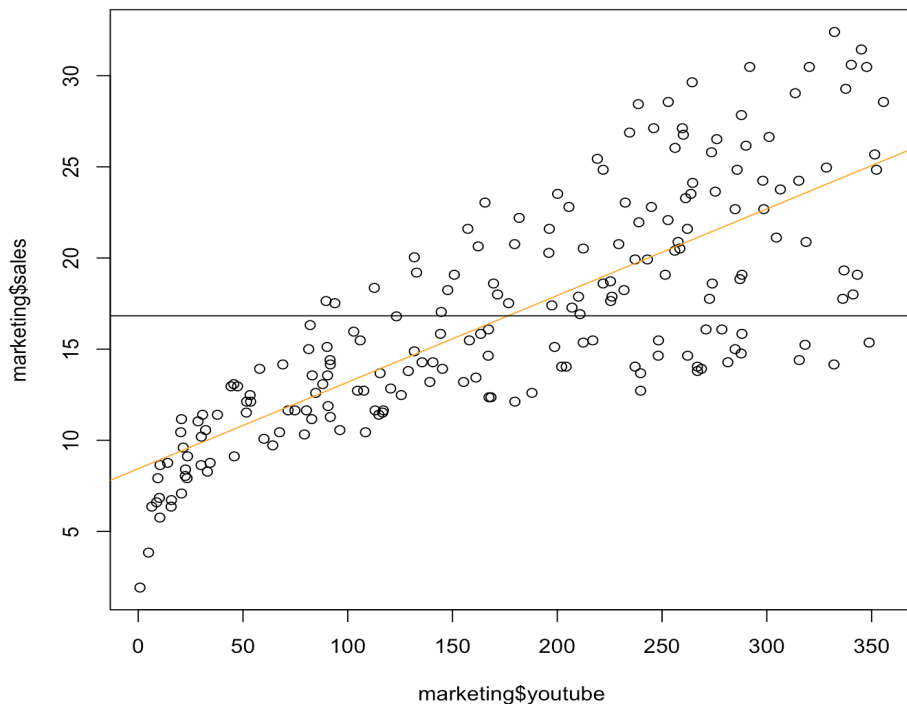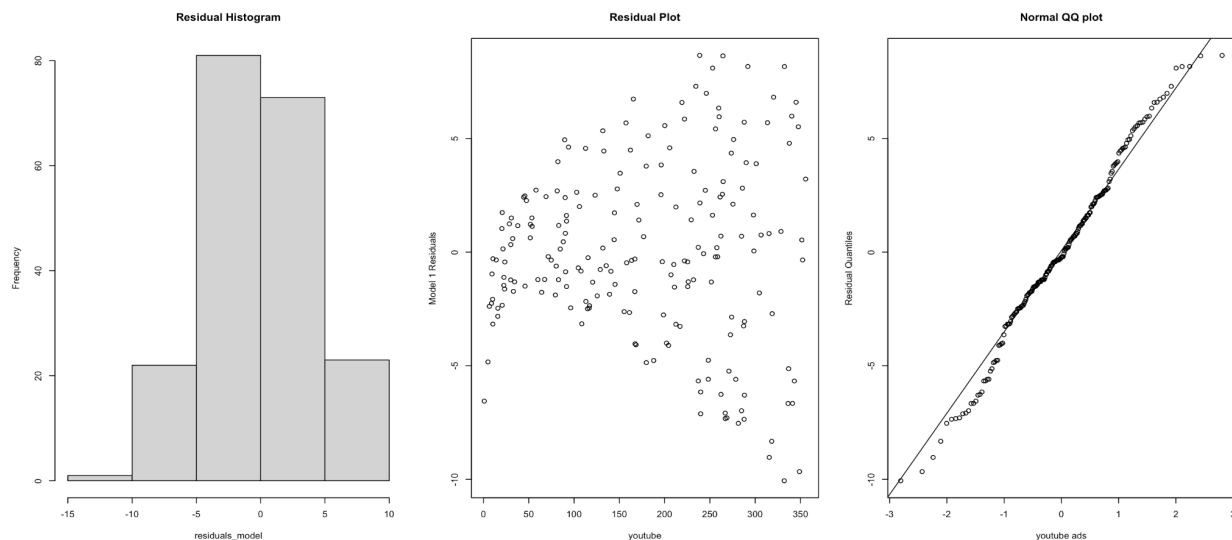During the analysis of the data set "marketing" many conclusions were found, specifically on the predictors included in the set and sales prediction. The results of the findings include a simple linear regression model as well as a multiple linear regression model that can be used to accurately predict sales.

