# Exploring the Interplay of Price, Cut, and Carat in Diamonds: A Data Analysis.

## 1. Abstract

In the world of diamonds, various attributes contribute to a gemstone's overall value and desirability. Among these attributes, the cut, carat, and price are pivotal factors influencing consumer preferences and market dynamics. This project aims to delve into the intricate relationship between the cut quality, carat weight, and price of diamonds using a comprehensive data analysis approach. By exploring patterns and correlations within these attributes, we seek to uncover insights that shed light on the factors driving diamond pricing and provide valuable information for jewelers, buyers, and enthusiasts to understand how different aspects come together to define a diamond's worth.

## 2. Introduction

Diamonds, renowned for their exquisite beauty and enduring allure, have captivated individuals for centuries. Their timeless elegance, coupled with their value as precious gemstones, makes the study of their attributes and market trends a captivating endeavor. In this data analysis project, we delve into the multifaceted world of diamonds by exploring the intricate relationship between their price, cut, and carat weight.

Understanding how these attributes interplay is not only of interest to gemologists and jewelry enthusiasts but also holds practical implications for traders, retailers, and consumers. Cut, referring to the craftsmanship and shape of a diamond, plays a pivotal role in its brilliance and overall visual appeal. Carat, a measure of a diamond's weight, is closely linked to its rarity and value. Price, the ultimate manifestation of a diamond's worth, is influenced by complex factors.

By meticulously dissecting the data, we aim to uncover patterns, correlations, and insights that illuminate the connections between these attributes. This analysis is not merely an exercise in data exploration; it is an exploration of the essence that gives diamonds their allure and monetary value.

Through this project, we seek to provide a nuanced perspective on the intricate world of diamonds, bridging the realms of science and art. By unraveling the relationships between price, cut, and carat, we hope to contribute to a deeper appreciation of the craftsmanship and allure that diamonds epitomize.

# 3. Dataset Description

### 3.1 Overview

The diamond dataset used in this project is a comprehensive collection of diamond attributes and their corresponding prices. This dataset is commonly utilized in the field of data analysis to understand the factors that influence the price of diamonds. Diamonds are valued not only for their intrinsic beauty but also for their rarity and the "Four Cs" criteria: cut, carat, color, and clarity. In this analysis, we will focus on exploring the relationship between the attributes of cut, carat, and price.

### 3.3 Attributes of Interest

The primary attributes of interest for this analysis are:

- **Cut:** This attribute refers to the quality of the diamond's cut, which affects its brilliance and overall appearance. The cut is categorized into levels such as "Ideal," "Premium," "Good," "Fair," and "Very Good."

- **Carat:** Carat is a unit of weight used to measure the size of diamonds. It's one of the most important factors influencing a diamond's price, as larger diamonds tend to be more valuable.

- **Price:** Price is the target variable in our analysis. It represents the monetary value assigned to each diamond based on its various attributes, including cut, carat, color, and clarity.

### 3.4 Dataset Structure

The dataset consists of several columns, each containing information about different attributes of the diamonds. Some of the key columns include:

- **Cut:** Categorical variable representing the quality of the diamond's cut.

- **Carat:** Numeric variable representing the weight of the diamond in carats.

- **Price:** Numeric variable representing the price of the diamond in a specific currency.

- **Other Columns:** The dataset also include additional columns such as color, clarity, depth,and its dimensions.

