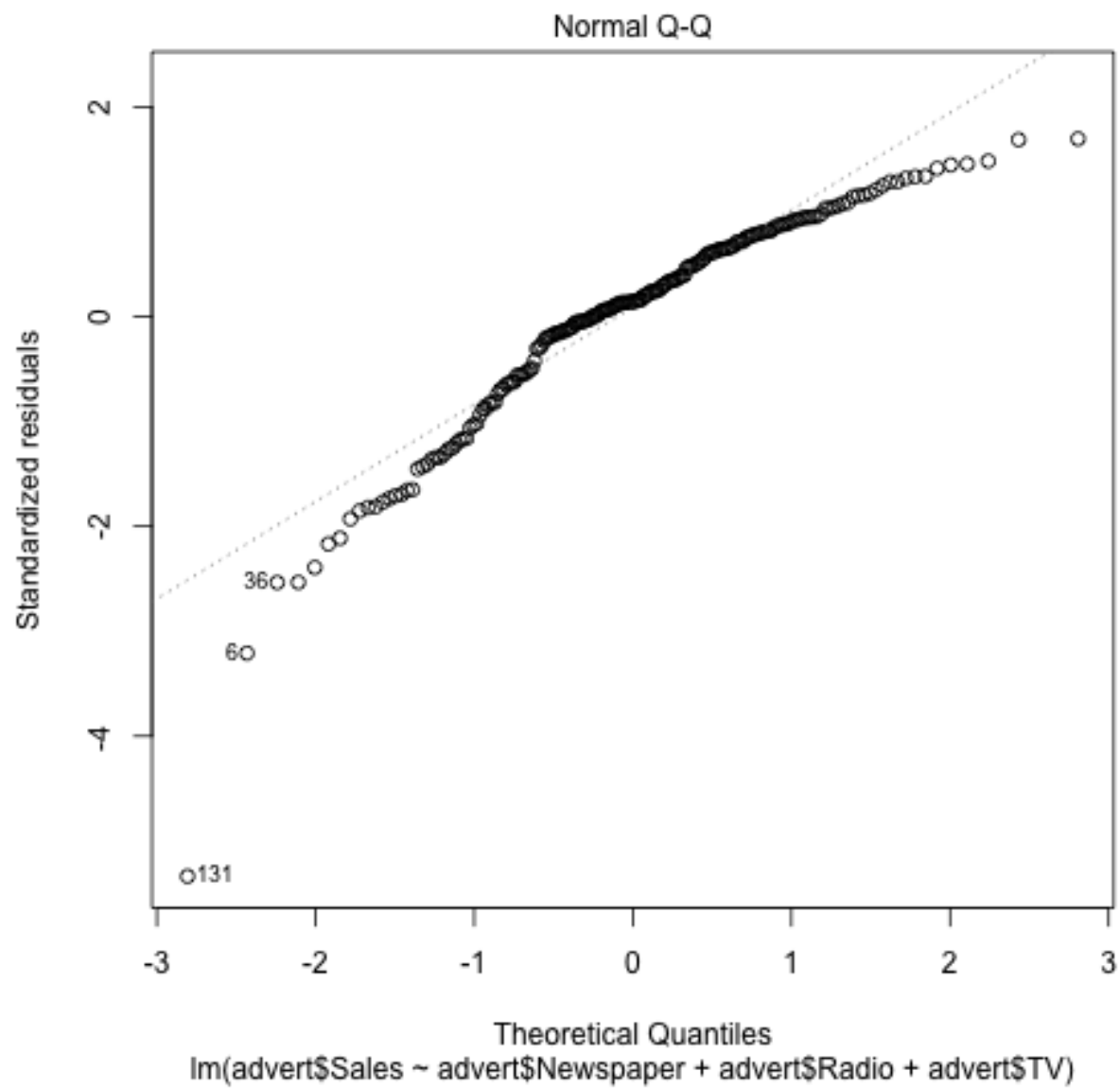
Figure 5: Figure 8: Scatterplot Matrix

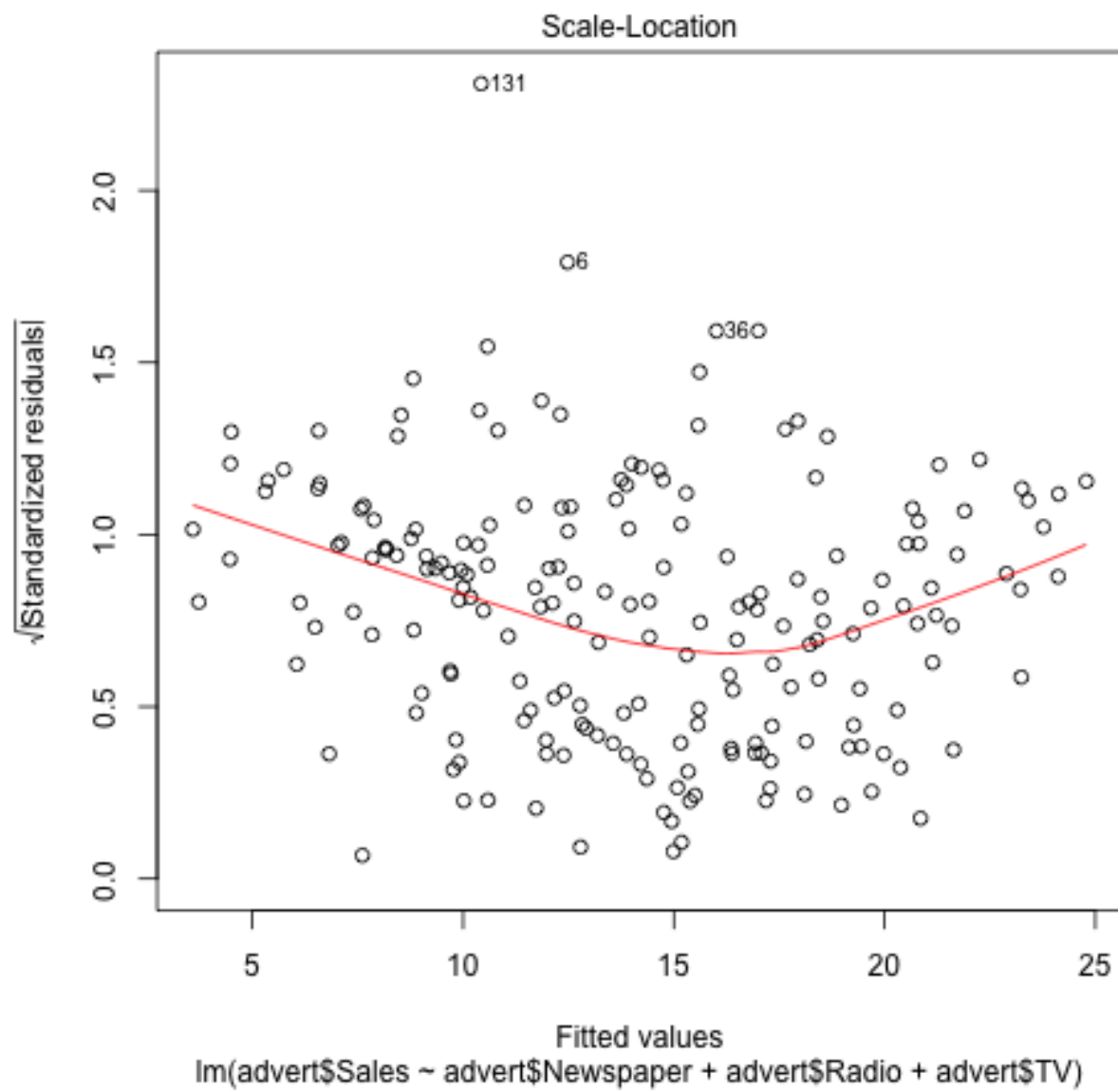
3 more plots can represent the relationship between Sales and the 3 media outlets:

| | 1 | 2 |
|---|-------------|-------------------|
| 1 | RSE | 1.68551037341474 |
| 2 | R^2 | 0.897210638178952 |
| 3 | F-statistic | 570.270703659094 |

Table 10: Quality Indices of the Multiple Linear Regression of Sales and Newspaper, TV and Radio







Conclusions

In conclusion, the multiple linear regression is more accurate than the single linear regressions. From the tables above (**F-statistic**) from multiple regression, at the very least one of the predictors can be used to predict **Sales**. Nevertheless, we also find that not all of the predictors are statistically significant (from the p-values). Hence, the prediction would be more accurate if the **newspaper** budget is not avoided based on its corresponding p-value. Other indicators, including Residual Standard Error, R^2 and F-statistic, also comment on the fit of the linear model - the smaller these values/statistics or approaching 0, the more accurately the model represents the data.