The most significant predictor was chosen to be used for the simple linear regression model. All three were tested, and the YouTube predictor was found to be the most significant. A model was made using each of the predictors so the accuracy of each model could be tested. Each model's adjusted R squared value and p-value was used to determine which model best fit the data. The model using Youtube as the predictor had the highest adjusted R squared value, which was 0.6099, and a p value of almost 0. Below is the scatter plot of Sales vs. YouTube Advertising:



According to the chosen model, the linear function for the expected value of sales is $\hat{y} = 8.439112 + 0.047537x$. That is the equation of the orange line present in the scatter plot. The intercept, 8.439112, means that without the predictor variable there is expected to be $8,439

in sales. As YouTube advertisements are increased by $1, there can be an expected $47 increase in sales. The regression standard error of the model is 0.3429913, meaning the difference between actual data points and ones predicted by the model is fairly small. The coefficient of determination value is 0.6118751, so about 61% of the variation in sales is able to be explained through the model, which indicates a relatively good fit. Because the slope value is positive and relatively significant, it can be concluded there is a positive relationship with moderate strength between the two variables. Below is the residual plot, histogram of residuals, and QQ plot:
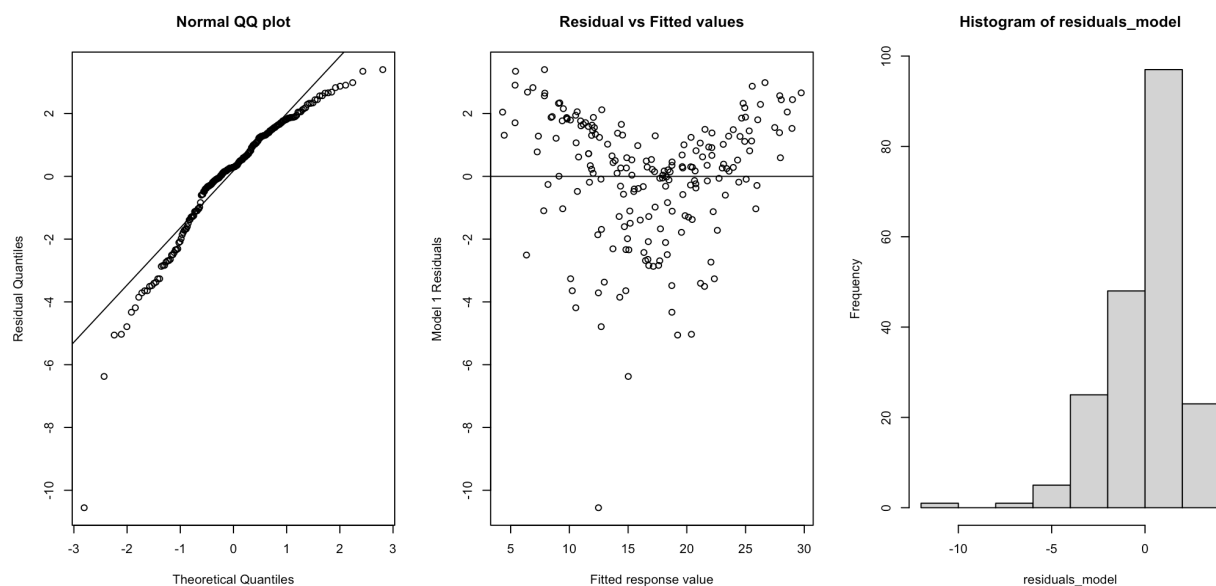


The residual plot and residual histogram do not suggest the model fits the data well at smaller values due to the clustering in the beginning of the plot. The normal QQ plot does suggest the residuals are roughly distributed aside from the beginning of the plot. These three models suggest that the model is more accurate for x values that are 200 and fits the model much less accurately in the lower bounds of the x values. Ultimately, they do still support that the model fits the majority of the data.