## 4. Data Analysis:

### 4.1 Descriptive Statistic

To begin, the summary statistics provide a brief overview of the key characteristics of the data set. These statistics are calculated to help understand the central tendency, variability, and distribution of the data.

```
##      price                cut            carat
##  Min.   :  326   Fair     : 1610   Min.   :0.2000
##  1st Qu.:  950   Good     : 4906   1st Qu.:0.4000
##  Median : 2401   Very Good:12082   Median :0.7000
##  Mean   : 3933   Premium  :13791   Mean   :0.7979
##  3rd Qu.: 5324   Ideal    :21551   3rd Qu.:1.0400
##  Max.   :18823                     Max.   :5.0100
```
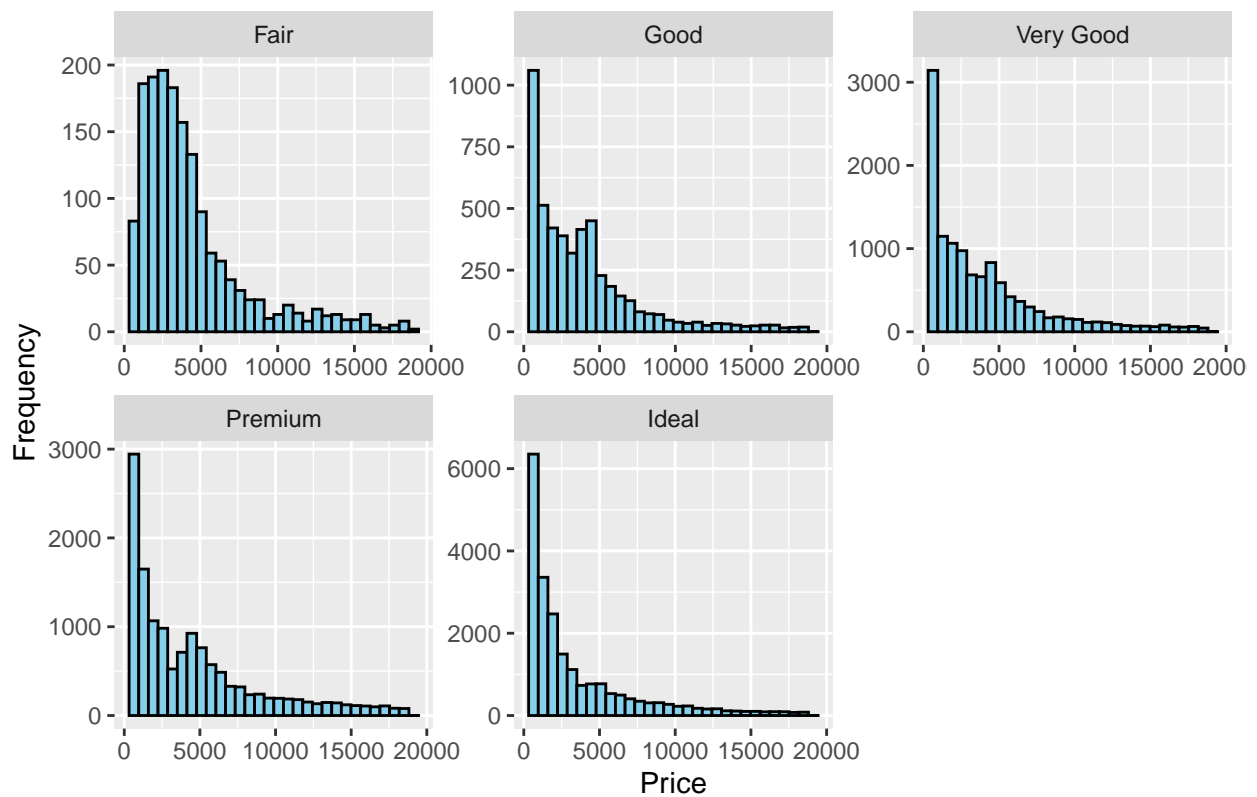
Based on the summary statistics provided above, we can make the following observations about the diamond dataset:

- Price ranges from $326 to $18823, with a median price of $2401 and a mean price of approximately $3933.

- Cut categories include Fair, Good, Very Good, Premium, and Ideal, with varying counts for each category.

