

Multiple Regression Analysis

Xiaoqian Zhu

Abstract

In this report we extend the scope of the previous HW that apply computational tools to reproduce the main results displayed in section 3.2 (page 71 to 82) of the book *An Introduction to Statistical Learning* by Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani. [Link to the book](#)

Introduction

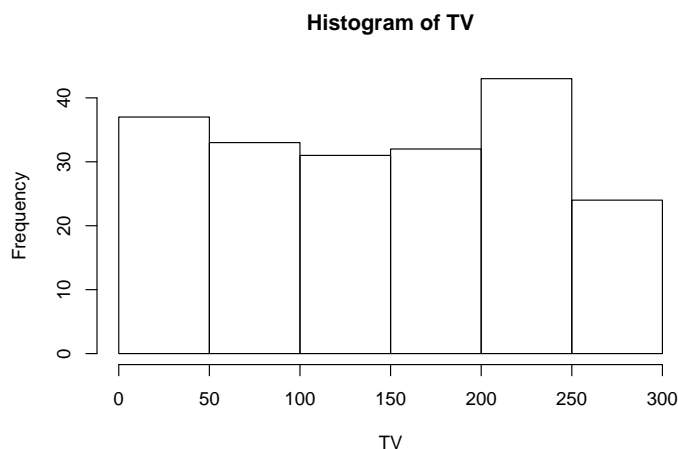
The overall goal is to provide advice on how to improve sales of the particular product. More specifically, the idea is to see whether there is an association between advertising on medium and sales, and if so, which medium would be the best to improve sales of the particular product. In this homework, we are specifically looking at how the advertisement on TV, radio and newspaper affect sales of the particular product; in other words, we are seeing advertising on TV, newspaper and radio as explanatory variables to explain the responses to sales. I will first check the association between each variable and sales; then, I am going to develop a multivariable regression model that can be used to estimate sales on the basis of advertising budget on TV, radio and newspaper using the advertising data set, which compiled by Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani.

Data

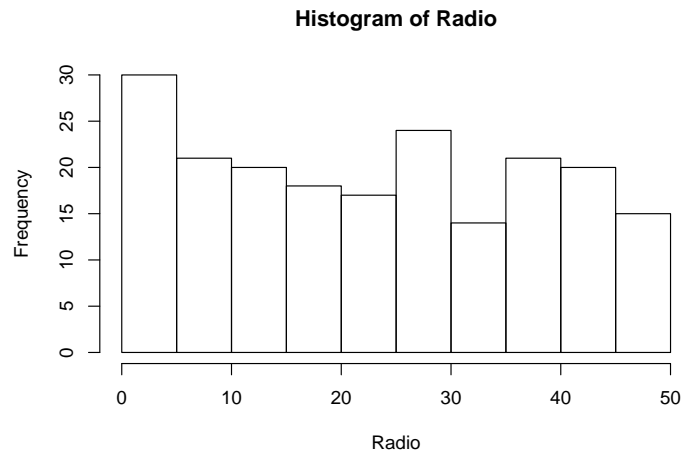
In this homework, we are working with Advertising data set ([Link for Advertising.csv](#)), which consists of *Sales* (in thousands of units) of a particular product in 200 different markets, along with advertising budgets (in thousands of dollars) for the product in each of those markets for three different media: *TV*, *Radio*, and *Newspaper*. In this homework, we are working on the three media and *sales* to find out whether they have an association.

We can easily see how *TV*, *Radio*, *Newspaper*, and *sales* distribute through histograms below:

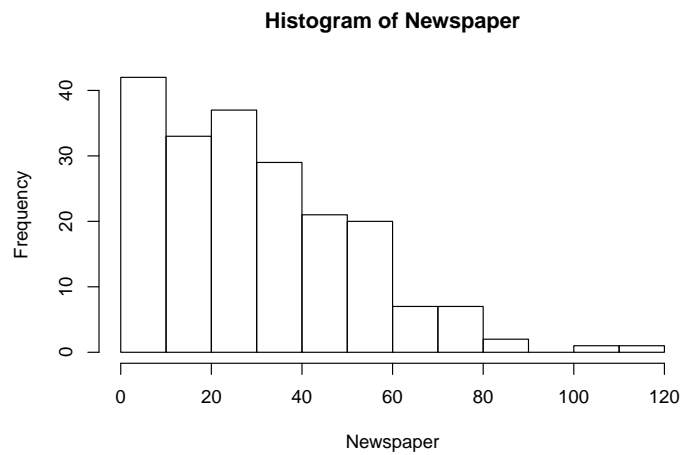
```
advertising <- read.csv('~/.stat159-hw/stat159-fall2016-hw03/data/Advertising.csv', header = TRUE)
hist(advertising$TV, xlab="TV", main="Histogram of TV")
```



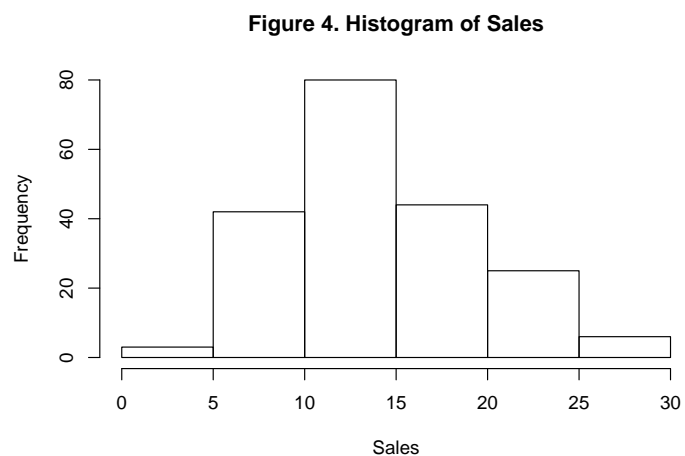
```
advertising <- read.csv('~/.stat159-hw/stat159-fall2016-hw03/data/Advertising.csv', header = TRUE)
hist(advertising$Radio, xlab="Radio", main="Histogram of Radio")
```



```
advertising <- read.csv('~/.stat159-hw/stat159-fall2016-hw03/data/Advertising.csv', header = TRUE)
hist(advertising$Newspaper, xlab="Newspaper", main="Histogram of Newspaper")
```



```
advertising <- read.csv('~/.stat159-hw/stat159-fall2016-hw03/data/Advertising.csv', header = TRUE)
hist(advertising$Sales, xlab="Sales", main="Figure 4. Histogram of Sales")
```



These histograms give us a broad view of how the Advertising dataset distribute.

Methodology

Looking at the data set, we are going to study the relationship between *Sales* and a group of variables, including, *TV*, *Radio*, and *Newspaper*. We start our analysis by setting up null and alternative hypothesis. The null hypothesis is, $H_0: \beta_1 = \beta_2 = \beta_3 = 0$, that there is no relationship between TV, radio and newspaper advertising budget and sales; the alternative hypothesis is, H_1 , that there exists one β_i that not equal to 0. To test the hypothesis, we apply a multiple linear regression model, $Sales = \beta_0 + \beta_1 TV + \beta_2 Radio + \beta_3 Newspaper + error$, on the data set, Advertising.csv, to get the estimates of β_0 , β_1 , β_2 , and β_3 , which are $\hat{\beta}_0$, $\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_3$. In this homework, we will also use the least squares model to our data for the regression analysis like the previous one.

Results

We first apply a simple linear regression on each explaining variable tables with Sales. Table 1 shows results from the simple single regression of sales on TV.

```
load("~/stat159-hw/stat159-fall2016-hw03/data/regression.RData")
library("xtable")
table_TV <- xtable(regression_sum_TV,
                   caption = 'Simple Regression of Sales on TV',
                   digits = 4)

row.names(table_TV) <- c('(Intercept)', 'TV')

print(table_TV, caption.placement = 'top', comment = getOption("xtable.comment", FALSE))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.0326	0.4578	15.3603	0.0000
TV	0.0475	0.0027	17.6676	0.0000

Table 2 shows the results from simple single regression of sales on radio.

