

# Introduction to Data Science

Ryan G

February 12, 2020

## Contents

<b>Simple Linear Regression</b>	<b>IntroDataSci</b>	<b>1</b>
Load Packages . . . . .		1
Question 1 . . . . .		2
.1 (a) Import the Data ATTACH . . . . .		2
.2 (b) Construct Scatter Plots . . . . .		3
.3 (c) Find the Correlation Coefficient LINEAR:REGRESSION . . . . .		12
.4 (d) Assess the accuracy of the parameter estimates MODELEVALUATION . .		13
Question 02 . . . . .		27
.1 (a) Upload the Auto Dataset and explore it. . . . .		27
.2 (b) Construct scatter plots to visualize the relationship between . . . . .		27
.3 Repeat the analysis in Q1 (c) to (i) using mpg and weight. . . . .		27
<b>Multiple Linear Regression</b>		<b>27</b>
Question 01 - Multiple Linear Regression . . . . .		28
.1 Load Packages . . . . .		28
.2 (a) Upload the data "Advertising.csv" and explore it. ATTACH . . . .		28
.3 . . . . .		30
.4 (b) Find the Covariance and Correlation Matrix of Sales, TV, Radio . . . . .		30
.5 (c) Construct the multiple linear regression model and find the . . . . .		32
.6 (e) Assess the overall accuracy of the model. . . . .		34
.7 (f) Calculate the predicted values and residuals . . . . .		36
.8 (g) Plot the residuals against the predicted values . . . . .		37
.9 (h) Plot the histogram of the residuals . . . . .		39
.10 (i) Comment on the residual plots . . . . .		44
.11 (j) Use the multivariate model for predictions . . . . .		46
Question 02: Non Linear Models: Use Advertising data set . . . . .		46
.1 (a) Add the Interaction Term TV*Radio and test the significance of . . . . .		46
.2 corplot . . . . .		47
.3 (b) Give the resulting model after considering this interaction . . . . .		49
.4 (c) Construct the Polynomial Regression Model of order 3 and test . . . . .		49
.5 (d) Give the resulting selected model . . . . .		53

Material of Tue 12 2019, week 2

## Load Packages

---

```
1  # Load Packages
2  if(require('pacman')){
3    library('pacman')
4  }else{
5    install.packages('pacman')
6    library('pacman')
7  }
8
9  pacman::p_load(caret, scales, ggplot2, rmarkdown, shiny, ISLR, class,
10 ↪ BiocManager,
11 ↪ corrplot, plotly, tidyverse, latex2exp, stringr,
12 ↪ reshape2, cowplot, ggpubr,
13 ↪ rstudioapi, wesanderson, RColorBrewer, colorspace,
14 ↪ gridExtra, grid, car,
15 ↪ boot, colourpicker, tree, ggtree, mise, rpart,
16 ↪ rpart.plot, knitr, MASS,
17 ↪ magrittr, EnvStats, tidyverse, tidyr, devtools, bookdown,
18 ↪ leaps, car, clipr,
19 ↪ tikzDevice, e1071)
20
21 mise()
22 set.seed(0932)
```

## Question 1

---

(a) Import the Data

ATTACH

```
1  setwd("/home/ryan/Dropbox/Notes/DataSci/IntroDataSci/Org-Babel/")
2  getwd()
3  adv <- read.csv(file =
4 ↪   "../data/83/4c42c3-8dd8-4fef-b402-7e341764d5e9/Advertising.csv",
5 ↪   header = TRUE, sep = ",")
```

## 1. Inspect the structure of the Data Set

```
1 head(adv)
```

```
##      TV Radio Newspaper Sales
## 1 230.1 37.8      69.2 22.1
## 2  44.5 39.3      45.1 10.4
## 3  17.2 45.9      69.3  9.3
## 4 151.5 41.3      58.5 18.5
## 5 180.8 10.8      58.4 12.9
## 6   8.7 48.9      75.0  7.2
```

outputoutputoutput

#valuevaluevalue+Bboth boutput

str(adv)

#+END\_SRC

#+BEGIN\_EXAMPLE

```
## 'data.frame':    200 obs. of  4 variables:
## $ TV          : num  230.1 44.5 17.2 151.5 180.8 ...
## $ Radio       : num  37.8 39.3 45.9 41.3 10.8 48.9 32.8 19.6 2.1 2.6 ...
## $ Newspaper: num  69.2 45.1 69.3 58.5 58.4 75 23.5 11.6 1 21.2 ...
## $ Sales      : num  22.1 10.4 9.3 18.5 12.9 7.2 11.8 13.2 4.8 10.6 ...
```

```
1 summary(adv)
```

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

### (b) Construct Scatter Plots

So I'd like to do a shiny ggplot here, however it's probably just as easy to use `tabset` by appending `{.tabset}` to the heading

```
1 # par(mfrow=c(2,2))
2 # plot(lm(y~x))
```

Multiple Plots may be fitted into one output using either the `par()` package or the `layout()` package, I personally prefer the `layout()` package, I think because in the past I had a bad experience with `par()`:

- In order to use `par`:

- `par( mfcol = c(ROW, COLS))`
- `par (mfcol = c(ROW, COLS))`
- \* That's not a typo, `mfrow` and `mfcol` are identical in this case

- In order to use `layout`:

- `layout(MATRIX)`
  - \* The matrix should be a grid, the plots will be fed to that grid in numerical order so for example:
    - `layout(matrix(1:3, nrow = 1))` will fit the plots to the following matrix in the order specified:
- $$\begin{matrix}
 & 1 & 2 & 3 \\
 \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}
 \end{matrix}$$

- **\*\***In order to use `'grid.layout()'`:

- `grid.arrange(plot1, plot2, ncol = 2))`
- \* This is the only one that will work with `ggplot2`

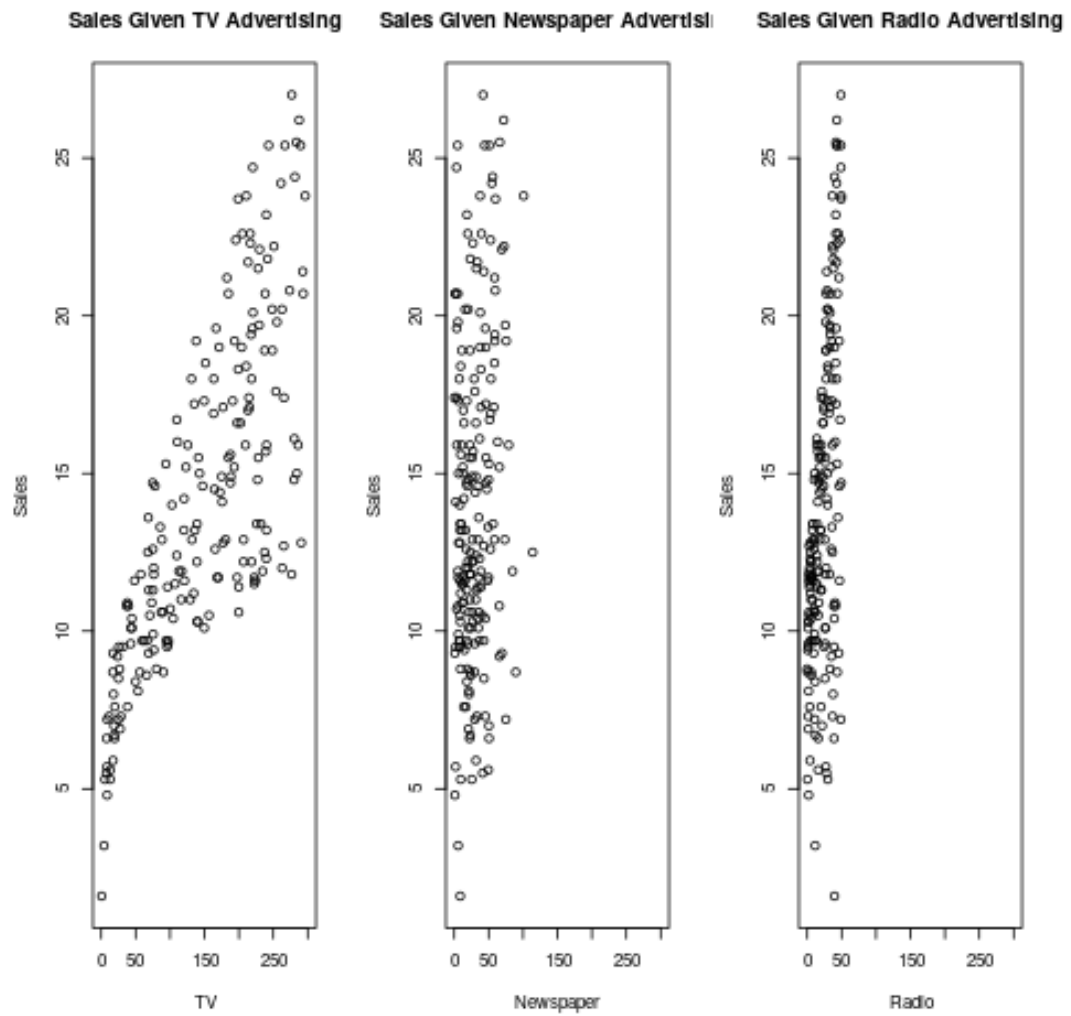
## 1. Multi Fit Base Plots

R:BASEPLOT

```

1  # Set the layout:
2  # Using `layout()` command:
3
4  layout(matrix(1:3, nrow =1))
5
6  # using `par()` command:
7
8  #par(mfrow=c(1,3)) # Specify the
9
10 # Set the plot Domain
11 pdom <- c(0, 300) #Plot Domain
12
13 #Generate the plots
14 plot(formula = Sales ~ TV, data = adv, xlim = pdom,
15       main = "Sales Given TV Advertising")
16 plot(formula = Sales ~ Newspaper, data = adv, xlim = pdom,
17       main = "Sales Given Newspaper Advertising")
18 plot(formula = Sales ~ Radio, data = adv, xlim = pdom,
19       main = "Sales Given Radio Advertising")

```



## 2. GGPlot

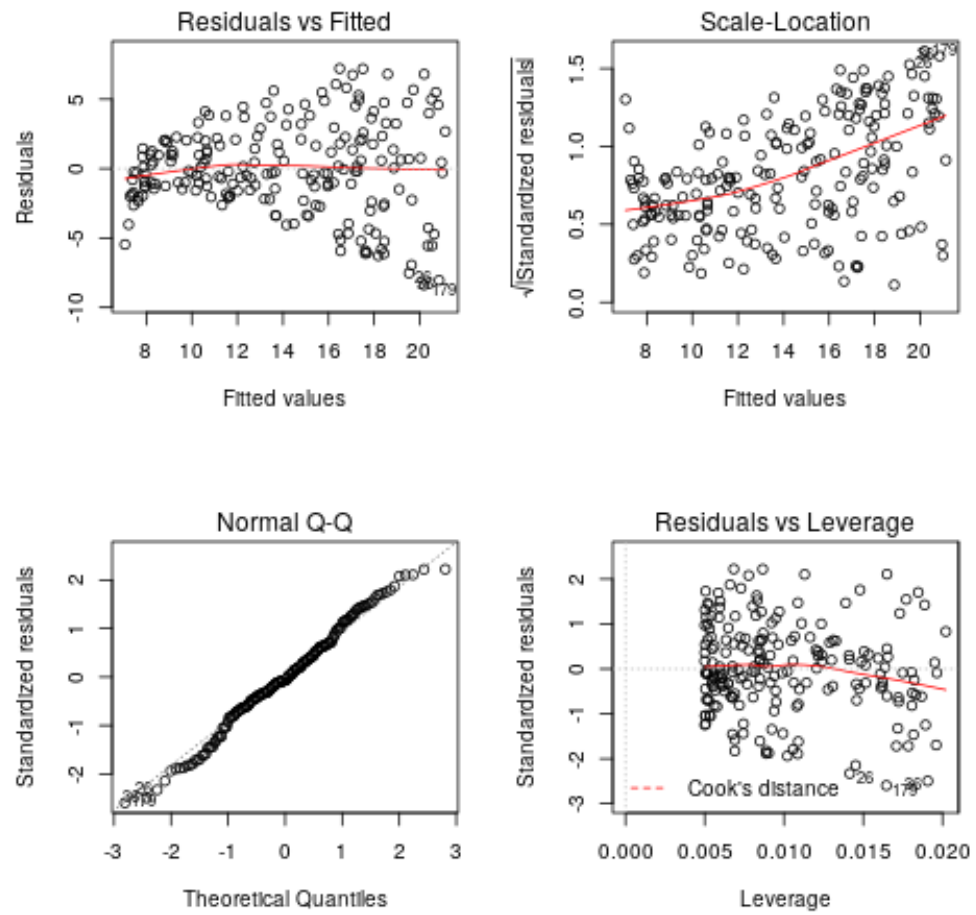
GGPLOT2:LINEAR:REGRESSION

(a) Television Advertising

```

1  adv$MeanAdvertising <- rowMeans(adv[,c(!(names(adv) ==
   ↪  "Sales"))])
2
3  AdvTVPlot <- ggplot(data = adv, aes(x = TV, y = Sales, col =
   ↪  MeanAdvertising)) +
4    geom_point() +
5    theme_bw() +
6    stat_smooth(method = 'lm', formula = y ~ poly(x, 2, raw =
   ↪  TRUE), se = FALSE) +
7    ##stat_smooth(method = 'lm', formula = y ~ log(x), se =
   ↪  FALSE) +
8    labs(col = "Mean Advertising", x= "TV Advertising")
9  print(AdvTVPlot)
10
11  if(knitr::is_html_output()){
12    ggplotly(knitr::is_latex_output())
13  } else {
14    AdvTVPlot
15  }