- Carat weight ranges from 0.2 to 5.01, with a median of 0.7 and an average of around 0.7979 carats.
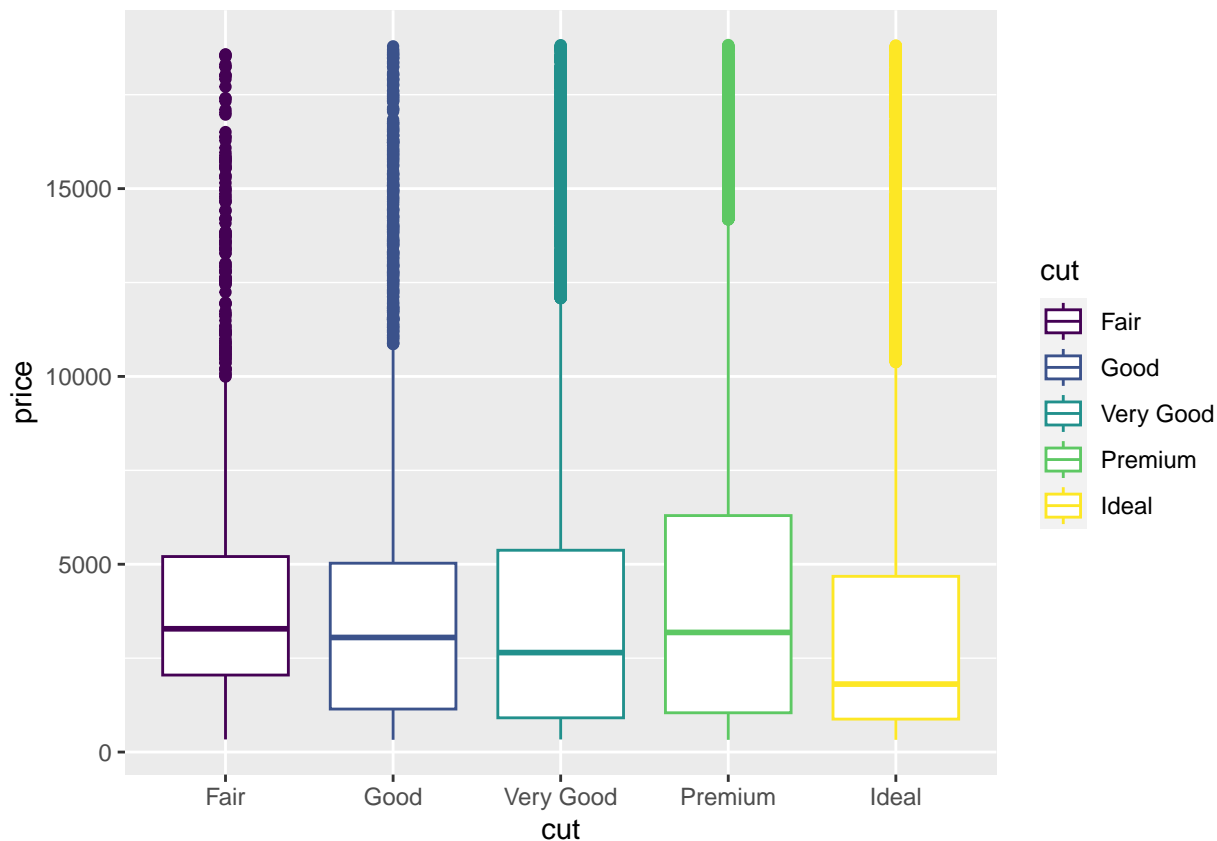
### 4.2 Exploring Price and Cut Relationships

The visualization below depicts the correlation between diamond cut and price through faceted histograms that highlight their respective skewness. As observed in the ensuing visualization, the price distribution for each diamond cut skews to the right, implying an abundance of diamonds within specific price ranges. This suggests that there are comparatively fewer diamonds with exceptionally high prices across all cut categories.
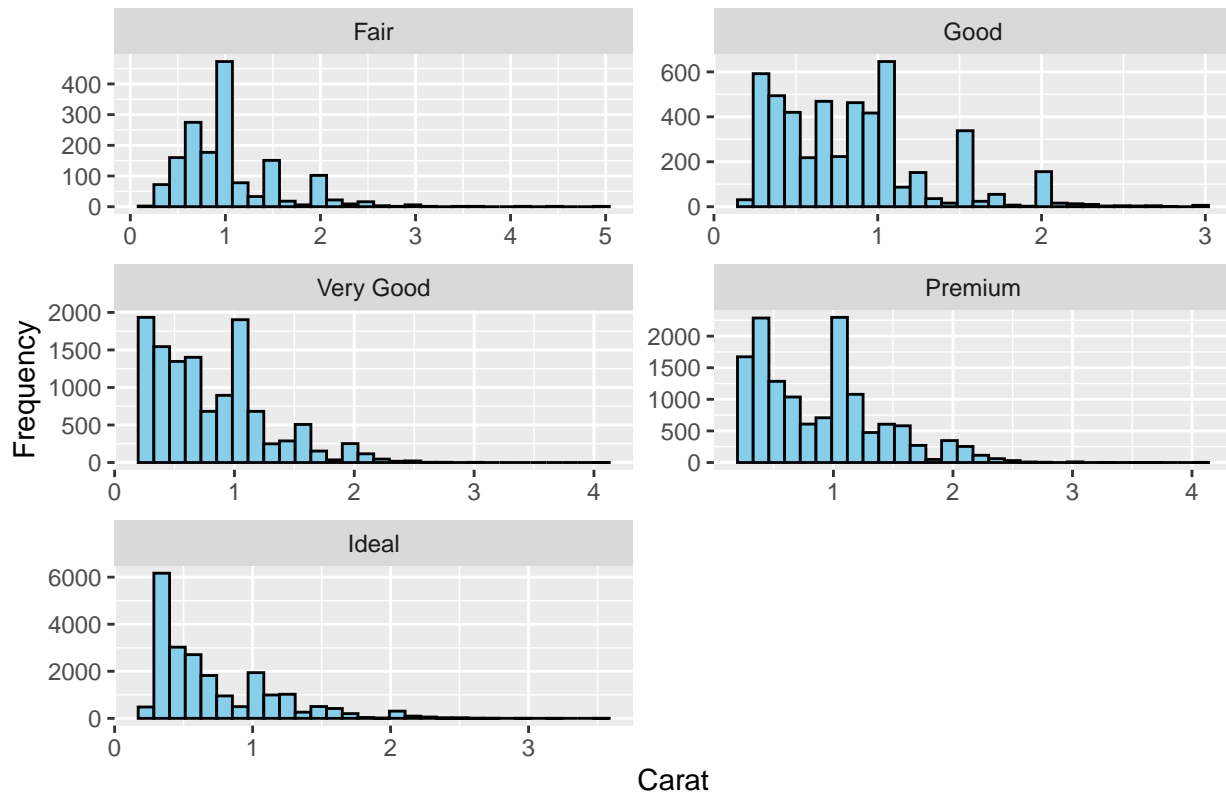
## Distribution of Diamond Prices by Cut



