

Color Writeup

Terrance Luangrath

2025-03-09

Data Description

The data set, `diamonds4.csv`, contains 5 variables that describes more than 1000 different diamonds for sale.

- **Carat:** *fill in description*
- **Clarity:** *fill in description*
- **Color:** how colorless a diamond is, the more colorless the diamonds, the rarer it tends to be
- **Cut:** *fill in description*
- **Price:** the price value of the diamond in USD

Variable Analysis: Color

The diamond color refers to how colorless a diamond is. From the Blue Nile website, color is the second most important of the 4Cs of diamond. The more colorless the diamond is the rarer it is. Diamond colors are classify in three main categories,

- Colorless Diamonds
- Near-Colorless Diamonds
- Faint Diamonds

from these three categories, diamond color are grades from the ranges of D (colorless) to K (faintly colored). However, in this dataset, there exists only diamond graded from D to J.

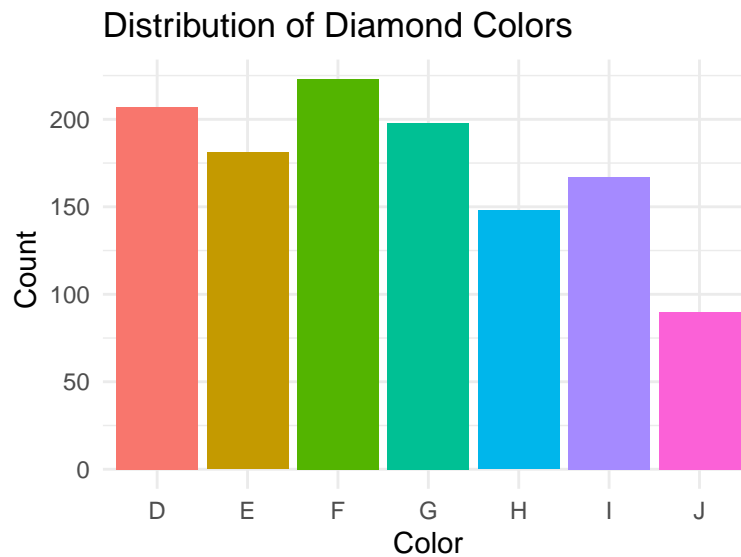


Figure 1: Diamond Color Distribution

Looking at Figure ??, the distribution of diamond color appears fairly consistent, with counts ranging between 150-200 for most color grades. However, color J seems to have around 100 diamonds in the data set.

Each diamond color falls into one of the three categories from Colorless to Faintly Colored.

- Colorless diamonds: D, E, F Color Diamonds
- Near-colorless diamonds: G, H, I, J Color Diamonds
- Faint color diamonds: K Color Diamonds

Since there is no K in the data set, it would be assumed there is no faint color diamonds in the data set. To facilitate the analysis, we will create a new column called `color_category` to group each `color` with their respected category.

Distribution of Diamond Color Groups

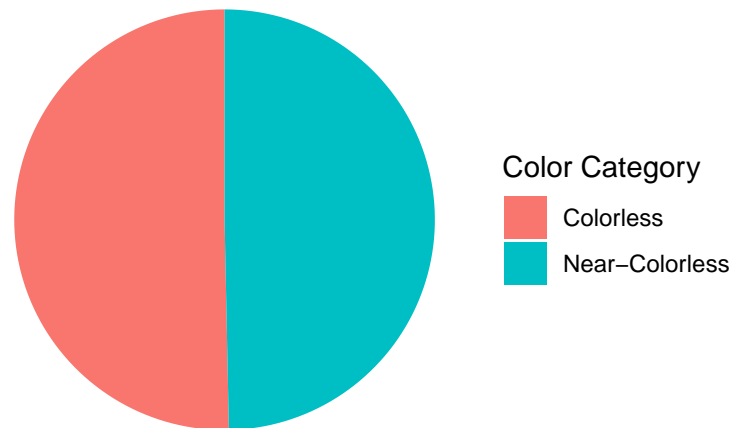


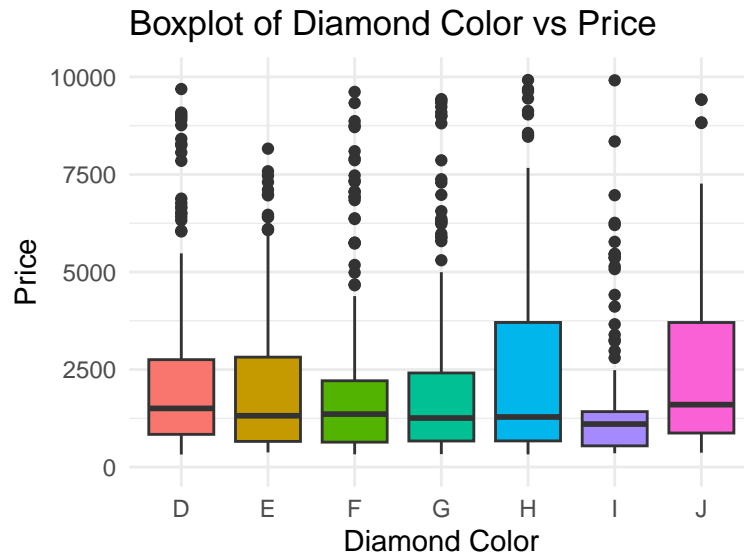
Figure 2: Diamond Color Category visualize through a Pie Chart

