

Stat Project 1

Trey Hamilton, Maxwell Levinson, Greg Madden, and Andrew Setaro

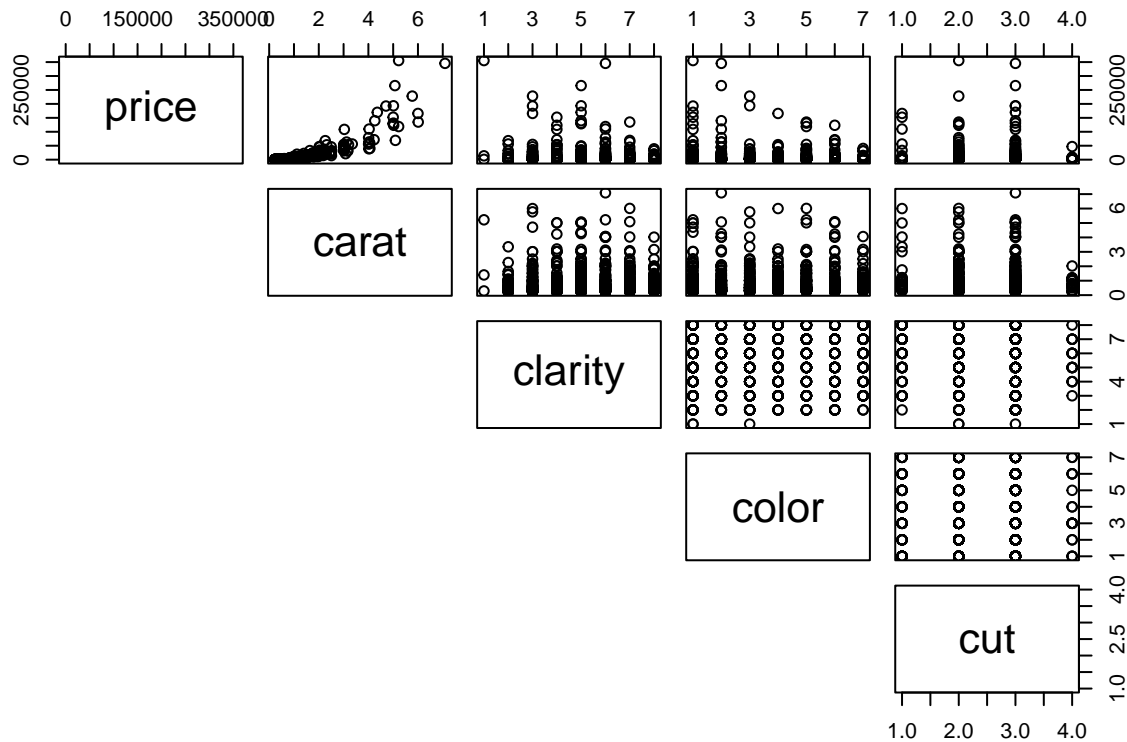
4/14/2022

We have been approached by Blue Nile to perform the following tasks:

1. Use data visualizations to explore how price is related to the other variables (carat, clarity, color, cut), as well as how the other variables may relate to each other. Address the various claims on the diamond education page on Blue Nile.
2. Fit an appropriate simple linear regression for price against carat.

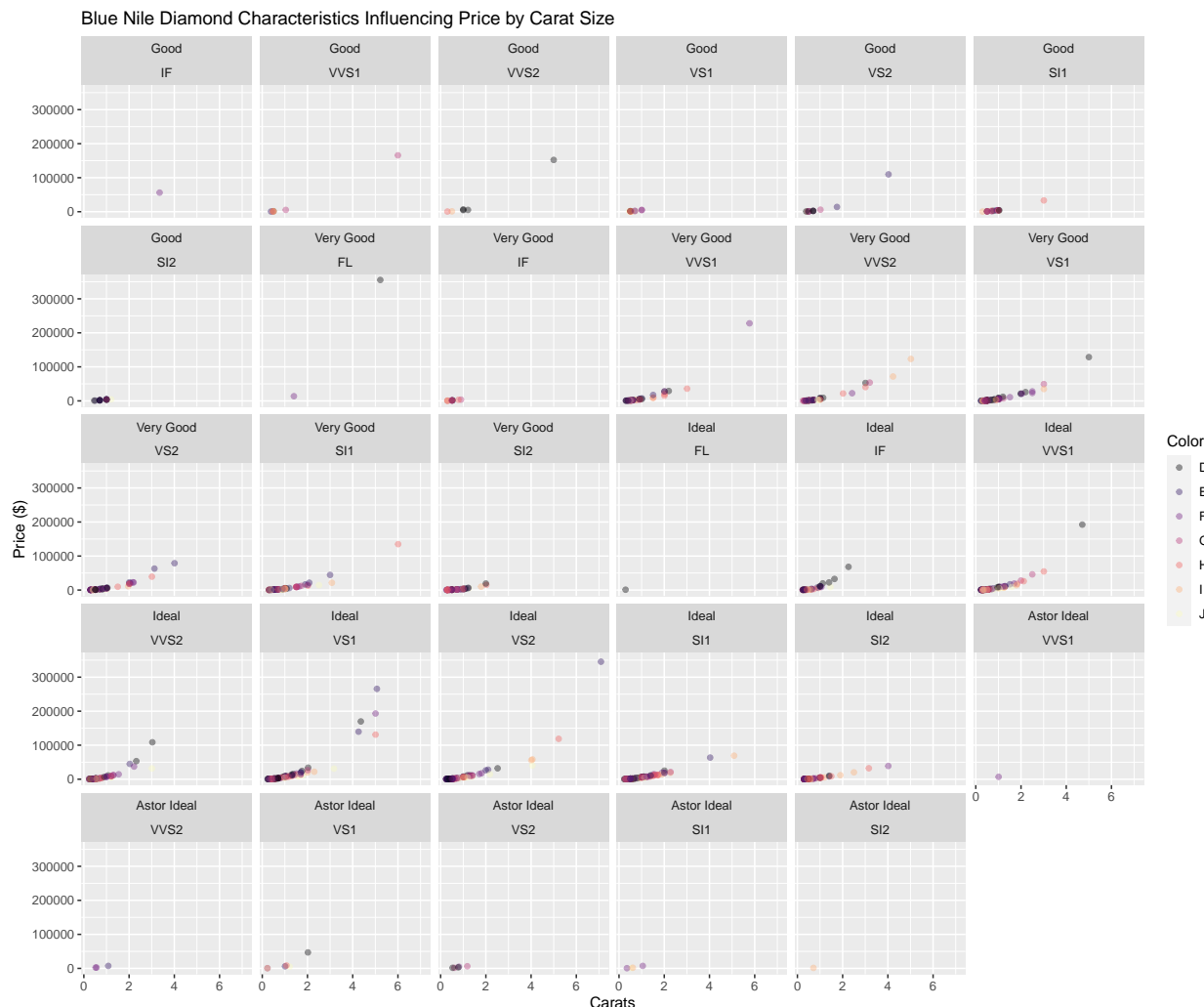
Description of the Data and Variables:

Plotting a pair-wise scatterplots for the data below:



Note that the x-axis above corresponds to increasing desirability of the factored categorical variables: clarity, color, and cut. In the above scatterplot matrix how price appears to have the clearest linear relationship with carats.