```

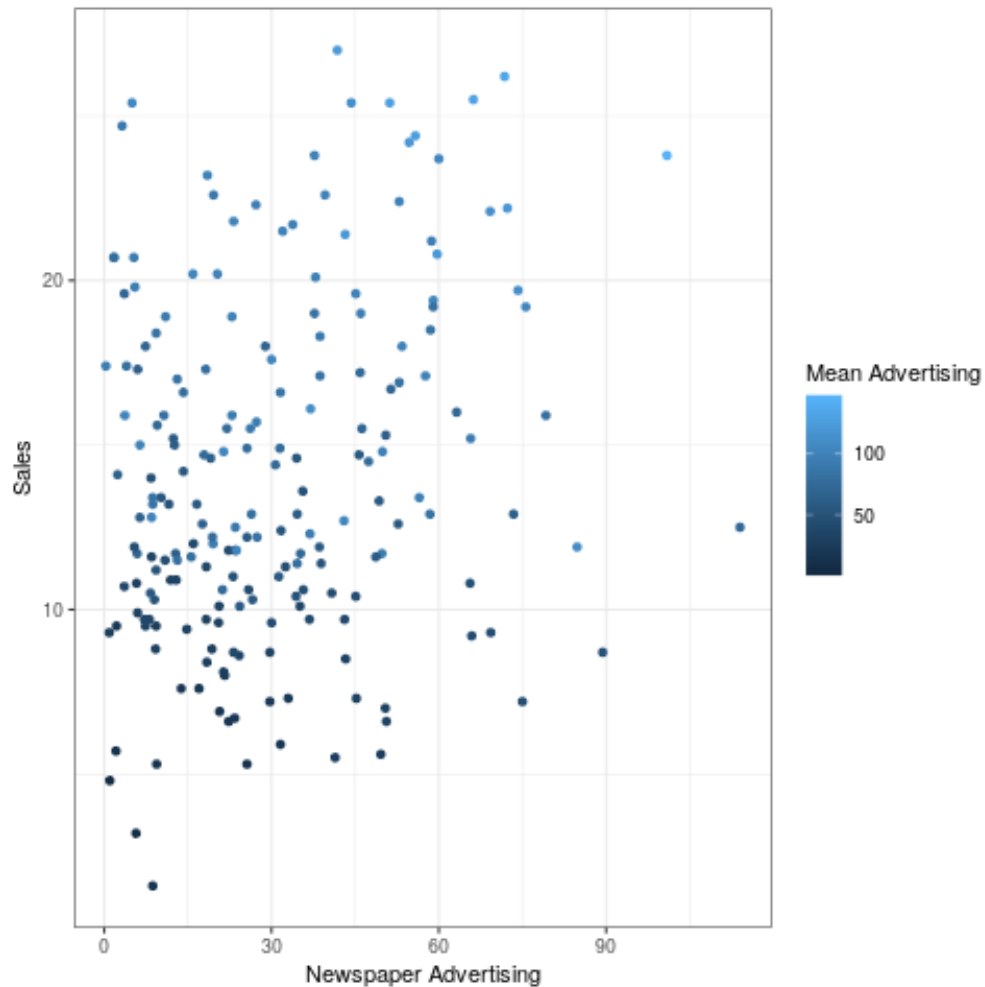


(b) Radio Advertising

```

1   AdvRadPlot <- ggplot(data = adv, aes(x = Radio, y = Sales,
    ↪   col = MeanAdvertising)) +
2     geom_point() +
3     theme_bw() +
4     labs(col = "Mean Advertising", x = "Radio Advertising") +
5     geom_smooth(method = 'lm')
6
7   # padv %>% ggplotly() plotly doesn't work with knitr/LaTeX so
    ↪   test the output and choose accordingly:
8   #This could be combined into an interactive graph by
    ↪   wrapping in ggplotly(padv)
9
10  if(knitr::is_html_output()){
11    AdvRadPlot %>% ggplotly()
12  } else {
13    AdvRadPlot
14  }

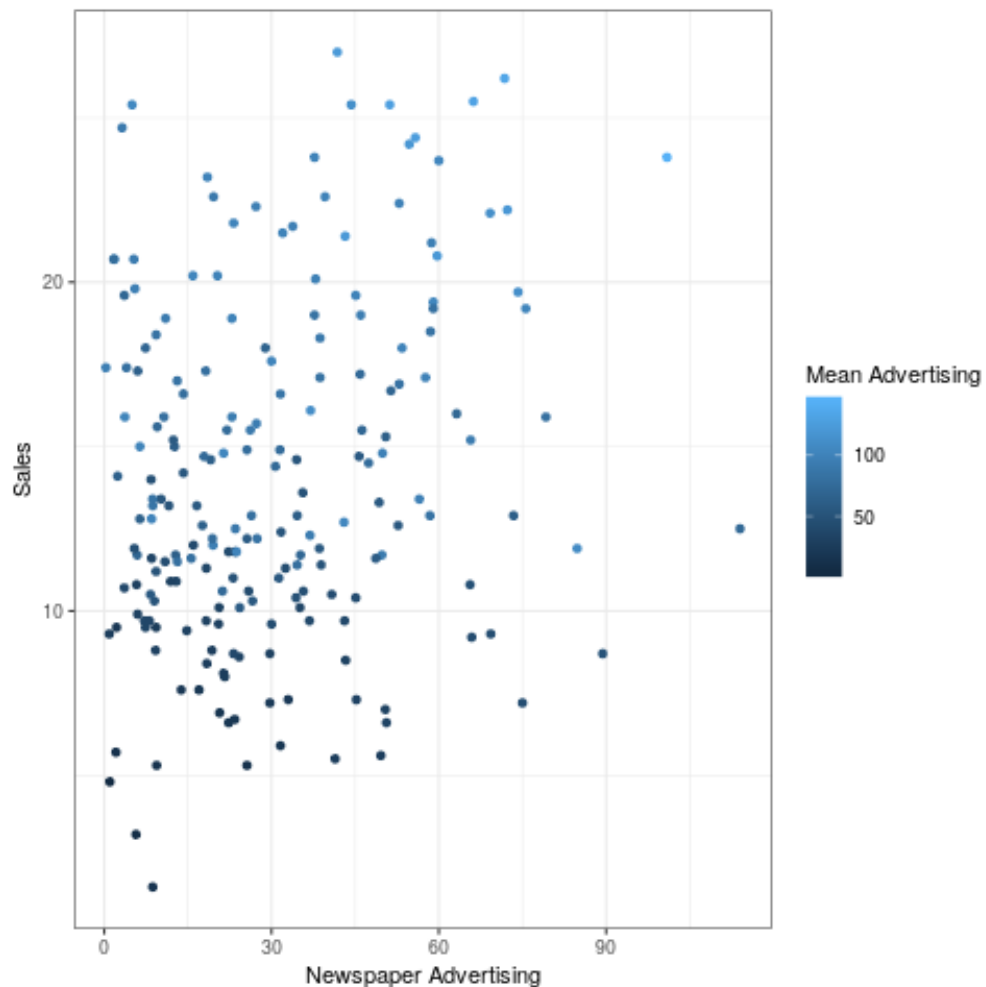
```





(c) Newspaper Advertising

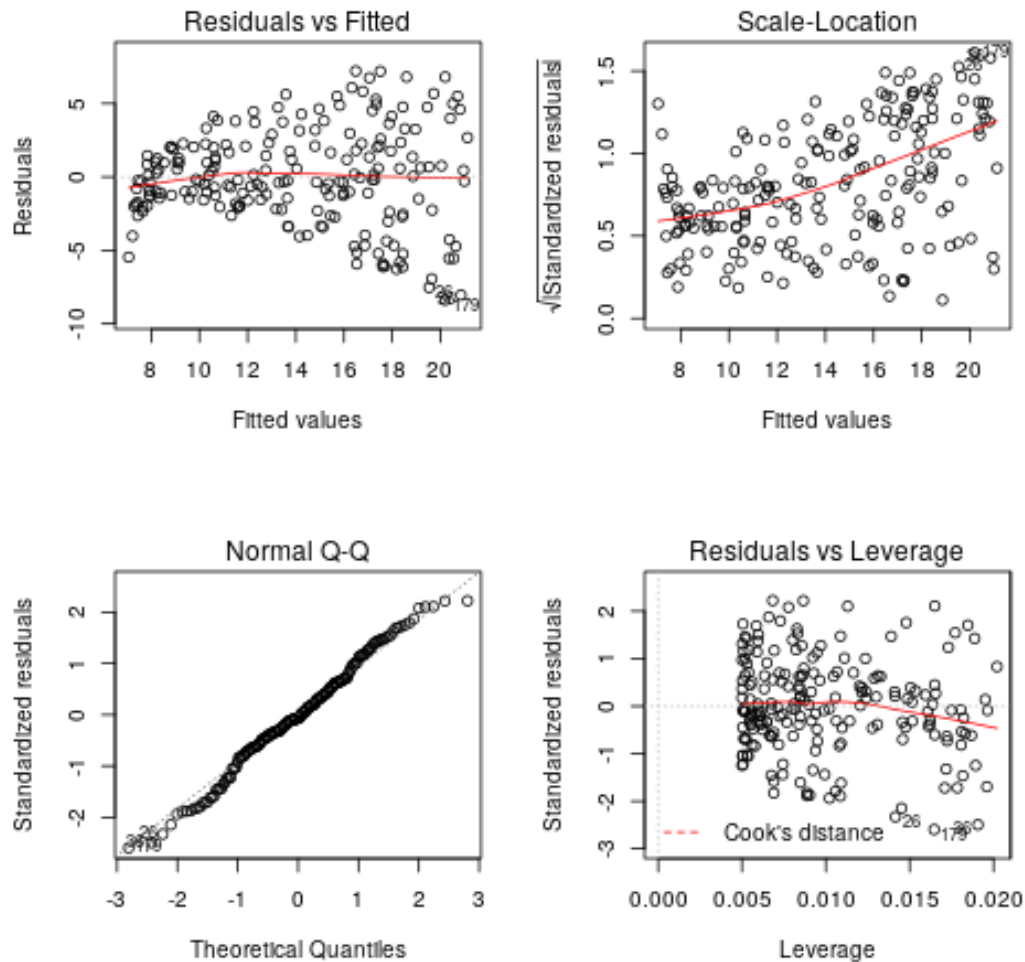
```
1 AdvNewsPlot <- ggplot(data = adv, aes(x = Newspaper, y =  
  ↳ Sales, col = MeanAdvertising)) +  
2   geom_point() +  
3   theme_bw() +  
4   labs(col = "Mean Advertising", x = "Newspaper Advertising")  
5  
6 # padv %>% ggplotly() plotly doesn't work with knitr/LaTeX so  
7 ↳ test the output and choose accordingly:  
8 #Thise could be combined into an interactive graph by wrapping  
9 ↳ in ggplotly(padv)  
10  
11 if(knitr::is_html_output()){  
12   AdvNewsPlot %>% ggplotly()  
13 } else {  
14   AdvNewsPlot  
15 }
```



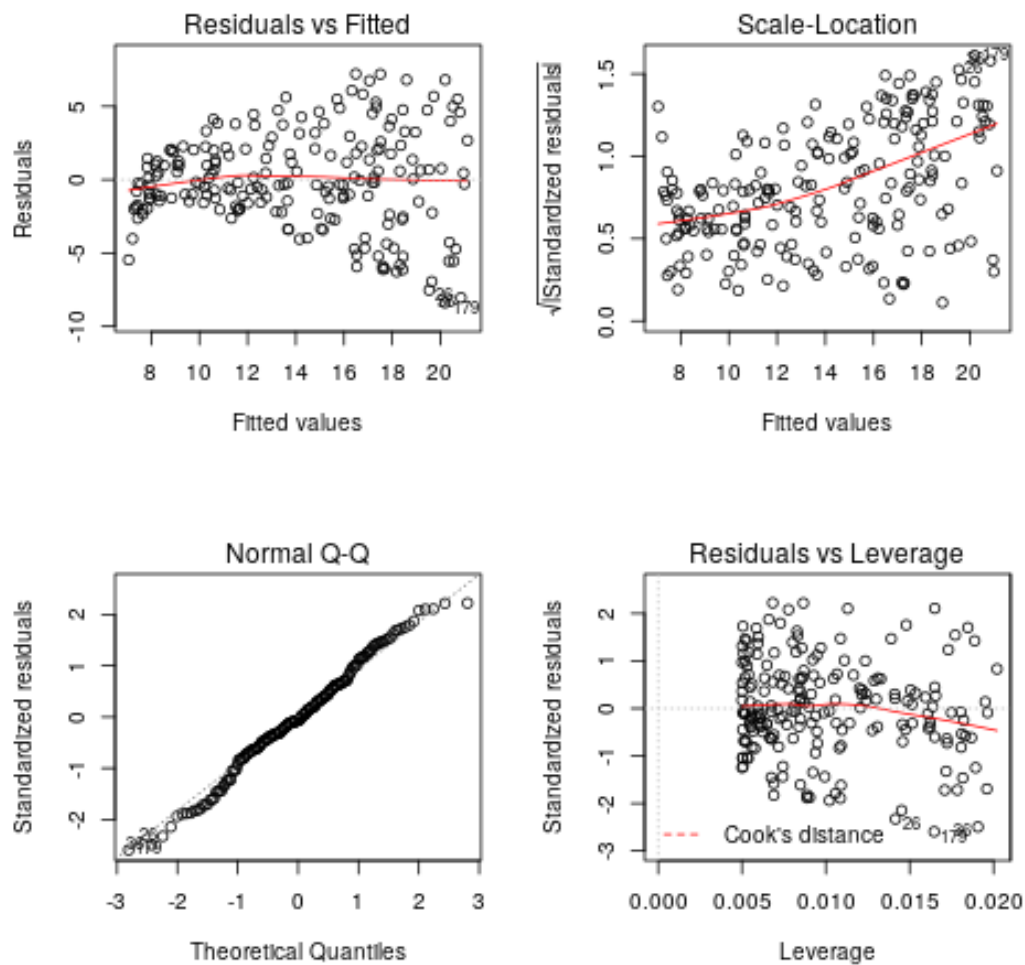
### 3. Base Plot

R:BASEPLOT

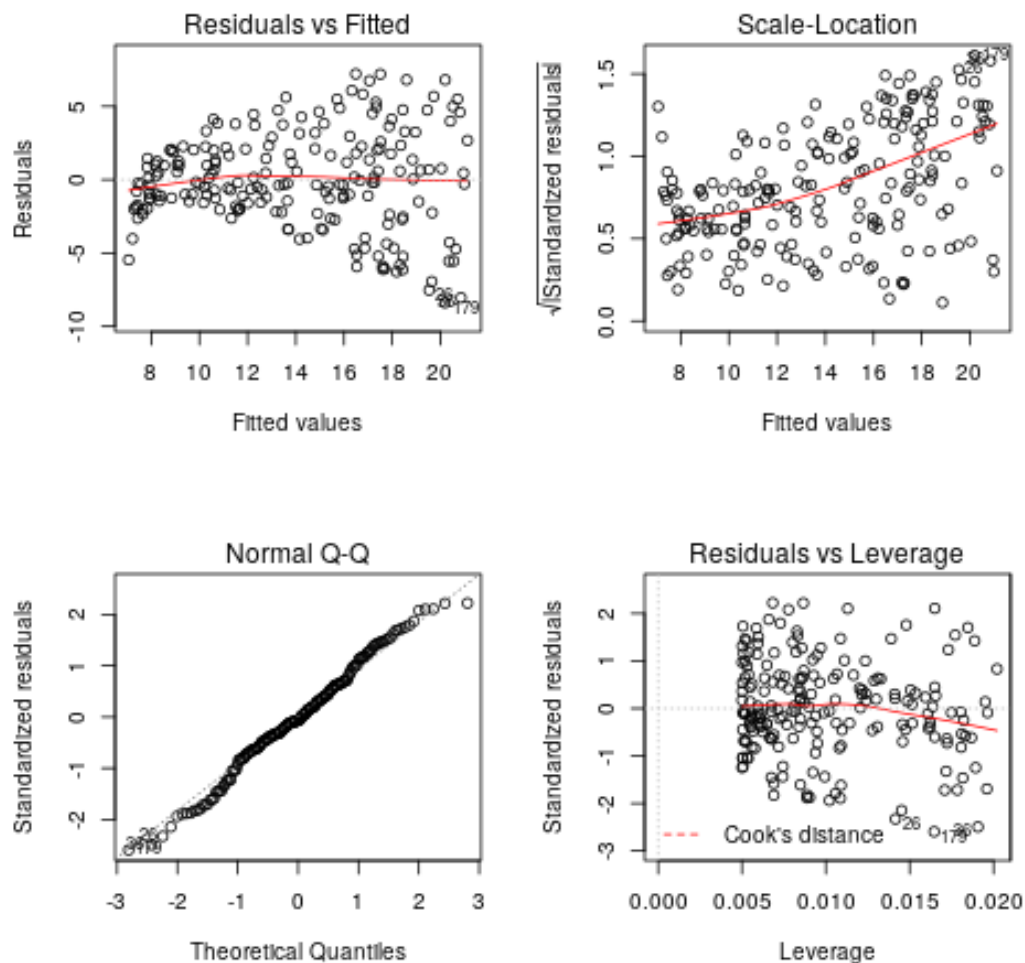
```
1 pdom <- c(0, 300) #Plot Domain
2 plot(formula = Sales ~ TV, data = adv, xlim = pdom,
3       main = "Sales Given TV Advertising")
```



```
1 plot(formula = Sales ~ Newspaper, data = adv, xlim = pdom,
2       main = "Sales Given Newspaper Advertising")
```



```
1 plot(formula = Sales ~ Radio, data = adv, xlim = pdom,
2       main = "Sales Given Radio Advertising")
```



### (c) Find the Correlation Coefficient

linear:regression

The correlation coefficient can be found by using the `cor` function, it is a measurement of the strength of a linear relationship ranging from -1, to 1, wherein a value of 0 would represent no relationship.

