UVA | SCHOOL *of* DATA SCIENCE

2023

# Price vs. Diamond Clarity: Studying Diamonds From Blue Nile

*Hayden French, Addison Gambhir, Rachel Holman, Isha Thukral*

# Table of Contents

# Summary of Findings

Diamond shopping for the first time this summer? Struggling to find that perfect stone in your price range? Researchers at the University of Virginia have conducted a comprehensive analysis of over a thousand real diamonds and have some tips for you about smarter diamond shopping. For example, don't sweat the cut. Out of all the properties studied, the cut of the diamond had the least impact on the price. It's much better to start by thinking about what size diamond is right for you. If you're shopping on a budget, those extra carats can add up fast- in fact, the study concluded a 44% increase in carat weight resulted in a doubling of price! One tip that might save you money is buying just shy of whole and half carat weights as many customers prefer those even numbers. You just might walk away with a similarly sized diamond at a cheaper price. After finding the appropriate size, you'll need to pick the right color. There's no trick here- the study shows that there's a direct relationship between the amount of color and the price. You'll have to decide for yourself if you're ok with a little coloration, or if you want to splurge on an absolutely colorless diamond. Once you have an idea of your preferred carat and color, you're in a great position! Online retailers such as Blue Nile have a huge selection and experts can help you narrow in on the right fit for you. We're confident that with these tips you can find a perfect gift for that special someone!

# Data Description

The data set that we will be working with describes more than 1,000 different diamonds that are for sale on [Blue Nile's Website](). This file contains a subset of the diamonds on Blue Nile as well as from other internet resources. The variables are:

- **Carat**: Weight of the diamond in carats (0.23 - 7.09)
- **Clarity**: Measurement of how clear the diamond is  (clarity types listed below)
  - I1:  *"Included"* Diamonds with obvious inclusions that impact beauty
  - SI2: *"Slightly Included"* Diamonds with inclusions detectable to keen unaided eye, especially when viewed from the side
  - SI1: *"Slightly Included"* Diamonds with inclusions noticeable at 10x magnification (best value)
  - VS2: *"Very Slightly Included"* Diamonds with minor inclusions that are somewhat easy to see at 10x magnification
  - VS1: *"Very Slightly Included"* Diamonds with minor inclusions that are difficult to see
  - VVS2: *"Very Very Slightly Included"* Diamonds with minuscule inclusions that are difficult even for trained eyes to see under 10x magnification
  - VVS1: *"Very Very Slightly Included"* Diamonds with minuscule inclusions that are difficult even for trained eyes to see under 10x magnification

- ○ IF: *"Internally Flawless"* Diamonds with no inclusions within the stone, only surface characteristics set the grade
  - ○ FL: *"Flawless"* Diamonds with no internal or external characteristics (rare)
- **Color**: Measurement of faint diamond color (color types listed below)
  - ○ D: Rarest and highest quality with a pure icy look
  - ○ E: Rarest and highest quality with a pure icy look
  - ○ F: Rarest and highest quality with a pure icy look
  - ○ G: No discernible color; great value for the quality
  - ○ H: No discernible color; great value for the quality
  - ○ I: No discernible color; great value for the quality
  - ○ J: No discernible color; great value for the quality
- **Cut**: Cut quality of Diamond (cut types listed below)
  - ○ Good: This cut represents roughly the top 25% of diamond cut quality. It reflects most of the light that enters, but not as much as a Very Good cut grade.
  - ○ Very Good: This cut represents roughly the top 15% of diamond cut quality. It reflects nearly as much light as the ideal cut, but for a lower price.
  - ○ Ideal: This rare cut represents roughly the top 3% of diamond cut quality. It reflects most of the light that enters the diamond.
  - ○ Astor Ideal: These diamonds are lab grown by Blue Nile™ and are certified by GemEx® to look perfect.
- **Price**: Price in U.S. Dollars ($322 - $355403)

In addition to the variables listed above, our team created one new variable based on the existing variables to further our analysis. We created a carat weight range variable which separated the diamonds into categories based on how many carats each contained. For example, all diamonds with less than 1 carat were categorized as "<1" on the carat weight range while those between 1 and less than 2 carats were marked "1 - 1.9" and so on. This helped to show how the average price changes within each carat weight range.

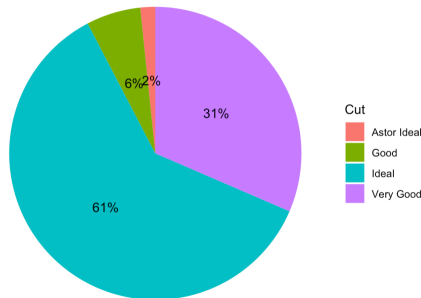# How Diamond Features Relate to One Another

Data visualization can be a game-changer when it comes to comprehending the vast and complex world of diamonds. We may learn a lot about diamond properties from Blue Nile's educational website, and by utilizing the power of visualizations, we can find insightful patterns and comprehend the dataset.