As for the multiple linear regression model, it was decided to drop the newspaper predictor. Several steps went into making that decision. First, a model was made with all three predictors. The model had an adjusted R squared value of 0.8956, so a large portion of the variability in sales is accounted for in the model. However, the newspaper predictor has a noticeably high p-value of 0.86 and its coefficient is almost 0 (-0.0010). This indicates that the predictor is not very significant to the model. YouTube and Facebook both have very low p-values, more significant coefficients, and three stars in the summary of the model. A model was created without the newspaper predictor and it had a higher adjusted R squared value, which

was 0.8956. This value being higher suggests the model is a better fit of the data. The linear equation for the model is $E(sales) = 3.50532 + 0.04575(youtube) + 0.18799(facebook)$. So, without advertising on YouTube or Facebook $3,505 in sales can be expected. For every dollar spent on YouTube advertisements, sales are expected to increase $45.75. For every dollar spent on Facebook advertisement, sales are expected to increase about $188. The regression standard error of the model is 2.018, meaning that the observed sale values are roughly $2,018 away from the predicted values. Both of these values lend support to the model. Below are the plots used for model evaluation and residual analysis:
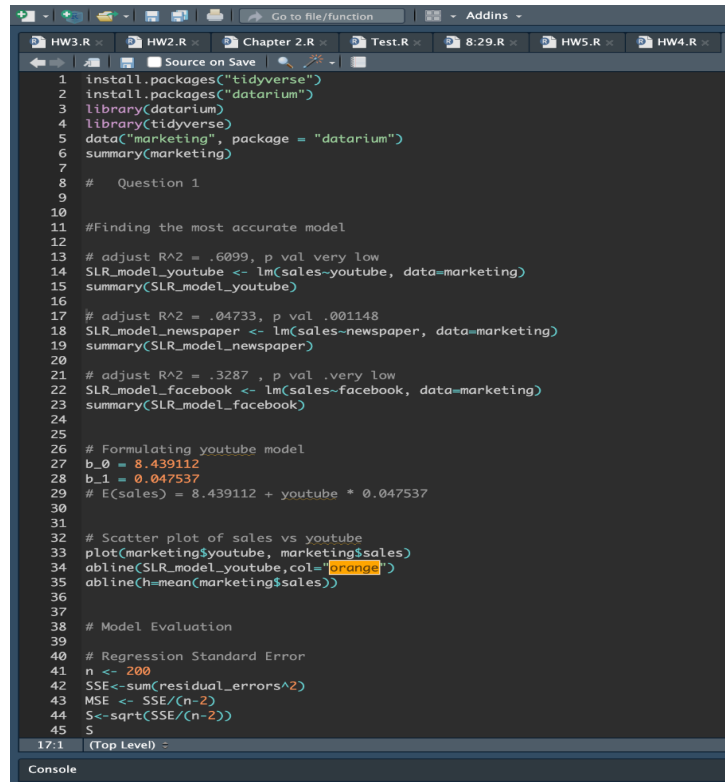


All of these plots ultimately support the model, but there are some flaws. The QQ plot indicates that the residuals are roughly normally distributed. The residual plot has some outliers and could have less clustering, but they are distributed random enough. This histogram is also a little skewed but is roughly normal.

Both the models fit the data well, but because the multiple linear regression model has a higher adjusted R squared and lower residual standard error it is reasonable to say it fits the data better. Using that model, we can be 95% confident that the expected sales are between $13,242 and $14,255 given $50.90 went to YouTube advertisements, $42.10 were used for advertisements on Facebook, and $87.60 were used to advertise in the newspaper. Likewise, we can be 95% confident that actual sales will be between $9,737 and $17,759 when there are  $50.90 going

toward YouTube advertisements, $42.10 are used for Facebook advertisements, and $87.60 are used to advertise in the newspaper.