The Pearson Correlation Coefficient tends to be used over other models and its value is determined by:

$$r_{xy} = \frac{\sum_{i=1}^n [x_i - \bar{x}] \times (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n [(x_i - \bar{x})^2]} \sqrt{\sum_{i=1}^n [(y_i - \bar{y})^2]}}$$

Some of the assumptions underlying the Correlation Coefficient are: <sup>1</sup>

- Independent Observations

---

<sup>1</sup>Correlation Coefficient

- Normally distributed observations (i.e. follows a bell curve)
- homoscedasticity <sup>2</sup>
  - This means equal variance of observations
    - \* i.e. all there is no pattern between the variables and the plot, the points should make a rectangle, not a triangle
- Normally distributed points
- the points must make a straight line not a curve

the correlation coefficient in this case can be found by using `cor(x = adv$TV, y = adv$Sales)` and provides that  $r \approx 0.78$ . This might not be a meaningful value however because the variance of the sales appears to increase as advertising increases, if that is overlooked however the Pearson correlation coefficient provides that the model is a reasonably strong positive linear model.

#### (d) Assess the accuracy of the parameter estimates

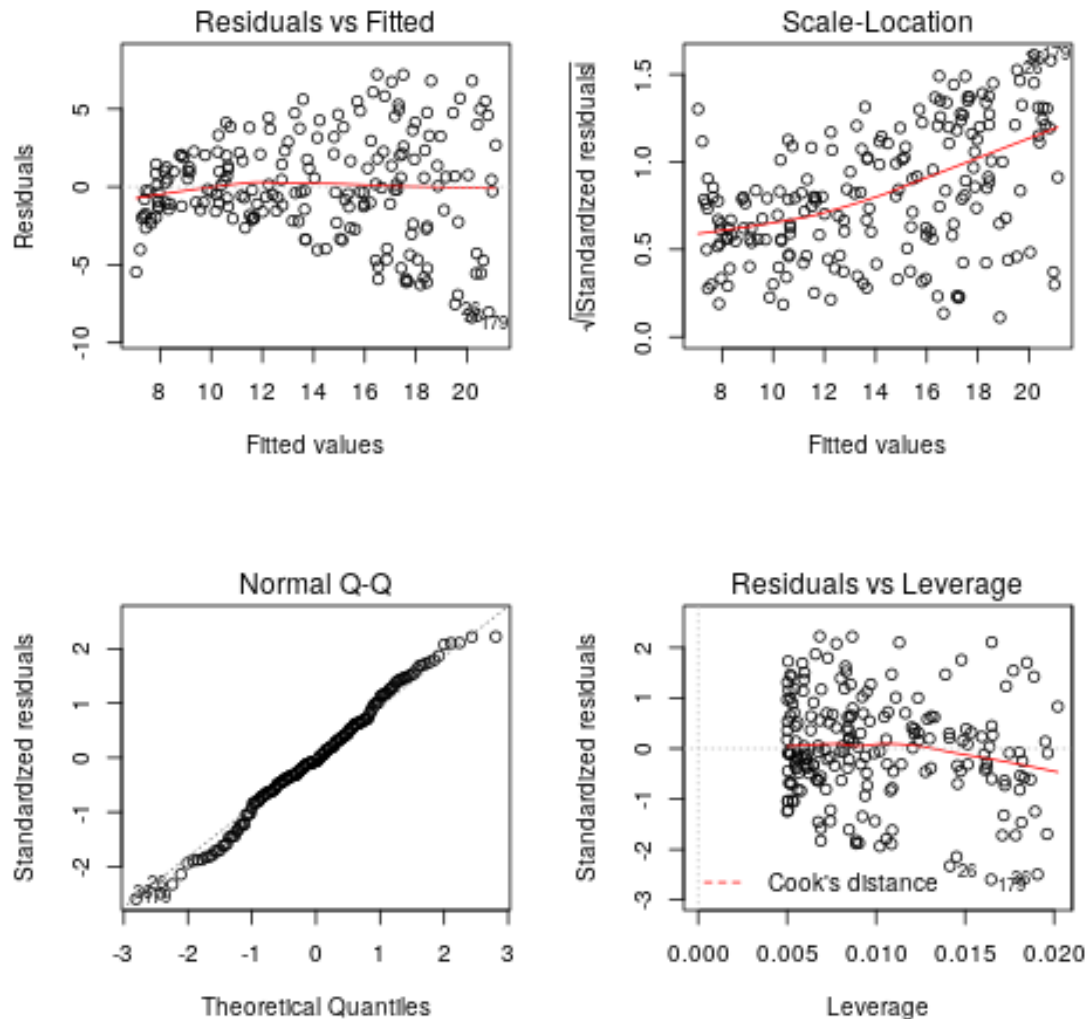
#### ModelEvaluation

The parameter estimates may be returned by summarising the model with `summary(lm)`

```
1  lmMod <- lm(formula = Sales ~ TV, data = adv)
2  lmSum <- summary(lmMod)
3  lmSum
```

---

<sup>2</sup>PennState University



```
##
## Call:
## lm(formula = Sales ~ TV, data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.3860 -1.9545 -0.1913  2.0671  7.2124
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.032594   0.457843  15.36   <2e-16 ***
## TV           0.047537   0.002691  17.67   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.259 on 198 degrees of freedom
```

```
## Multiple R-squared:  0.6119, Adjusted R-squared:  0.6099
## F-statistic: 312.1 on 1 and 198 DF,  p-value: < 2.2e-16
```

```
1 lmSum$coefficients
```

```
              Estimate Std. Error t value    Pr(>|t|)
(Intercept) 7.03259355 0.457842940 15.36028 1.40630e-35
TV           0.04753664 0.002690607 17.66763 1.46739e-42
```

```
##              Estimate Std. Error t value    Pr(>|t|)
## (Intercept) 7.03259355 0.457842940 15.36028 1.40630e-35
## TV          0.04753664 0.002690607 17.66763 1.46739e-42
```

```
1 lmMod2 <- lm(formula = Sales ~ TV, data = adv)
```

In this case we have:

- a slope of  $\beta_1 \approx 0.048 \pm 0.0027$
- an Intercept of  $\beta_0 \approx 7 \pm 0.46$

The standard deviation of a statistic used as an estimator of a population parameter is often referred to as the **standard error of the estimator (S.E.)**; it is the  $\pm$  values specified above:

- Standard Error of Slope Coefficient  $\sigma_{\beta_1} = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SS_x}} = 0.00027$
- Standard Error of Intercept Coefficient  $\sigma_{\beta_0} = \frac{s}{\sqrt{SS_x}} = 0.46$

Where:

- $s$  is the sample standard deviation (OF WHAT?)
- $SS_x = \sum_{i=1}^n [x_i^2] - n \cdot (\bar{x})^2$
- $s$  is the sample standard deviation of  $x$
- because the sample standard deviation of  $x$  predicts the deviation of  $y$  anyway

You may also have the standard deviation of the residuals (the distance along the y-axis of a point from the regression line), this is known as the **Residual Standard Error** and is calculated via the *Ordinary Least Squares Method*<sup>3</sup>, it is given by:

<sup>3</sup>i.e. choosing  $\beta_0$  and  $\beta_1$  to minimise ( $\mathbf{RSS} = \sum_{i=1}^n [(y_i - \hat{y}_i)^2]$ )

$$\sigma_{\varepsilon} = S.E. = \sqrt{\frac{RSS}{N}} \quad (1)$$

(2)

$$= \sqrt{\frac{\sum_{i=1}^n [(y_i - \hat{y}_i)^2]}{N}} \quad (3)$$

Which you'll notice is identical to the **RMSE**.

so by the empirical method  $2 \times S.E.$  would represent a 95% confidence interval (rather than prediction interval) of the expected  $y$ -values. Drawing such a confidence interval:

```
1 paramint <- confint(object = lm(adv$Sales ~ adv$TV), level = 0.95) %>%
  ↪ signif(2)
2 paramint
```

```
##                2.5 % 97.5 %
## (Intercept) 6.100  7.900
## adv$TV      0.042  0.053
```

So drawing from this we could expect, with only a 5% probability of incorrectly rejecting the null hypothesis that there is no relationship, that in the absence of advertising, the TV sales to fall between 6.1 and 7.9.

With the same degree and type of certainty it could also be included that for every \$1000 increase in advertising, the tv sales will increase by between 42 and 53.

- (f) Test the significance of the slope of the linear model If it is appropriate to fit a linear model to data, then we can test for correlation between the data points by considering whether or not the slope value is non-zero  $\beta_1 \neq 0$ , this is because a zero coefficient would be such that the model would predict  $Y = C + \varepsilon$ , this means that  $X$  is not a feature/predictor of  $Y$ , however  $Y$  may still be a function of (or rather response variable of) other values other factors that are 'behind the scenes'.<sup>4</sup>

So our hypotheses would be:

$$H_0 : \beta_1 = 0 \quad (\text{The null hypothesis is that nothings related}) \quad (4)$$

$$H_1 : \beta_1 = 1 \quad (5)$$

So our interest is to determine how far from 0 our expected  $\beta_1$  value needs to be from 0 for us to conclude

The expected value of  $\beta_1$  is so far from zero we can conclude that it it's not zero at some significance level <sup>†</sup>"

<sup>4</sup>Refer to page 67 of the text book, section [3.1.2]



#+BEGINQUOTE † at some low probability of incorrectly rejecting the null hypothesis

#+ENDQUOTE

The problem is defining how far from zero is far enough, for this we use the expected distance from the regression line, the standard error from above, a value observed too many standard deviations to the right of the mean are not very likely to occur.

- (a) Choosing a Parametric method A statistical method that relies on an underlying assumption of the statistical distribution of the data is known as a parametric method, in this case, it is a fundamental assumption of **Ordinary Least Squares** Linear regression that the data is normally distributed.<sup>5</sup>

This is a situation where we use the *Student's t-test* because this is a sample, and the population standard deviation ( $\sigma$ ) is not known and hence the confidence interval for the mean must be made broader in order to account for the fact that the sample standard deviation  $s$  is being used to estimate  $\sigma$

because the sampled population is normally distributed, the sampling distribution of  $\bar{x}$  will be normally distributed<sup>6</sup> (regardless of sample size) and centred about  $\mu$  with a standard deviation of  $\frac{\sigma}{\sqrt{n}}$ . If the population was non-normal the sampling distribution will be approximately normal for  $n \geq 30$ .

Because  $\frac{\sigma}{\sqrt{n}}$  is the standard deviation of the the sample mean  $\bar{x}$  it is referred to as the **Standard Error of the mean**<sup>7</sup>, so we could calculate the critical value along the standard normal distribution corresponding to the the sampling distribution in order to determine probabilities, however,  $\sigma$  is unknown and using  $s$  instead will not create a normal distribution, the distribution it creates is Gosset's **Student's t-distribution**<sup>8</sup>:

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad (6)$$

So in this case our test statistic will be:

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} \quad (7)$$

In order to perform this test in R we can use `qnorm` and `qt` to return critical values, `t.test` will perform a hypothesis test directly from input data but that's not suitable here.

```
1 tcritval <- qt(p = 0.05, df = nrow(adv)-2 )
2 tcritval %>% signif(2)
```

<sup>5</sup>Mendenhall, *Introduction to Probability & Statistics* p.254 [7.4]

<sup>6</sup>By the Central Limit Theorem

<sup>7</sup>Mendenhall, *Introduction to Probability & Statistics* p.254 [7.4]

<sup>8</sup>Mendenhall, *Introduction to Probability & Statistics* p.254 [7.4]

```
## [1] -1.7
```

So the critical t-value is -1.7 and from the summary call from before we have that the t-statistic is 17, which far exceeds this, as a matter of fact further over to the right the p-value is reported at  $\alpha = 10^{-16}$ .

In practice you'd just read off the  $p$ -values and pick the ones with \* to the right of them, the more \* the more significance.

Hence we reject the hypothesis that no relationship exists at an extremely low probability of incorrectly doing so (i.e. low probability of committing type 1 error).

2. (g) Plot the straight line within the scatter plot and comment

- (a) Base Plot In order to plot this inside base packages, feed the model object, i.e.  $\text{lm}(Y \sim X)$  inside a call to `abline()` in order to plot the model over the top of the base plot, so all together it might look like: <sup>9</sup>

```
Form <- Sales ~ TV
Lmodel <- lm(formula = Form, data = adv, na.action = na.exclude)
plot(Form, data = adv)
abline(Lmodel)
```

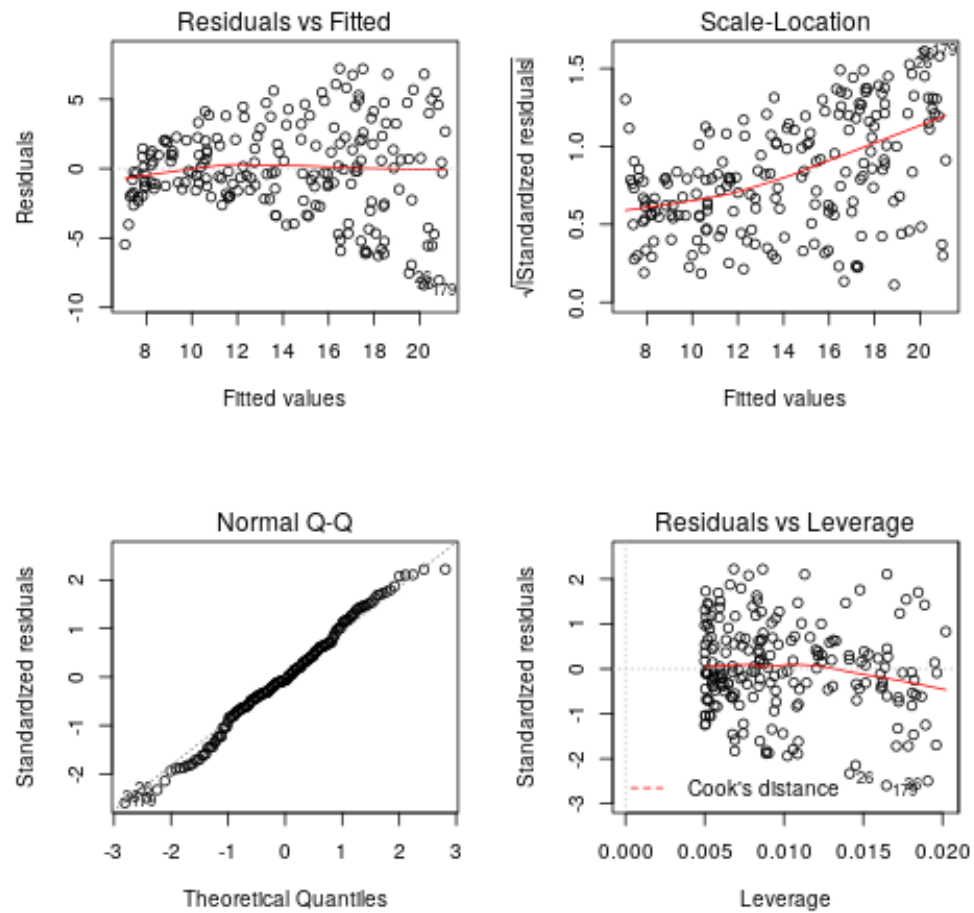
Or you could do it like this even, but I think the way above is better syntax because it will behave better with 'predict' function and follows tidyverse syntax

```
Lmodel <- lm(adv$Sales ~ adv$TV)
plot(x = adv$TV, y = adv$Sales)
abline(a = Lmodel$coefficients[1], b = Lmodel$coefficients[2])
```

```
1 plot(formula = Sales ~ TV, data = adv, xlim = pdom,
2       main = "Sales Given TV Advertising")
3
4 abline(lmMod)
```