Speaking about price, knowing how it relates to other characteristics is essential to understanding the diamond market. By exploring data visually, we can learn more about this connection and find insightful information. We can determine how various qualities affect diamond pricing by making scatter plots using price as one of the variables and looking at its link with cut,
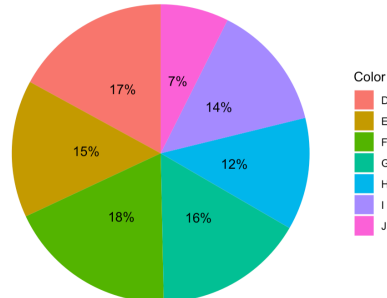
color, and clarity. For instance, we might note that diamonds with higher cut grades typically sell for more money or that a diamond's clarity grade significantly affects its price. These visualizations provide us a better understanding of how price and various features interact, providing insightful advice for both consumers and business experts.

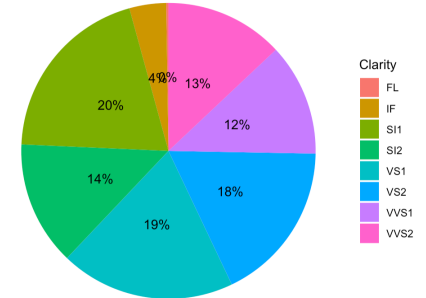## *Visual 1: Pie Chart Insights for Exploring Cut, Clarity, and Color*



Understanding the distribution and importance of various qualities is crucial in the world of diamonds. We can learn a lot about the distributions and rankings of critical characteristics like cut, clarity, and color within our sunset of the Blue Nile diamond collection by presenting data as pie charts.

Beginning with the cut pie chart. The graph demonstrates that the majority of diamonds, or 61% of the dataset, fall into the "Ideal" cut category. This implies that the vast majority of the diamonds given by Blue Nile are exceptionally precise and well-crafted, fulfilling the highest industry standards for cut quality.

Moving on to color, the pie chart shows which color grades are most common in the dataset. With 18% of the diamonds, the "F" color grade has the highest percentage. It's vital to keep in mind that colorless diamonds (usually in the D-F range) are regarded as the most valuable and sought in the diamond market due to their scarcity and outstanding capacity to reflect light.

For clarity, let's examine the rightmost pie chart now. The distribution of clarity grades among the diamonds in the dataset is shown in the chart. We see that there is a fairly even distribution of charity types in our data set, apart from the FL "flawless" and IF "internally flawless" categories. It's important to keep in mind that diamonds with higher clarity classifications, like VS2 and VVS2, usually have less obvious defects or flaws, boosting their general appeal and value.

We can create predictions about the pie chart proportions of diamonds available for sale using the rankings, pricing data, and outside industry expertise. Given the high percentage of "Ideal" cut diamonds, it is safe to infer that Blue Nile has a large assortment of well-cut diamonds to satisfy the

needs of clients looking for great brightness and fire. It also shows that Blue Nile prioritizes providing diamonds with top-tier color quality, coinciding with the industry's emphasis on colorless diamonds as the most desirable and precious, given the prominence of higher color grades like "F" and "D" Last but not least, the distribution of clarity grades shows that Blue Nile concentrates on offering diamonds with average to high clarity ratings, making sure that consumers have access to diamonds with few inclusions or other flaws.

## *Visual 2: Scatterplot for Clarity*



Scatterplot of Diamonds Prices by Carat and Clarity

Scatter plots are yet another effective visualization technique for investigating the dataset. They allow us to evaluate the relationship between two variables, such as carat and price. By plotting price on the y-axis and carat on the x-axis, we can examine how these two parameters are linked. Using scatter plots, like the one above, we can gain insights into the role of different clarity types in pricing carat weight. We can observe above how clarity types affect the pricing of diamonds. Notably, there is an outlier within the graph corresponding to the 'FL' grade, which represents flawless diamonds. These diamonds are extremely rare within our dataset, driving their prices significantly higher, as indicated by the steeper slope associated with this grade. However, apart from the 'FL' grade, the scatter plot suggests that the other clarity measurements do not have a significant impact on the price of diamonds below two carat weight.

Similar plots were made to show how other variables affect pricing of diamonds and they serve the dual purpose of addressing many of the claims found on Blue Nile's website so we will explore those in the following section.

# Addressing Blue Nile Claims

One of our goals throughout this analysis is to explore how price is related to the other variables, but we are also interested in addressing claims made on the [diamond education page on Blue Nile](#). Such claims are as follows:

1) The cut of the diamond can be the biggest factor in the price tag ([as claimed by Blue Nile](#)). "The Ideal cut diamond, and the super-ideal Astor by Blue Nile™, are the most expensive diamond cuts because they optimize light performance and create the most impressive sparkle."
2) "Buy shy" to save money. Select a carat weight slightly below the whole and half-carat marks. For example, instead of a 2.00 carat diamond, consider buying a 1.90 carat weight. This will save a considerable amount of money, and the slight difference will never be noticed."
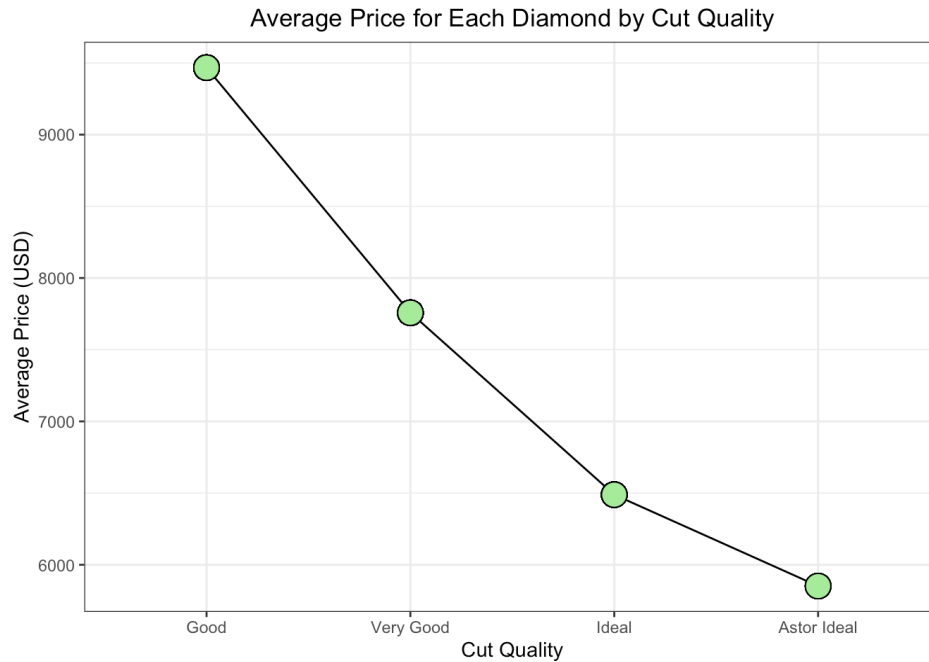3) "The absence of color in a diamond is the rarest and therefore, the most expensive."

We address each of these claims using visualizations in the following sections.

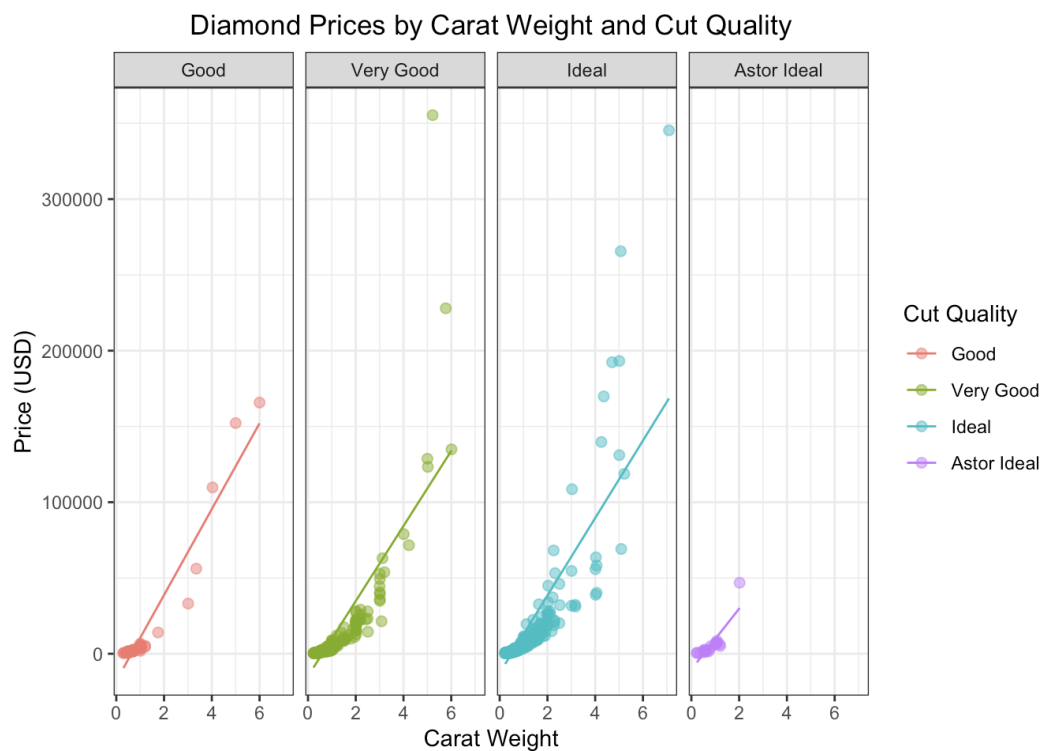## *Claim 1: Cut is the Biggest Factor in Price Tag*

One of the most crucial things to think about when buying a diamond is its cut. The cut, or shape and facets of a diamond, is what ultimately defines its brilliance and shine. A higher grade cut is thought to attract a higher price since it enhances light performance and produces a stunning glitter.

But as we examine the data and go further into the world of diamonds, a new viewpoint starts to take shape. While the cut undoubtedly contributes significantly to a diamond's overall beauty, the price may not be primarily determined by the cut. We can learn about the subtleties and hidden complexities of diamond price by investigating extensive data.

Below we utilize a line graph to show the average price for the diamonds grouped by cut quality:

Average Price for Each Diamond by Cut Quality

This line graph above shows that the average price is much higher for diamonds of "Good" cut quality than those with "Very Good", "Ideal" or even "Astor Ideal" cut qualities. This affirms the claim that cut quality is a big factor in determining price, but the correlation is the opposite of what one might expect. To understand why, we will evaluate another graph.
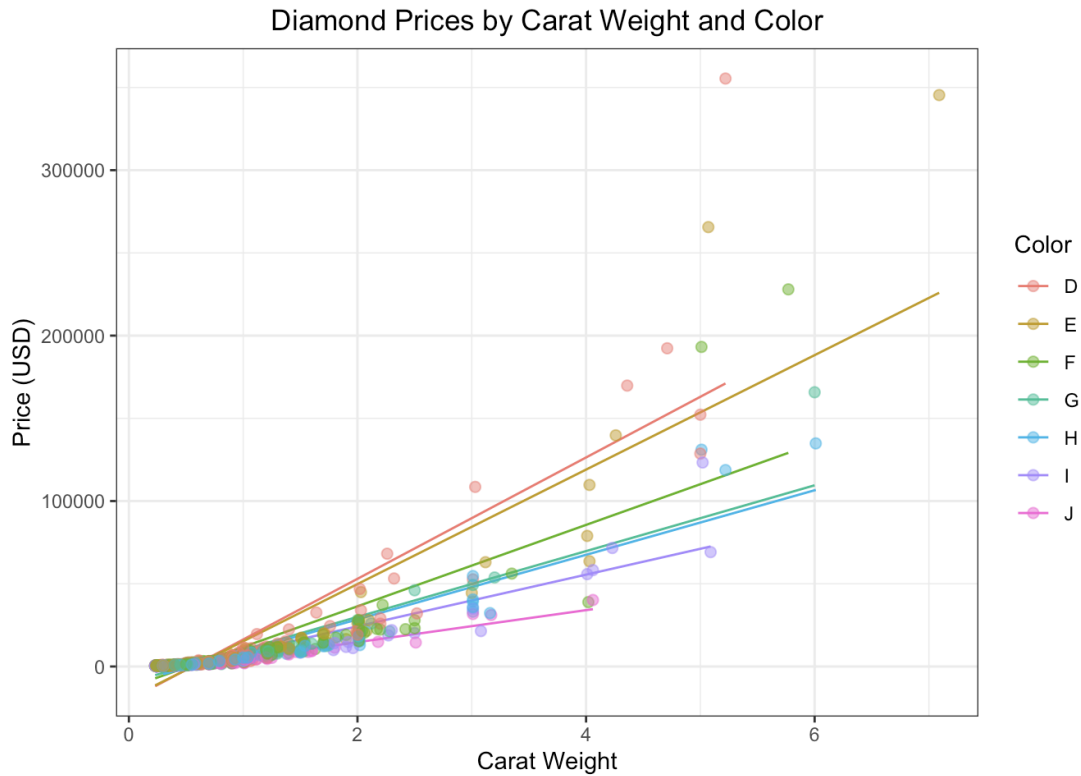


Diamond Prices by Carat Weight and Cut Quality

The above visualization represents the relationship of price against carat size and grouped by cut quality. While Blue Nile asserts that cut is the most important factor, our analysis indicates otherwise. Looking at the linear relationship between price and carat weight for each cut type (shown by the lines on each scatterplot), we see that the slopes are approximately parallel. In other words, the price per carat ratio is approximately the same despite the differing cut types. The discrepancy between this finding and that from the previous plot may be explained by the fact that there are far fewer diamonds in our data set with "Good" cut quality than "Very Good" which in turn has fewer observations than "Ideal". Because there are more lower carat and therefore lower price diamonds of better quality, our averages are skewed. Our findings demonstrate that while a well-cut diamond may be visually appealing, it does not necessarily command a higher price if its carat weight is lower. This suggests that customers prioritize the size and weight of a diamond over its cut when making purchasing decisions, rendering Blue Nile's claim about the significance of cut questionable.
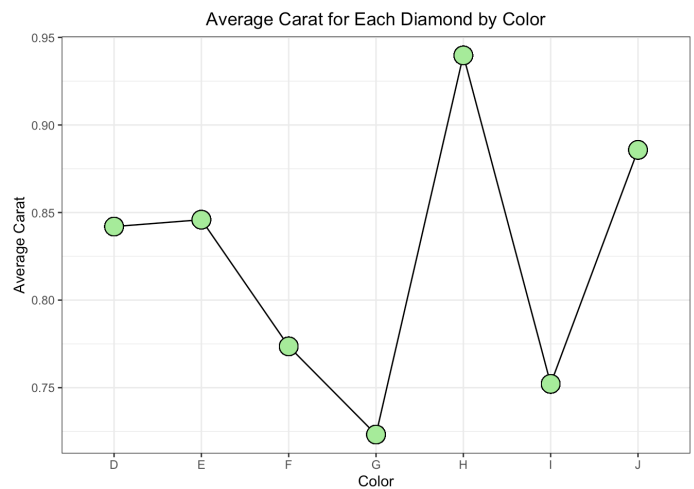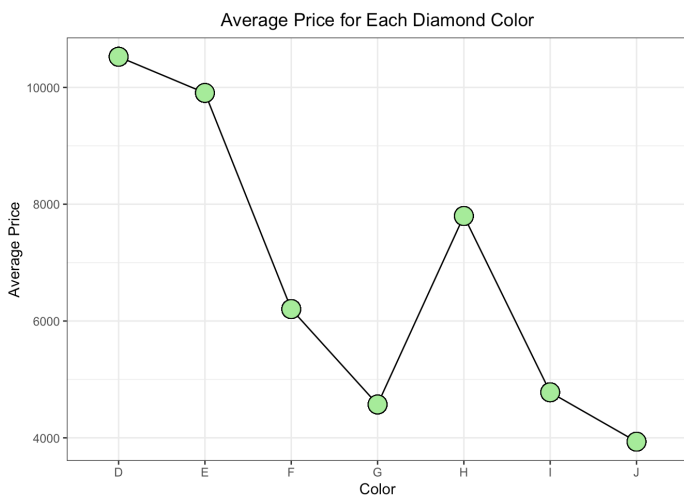
## Claim 2: The Absence of Color in a Diamond is the Most Expensive