Examining the box plot below, a clear trend emerges: diamonds with premium cuts command higher prices, while those with ideal cuts are associated with lower prices. This simple yet insightful observation underscores the influence of cut quality on diamond pricing.
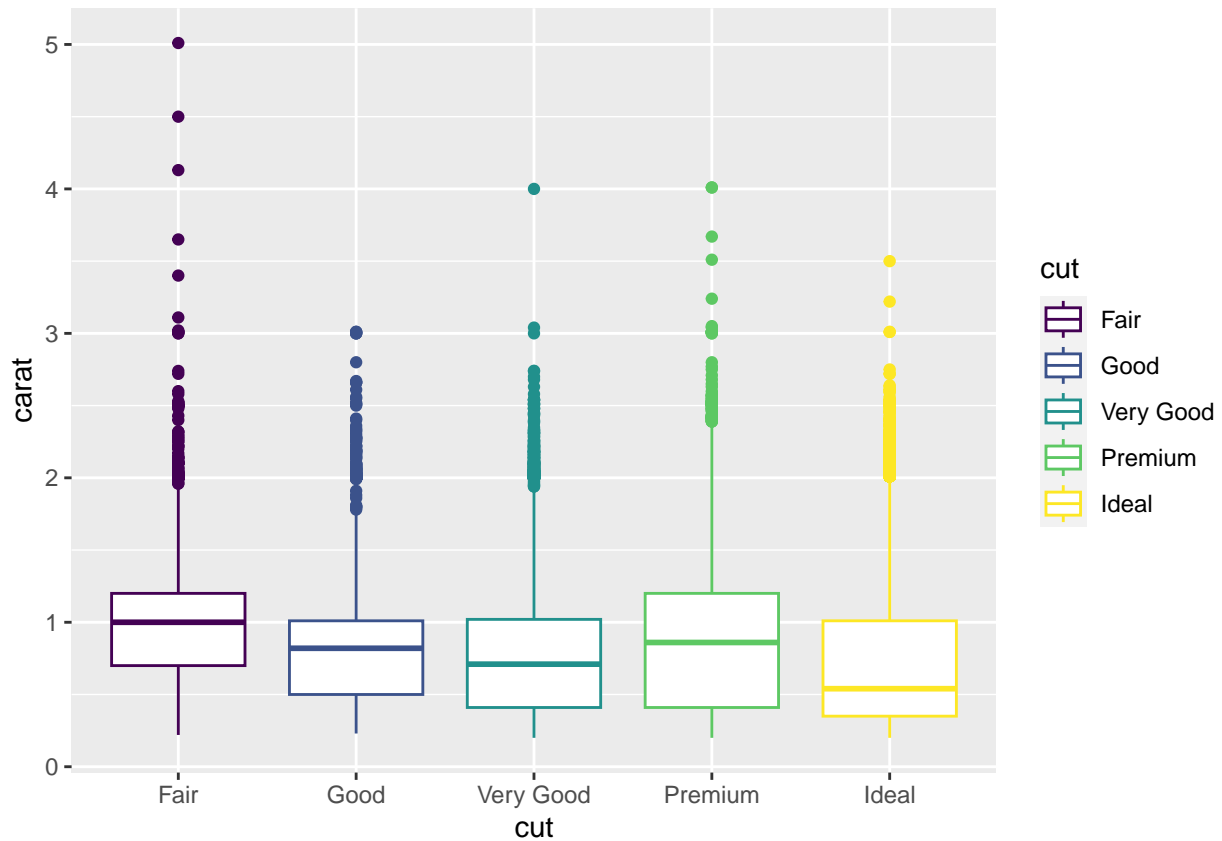
### 4.3 Exploring Carat and Cut Relationships