Plotting the response variable of interest (price) against carat, faceted by cut and clarity, with color indicating different diamond color characteristics:



In the above plot, you can see that certain cut/clarity combinations (e.g., Good cut and VVS2 clarity) have different slopes in terms of their price ~ carat relationships. For example, ideal cut with FL (Flawless) clarity appears to have a higher price per additional carat size than the Ideal Cut with SI2 (Slightly Included). In addition, more desirable colors (i.e. D-F) appear to cluster at the lower price ranges.

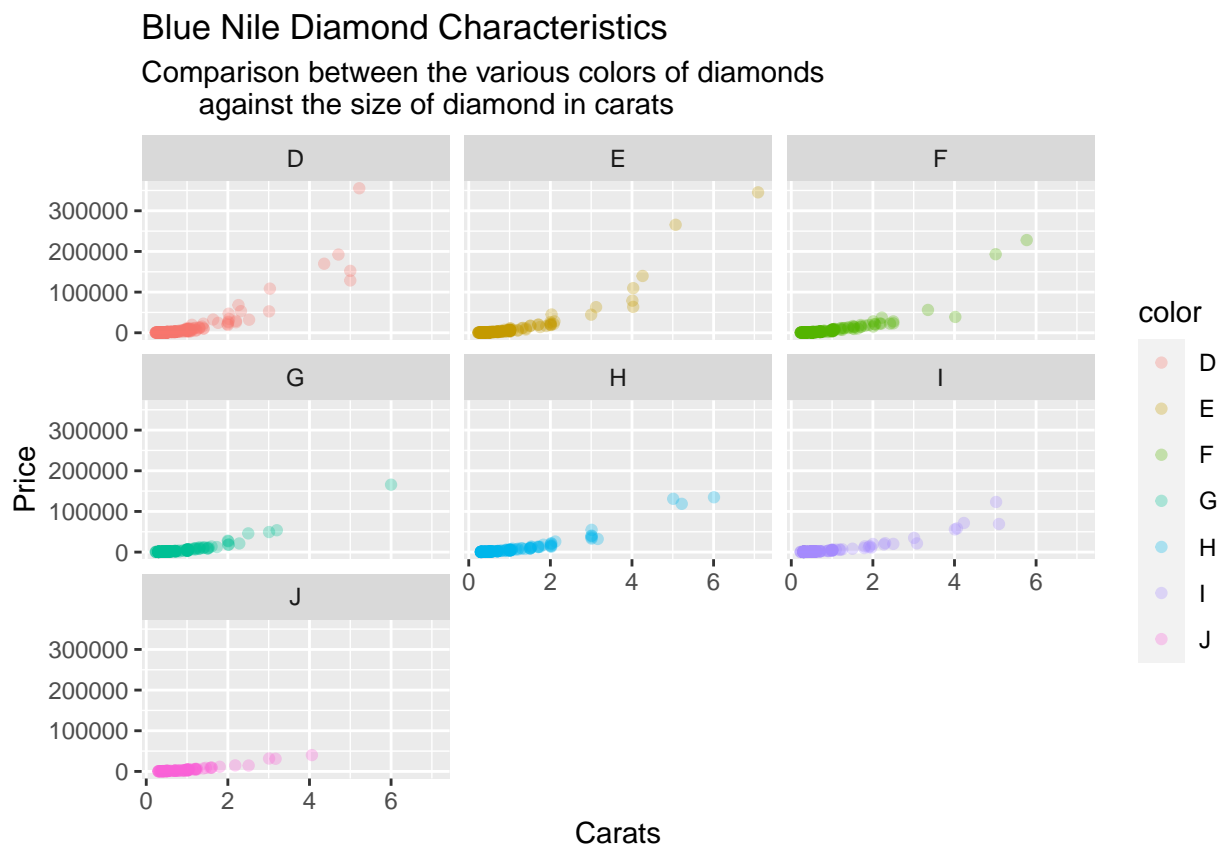
- Color: <https://www.bluenile.com/education/diamonds/color> + “Did you know that it’s very rare to find a diamond that doesn’t have any color at all? In fact, diamonds are found in almost any naturally occurring color, including gray, white, yellow, green, brown, and pink. The absence of color in a diamond is the rarest and therefore, the most expensive. While the majority of our customers choose a D or E color grade, many go with a beautiful near colorless grade to make the most of their budget and allocate more on the best cut that they can afford (which gives them more sparkle).”

Deciding whether or not you want to spend more on diamond color grade is partly related to the size and

shape of the diamond that you are considering, and your setting preference. You can save by knowing how color affects these attributes.”

Assertion above is that when choosing a specific color grade of a diamond, it is partly related to the size of the diamond so lets compare the various sizes of diamonds across all of the color ranges and see what assumptions we can make: `white_check_mark` eyes raised_hands

```
Data %>%
  ggplot(aes(x=carat, y = price, color = color)) +
  geom_point(alpha = 0.3) +
  labs(x="Carats", y="Price",
       title = "Blue Nile Diamond Characteristics",
       subtitle = "Comparison between the various colors of diamonds
                  against the size of diamond in carats")+
  facet_wrap(~color)
```