Diamond color, one of the four Cs (Cut, Color, Clarity, and Carat), refers to the degree of colorlessness in a diamond, which is evaluated using a standardized diamond color chart. Diamonds with higher color grades are closer to being colorless, while those with lower grades often exhibit noticeable color tints as indicated on the diamond color scale. A colorless diamond appears transparent, while diamonds with lower color grades may have a warm hue. According to BlUe Nile: Diamond color is crucial when purchasing a diamond, just like the other aspects of the four Cs. In the graphs and analysis below, we investigate the validity of this claim.

Our findings support Blue Nile's claim that color is indeed a significant factor influencing diamond pricing. Diamonds are assigned letter ratings from D to J, with D representing a pure and colorless diamond, and the color intensity increasing as the letters progress. Our analysis revealed a clear linear relationship between color quality and price per carat.
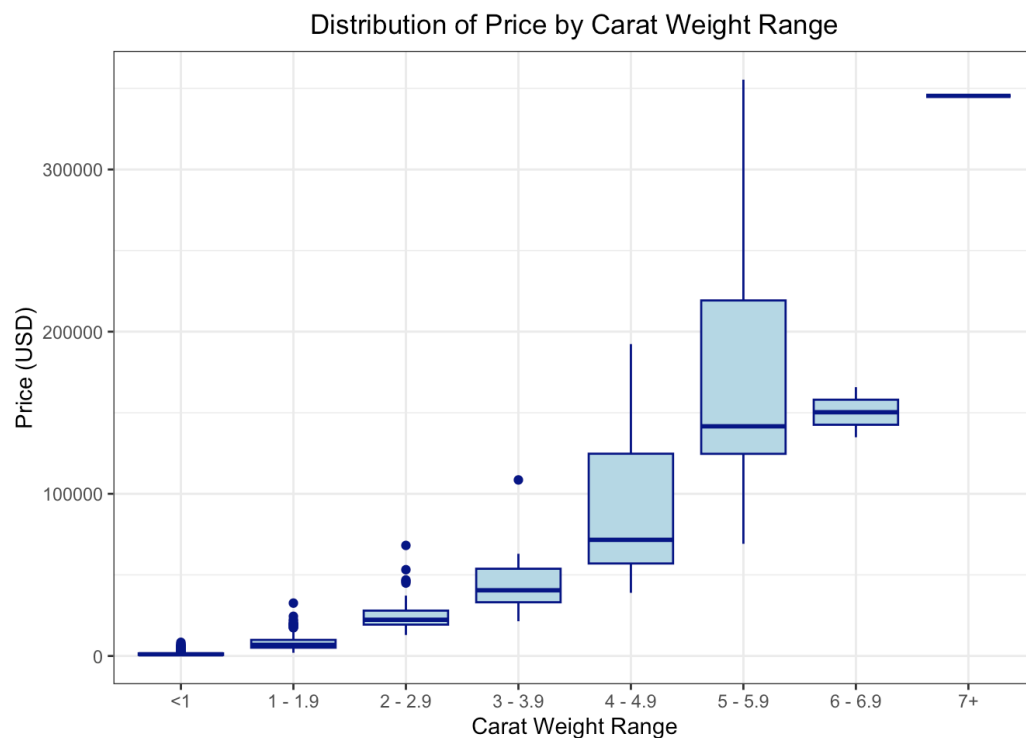
Diamond Prices by Carat Weight and Color

The plot of our data clearly illustrates that diamonds with higher color grades, closer to the beginning of the alphabet, command higher prices per carat compared to those with lower color grades. A diamond rated as D, characterized by its colorlessness, fetches a higher price than diamonds rated further down the color scale. This correlation suggests that the absence of color contributes to the perceived value and desirability of a diamond, aligning with Blue Nile's claim that diamonds with minimal color are more expensive.