```
load("~/stat159-hw/stat159-fall2016-hw03/data/regression.RData")
library("xtable")
table_Radio <- xtable(regression_sum_Radio, caption = 'Simple Regression of Sales on Radio', digits = 4)

row.names(table_Radio) <- c('(Intercept)', 'Radio')

print(table_Radio, caption.placement = 'top', comment = getOption("xtable.comment", FALSE))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.3116	0.5629	16.5422	0.0000
Radio	0.2025	0.0204	9.9208	0.0000

Table 3 shows the results from simple single regression of sales on newspaper.

```
load("~/stat159-hw/stat159-fall2016-hw03/data/regression.RData")
library("xtable")
table_Newspaper <- xtable(regression_sum_Newspaper, caption = 'Simple Regression of Sales on Newspaper',
                           digits = 4)

row.names(table_Radio) <- c('(Intercept)', 'Newspaper')
```

```
print(table_Newspaper, caption.placement = 'top',comment = getOption("xtable.comment", FALSE))
```

Table 3: Simple Regression of Sales on Newspaper

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	12.3514	0.6214	19.8761	0.0000
Newspaper	0.0547	0.0166	3.2996	0.0011

Then, we apply the a multiple linear regression of sales on TV, radio and regression. Table 4 shows the results from this multiple linear regression.

```
load("~/stat159-hw/stat159-fall2016-hw03/data/regression.RData")
library("xtable")
table_multi <- xtable(regression_sum_multi,caption = 'Multiple Regression Table',digits = 4)

row.names(table_multi) <- c('(Intercept)', 'TV', 'Radio', 'Newspaper')

print(table_multi, caption.placement = 'top',comment = getOption("xtable.comment", FALSE))
```

Table 4: Multiple Regression Table

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.9389	0.3119	9.4223	0.0000
TV	0.0458	0.0014	32.8086	0.0000
Radio	0.1885	0.0086	21.8935	0.0000
Newspaper	-0.0010	0.0059	-0.1767	0.8599

Moreover, in this homework, we also compute the correlation matrix among TV, radio, newspaper, and sales. Table 5 shows the correlation matrix for TV, radio, newspaper, and sales for the Advertising data.

```
load("~/stat159-hw/stat159-fall2016-hw03/data/correlation-matrix.RData")
library("xtable")
cor_matrix <- xtable(correlation_matrix, caption = 'Correlation matrix for TV, radio, newspaper, and sales',
print(cor_matrix, caption.placement = 'top',comment = getOption("xtable.comment", FALSE))
```

Table 5: Correlation matrix for TV, radio, newspaper, and sales for the Advertising data

	X	TV	Radio	Newspaper	Sales
X	1.0000	0.0177	-0.1107	-0.1549	-0.0516
TV	0.0177	1.0000	0.0548	0.0566	0.7822
Radio	-0.1107	0.0548	1.0000	0.3541	0.5762
Newspaper	-0.1549	0.0566	0.3541	1.0000	0.2283
Sales	-0.0516	0.7822	0.5762	0.2283	1.0000

When we perform multiple linear regression, we usually are intrested in answering a few important questions.

1. Is at least one of the predictors useful in predicting the response?
2. Do all predictors help to explain the response, or is only a subset of the predictors useful?
3. How well does the model fit the data?
4. How accurate is the prediction?

We now address each of these questions in turn.

One: Is at least one of the predictors useful in predicting the response?