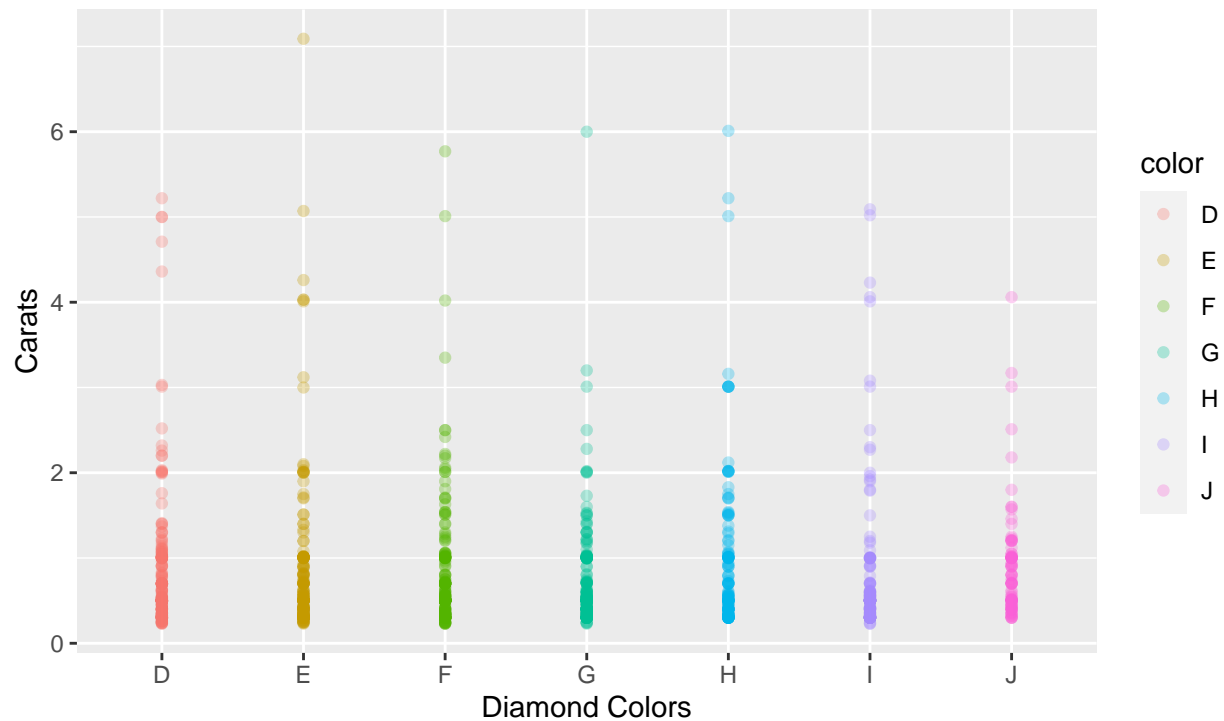


The color of a diamond is more visible in larger diamonds. As shown above, when choosing the size of the diamond, color is a relative factor as color groups ‘D’ and ‘E’ which are considered to be colorless and provide the closest look to ‘icy’ are larger in size.

```
Data %>%
  ggplot(aes(x=color, y = carat, color = color)) +
  geom_point(alpha = 0.3) +
  labs(x="Diamond Colors", y="Carats",
       title = "Blue Nile Diamond Characteristics",
       subtitle = "Comparison between the various colors of diamonds
                  against the size of diamond in carats")
```

Blue Nile Diamond Characteristics

Comparison between the various colors of diamonds
against the size of diamond in carats

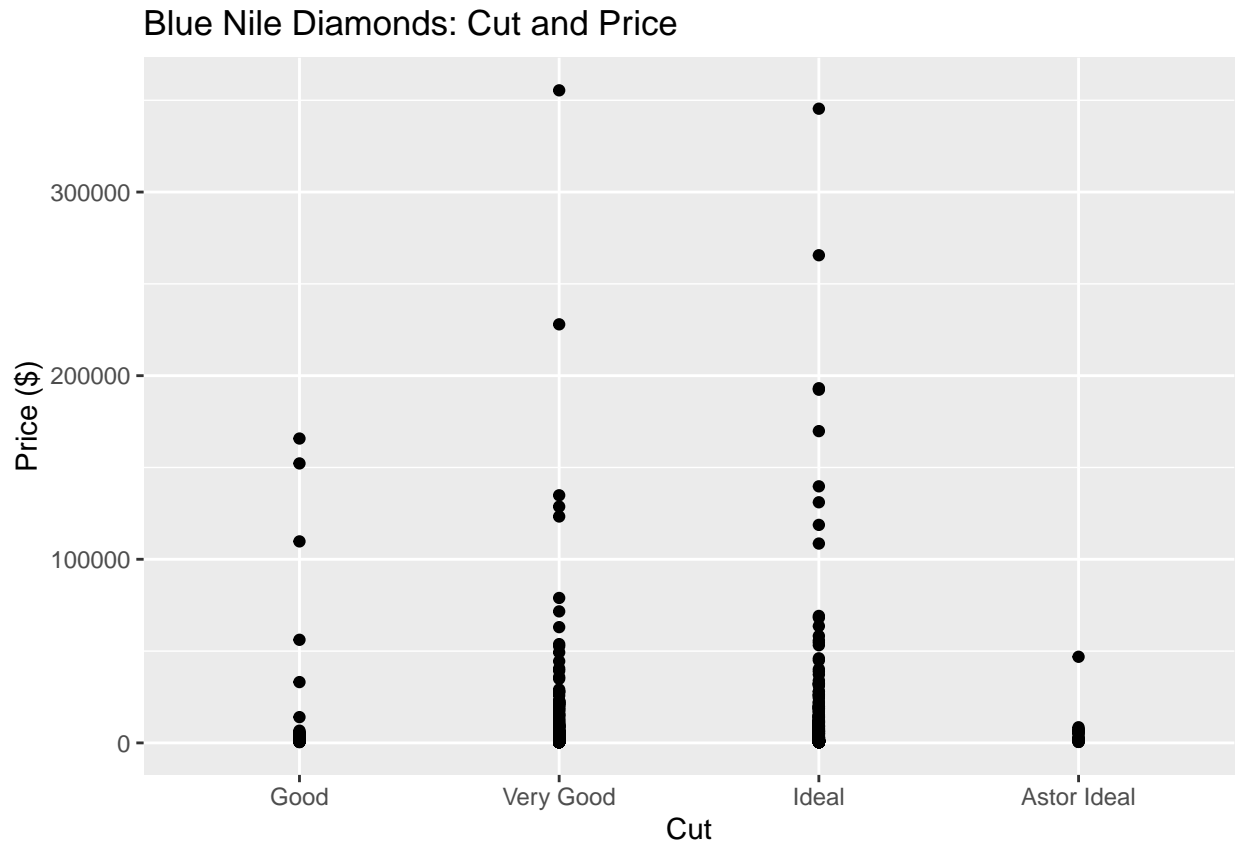


The color of a diamond is more visible in larger diamonds. As shown above, when choosing the size of the diamond, color is a relative factor as color groups 'D' and 'E' which are considered to be colorless and provide the closest look to 'icy' are larger in size.

##Addressing various claims on the diamond education page on Blue Nile:

- Cut: <https://www.bluenile.com/education/diamonds/cut> + "A diamond's cut refers to how well-proportioned the dimensions of a diamond are, and how these surfaces, or facets, are positioned to create sparkle and brilliance. For example, what is the ratio of the diamond's diameter in comparison to its depth? These small, yet essential, factors determine the diamond's beauty and price."