Looking at ??, the distribution of diamond colors seems to be nearly even, with `Colorless` containing slightly higher amount of diamonds. With a well spread of diamond colors, further analysis is needed to understand the claim and how diamond color would be correlates with the other factors in the data set.

Bivariate Analysis

Bivariate Analysis on Prices

```
## # A tibble: 7 x 6
##   color   min    q1 median    q3    max
##   <chr> <int> <dbl> <dbl> <dbl> <int>
## 1 D       322  882   1781 6048. 355403
## 2 E       376  705   1602 5138  345397
## 3 F       328  704   1485 4372  227960
## 4 G       332  689.   1374 3872.  165766
## 5 H       326  728.   1575 6135  134856
## 6 I       354  574.   1212 1680.  123311
## 7 J       369  958.   1803 4504   40184
```

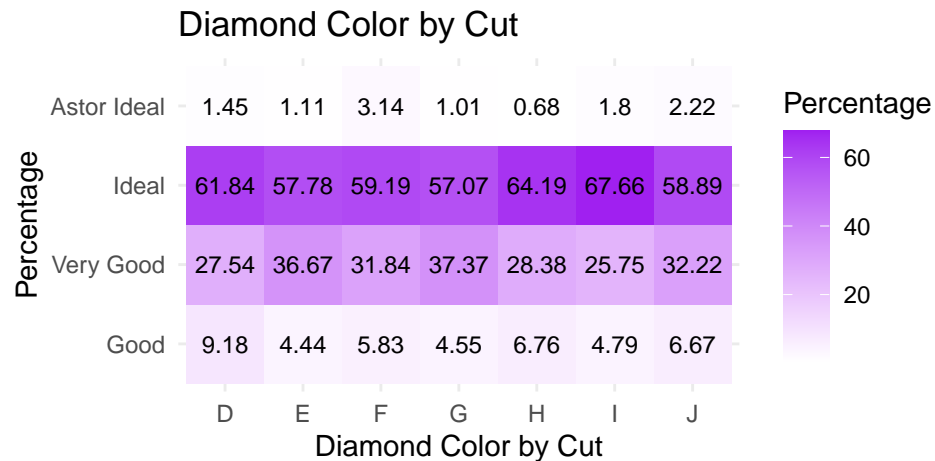


On the Blue Nile site, it mentions that color is the second most important among the 4C's for the price of the diamond. Reviewing the box plot, the colorless diamonds (D, E, and F) tends to have higher prices than the near-colorless diamond (G, H, I, and J). However, some colorless diamond have higher median values than the near-colorless diamond:

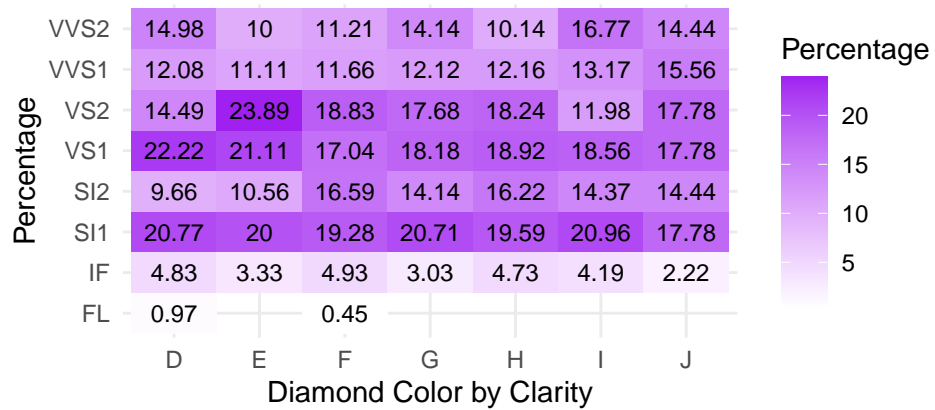
- H have higher median values than the colorless diamond F
- J (the lowest color grade) has a higher median value than all the colorless diamond.

This suggests that color does play a significant factors in the diamond price, but it's not the most important factors to the diamond prices. Other characteristics, such as `cut`, `clarity`, and `carat`, needs to be considered to help prove the claim.

Color vs Other Diamond Characteristics



Diamond Color by Clarity



Diamond Color by Carat