---

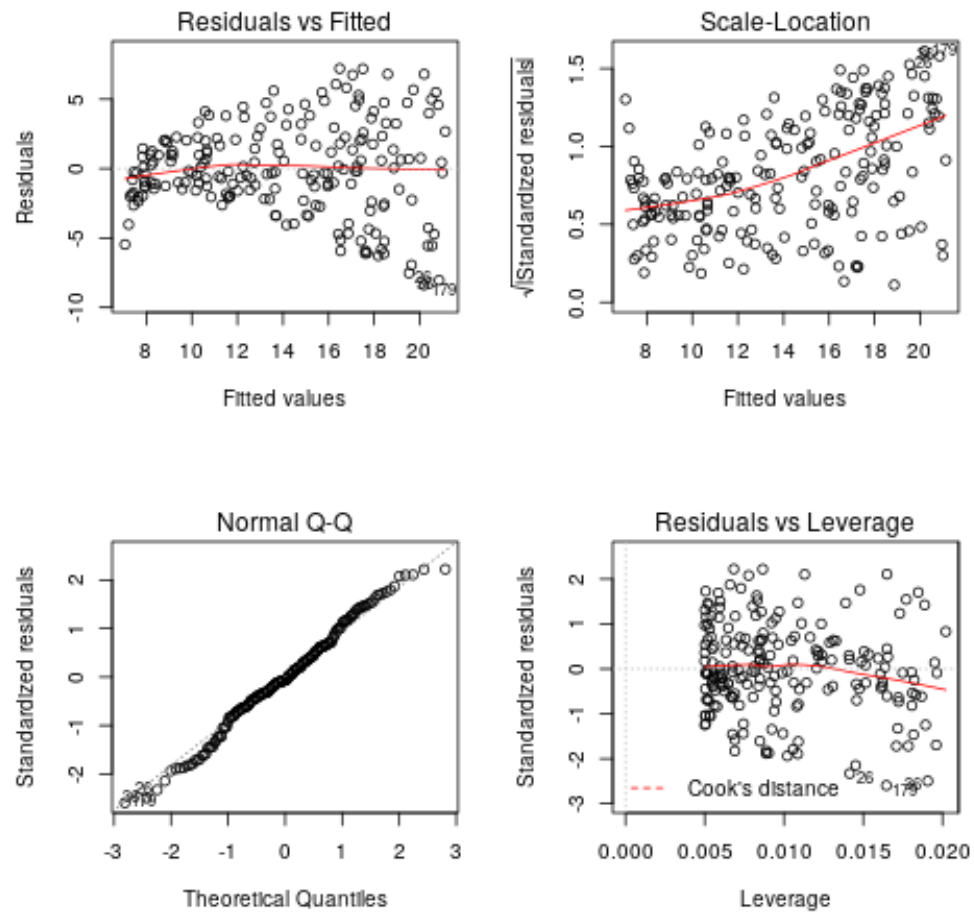
<sup>9</sup> `na.exclude` will pad values extracted so lengths are the same, `na.omit` will not



(b) GGplot

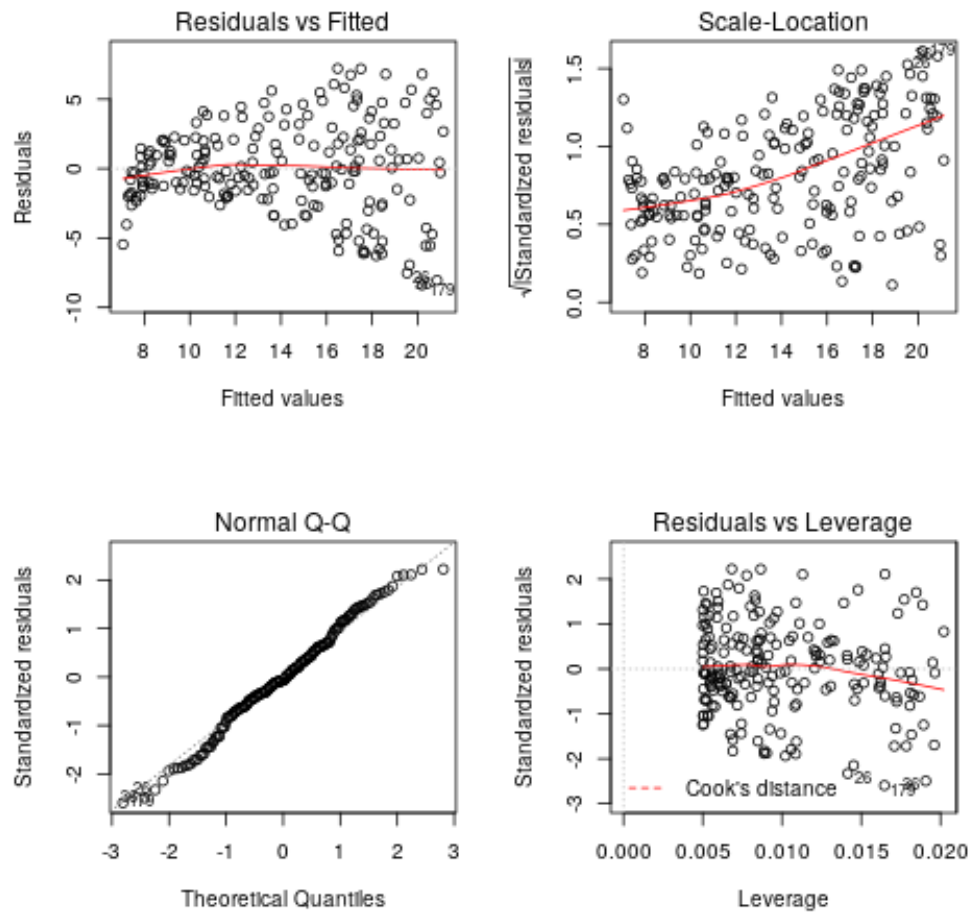
LINEAR:REGRESSION:GGPLOT2

```
1 AdvTVPlot <- ggplot(data = adv, aes(x = TV, y = Sales, col =
  ↳ MeanAdvertising)) +
2   geom_point() +
3   theme_bw() +
4   stat_smooth(method = 'lm', formula = y ~ x, se = FALSE)
5
6 AdvTVPlot
```



If we needed to feed ggplot a specific model we could do that like this, but it's a whole thing to do and you'd probably rather not do it this way, but if you really really need to

```
1 AdvTVPlot <- ggplot(data = adv, aes(x = TV, y = Sales, col =
  ↳ MeanAdvertising)) +
2   geom_point() +
3   theme_bw() +
4   stat_smooth(
5     method = "lm",
6     mapping = aes( y = predict(lmMod)
7                     )
8   )
9
10 AdvTVPlot
```



3. (h) Assess the overall accuracy of the model The model can be assessed by considering the:

- Coefficient of determination  $R^2$  which is the proportion of variance in the data that is explained by the model
- Only in the case of simple linear regression is  $R^2 = (r)^2$
- The Residual Standard Error is the standard deviation of the residuals, i.e. it is the expected distance between each point to the regression line, taken along the  $y$ -axis.

(a) Terminology The textbook makes, in my opinion, a mistake in that it refers to the the *Root Mean Square Error (RMSE)* as the *Residual Standard Error (RSE)* <sup>10</sup>, this is true, the standard error of the residuals ( $\varepsilon$ ) would be the RMSE, so we would have  $RMSE = \sigma_{\varepsilon}$ , that's fine.

<sup>10</sup>Refer to Page 69 of the TB for RMSE definition, the TB divides by DF which is probably more correct that dividing by sample size.

The issue is there is another common term used called the *Relative Squared Error* (**RSE**) is often used <sup>11</sup> and so this is hence ambiguous, hence forth I will:

- Refer to the Standard Error of the residuals ( $\sigma_{varepsilon}$ ) as **RMSE**:

$$- \text{RMSE} = \sqrt{\frac{\sum \varepsilon^2}{n}}$$

- Refer to the Relative Standard Error as **RSE**

$$- \text{RSE} = \frac{\sigma_{\varepsilon}^2}{\sigma_y^2} = \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

- \* The advantage to the RSE is that it can be compared between models with different units, whereas the RMSE cannot, just another tool in the belt I suppose.

(b) Root Mean Square Error LOSSFUNCTION Recall that the model was of the form  $Y = \beta_1 X + \beta_0 + \varepsilon$ , the **RMSE** (*Root Mean Square Error*) is the standard deviation of  $\varepsilon$  as measured along the  $Y$ -axis:

$$\sigma_{\varepsilon} = \sqrt{\frac{\sum_{i=1}^n [(y_i - \hat{y}_i)^2]}{N}} \quad (8)$$

This value can be returned from R by investigating the anova table:

```
1 anova(lmMod)
```

```
## Analysis of Variance Table
##
## Response: Sales
##           Df Sum Sq Mean Sq F value    Pr(>F)
## TV           1 3314.6  3314.6    312.14 < 2.2e-16 ***
## Residuals 198 2102.5    10.6
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From the ANOVA table it can be seen that the average squared residual is 10.6

---

<sup>11</sup>[An Introduction to Data Science : Model Evaluation - Regression](#)

$$\text{mean}(\varepsilon^2) = 10.6 \quad (9)$$

$$\Rightarrow \frac{1}{n} \cdot \sum_{i=1}^n [\varepsilon_i] = 10.6 \quad (10)$$

$$\Rightarrow \frac{1}{n} \cdot \sum_{i=1}^n [(\hat{y}_i - y_i)^2] = 10.6 \quad (11)$$

$$\Rightarrow \sqrt{\frac{1}{n} \cdot \sum_{i=1}^n [(\hat{y}_i - y_i)^2]} = 3.2 \quad (12)$$

$$(13)$$

$$\Rightarrow \sigma_\varepsilon = 3.2 \quad (14)$$

Thus we may conclude that we expect the model to predict the sales within  $\pm 3.2$  units, which is quite predictive and hence useful.

- (c) Coefficient of Determination The coefficient of determination is the proportion of variation within the model that is explained by the model:

$$R^2 = \frac{TSS - RSS}{TSS} \quad (15)$$

$$= \frac{3315}{3315 + 2103} \quad (16)$$

$$(17)$$

$$= 0.612 \quad (18)$$

In practice we would simply extract the coefficient of determination ( $R^2$ ) from the model-summary:

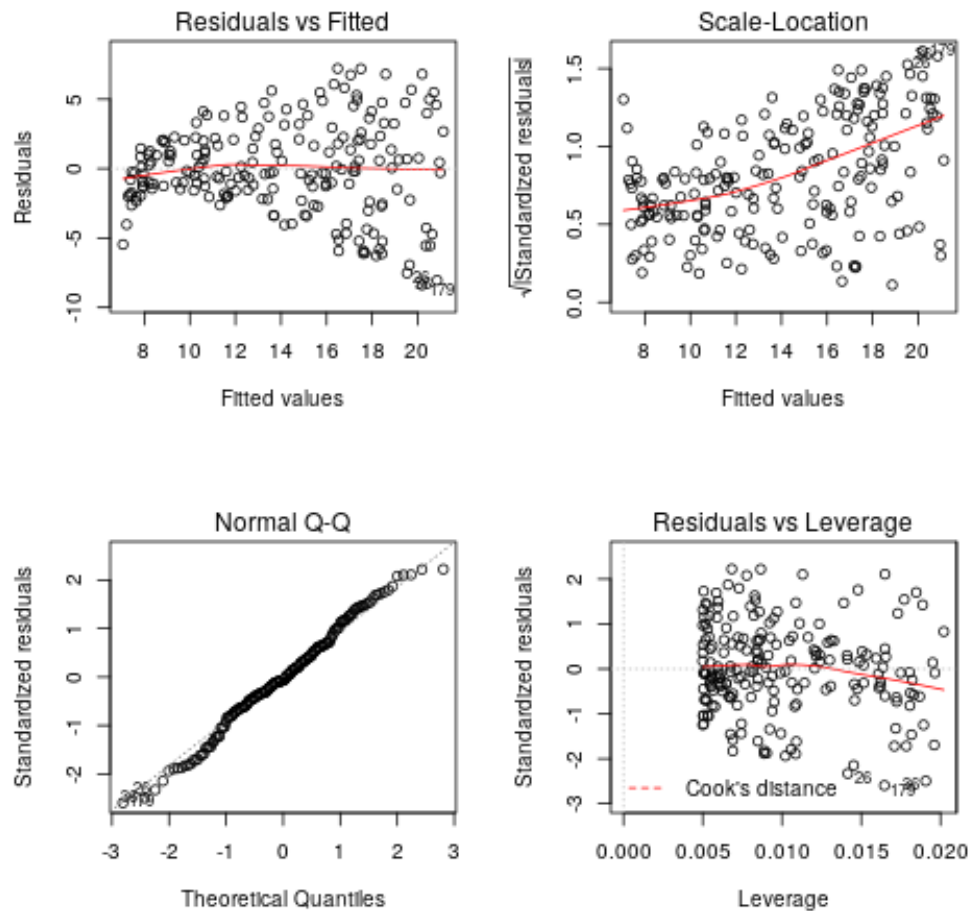
```
1 lmSum$r.squared %>% round(3) %>% percent()
```

```
## [1] "61.2%"
```

This value suggests that a reasonable amount of the variation is explained by the model, but perhaps a non-linear model could explain more of the variance. (be careful a significant coefficient of determination doesn't necessarily mean that the slope is significantly different from 0)

- (d) Residual Analysis

```
1 layout(matrix(1:4, nrow = 2))
2 plot(lmMod)
```



- The residual plot does not appear to be normally distributed, there is a slight logarithmic trend, this violates assumptions of the linear model undermining the predictive capacity of the model in this case.
- The variance is also non-constant, for a linear model to be used in must be homoscedastic (i.e. constant variance), this is not the case implying that the assumptions of the linear model have been violated and hence this model may not be appropriate <sup>12</sup>
- the standardised residuals should be normally distributed with a mean of 0 and standard deviation of 1, whilst the standard deviation appears acceptable, the standardised residuals are centred around  $\approx 3/4$  with a positive upward slope violating the assumption of normality.
- The normal Q-Q plot is a straight line so actually the data is probably normally distributed, the only issue is the heteroscedasticity of the data.

<sup>12</sup>refer to page 96 of the TB, log or exp transforming may be appropriate here, the data is not homoscedastic and is hence said to be heteroscedastic. ##### How to use predict



- The Cook's Distance plot suggests that there are some points with a high amount of leverage, so perhaps there are some outliers or perhaps the increasing variance is undermining the appropriateness of the model.
4. (i) Use the model to make predictions When making predictions is important to ensure that the names of a data frame are syntactically correct, otherwise you will have a bad day trying to get predict to work and ggplot2 to work because specifying the data frame names in a formula will be difficult, make sure that names are always syntactically valid.

what is important is you create your model with the correct syntax, if you create your model like this:

```
mymodelWRONG <- lm(adv$Sales ~ adv$TV)
```

you won't be able to predict data like this:

```
predict(object = lmMod, newdata = data.frame("TV" = 300))
```

you'll just get an error that says 'newdata' had 1 row but variables found have 200 rows, you have to give the variables corresponding names so that the model object can save them for later and make the connection, for instance, if inspect the terms from above you will get:

```
mymodelWRONG[["terms"]]
```

which outputs, at the tail end:

```
adv$Sales    adv$TV
"numeric" "numeric"
```

where as if you create the model like this:

```
lmModCORRECT <- lm(formula = Sales ~ TV, data = adv)
predict(object = lmModCORRECT, newdata = data.frame("TV" = 300))
```

and inspect the terms with:

```
lmModCORRECT[["terms"]]
```

you will get this as output

```
Sales      TV
"numeric" "numeric"
```

where Sales and TV are the outputs of `names(adv)` and so I can use that when I use `predict`. You should not use `attach` it will cause problems later, however, it can be nice to use `attach` just before a `predict` call to get auto completed names and then remove `attach` and re-execute the script .

So always use the `lm(formula = Y~X, data = myDF)` because it works the best; you have to use the same syntax/format when using `predict` or `ggplot` anyway so there's no reason not to use the same syntax throughout anyway.

Also the lecturer said to use lists, I reckon use data frames because that way your `newdata` matches the input data one-to-one, moreover:

- It makes it far simpler to assign names, because again, the input/output data will all be the same format
- when creating *Lasso Regression Models* you have to use matrices as input data and it's easier to set your workflow up to go from dataframe to matrix (You have to do this in predictive modelling)

#### (a) Predict the Data

##### i. One Point

```
1 input = 3
2 output <- predict(object = lmMod, newdata = data.frame("TV"
  ↳ = 3))
3 predDatasingl <- data.frame(input, output)
4 names(predDatasingl) <- names(adv[c(1,4)])
5
6 print(predDatasingl)
```

```
##    TV    Sales
## 1   3  7.175203
```

##### ii. Multiple points

```
1 input <- seq(from = 100, to = 900, by = 100)
2 output <- predict(object = lmMod, newdata = data.frame("TV"
  ↳ = input))
3
4 predDF <- data.frame(input, output)
5 names(predDF) <- names(adv[c(1,4)])
6 predDF
```

```
##    TV    Sales
## 1 100 11.78626
## 2 200 16.53992
## 3 300 21.29359
## 4 400 26.04725
## 5 500 30.80091
```

```
## 6 600 35.55458
## 7 700 40.30824
## 8 800 45.06191
## 9 900 49.81557
```

## Question 02

---

**(a) Upload the Auto Dataset and explore it.**

**(b) Construct scatter plots to visualize the relationship between**

mpg and displacement, weight and acceleration: :CUSTOM\_ID: b-construct-scatter-plots-to-visualize-the-relationship-between-mpg-and-displacement-weight-and-acceleration

**Repeat the analysis in Q1 (c) to (i) using mpg and weight.**

## Multiple Linear Regression

Material of Tue 19 March 2019, week 3

## Question 01 - Multiple Linear Regression

### Load Packages

```
1  # Load Packages
2  if(require('pacman')){
3    library('pacman')
4  }else{
5    install.packages('pacman')
6    library('pacman')
7  }
8
9  pacman::p_load(caret, scales, ggplot2, rmarkdown, shiny, ISLR, class,
10 ↪ BiocManager,
11 ↪ corrplot, plotly, tidyverse, latex2exp, stringr,
12 ↪ reshape2, cowplot, ggpubr,
13 ↪ rstudioapi, wesanderson, RColorBrewer, colorspace,
14 ↪ gridExtra, grid, car,
15 ↪ boot, colourpicker, tree, ggtree, mise, rpart,
16 ↪ rpart.plot, knitr, MASS,
17 ↪ magrittr, EnvStats, tidyverse, tidyr, devtools, bookdown,
18 ↪ leaps, car, clipr,
19 ↪ tikzDevice, e1071)
20
21 mise()
22 set.seed(0932)
```

(a) Upload the data "Advertising.csv" and explore it.

ATTACH

First import the data and investigate it:

```
1  adv <- read.csv(file = "../data/83/4c42c3-8dd8-4fef-b402-7e341764d_
↪ 5e9/Advertising.csv", header = TRUE, sep =
↪ ",")
```

```
1  head(adv)
```

```
##      TV Radio Newspaper Sales
## 1 230.1  37.8      69.2  22.1
## 2  44.5  39.3      45.1  10.4
```

```
## 3  17.2  45.9      69.3   9.3
## 4 151.5  41.3      58.5  18.5
## 5 180.8  10.8      58.4  12.9
## 6   8.7  48.9      75.0   7.2
```

```
1 writeLines("\n")
```

```
1 print("***Dimensions***")
```

```
## [1] "***Dimensions***"
```

```
1 writeLines("\n")
```

```
1 dim(adv)
```

```
## [1] 200   4
```

```
1 writeLines("\n")
```

```
1 print("***Summary***")
```

```
## [1] "***Summary***"
```

```
1 writeLines("\n")
```

```
1 summary(adv)
```

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

```
1 writeLines("\n")
```

```
1 print("***Structure***")
```

```
## [1] "***Structure***"
```

```
1 writeLines("\n")
```

```
1 str(adv)
```

```
## 'data.frame':    200 obs. of  4 variables:
## $ TV          : num  230.1 44.5 17.2 151.5 180.8 ...
## $ Radio       : num  37.8 39.3 45.9 41.3 10.8 48.9 32.8 19.6 2.1 2.6 ...
## $ Newspaper: num  69.2 45.1 69.3 58.5 58.4 75 23.5 11.6 1 21.2 ...
## $ Sales      : num  22.1 10.4 9.3 18.5 12.9 7.2 11.8 13.2 4.8 10.6 ...
```

```
1 writeLines("\n")
```

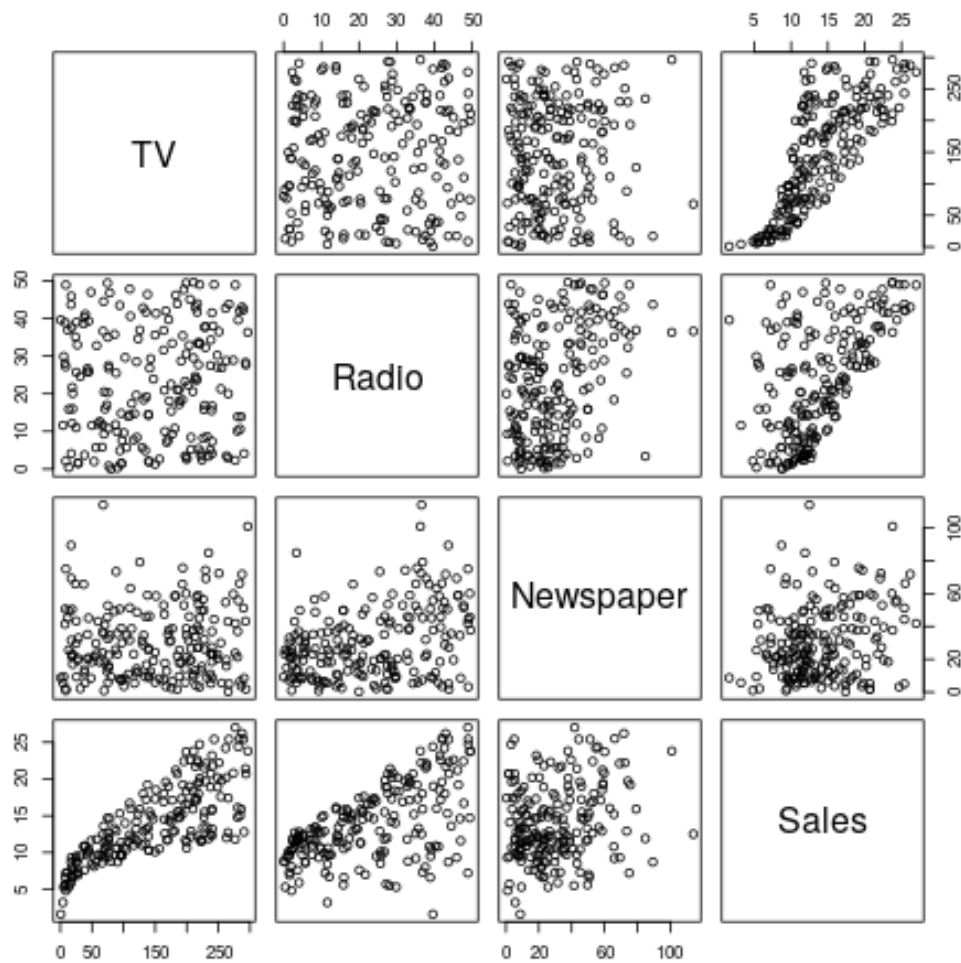
From this we can tell that there is one output, with 3 input values and 200 Observations.

## (b) Find the Covariance and Correlation Matrix of Sales, TV, Radio

and Newspaper. :CUSTOM<sub>ID</sub>: b-find-the-covariance-and-correlation-matrix-of-sales-tv-radio-and-newspaper.  
:CLASS: tabset

### 1. Base Packages

```
1 pairs(x = adv)
```

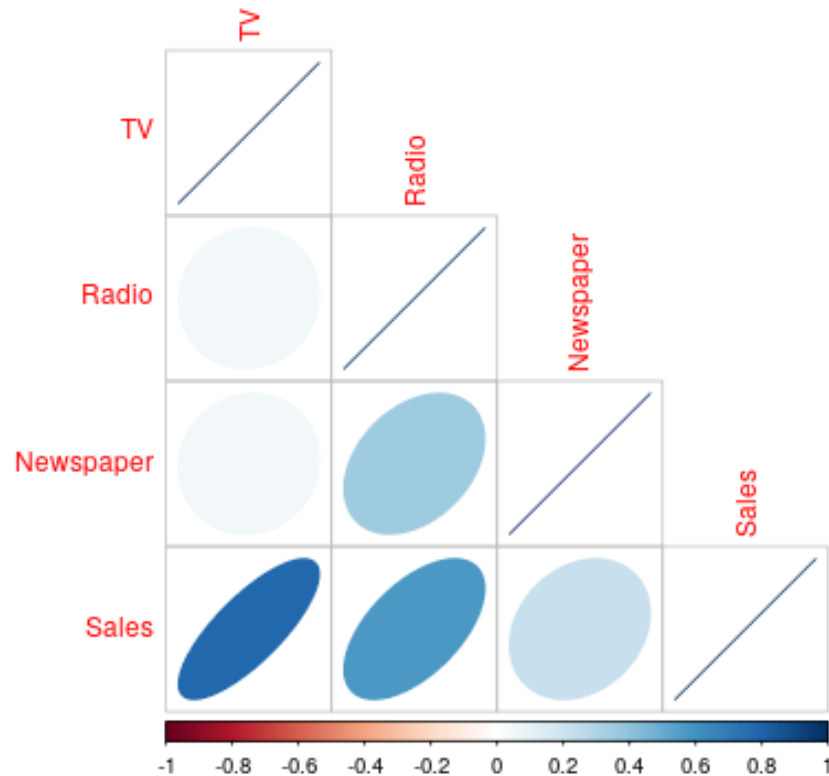


2. `corrplot` In order to use `corrplot` first create a correlation matrix using `cor(adv)` then feed that matrix to `corrplot` with the command `corrplot(cor(adv))`

```

1 # coriris <- cor(iris[,!(names(iris) == "Species")])
2 # corrplot(method = 'ellipse', type = 'lower', corr = coriris)
3
4 corMat <- cor(x = adv)
5 corrplot(corr = corMat, method = "ellipse", type = "lower")

```



</div>

From this we can see that there is a significant amount of correlation between Radio and newspaper (more so even than newspaper and sales), we should consider this variable interaction when deciding upon our model

### (c) Construct the multiple linear regression model and find the

least square estimates of the model :CUSTOM<sub>ID</sub>: c-construct-the-multiple-linear-regression-model-and-find-the-least-square-estimates-of-the-model

A multiple linear regression would give the model:

```
1 advModMult <- lm(formula = Sales ~ TV + Radio + Newspaper, data = adv)
2 advModMult
```



```
##
## Call:
## lm(formula = Sales ~ TV + Radio + Newspaper, data = adv)
##
## Coefficients:
## (Intercept)          TV          Radio    Newspaper
##    2.938889    0.045765    0.188530   -0.001037
```

This gives that the appropriate linear model is:

$$Y_{Sales} = 0.0458 \times TV + 0.19 \times Radio - 0.001 \times Newspaper$$

the fact that Newspaper has a negative coefficient despite being positively correlated with sales is indicative of the weak effect newspaper advertising has on sales as well as the interaction between newspaper and sales.

from this it can be `###` (d) Test the significance of the parameters and find the resulting model to model Sales in terms of advertising modes, TV, Radio and Newspaper. First Summarise the Model:

```
1 advModMult %>% summary()
```

```
##
## Call:
## lm(formula = Sales ~ TV + Radio + Newspaper, data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.8277 -0.8908  0.2418  1.1893  2.8292
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.938889   0.311908   9.422  <2e-16 ***
## TV           0.045765   0.001395  32.809  <2e-16 ***
## Radio        0.188530   0.008611  21.893  <2e-16 ***
## Newspaper   -0.001037   0.005871  -0.177    0.86
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.686 on 196 degrees of freedom
## Multiple R-squared:  0.8972, Adjusted R-squared:  0.8956
## F-statistic: 570.3 on 3 and 196 DF,  p-value: < 2.2e-16
```

A summary of the model provides that given the hypothesis test:

$$H_0 : \beta_i = 0 H_a : \beta_i \neq 0 \quad \forall i \in \mathbb{N}$$

There would be an extremely low probability of incorrectly rejecting the null hypothesis that the given a linear model the coefficients would be zero, hence it is accepted that the coefficients are non-zero except for newspaper advertising.

There is a high probability of incorrectly rejecting the null-hypothesis and hence that should not be rejected and the newspaper advertising should not be seen as significant predictor for sales in this linear model.

It is hence appropriate to remove, via backwards selection, Newspaper from the model, which gives:

```
1  ## Calculate Power??
2
3  # n <- length(adv)
4  # sigma <- sd(adv$Newspaper)
5  # sem = sigma/sqrt(n)
6  # alpha <- 0.05
7  # mu0 <- 0
8  # q <- qnorm(p = 0.005, mean = mu0, sd = sem);q
9  # mu <- q # assumed actual mean value
10 # pnorm(q, mean = mu, sd = sem, lower.tail = FALSE)
```

```
1  advModMultb1 <- lm(formula = Sales ~ TV + Radio, data = adv)
2  advModMultb1
```

```
##
## Call:
## lm(formula = Sales ~ TV + Radio, data = adv)
##
## Coefficients:
## (Intercept)          TV          Radio
##      2.92110      0.04575      0.18799
```

```
1  advModMultb1.sum <- summary(advModMultb1)
```

$$\text{Sales} = 0.045 \times \text{TV} + 0.19 \times \text{Radio} + 2.9211$$

All parameters are highly significant and hence the model is deemed adequate, the model explains 100% of the variation, this can be found by appending `$r.squared` to the model object.

### (e) Assess the overall accuracy of the model.

In order to assess the model, consider the anova table and derive the RMSE and  $R^2$  values:

```
1  advModMultb1.anova <- anova(advModMultb1)
2  advModMultb1.anova
```

```
## Analysis of Variance Table
##
## Response: Sales
##           Df Sum Sq Mean Sq F value    Pr(>F)
## TV          1 3314.6   3314.6 1172.50 < 2.2e-16 ***
## Radio        1 1545.6   1545.6  546.74 < 2.2e-16 ***
## Residuals 197   556.9     2.8
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From this we can see that the  $F$  statistic is associated with a very low p-value, these are the exact same p-values from the summary call.