# Appendix

```r
1   install.packages("tidyverse")
2   install.packages("datarium")
3   library(datarium)
4   library(tidyverse)
5   data("marketing", package = "datarium")
6   summary(marketing)
7
8   #   Question 1
9
10
11  #Finding the most accurate model
12
13  # adjust R^2 = .6099, p val very low
14  SLR_model_youtube <- lm(sales~youtube, data=marketing)
15  summary(SLR_model_youtube)
16
17  # adjust R^2 = .04733, p val .001148
18  SLR_model_newspaper <- lm(sales~newspaper, data=marketing)
19  summary(SLR_model_newspaper)
20
21  # adjust R^2 = .3287 , p val .very low
22  SLR_model_facebook <- lm(sales~facebook, data=marketing)
23  summary(SLR_model_facebook)
24
25
26  # Formulating youtube model
27  b_0 = 8.439112
28  b_1 = 0.047537
29  # E(sales) = 8.439112 + youtube * 0.047537
30
31
32  # Scatter plot of sales vs youtube
33  plot(marketing$youtube, marketing$sales)
34  abline(SLR_model_youtube,col="orange")
35  abline(h=mean(marketing$sales))
36
37
38  # Model Evaluation
39
40  # Regression Standard Error
41  n <- 200
42  SSE<-sum(residual_errors^2)
43  MSE <- SSE/(n-2)
44  S<-sqrt(SSE/(n-2))
45  S
```

```r
38    # Model Evaluation
39
40    # Regression Standard Error
41    n <- 200
42    SSE<-sum(residual_errors^2)
43    MSE <- SSE/(n-2)
44    S<-sqrt(SSE/(n-2))
45    S
46
47    # Coefficient of Determination
48    y_bar <- mean(marketing$sales)
49    fitted_values<-fitted.values(SLR_model_youtube)
50    SSR<-sum((fitted_values-y_bar)^2)
51    SST<-sum((marketing$sales-y_bar)^2)
52    R_squrd<-SSR/SST
53    R_squrd
54
55
56    # Residual Plot Against YouTube
57    residuals_model <- residuals(SLR_model_youtube)
58    plot(marketing$youtube, residuals_model, xlab="youtube",
59          ylab="Model 1 Residuals", col="black", main = "Residual Plot")
60
61    # Histogram of Residuals
62    hist(residuals_model, breaks=6, main = "Residual Histogram")
63
64    # QQ Plot
65    qqnorm(residuals_model, main="Normal QQ plot", xlab="youtube ads",
66          ylab="Residual Quantiles")
67    qqline(residuals_model)
68    par(mfrow=c(1,3))
69
70    #   Question 2
71
72
73    #Fitting the model with all 3 predictors
74    MLR_model <- lm(sales~youtube+facebook+newspaper, data = marketing)
75
76    # p val of newspaper is high at .86 and its coefficient is almost 0 (-0.001037)
77    # has 0 stars as well
78    # has lower adjusted r squared than model with dropped predictor
79    summary(MLR_model)
80
81    MLR_model2 <- lm(sales~youtube+facebook, data = marketing)
82
```

17:1   (Top Level) ≑

Console

```r
81    MLR_model2 <- lm(sales~youtube+facebook, data = marketing)
82
83    # Model Evaluation
84    summary(MLR_model2)
85    # s = 2.018; observed sales is roughly $2,018 away from the predicted values
86    # adj r squared = 0.8962
87
88    # Residual Analysis
89    residuals_model<-residuals(MLR_model2)
90    fitted_model<-predict(MLR_model2)
91
92    plot(fitted_model, residuals_model, xlab="Fitted response value",
93          ylab="Model 1 Residuals", col="black", main= "Residual vs Fitted values")
94    abline(h=0)
95
96    hist(residuals_model, breaks=6)
97    # Adds box to histogram
98    box()
99    qqnorm(residuals_model, main="Normal QQ plot", xlab="Theoretical Quantiles",
100          ylab="Residual Quantiles")
101    qqline(residuals_model)
102
103
104   #   Question 3
105   newdat <- data.frame(youtube = 50.9, facebook = 42.1, newspaper = 87.6)
106   CI <- predict(MLR_model2, newdat, se.fit = TRUE, interval = "confidence",
107                 level = 0.95)
108   CI
109   PI <- predict(MLR_model2, newdat, se.fit = TRUE, interval = "prediction",
110                 level = 0.95)
111   PI
```

110:28  (Top Level) ≑

Console

This is my own unaided work

Daria Casey