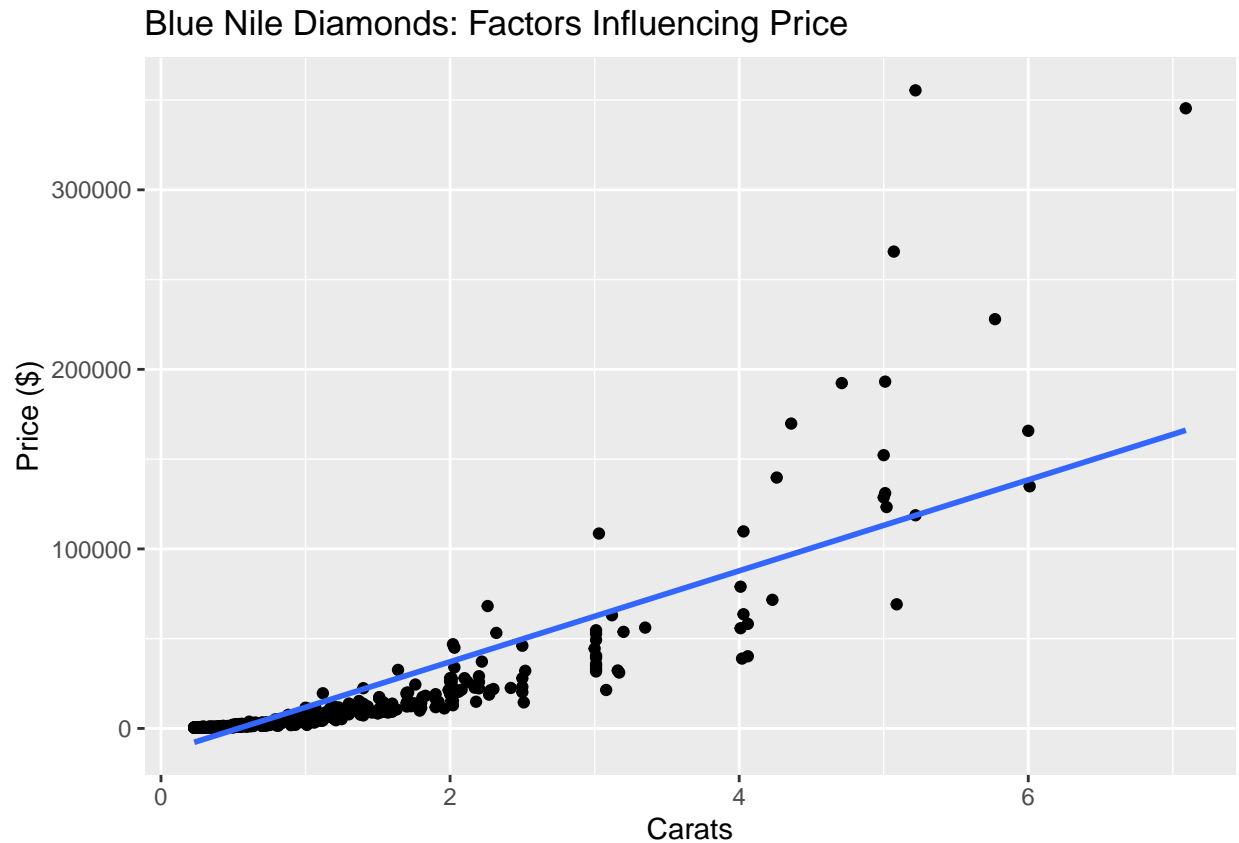
Assertion above is that better cuts correlate with higher price. Let's check the scatterplot to see if that bears out in the data:



As shown above, increasing quality of diamond cut by itself does not seem to have a linear relationship with price, contrary to the Blue Nile's claim.

Description of how we fitted the regression of price against carat:

First checking a scatterplot for Price ~ Carat:



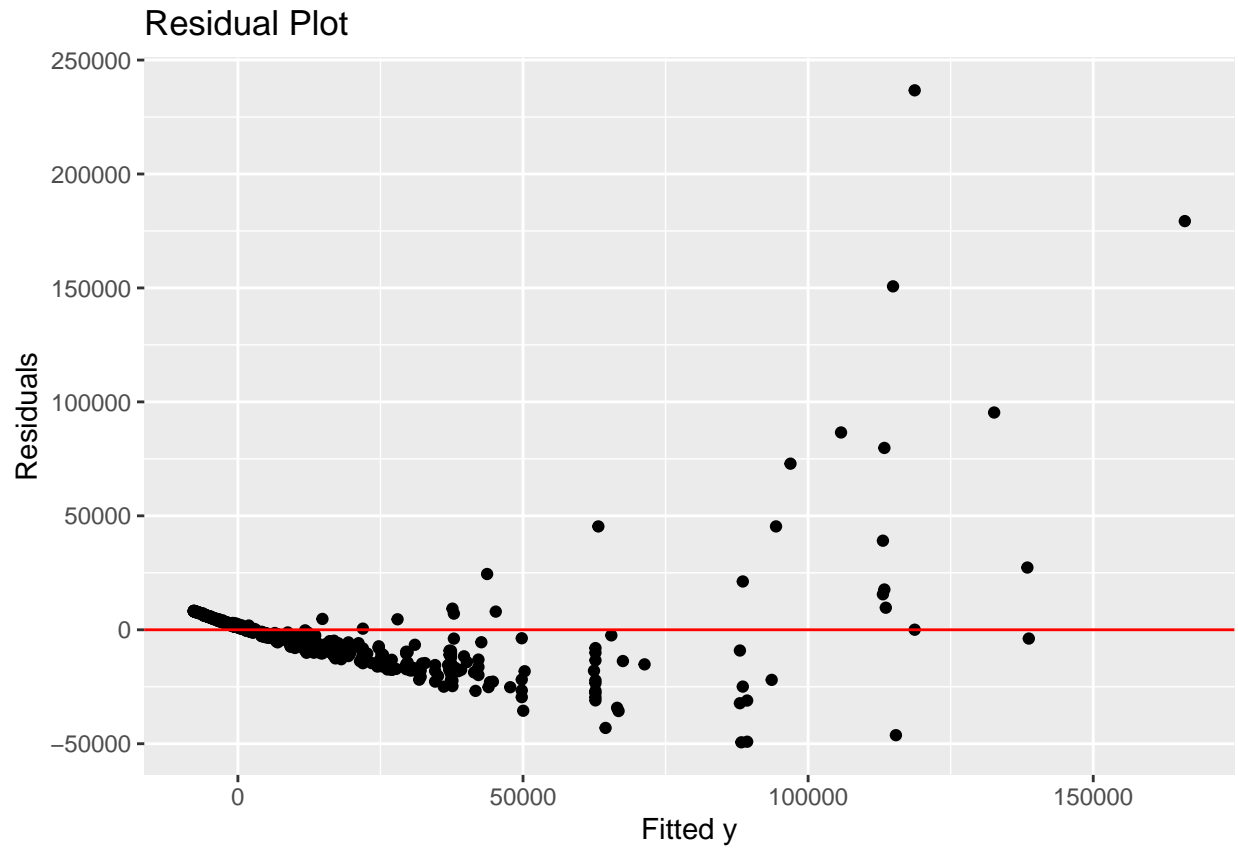
Note that the relationship between Price~carats is roughly linear, however the variance of price over carats does not appear to be constant.