1. RMSE The RMSE value is the Root mean square error, it is the standard error of the residuals (recall that the standard error is the standard deviation of a model parameter), so the RMSE is basically the standard deviation of the residuals of the model (as measured along the  $y$ -axis).

$$\text{RMSE} = \sqrt{\frac{\sum \varepsilon^2}{n}}$$

```
1 rmse <- function(model){
2   # (sum(model$residuals**2)/length(model$residuals))**0.5
3   sd(advModMultb1$residuals)
4 }
5
6 rse <- function(model){
7   var(model$residuals)/var(model$residuals - model$fitted.values)
8 }
9
10 data.frame("RMSE" = rmse(advModMultb1) , "RSE" = rse(advModMultb1) )
```

```
##           RMSE           RSE
## 1 1.672891 0.1028057
```

So the expected error of the model is  $\pm 1.7$  units of sale, we could use this to create a confidence interval by multiplying by the corresponding *Student's t-statistic* and saying that we would expect an observed value to lie within  $1.96 * \text{S.E}$  of the model 95% of the time (is this correct or is it the expected mean or something because of the Central Limit Theorem?).

2. Coefficient of Determination The coefficient of determination is 89.7%, this can be returned by extracting the value from the object by appending `$r.squared` to the model summary (i.e. `summary(lm( Y ~ X1 + X2 ))$r.squared`).

The  $R^2$  value will be very nearly the same between the initial model and the backwards selected model, however given that the initial model had non-significant predictors it could be considered as over-parameterized, i.e. it violates [Occam's Razor](#).

The coefficient of determination is the ratio of the total variance of the data that is explained by the model, in this case it could be determined by:

$$R^2 = \frac{TSS - SSE}{TSS} = \frac{3314.6 + 1546 + 0.1}{3315 + 1545 + 0.1 + 556.8} = 89\%$$

3. Coefficient of Determination from ANOVA The  $R^2$  value is derived from the ANOVA table thusly:

```
1 advModMultb1.anova <- anova(advModMultb1)
2 TSS_Multb1 <- advModMultb1.anova$`Sum Sq` %>% sum()
3 RSS_Multb1 <- advModMultb1$residuals^2 %>% sum()
4
5 ((TSS_Multb1- RSS_Multb1)/(TSS_Multb1)) %>% signif(3) %>% percent()
  ↳ # Requires `scales` package
```

```
## [1] "89.7%"
```

4. Residual analysis When determining the accuracy or performance of a model, always turn your mind to residual analysis (i.e. how normally are they distributed), this will be performed below.

#### (f) Calculate the predicted values and residuals

The residuals and fitted values may be returned by extracting them from the model (you could always derive them from first principles but you should use the tools at your disposal, it is quicker, less error prone and makes more readable code)

```
1 ResDF <- data.frame("Input" = advModMultb1$fitted.values -
  ↳ advModMultb1$residuals, "Output" = advModMultb1$fitted.values,
  ↳ "Error" = advModMultb1$residuals)
2 ResDF %>% head()
```

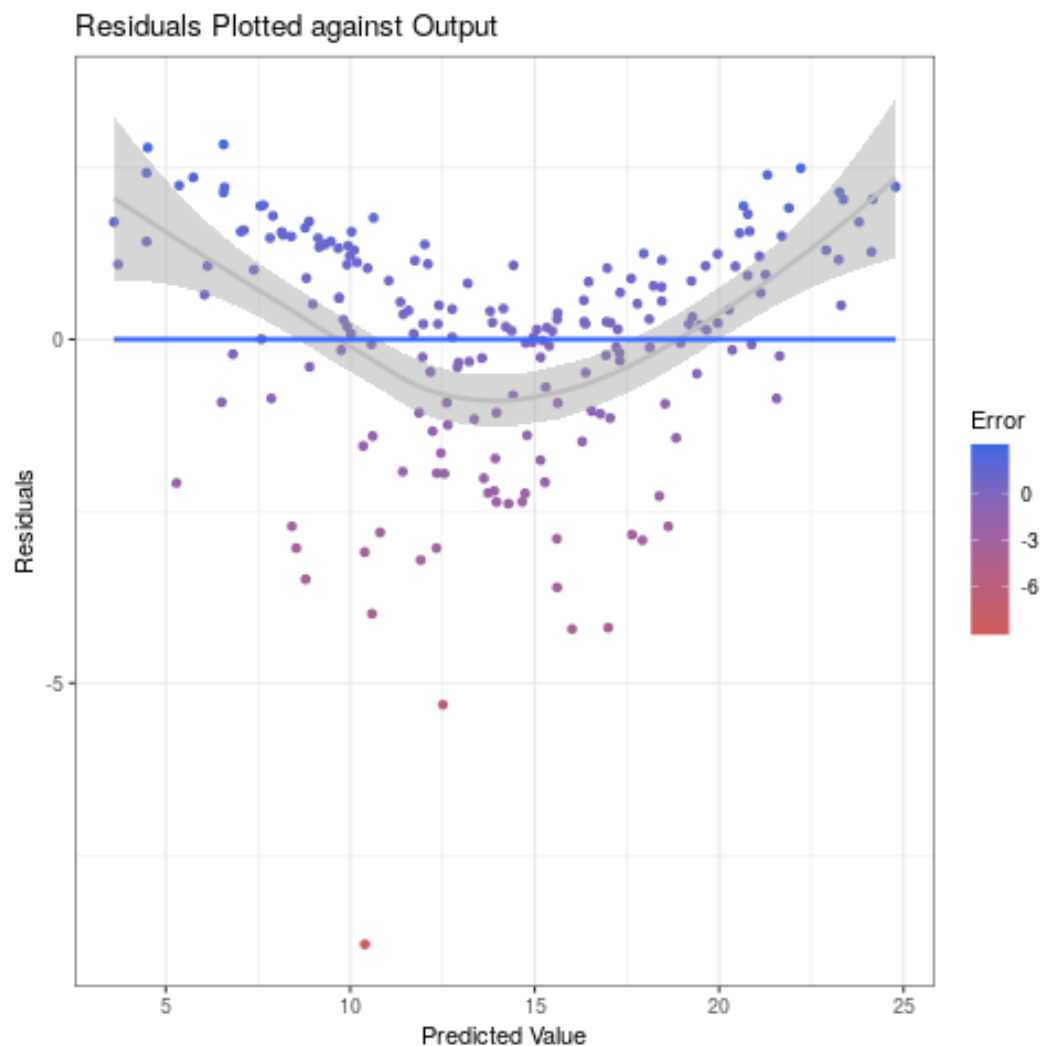
```
##      Input  Output      Error
## 1 19.01093 20.55546  1.5445354
## 2 14.29072 12.34536 -1.9453623
## 3 15.37404 12.33702 -3.0370177
## 4 16.73423 17.61712  0.8828840
## 5 13.54782 13.22391 -0.3239081
## 6 17.82417 12.51208 -5.3120845
```

Predict will also return fitted values if no newdata is specified.

(g) Plot the residuals against the predicted values

1. ggplot

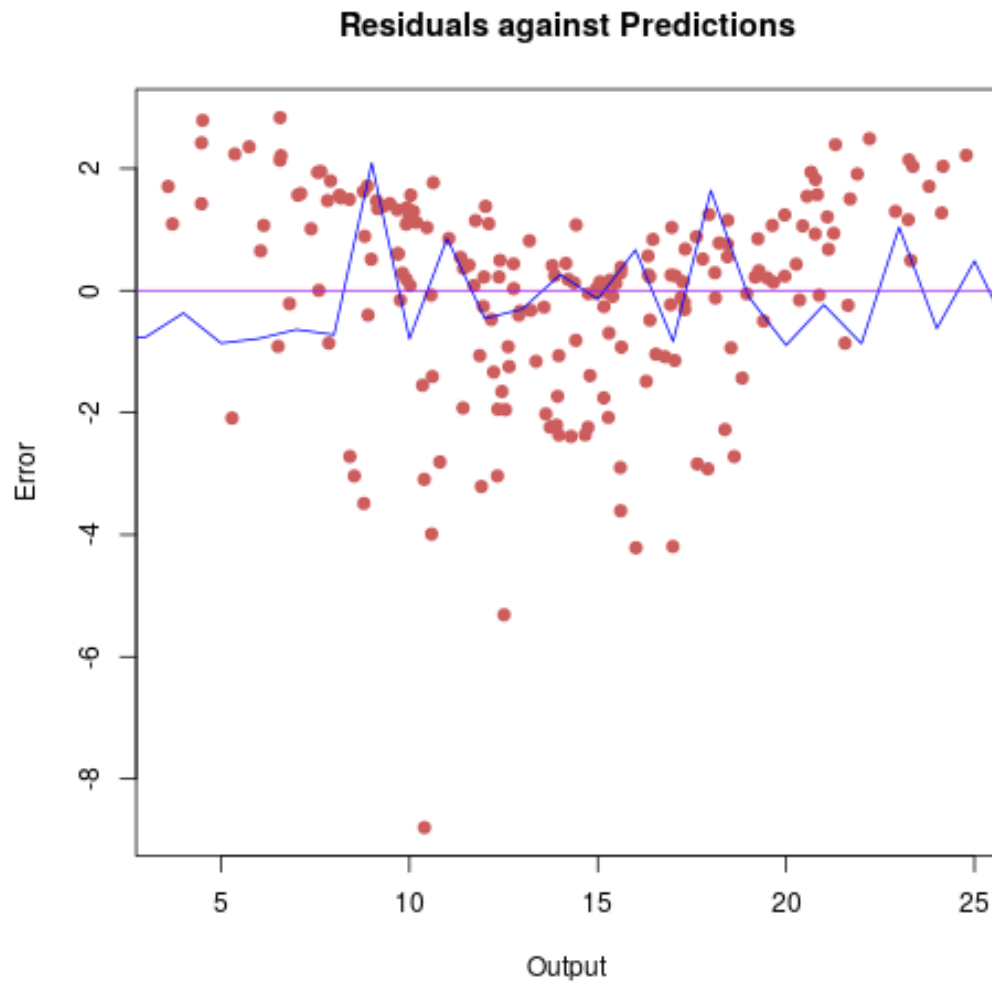
```
1 ggplot(data = ResDF, aes(x = Output, y = Error, col = Error )) +  
2   geom_point() +  
3   theme_bw() +  
4   stat_smooth(col = "grey")+  
5   stat_smooth(method = "lm", se = 0, ) +  
6   scale_color_gradient(low = "indianred", high = "royalblue") +  
7   labs(y = "Residuals", x = "Predicted Value", title = "Residuals  
   ↪ Plotted against Output")
```



```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

## 2. base

```
1 plot(Error ~ Output, data = ResDF, pch = 19, col = "IndianRed",  
  ↪   main = "Residuals against Predictions")  
2 smooth <- loess(formula = Error ~ Output, data = ResDF, span = 0.8)  
3 predict(smooth) %>% lines(col = "Blue")  
4 abline(lm(Error ~ Output, data = ResDF), col = "Purple")
```



</div>

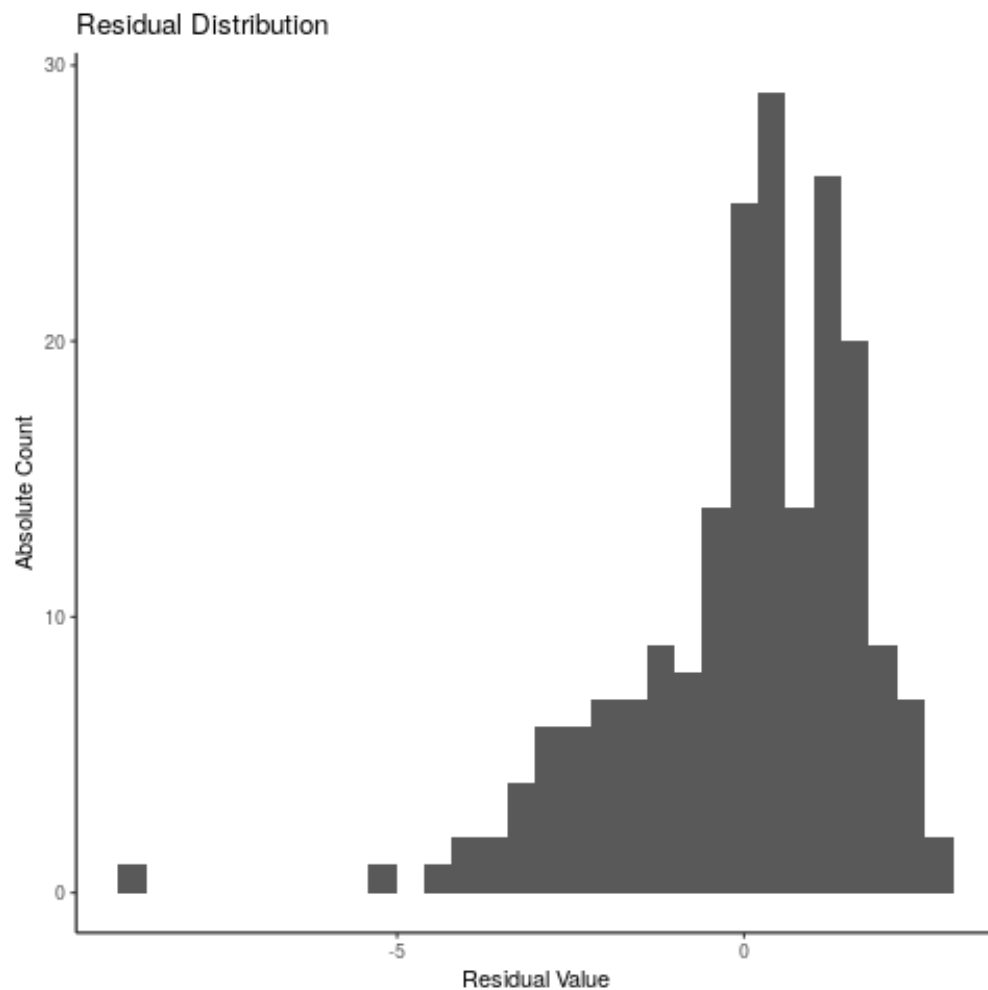
From this it can be seen that the residuals are centred around zero, but they are not normally distributed, the model performance appears to fail normality assumptions for values  $\in (10, 20)$

## (h) Plot the histogram of the residuals

1. ggplot

(a) Absolute Count

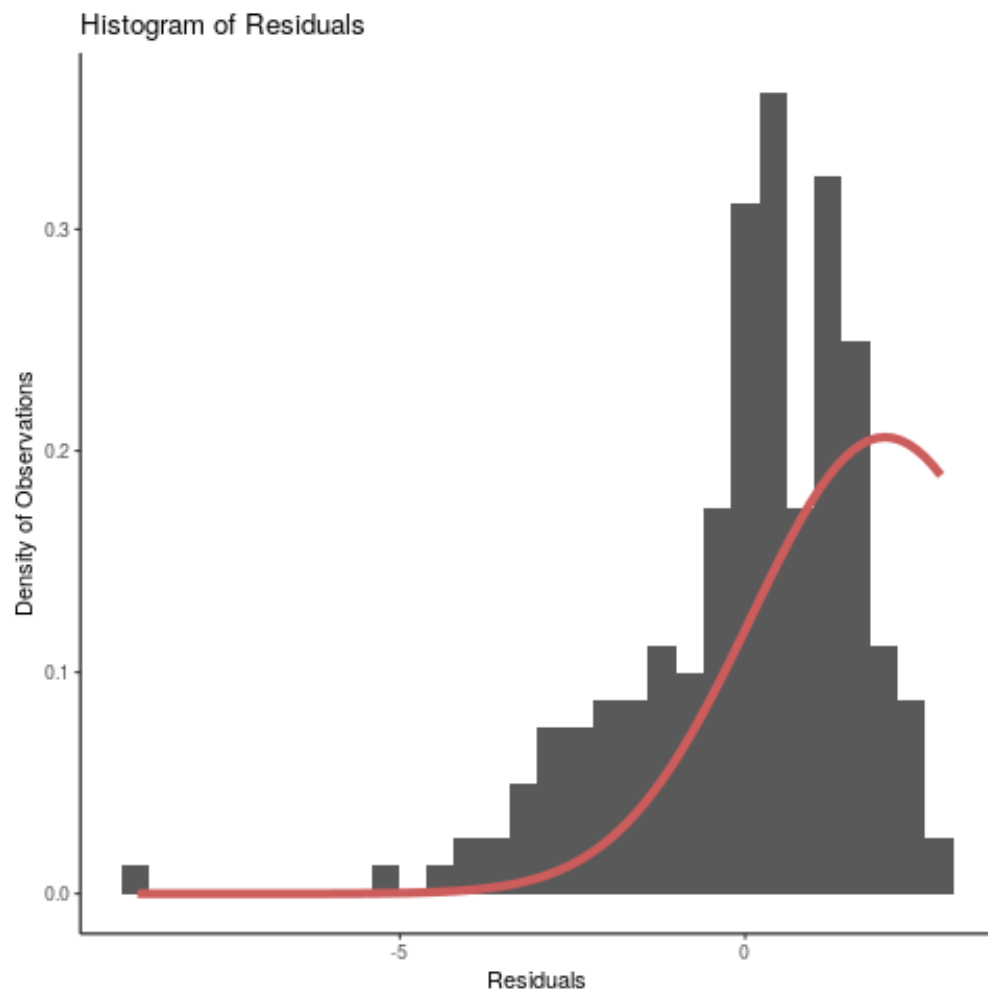
```
1 ggplot(data = ResDF, aes(x = Error, col = Output)) +  
2   geom_histogram() +  
3   theme_classic() +  
4   labs(y = "Absolute Count", x = "Residual Value", title =  
        ↪ "Residual Distribution")
```



## `stat\_bin()` using `bins = 30`. Pick better value with `binwidth`.