The visualization below depicts the correlation between diamond cut and carat through faceted histograms that highlight their respective skewness. As observed in the ensuing visualization, the carat distribution for each diamond cut slightly skews to the right, implying an abundance of diamonds within specific carat ranges. This suggests that there are comparatively fewer diamonds with exceptionally high carat across all cut categories.
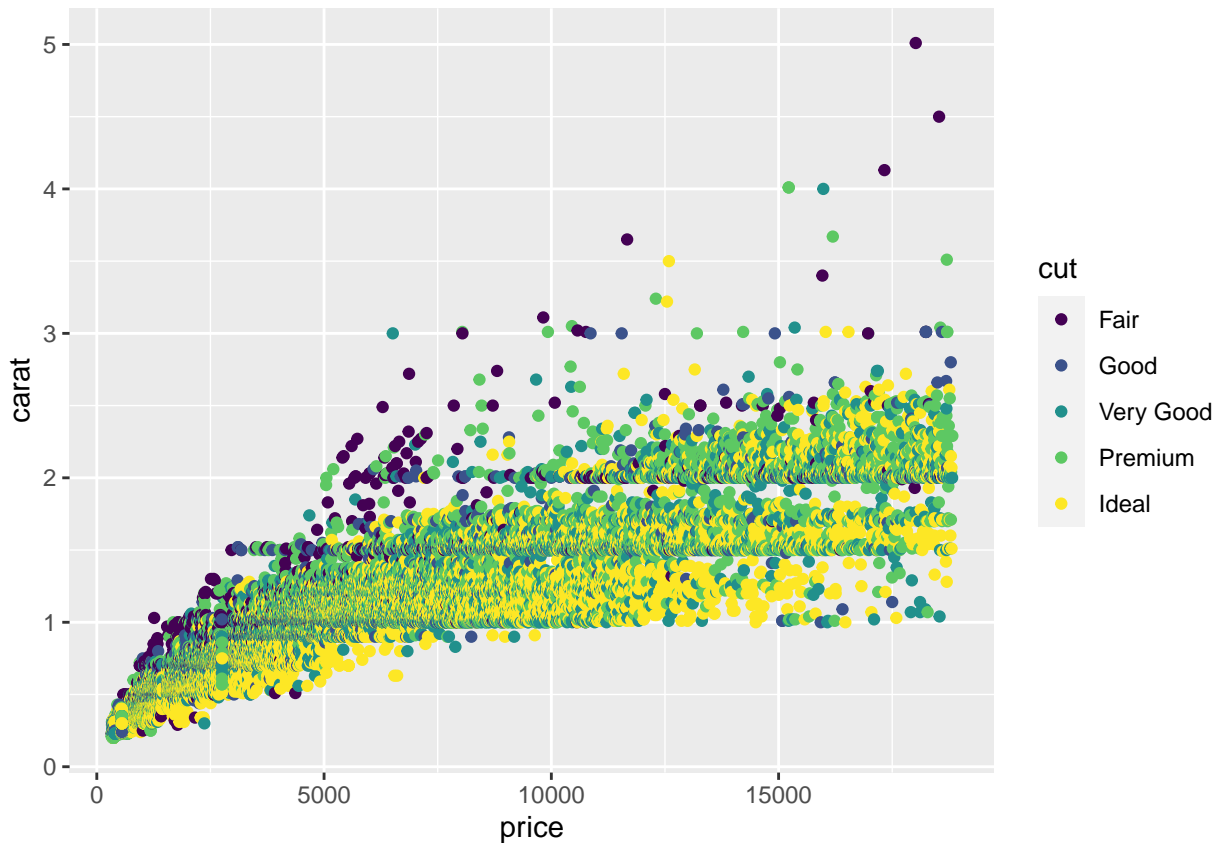
Distribution of Carat Values by Cut