Fitting a preliminary simple linear regression model:

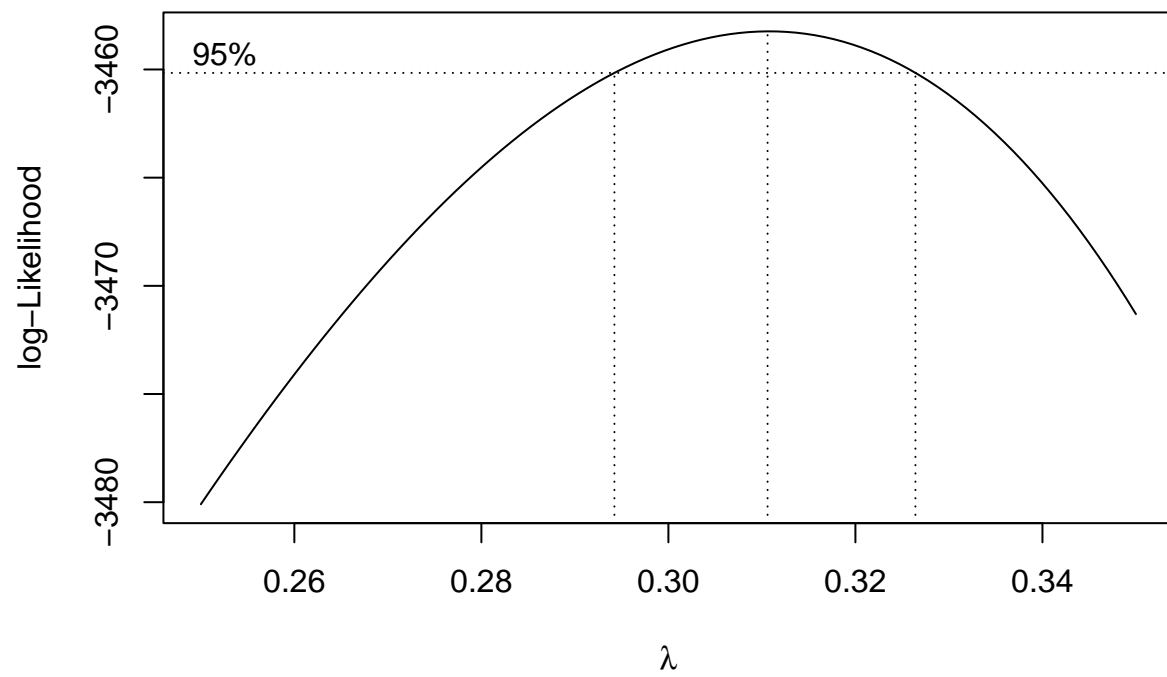
```
##
## Call:
## lm(formula = price ~ carat, data = Data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -49375  -5048   1867   4965  236711
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept) -13550.9      559.7  -24.21 <0.0000000000000002 ***
## carat       25333.9      494.4   51.24 <0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 13560 on 1212 degrees of freedom
## Multiple R-squared:  0.6842, Adjusted R-squared:  0.6839
## F-statistic: 2625 on 1 and 1212 DF, p-value: < 0.00000000000000022
```

Plotting the residuals.

Constant variance and mean of error = 0 assumptions do not appear to be met.



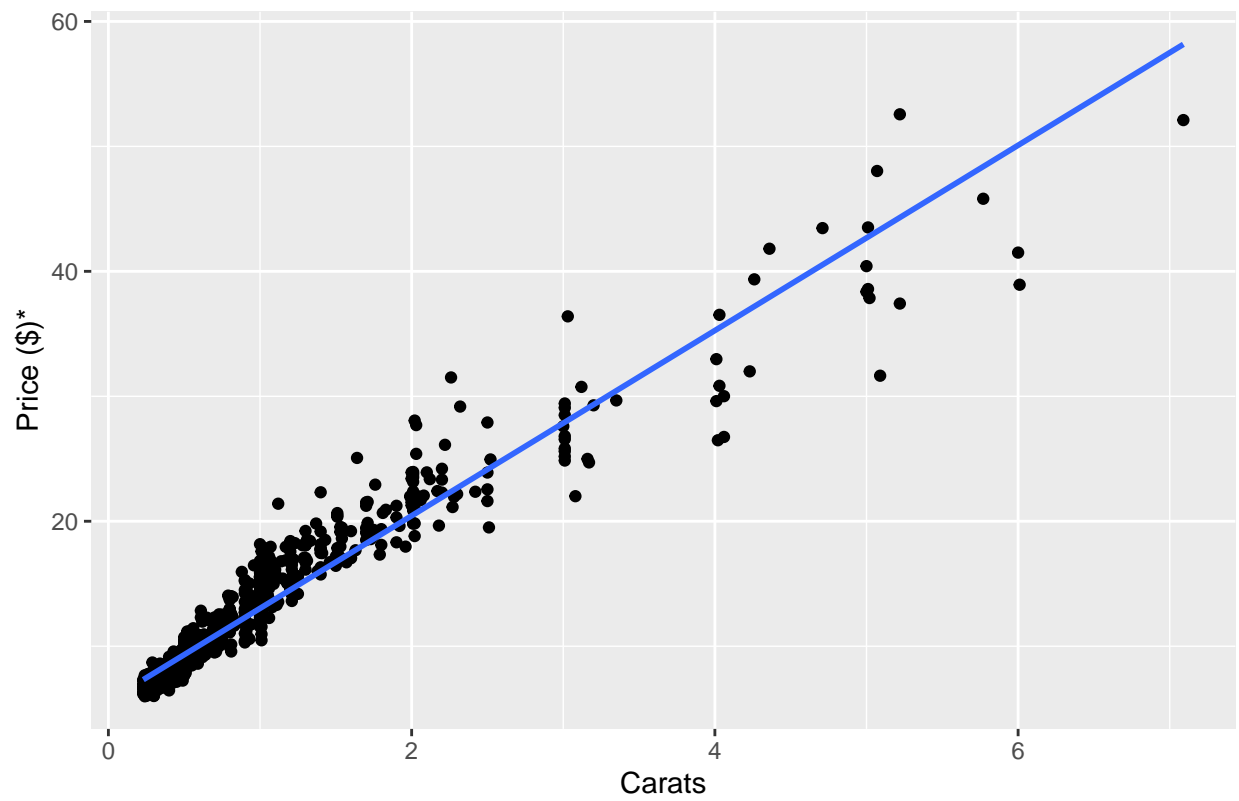
***Variance is not constant so attempting to transform y first. Will start with boxcox plot to see what the optimal lambda may be.