```
1 (setq org-complete-tags-always-offer-all-agenda-tags t)  
2 (setq org-fast-tag-selection-single-key nil)
```

```
1 df <- data.frame(x = rnorm(1000, 2, 2))
2
3 # overlay histogram and normal density
4 ggplot(ResDF, aes(x=Error)) +
5   geom_histogram(aes(y = stat(density))) +
6   stat_function(
7     fun = dnorm,
8     args = list(mean = mean(df$x), sd = sd(df$x)),
9     lwd = 2,
10    col = 'IndianRed'
11  ) +
12  theme_classic() +
13  labs(title = "Histogram of Residuals", y = "Density of
  ↳ Observations", x= "Residuals")
```



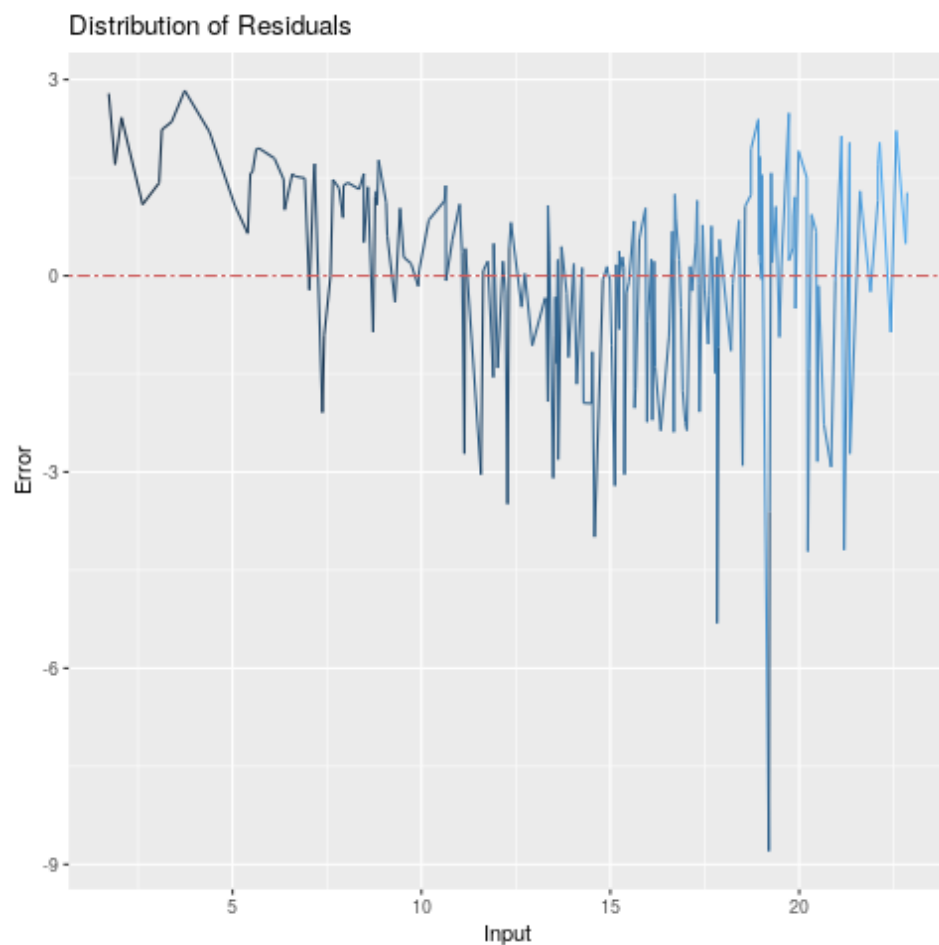
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



(c) White Noise We can visualise the residuals as white noise:

i. Our Residuals

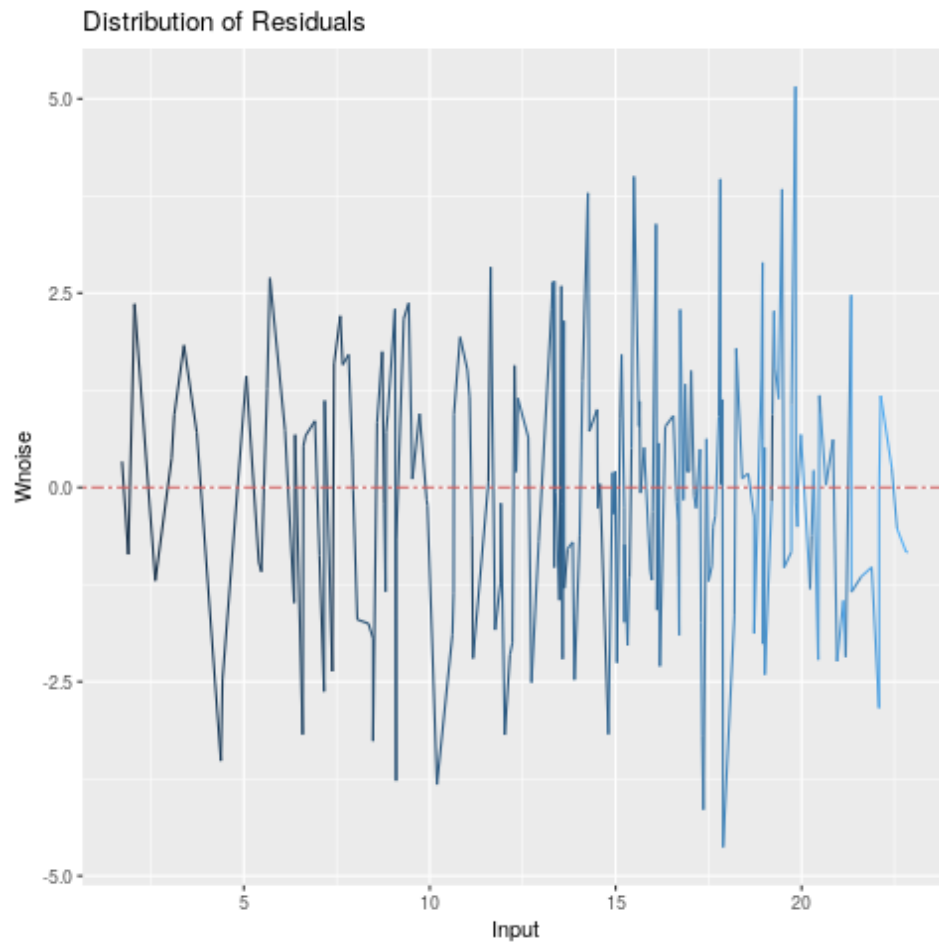
```
1 library(tidyverse)
2 ## Put some White Noise on the ResDF
3 ResDF <- cbind(ResDF, "Wnoise" = rnorm(nrow(ResDF), mean =
  ↪ 0, sd = sd(ResDF$Error)))
4 head(ResDF)
5
6
7 ## Our Residuals
8 ggplot(ResDF, aes(x = Input, y = Error, col = Output)) +
9   geom_line() +
10   geom_abline(slope = 0, intercept = 0, lty = "twodash",
  ↪   col = "IndianRed") +
11   theme(legend.position = "none") +
12   labs(title = "Distribution of Residuals")
```



```

1  ## White Noise
2
3  ggplot(ResDF, aes(x = Input, y = Wnoise, col = Output)) +
4    geom_line() +
5    geom_abline(slope = 0, intercept = 0, lty = "twodash",
6               ↪ col = "IndianRed") +
7    theme(legend.position = "none") +
8    labs(title = "Distribution of Residuals")

```



This clearly shows that the Residuals are non-normal.

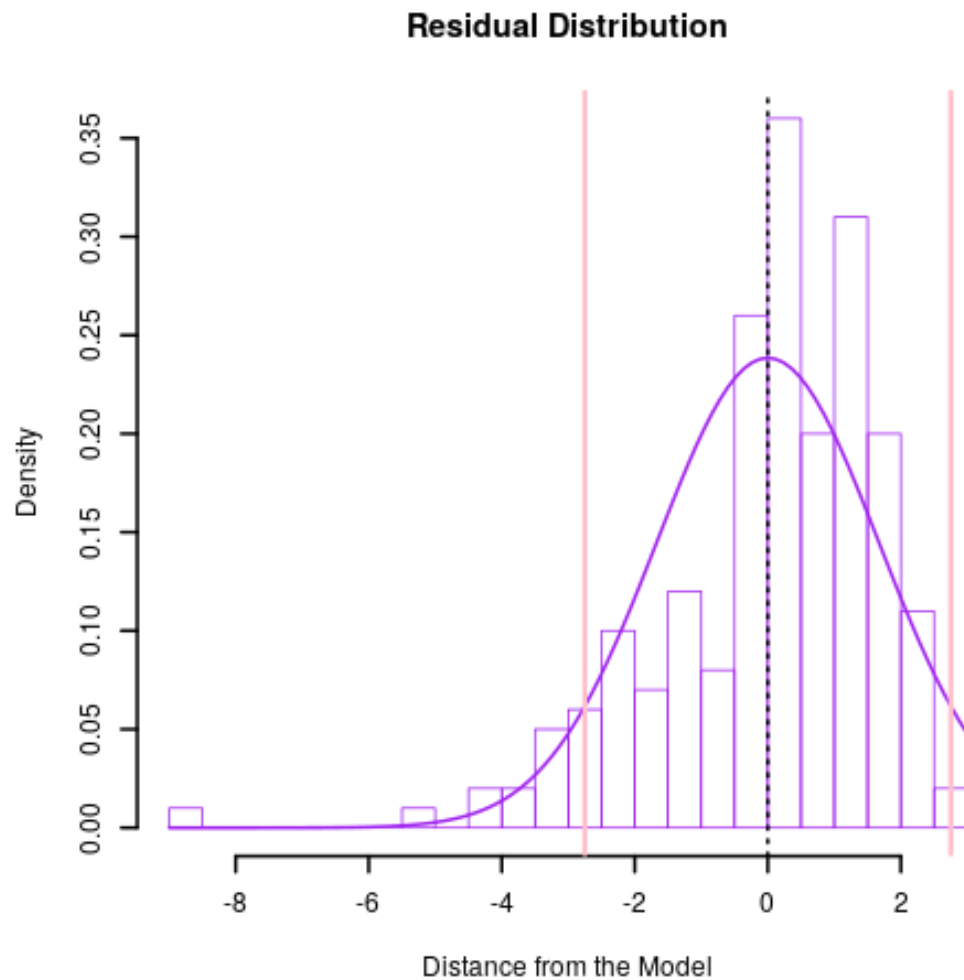
2. Base

GGPLOT2:HISTOGRAM:BASEPLOT

```

1 hist(ResDF$Error
2     , breaks = 30,
3     prob=TRUE,
4     lwd=2,
5     main = "Residual Distribution",
6     xlab = "Distance from the Model", border = "purple"
7     )
8
9     # Overlay the Normal Dist Curve
10 x <- 1:100    # Stupid base package wants some f(x), this is hacky
11             → but easier than stuffing around with lines or defining a
12             → function
13 curve(dnorm(x, mean(ResDF$Error), sd(ResDF$Error)), add=TRUE,
14       → col="purple", lwd=2) # Draws the actual density function
15
16 # lines(density(possiblevals_conf), col='purple', lwd=2) # Draws
17 → the observed density function
18 lwr_conf <- qnorm(p = 0.05, mean(ResDF$Error), sd(ResDF$Error))
19 upr_conf <- qnorm(p = 1-0.05, mean(ResDF$Error), sd(ResDF$Error))
20 abline(v=upr_conf, col='pink', lwd=3)
21 abline(v=lwr_conf, col='pink', lwd=3)
22 abline(v=mean(ResDF$Error), lwd=2, lty='dotted')

```