Also, the boxplot below, displays various diamond cuts—Fair, Good, Very Good, Premium, and Ideal—plotted against carat weight, including outlier data points. We can also see that carat and cut hold a negative correlation. Diamonds with Good, Very Good, and Ideal cuts generally have a carat weight of less than 1, as indicated by a median value of <1. While a few Fair and Premium cut diamonds slightly exceed 1 carat in weight, their medians remain at <=1 carat. The Premium cut exhibits the most expansive range between the 1st and 3rd quartiles. Larger diamonds are predominantly associated with the Fair cut.

The scatter plot reveals a robust positive correlation between carat and price, with a predominant concentration of diamonds with low carat values along the x-axis. It is evident that diamonds with lower carat weights tend to have correspondingly lower prices. As carat size escalates, there is a noticeable upward trend in diamond prices.

## 5. Regression Analysis

In the realm of data analysis, accurate regression models are vital for predictive modeling. We employ the lm function, a key tool in statistical programming, to build a robust regression model for our project. This model, crafted from carefully selected features (price, cut, carat), predicts diamond prices effectively.

Regression models unveil relationships between variables through mathematical representations. The lm function helps us intricately understand how chosen features (price, cut, carat) influence diamond prices. This predictive tool becomes invaluable, offering insight into future diamond prices by considering these amalgamated features. Below is the summary of the model built for this project.

```
##
## Call:
## lm(formula = price ~ carat + cut, data = diamond_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17540.7   -791.6    -37.6    522.1  12721.4
##
## Coefficients:
##               Estimate Std. Error  t value Pr(>|t|)
```

```
## (Intercept) -2701.38      15.43 -175.061  < 2e-16 ***
## carat        7871.08      13.98  563.040  < 2e-16 ***
## cut.L        1239.80      26.10   47.502  < 2e-16 ***
## cut.Q        -528.60      23.13  -22.851  < 2e-16 ***
## cut.C         367.91      20.21   18.201  < 2e-16 ***
## cut^4          74.59      16.24    4.593 4.37e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1511 on 53934 degrees of freedom
## Multiple R-squared:  0.8565, Adjusted R-squared:  0.8565
## F-statistic: 6.437e+04 on 5 and 53934 DF,  p-value: < 2.2e-16
```

**Model Summary:**

- The model's residual standard error is 1511, which represents the average magnitude of the differences between observed and predicted values.

- The multiple R-squared value of 0.8565 indicates that around 85.65% of the variability in diamond prices is explained by the model's predictors (carat and cut).

- The adjusted R-squared, which takes into account the number of predictors, remains the same as the multiple R-squared, indicating that the added predictors (cut levels) contribute meaningfully to explaining price variability.

- The F-statistic is 6.437e+04, and its associated p-value is practically zero ($< 2.2$e-16), suggesting that the overall model is highly significant.