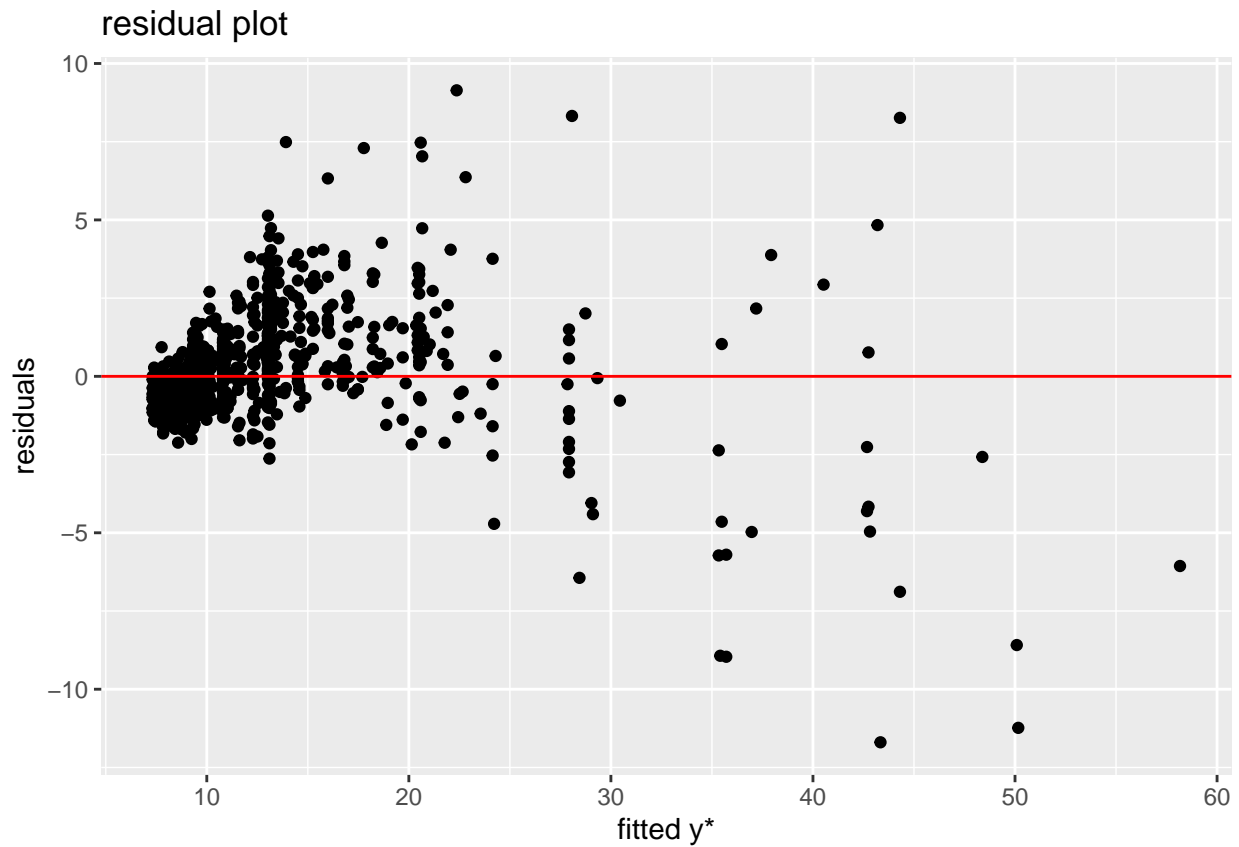
It appears that a lambda of 0.31 would be appropriate.

Replotting the scatter plot with the transformed y variable.

Blue Nile Diamonds: Factors Influencing Price



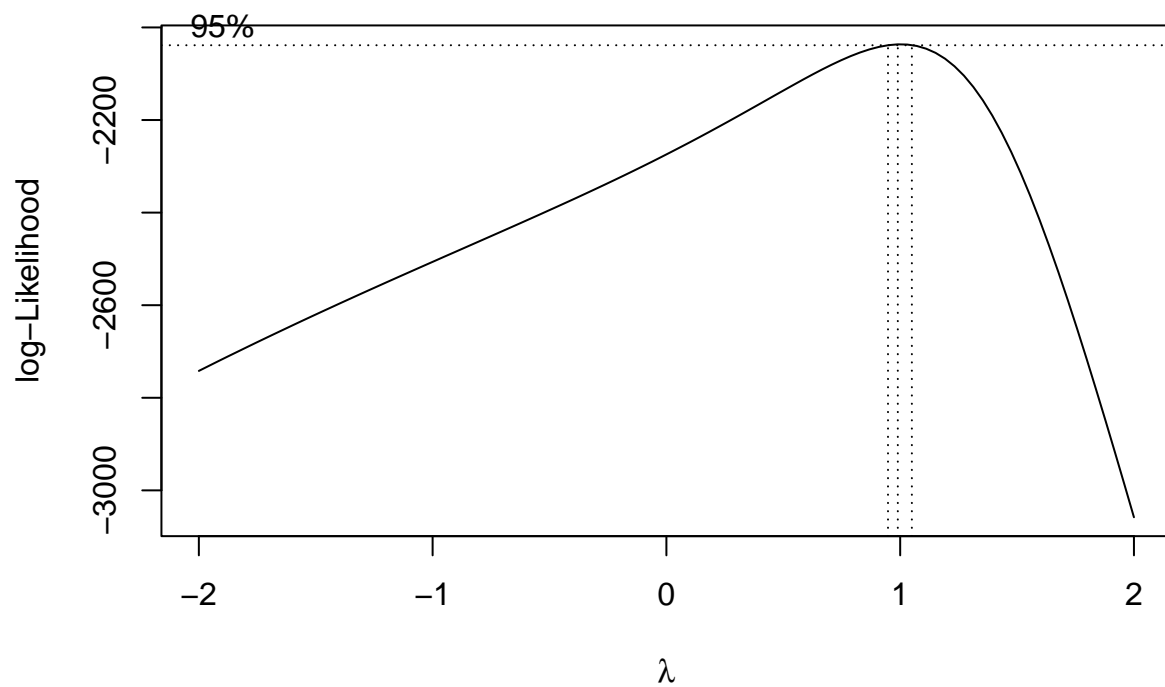
New residual plot:



Note: variance looks better but mean of errors still not equal to zero over x so will attempt to transform the x variable. Given the curved appearance, will try a square root transformation.

Confirming that the box cox plot looks better.