Average Price for Each Diamond Color



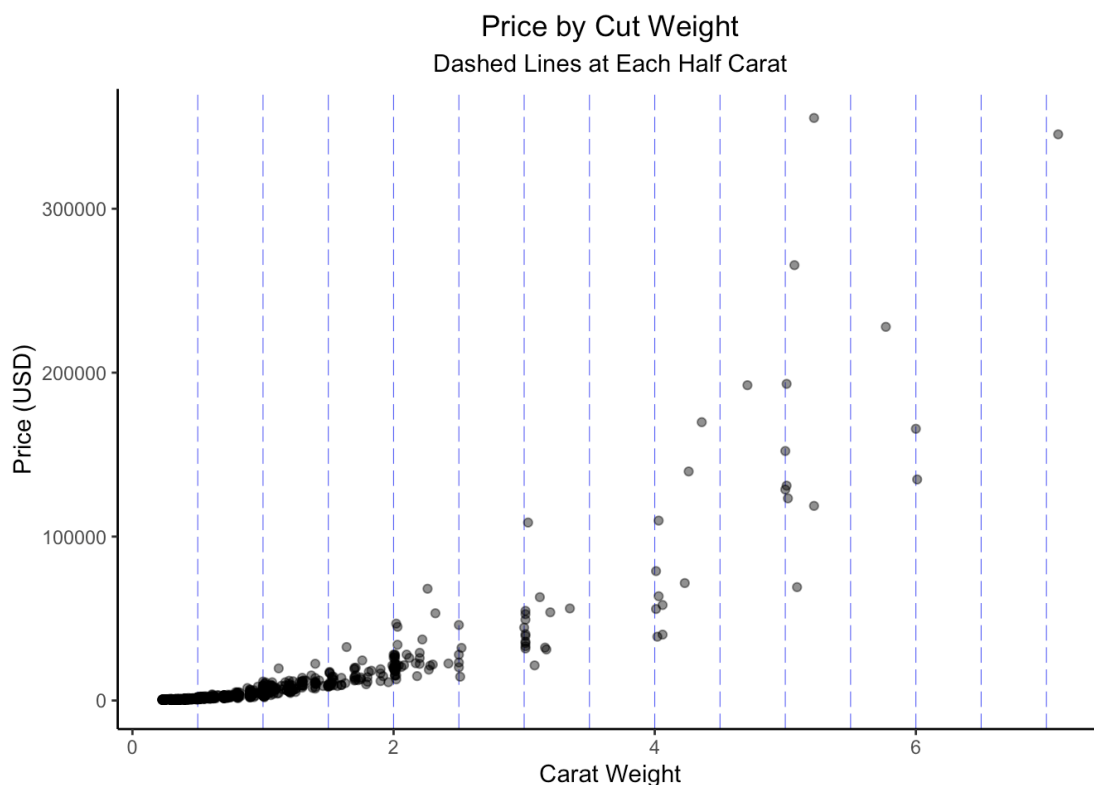Average Carat for Each Diamond by Color

The color D means the diamond is pure, icy, and lacking color. As the color values sequence up through the alphabet, the amount of color present increases. It is clear from the plot above that the average price per carat is higher for diamonds with less color than those with more. This affirms Blue Nile's claim that the absence of color makes a diamond more expensive. We do notice a spike in the average price for diamonds with the color H, but this may be partially attributed to the fact that the H colored diamonds in this data set were, on average, larger carat weights than the other colored diamonds.

## *Claim 3: Buying Shy of Whole and Half Carat Weights Will Save Lots of Money*

Many customers place a high premium on maximizing value when it comes to diamond purchases without compromising on quality. One tactic that has gained favor is the idea of "buying shy." This strategy involves choosing a carat weight that is just below the whole and half-carat markers, which may result in savings without a discernible visual difference. The psychological pricing markers frequently associated with whole and half carat weights, which are seen as significant milestones in diamond sizes, are likely to be to blame for the price gap. However, a 2.00 and 1.90 carat diamond's 0.10 carat difference is normally undetectable to the unaided eye. Blue Nile claims that "buying shy" is a valid strategy, since, according to Blue Nile: "it will save a considerable amount of money, and the slight difference will never be noticed".



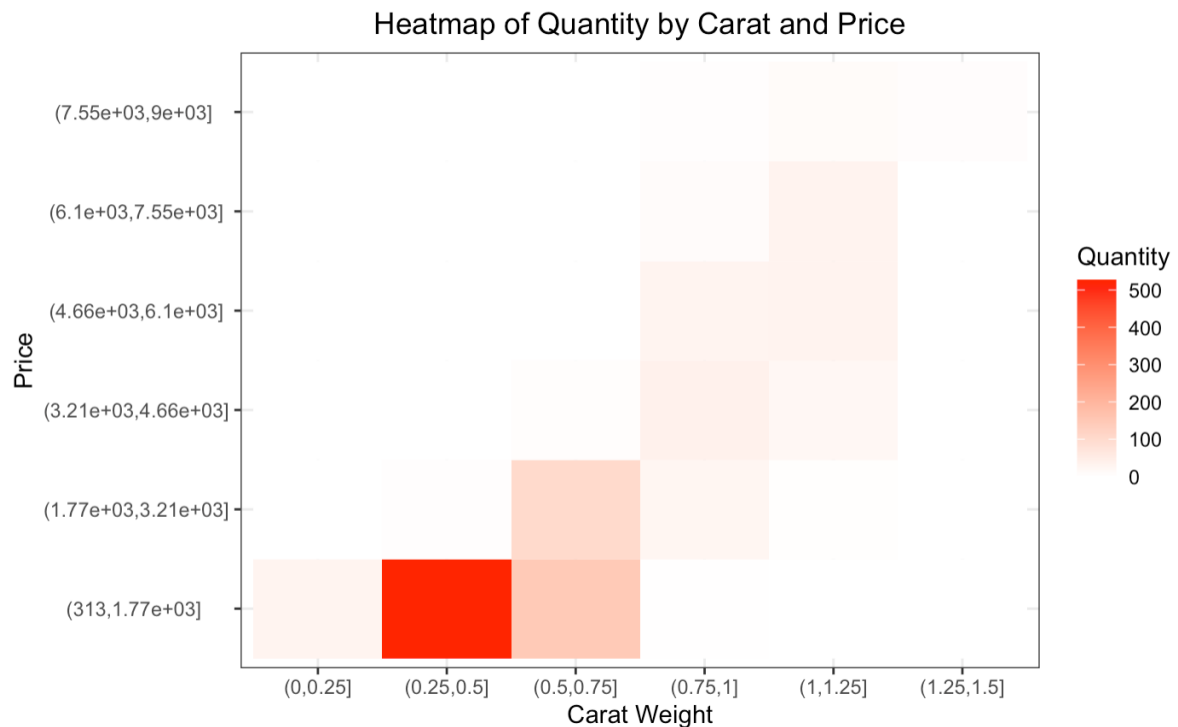Distribution of Price by Carat Weight Range

Looking at the boxplot above, we notice that the prices of diamonds seem to increase exponentially as we move up through the carat weight ranges. By comparing the median lines of each box, we do see that price jumps quite a significant amount for each additional whole size primarily once you are looking at diamonds larger than 4 carats. To further explore, we will visualize it through a scatter plot below.