Recall that in the simple linear regression setting, in order to determine whether there is a relationship between the response and the predictor we can simply check whether $\beta_0 = 0$. Now, we are performing a multiple linear regression with 3 predictors, thus we need to ask whether all of the regression coefficients are 0. In this case, whether $\beta_1 = \beta_2 = \beta_3 = 0$. As in the simple linear regression setting, we use a hypothesis test to answer this question. We test the null hypothesis versus the alternative (which I specifically listed in introduction). Such hypothesis test is performed by computing the F-statistic. When there is no relationship between the response and predictors, one would expect the F-statistic to take a value close to 1. On the other hand, if alternative hypothesis is true, then we expect F-statistic to be greater than 1.

```
load("~/stat159-hw/stat159-fall2016-hw03/data/regression.RData")
library(xtable)
Quantity <- c('RSS', 'R2', 'F-stat')
Value <- c(regression_sum_multi$sigma,
            regression_sum_multi$adj.r.squared,
            regression_sum_multi$fstatistic[1])
results2 <- data.frame(Quantity, Value)

table2 <- xtable(results2, caption = 'More information about the least squares model for the multiple r
print(table2, caption.placement = 'top', comment = getOption("xtable.comment", FALSE))
```

Table 6: More information about the least squares model for the multiple regression

	Quantity	Value
1	RSS	1.69
2	R2	0.90
3	F-stat	570.27

The The F-statistic for the multiple linear regression model obtained by regressing *sales* onto *radio*, *TV*, and *newspaper* is shown in Table 6. In this example, the F-statistic = 570. Since this is far larger than 1, it provides compelling evidence against the null hypothesis. In other words, the large F-statistic suggests that at least one of the advertising media must be related to *sales*.

Two: Do all predictors help to explain the response, or is only a subset of the predictors useful?

It is possible that all of the predictors are associated with the response, but it is more often the case that the response is only related to a subset of the predictors. Then, we would like to perform a variable selection that fit a single model involving each predictor to determine which predictors are associated with the response. However, variable selection in this example needs to run 8 tests, which is too tedious. Thus, we can look at the estimated coefficients in the multiple linear regression. The P-value tells whether the predictor has a statistically significant association with *sales*. If the p-value is small, then we would be able to believe that there is a relationship between the single predictor and *sales*. On the other hand, if the p-value is large, then we have reason to believe that there is no relationship between two. Looking at Tables of single linear regressions, we can see that P-values of the estimated coefficient of *TV* and *Radio* are both 0, which means that they are both statistically significant at 1% level, so we reject the null hypothesis that they do not have a relationship with *Sales*. However, p-value of *Newspaper* is too high, that we are not able to reject the null. Then, *Newspaper* fails to explain the response. In this case, only *TV* and *Radio* are the useful predictors.

Three: How well does the model fit the data?

Two of the most common numerical measures of model fit are RSE and R^2 . An R^2 values close to 1 indicates that the model explains a large portion of the variance in the response variable. We saw in Table 6, for the Advertising data set, the model that uses all there advertising media to predict sales has an R^2 value of 0.8972 .This is very close to 1, which means that our explanatory variables work well to explain changes in sales. Furthermore, RSE is the estimate of the standard deviation of error, and when it is small, it indicates a well-fitted model on data. In this multiple regression model, the RSE is 1.6855 .Thus, we can also tell from RSE that is a well-fitted model on advertising data set.

Four: How accurate is the prediction?

There are uncertainty associated with the estimated prediction of sales, to answer how accurate is the prediction, we can simply look at the standard error of each estimated coefficients in Table 4. The standard of TV, Radio and newspaper are all very small, which means that we can be confidently sure that the estimation of explanatory variables are close to the true model. Furthermore, the confidence can also be shown by R^2 like the previous question, that the R^2 value is close to 1, which means that the prediction model is well-fitted.

Conclusions

To be honest, it is a very long homework, that we try to recreate the study that was done in chapter 3.2 of *An Introduction of Statistical Learning* by Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani. We first extend our analysis in previous homework in order to accomodate two more predictors, *Radio* and *Newspaper* by running two more simple linear regression. Later, we perform a multiple linear regression for better prediction on sales. From the results, we can conclude that the estimated coefficients of *TV* and *Radio* are statistically significant, while *Newspaper* seems to do not have a relationship towards sales, according to p-value. \$1000 increase in TV advertising budget would lead to average increase in sales by 46 units, and \$1000 increase in Radio advertisng budget would lead to average increase in sales by 189 units. Such multiple linear regression has a R^2 close to 1 and a small RSE, thus, this model is well-fitted. The standard error of each variable shows that the estimation is close to the true model.