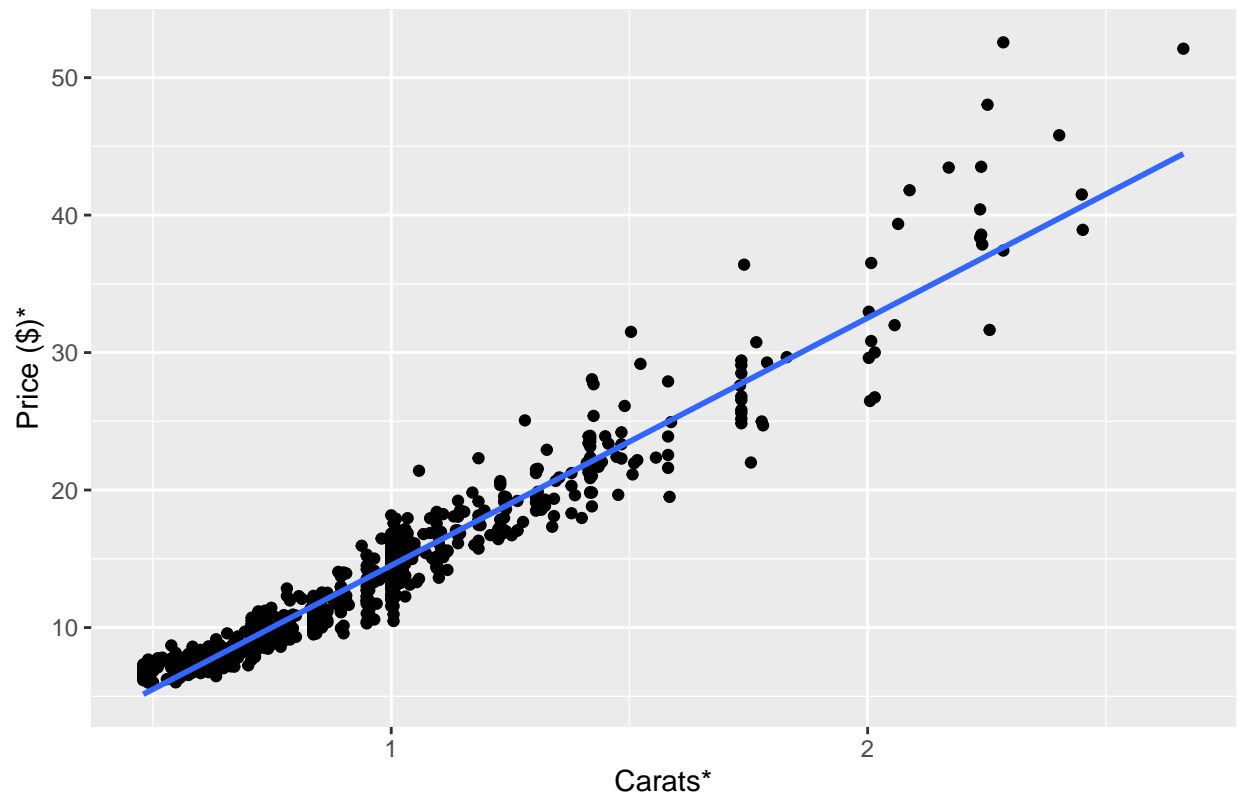
Boxcox plot:



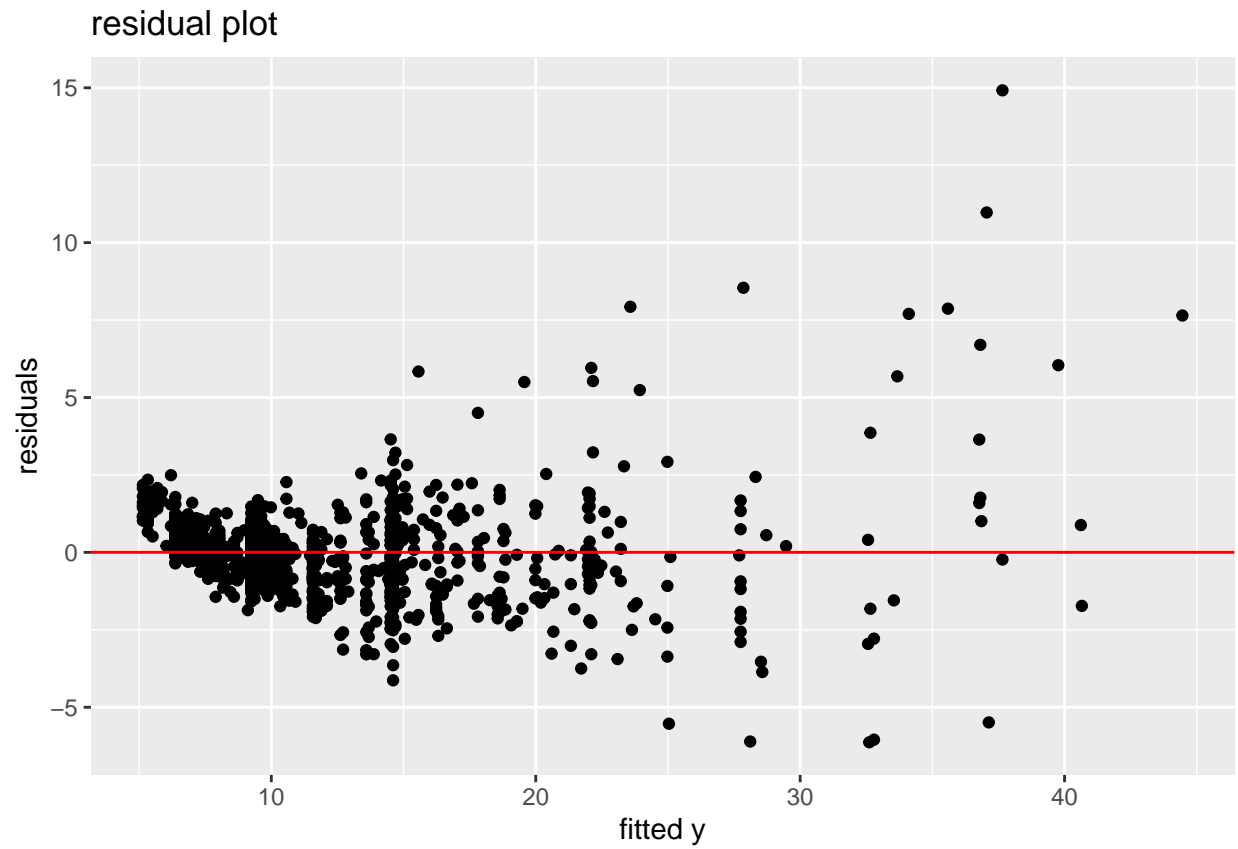
Now I see that confidence interval includes 1. Next step is to consider whether or not to transform x .

Plotting the scatterplot using the transformed x (\sqrt{x}) and y ($y^{0.31}$) variables.