## Price by Cut Weight
### Dashed Lines at Each Half Carat



Looking at this distribution with dashed marker lines at the whole and half carat marks, we do see clustering around the markers. This shows that many customers do tend to buy at the whole and half carat weight values. We can more easily see here that the price to carat ratio is relatively consistent and linear for diamonds less than 3 carats, but once the carat weight hits 4 carats and beyond the savings for buying shy are more significant. Based on this visualization, we are tempted to affirm Blue Nile's claim conditionally: buying shy does save money, but it is only saving a considerable amount when buying shy of carat weights larger than 4 carats.

Furthermore the heatmap below provides insightful information on the distribution of diamonds across various carat weights and price points. Using this, we can examine whether the quantity of diamonds purchased inside one of the most common carat ranges in the dataset—precisely, 0.25–1.5 weight—aligns with the idea of "buying shy" by concentrating on this range. The visualization indicates a larger quantity of diamonds are purchased at carat sizes slightly below one

carrot; however, the prices don't seem to be a driving factor in the quantity, and the price seems to be steady.
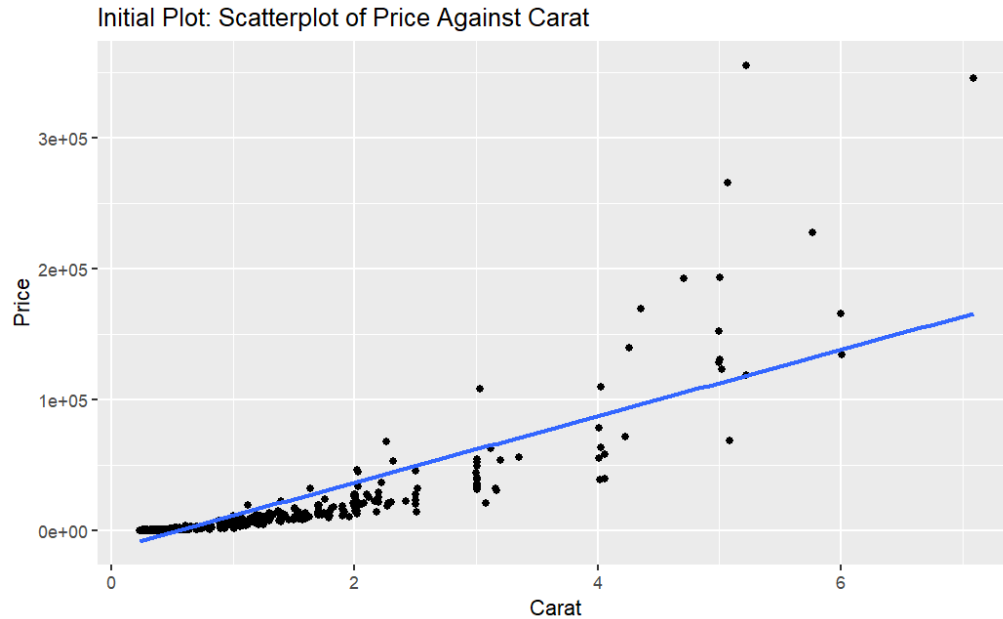


Heatmap of Quantity by Carat and Price

These findings challenge the notion that customers can achieve a noticeable cost reduction by purchasing diamonds slightly below the whole-carat weight. Instead, our analysis suggests that the price per carat remains relatively constant, indicating that customers might not be able to achieve significant cost savings by adopting the "buy shy" approach unless purchasing diamonds larger than 4 carats.
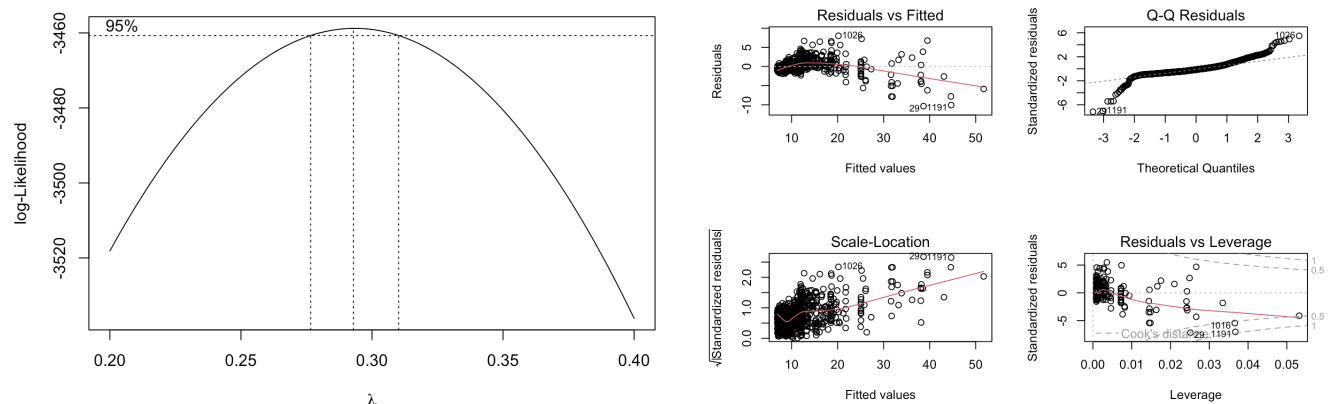
# Modeling