```
1 lwr_conf <- qnorm(p = 0.05, mean(ResDF$Error), sd(ResDF$Error))
2 upr_conf <- qnorm(p = 1-0.05, mean(ResDF$Error), sd(ResDF$Error))
```

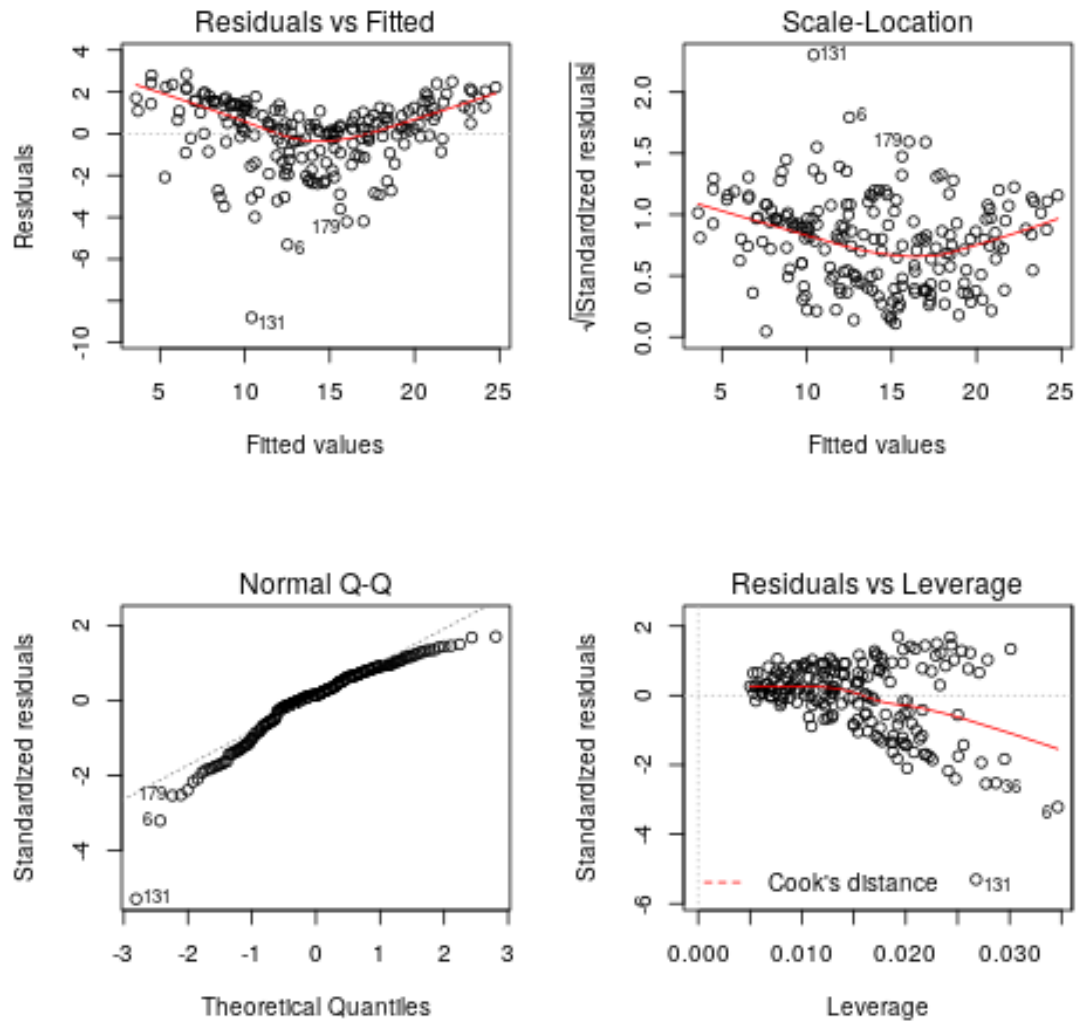
</div>

It hence appears that the Residuals skewed left, which suggests an upper bound on the residual value, the histogram is sufficiently non-normal to reject the assumption that the residuals are normally distributed

#### (i) Comment on the residual plots

The Residuals can be analysed by plotting the model:

```
1 layout(matrix(data = 1:4, nrow = 2))
2 plot(advModMultb1)
```



This is an inappropriate model because the residuals are non-normally distributed.

### (j) Use the multivariate model for predictions

```
1 newdata = data.frame("TV" = c(3, 1, 2), "Radio" = c(4, 5, 6))
2 Mod_Adv_Mult_b1 <- predict(object = advModMultb1, newdata)
3
4 mypreds <- data.frame("Input" = newdata, "Output" = Mod_Adv_Mult_b1)
5 mypreds # Careful with the names, they workout as Input.name in this
  ↪ case.
```

```
##   Input.TV Input.Radio   Output
## 1         3           4 3.810341
## 2         1           5 3.906826
## 3         2           6 4.140575
```

## Question 02: Non Linear Models: Use Advertising data set

---

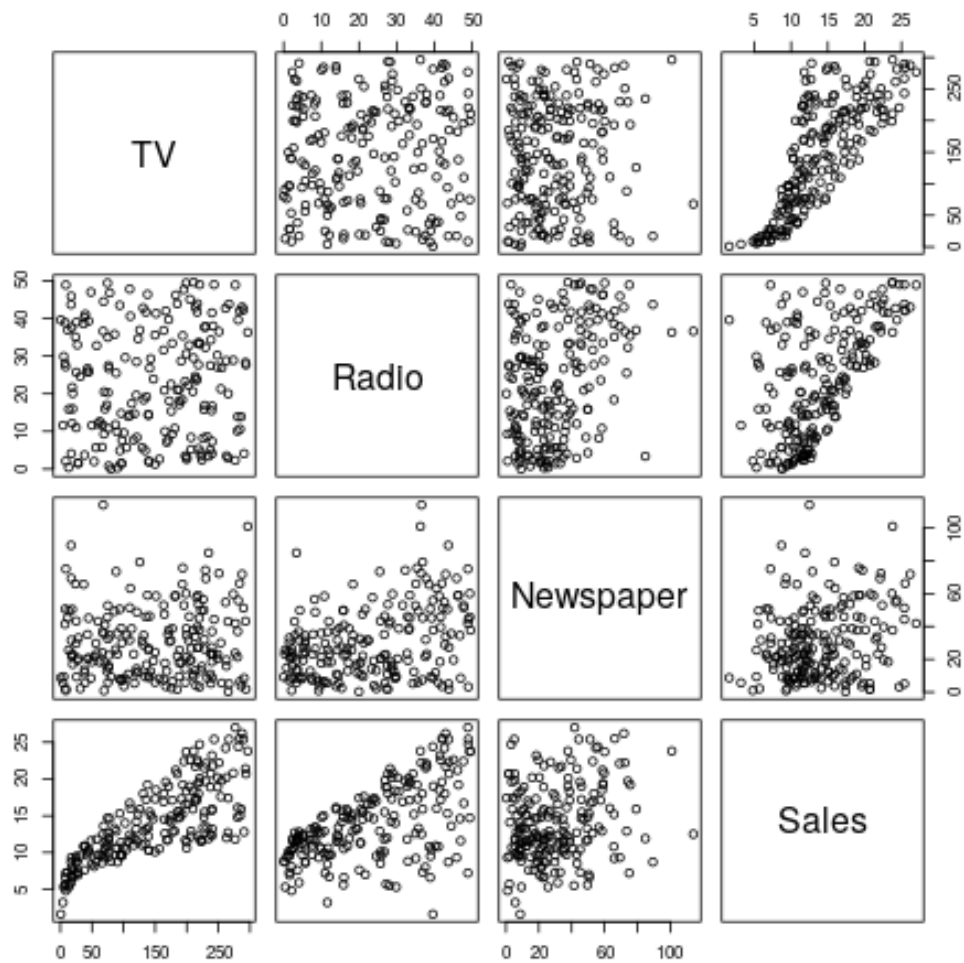
### (a) Add the Interaction Term TV\*Radio and test the significance of

the interaction term `{.tabset} :CUSTOM_ID: a-add-the-interaction-term-tvradio-and-test-the-significance-of-the-interaction-term-.tabset`

reconsider the correlationplots:

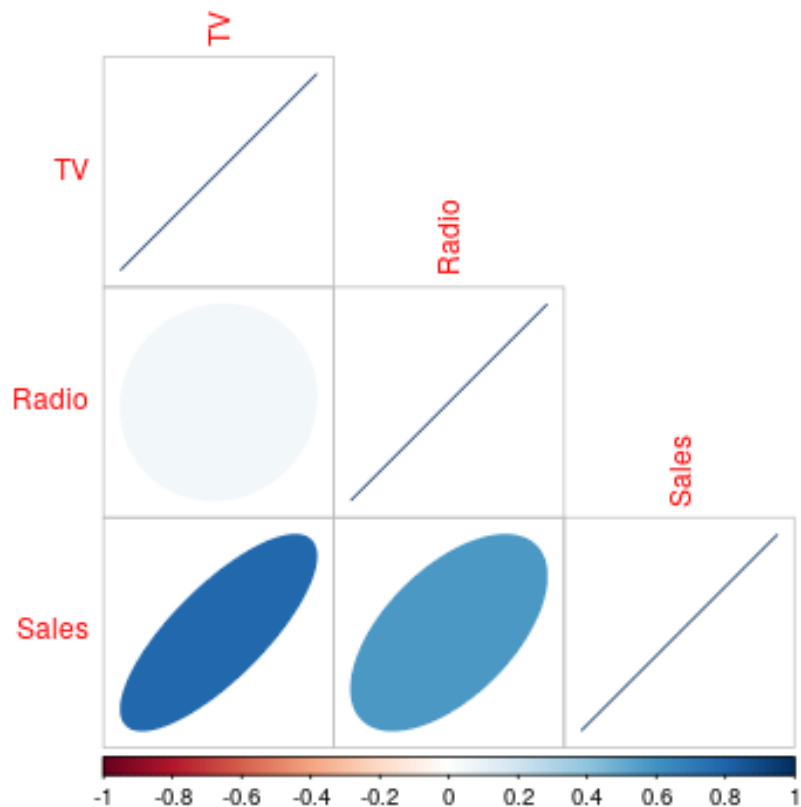
#### 1. Base

```
1 pairs(x = adv)
```



corrplot

```
1 adv[, names(adv)!="Newspaper"] %>% cor() %>% corrplot(method =  
  ↪ "ellipse", type = "lower")
```



</div>

From this we can determine that there is no interaction between TV and radio (however previously there was interaction between newspaper and Radio, we should consider that next)

1. Create the MLReg

```
1 int_mod_adv <- lm(formula = Sales ~ TV * Radio + TV + Radio, data =
  ↳ adv)
2 int_mod_adv
```

##

## Call:



```
## lm(formula = Sales ~ TV * Radio + TV + Radio, data = adv)
##
## Coefficients:
## (Intercept)          TV          Radio      TV:Radio
##      6.750220      0.019101      0.028860      0.001086
```

```
1 summary(int_mod_adv)
```

```
##
## Call:
## lm(formula = Sales ~ TV * Radio + TV + Radio, data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.3366 -0.4028  0.1831  0.5948  1.5246
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.750e+00  2.479e-01  27.233  <2e-16 ***
## TV           1.910e-02  1.504e-03  12.699  <2e-16 ***
## Radio        2.886e-02  8.905e-03   3.241   0.0014 **
## TV:Radio     1.086e-03  5.242e-05  20.727  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9435 on 196 degrees of freedom
## Multiple R-squared:  0.9678, Adjusted R-squared:  0.9673
## F-statistic: 1963 on 3 and 196 DF, p-value: < 2.2e-16
```

### (b) Give the resulting model after considering this interaction

term. :CUSTOM<sub>ID</sub>: b-give-the-resulting-model-after-considering-this-interaction-term.

The model will be of the form:

$$\text{Sales} = 0.0019 \times \text{TV} \times \text{Radio} + 0.002886 \times \text{TV} + 0.001086 \times \text{Radio}$$

Note that all the terms of the model are significant, hence we deem the model as adequate.

### (c) Construct the Polynomial Regression Model of order 3 and test

the model significance :CUSTOM<sub>ID</sub>: c-construct-the-polynomial-regression-model-of-order-3-and-test-the-model-significance

When creating polynomial models there are raw and orthogonal polynomials,

- A raw polynomial will be the the standard System of equations that you get by minimising the RSS

- The problem with this is that the values will be correlated with each other
- An orthogonal polynomial is a transformed but equivalent polynomial
  - The advantage to this is that the p-value's will be more meaningful because the coefficients aren't correlated with each other
  - The disadvantage is that the values are not directly related to our data and are not hence meaningful.

```

1  # Orthogonal (fixes the correlation between the coefficients))
2
3  polymodel_Orth <- lm(Sales ~ poly(x = TV, degree = 3, raw = FALSE),
   ↪ data = adv)
4  summary(polymodel_Orth)

```

```

##
## Call:
## lm(formula = Sales ~ poly(x = TV, degree = 3, raw = FALSE), data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.9734 -1.8900 -0.0897  2.0189  7.3765
##
## Coefficients:
##                                Estimate Std. Error t value
## (Intercept)                   14.0225     0.2286  61.353
## poly(x = TV, degree = 3, raw = FALSE)1  57.5727     3.2322  17.812
## poly(x = TV, degree = 3, raw = FALSE)2  -6.2288     3.2322  -1.927
## poly(x = TV, degree = 3, raw = FALSE)3   4.0074     3.2322   1.240
##                                Pr(>|t|)
## (Intercept)                   <2e-16 ***
## poly(x = TV, degree = 3, raw = FALSE)1  <2e-16 ***
## poly(x = TV, degree = 3, raw = FALSE)2   0.0554 .
## poly(x = TV, degree = 3, raw = FALSE)3   0.2165
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.232 on 196 degrees of freedom
## Multiple R-squared:  0.622, Adjusted R-squared:  0.6162
## F-statistic: 107.5 on 3 and 196 DF, p-value: < 2.2e-16

```

```

1  # Pure/Raw (has the advantage that you can interpret the coefficients,
   ↳ but the coefficients will depend on each other and be hence
   ↳ correlated.)
2
3  polymodel <- lm(Sales ~ I(TV*TV*TV) + I(TV*TV) + (TV), data = adv)
4  summary(polymodel)      # I is used to inhibit the interpretation of * as
   ↳ relating to the model,

```

```

##
## Call:
## lm(formula = Sales ~ I(TV * TV * TV) + I(TV * TV) + (TV), data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.9734 -1.8900 -0.0897  2.0189  7.3765
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5.420e+00  8.641e-01   6.272 2.23e-09 ***
## I(TV * TV * TV) 5.572e-07  4.494e-07   1.240 0.216519
## I(TV * TV)     -3.152e-04  2.022e-04  -1.559 0.120559
## TV             9.643e-02  2.580e-02   3.738 0.000243 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.232 on 196 degrees of freedom
## Multiple R-squared:  0.622, Adjusted R-squared:  0.6162
## F-statistic: 107.5 on 3 and 196 DF,  p-value: < 2.2e-16

```

```

1                                     # Instead of representing interaction it
   ↳ represents TV^3
2
3  polymodel_Raw<- lm(Sales ~ poly(x = TV, degree = 3, raw = TRUE), data =
   ↳ adv)
4  summary(polymodel_Raw)

```

```

##
## Call:
## lm(formula = Sales ~ poly(x = TV, degree = 3, raw = TRUE), data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.9734 -1.8900 -0.0897  2.0189  7.3765

```

```
##
## Coefficients:
##
##              Estimate Std. Error t value
## (Intercept)      5.420e+00  8.641e-01   6.272
## poly(x = TV, degree = 3, raw = TRUE)1  9.643e-02  2.580e-02   3.738
## poly(x = TV, degree = 3, raw = TRUE)2 -3.152e-04  2.022e-04  -1.559
## poly(x = TV, degree = 3, raw = TRUE)3  5.572e-07  4.494e-07   1.240
##
##              Pr(>|t|)
## (Intercept)      2.23e-09 ***
## poly(x = TV, degree = 3, raw = TRUE)1  0.000243 ***
## poly(x = TV, degree = 3, raw = TRUE)2  0.120559
## poly(x = TV, degree = 3, raw = TRUE)3  0.216519
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.232 on 196 degrees of freedom
## Multiple R-squared:  0.622, Adjusted R-squared:  0.6162
## F-statistic: 107.5 on 3 and 196 DF,  p-value: < 2.2e-16
```

The 3rd degree coefficient is not significant, hence we consider the 2nd degree:

```
1 quadmod <- lm(formula = Sales ~ I(TV*TV) + TV, data = adv)
2 summary(quadmod)
```

```
##
## Call:
## lm(formula = Sales ~ I(TV * TV) + TV, data = adv)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.6844 -1.7843 -0.1562  2.0088  7.5097
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.114e+00  6.592e-01   9.275  < 2e-16 ***
## I(TV * TV)  -6.847e-05  3.558e-05  -1.924   0.0557 .
## TV           6.727e-02  1.059e-02   6.349 1.46e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.237 on 197 degrees of freedom
## Multiple R-squared:  0.619, Adjusted R-squared:  0.6152
## F-statistic: 160.1 on 2 and 197 DF,  p-value: < 2.2e-16
```

The quadratic term is not sufficiently significant so the model is rejected

**(d) Give the resulting selected model**

The model selected is the multiple linear regression with the interaction from before:

$$\text{Sales} = 0.0019 \times \text{TV} \times + \text{Radio} \times 0.002886 + \text{TV} \times \text{Radio} \times 0.001086$$

This was converted from 'md' to 'org' using 'pandoc -f gfm' at time: 2020-02-08T05-00-13