Blue Nile Diamonds: Factors Influencing Price (SLR Model)



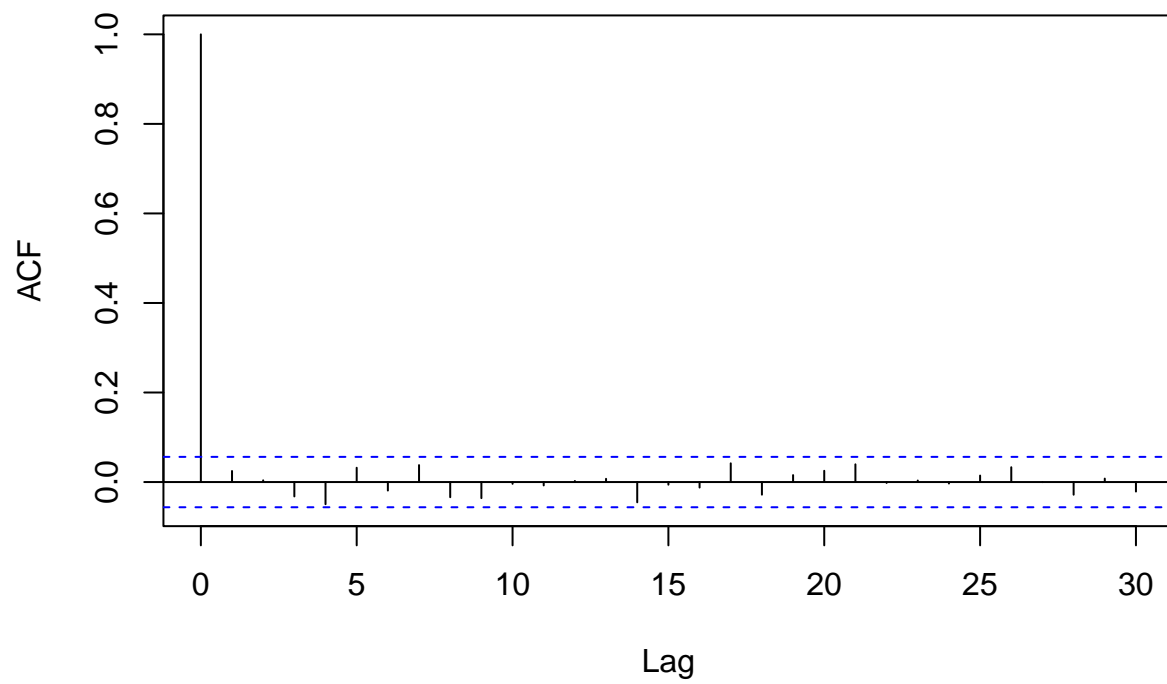
Fitting new model and creating another residual plot:



Not a perfect fit but overall improved and adequate for prediction.

ACF Plot:

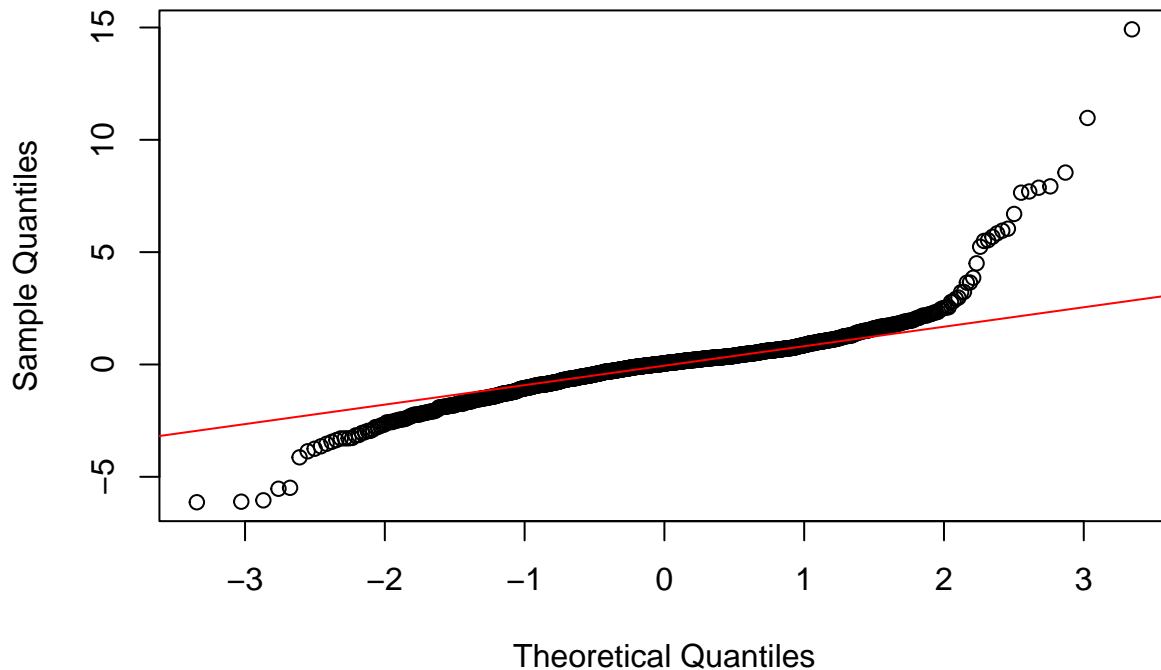
ACF Plot of Residuals with xstar and ystar



Errors do not appear correlated to each other.

QQ Plot:

Normal Q-Q Plot



Normality assumption is not met but this may be the least important.

Summarizing the final model below:

```
##
## Call:
## lm(formula = ystar ~ xstar, data = Data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.1320 -0.6377  0.0373  0.5315 14.9172
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  -3.4936     0.1137   -30.73 <0.0000000000000002 ***
## xstar         18.0085     0.1261   142.87 <0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.436 on 1212 degrees of freedom
## Multiple R-squared:  0.9439, Adjusted R-squared:  0.9439
## F-statistic: 2.041e+04 on 1 and 1212 DF, p-value: < 0.00000000000000022
```

Conclusions:

In conclusion, the blue nile data contain various characteristics for individual diamonds including carat (weight of the diamond), clarity, color, and cut; clearer, colorless, ideal cut, and heavier diamonds being more desirable and thus correlating with higher price. Certain claims by the website such as that the quality of a diamond's cut determines the price - do not appear to be supported by the data.

We fitted an appropriate simple linear regression for price against carat. Doing so required transformation of both the x and y variables to approximate the assumptions required for linear regression.

We can summarize our final regression equation for diamond price (y) as the following:

$$y^* = 18.0085x^* - 3.4936$$

where $x^* = \sqrt{x}$ and $y^* = y^{0.31}$

As you can see in the model summary printed above, the t test for slope P-value for xstar (<0.0000000000000002) indicates that we can reject the null hypothesis that the slope is equal to zero; in other words, carat is a significant predictor of diamond price.

Executive Summary:

Purpose: This project uses various data visualizations to explore how price is related to the other variables such as (carat, clarity, color, cut), as well as how the other variables may relate to each other.

Problem: When considering the rarity of a diamond, it is important to consider the 4Cs of Diamonds - cut, color, clarity, and carat weight. We wanted to determine which out of these characteristics are the most influential on the actual price of the diamond and how they all relate to one another.

Problem Analysis: With the data provided from Blue Nile, we should be able to check if the various claims found on Blue Nile's website are accurate.

Results of Analysis: The results of our project conclude that the website has several incorrect claims based on several of the models and data visualizations that were fitted. We found that after modeling the data in various ways, the most important factor is carat weight rather than the diamonds cut.

Recommendations: We would recommend that buyers consider the size, carat weight, of the diamond as the most influential factor that contributes to the price when looking to purchase a diamond. With this knowledge, buyers will have a better understanding on what characteristics have a more determinant role and will be better suited to find the best diamond for them.