Now that we have explored how each of the variables in our data set impacts the price of a diamond and elaborated on some of the claims found on Blue Nile's website, we intend to create a simple linear model quantifying the relationship between carat weight and price for a diamond. In order to account for and quantify variation in price based on diamond carat weight we first have to adjust our variables to ensure the model has constant variance and a linear pattern. Both our process for fitting this model and a contextually useful interpretation of the results are detailed in the following section.
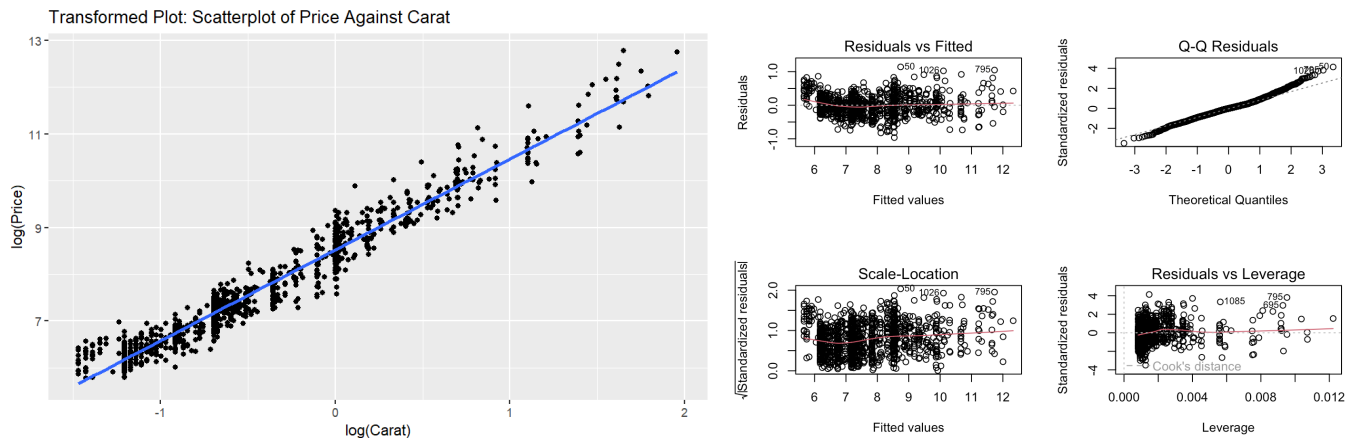
## *Linear Regression Fitting*

Upon plotting the initial scatterplot of price against carat it was clear that several assumptions for linear regression were not met.



In particular, we observed that our residuals did not have mean zero or constant variance. So, we knew that we would have to transform both our x and y variables (carat and price, respectively). For the y transformation, we generated a Box-Cox plot which suggested a lambda value around 0.3. However, after transforming with this lambda value we noticed from the residual plot that we still had increasing variance. This was concerning so we decided to also try a log transformation for y. Our reasoning was that we knew from the Box-Cox plot that our lambda would be less than 1, and a log transformation would allow us to better interpret our results. The log transformation additionally better stabilized the variance as shown below, so that was the transformation we went with.

We also needed to transform x, and a log transformation seemed to best fit the scatterplot we generated. After applying a log transformation to both x and y, we re-examined the residual plot and determined that we now met the assumptions for linear regression. Our errors were evenly scattered on each side of our regression line, and our variance was much more stable as is clear from the scatter and residual plots below.



Our final regression line corresponded with a R^2 value of .954, and had the following equation:

$$\hat{y}^* = 8.521 + 1.944x^*_{(y^*=\log(y)\,,\,x^*=\log(x))}$$

Taking into account the transformations performed, we can interpret this result as follows: For a 44% increase in carat size, the price of the diamond approximately doubles.

# Closing Remarks

In this report, we analyzed the veracity of several claims made by Blue Nile regarding their diamond collection. The first claim was that the cut of the diamond can be the biggest factor in the price tag. We found this claim to be mostly false. When examining the mean price by cut, we actually found cleaner cuts to have cheaper prices. The reasons for this were confounding variables, specifically the carat size, which shows that cut is not the main contributor to price. The second claim examined was that individuals looking for a deal should "buy shy" and select a carat weight just below the whole or half-carat marks. We found this claim to be plausible. From our regression analysis, it is certainly true that buying less carats will be cheaper. However, we were unable to conclude that the rate of change in the price was specifically affected around the whole and half carat marks. So we recommend that buyers take this strategy into consideration, but understand it might only come into play on a case-by-case basis. Finally, we looked at the claim that the less color in a diamond, the more expensive. We were able to decisively confirm this claim. Out of the entire color range examined, each step up in colorlessness corresponded to an increase in price. We hope these results will result in more well-informed consumers who can find the perfect diamond at the right price.