**Coefficients:**

The intercept of -2701.38 represents the estimated baseline price for diamonds with zero carat weight and cut quality, though this is impractical for diamonds. The carat coefficient of 7871.08 suggests that each additional carat increases the predicted diamond price by \$7871.08, assuming cut quality remains constant. The coefficients for different cut levels (L, Q, C, ^4) reveal their impact on price compared to the reference level, "Ideal." For example, an "L" cut elevates price by 1239.80, while "Q" and "C" cuts decrease price. The '^4' level indicates a nonlinear relationship between this cut and price.
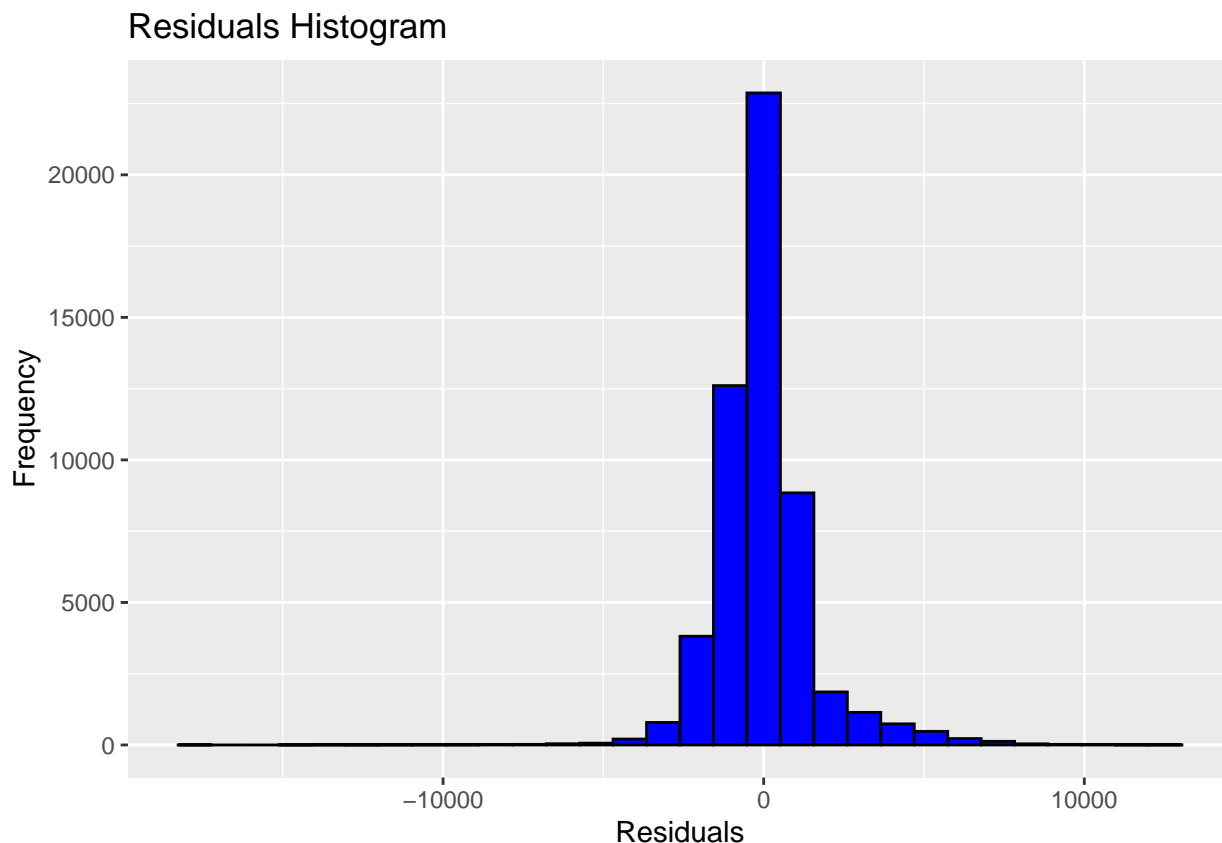
**Residuals:**

The residuals, which represent the differences between the actual and predicted prices, have a distribution with a minimum of -17540.7 and a maximum of 12721.4. The majority of residuals are within the range of -791.6 to 522.1, indicating that the model generally captures the price trends well.

In summary, the model suggests that carat weight and cut quality are strong predictors of diamond prices. The coefficients provide insights into the magnitude and direction of these

relationships, considering the different cut levels as well. The model appears to fit the data quite well, as indicated by the R-squared values and the significant F-statistic.

## 5.2 Evaluating Models Performance

### Residuals Histogram



## 5.1 Generating Predictions from the Model

Subsequently, we move forward to conduct predictions using the generated model. Our objective is to estimate the prices for diamonds with ideal, premium, and fair cuts, each weighing 5.0 carats. The prediction outcomes are presented below.

```
##        1        2        3
## 37280.86 36919.02 35479.94
```

Analyzing the predictions further, we observe that an Ideal cut diamond weighing 5.0 carats is estimated to cost $37,280.86. Comparatively, a Premium cut diamond of the same weight is projected to be priced at $36,919.20. In contrast, a Fair cut diamond with a 5.0-carat weight is anticipated to have a lower price of $35,479.94.

These predictions reflect the model's understanding of how different cut qualities influence the prices of diamonds of equal weight. The price variations suggest that cut quality significantly impacts diamond valuation, with an Ideal cut commanding the highest price, followed by Premium and Fair cuts. This insight is consistent with industry norms, where cut quality is pivotal in determining a diamond's market value.