

# Exploratory Data

Kylie A. Bem

Northeastern Uni  
Khoury College of Compu



Northeastern U

## Learning goals

- What is EDA?
- Variation and covariation
- “Interesting” visualizations

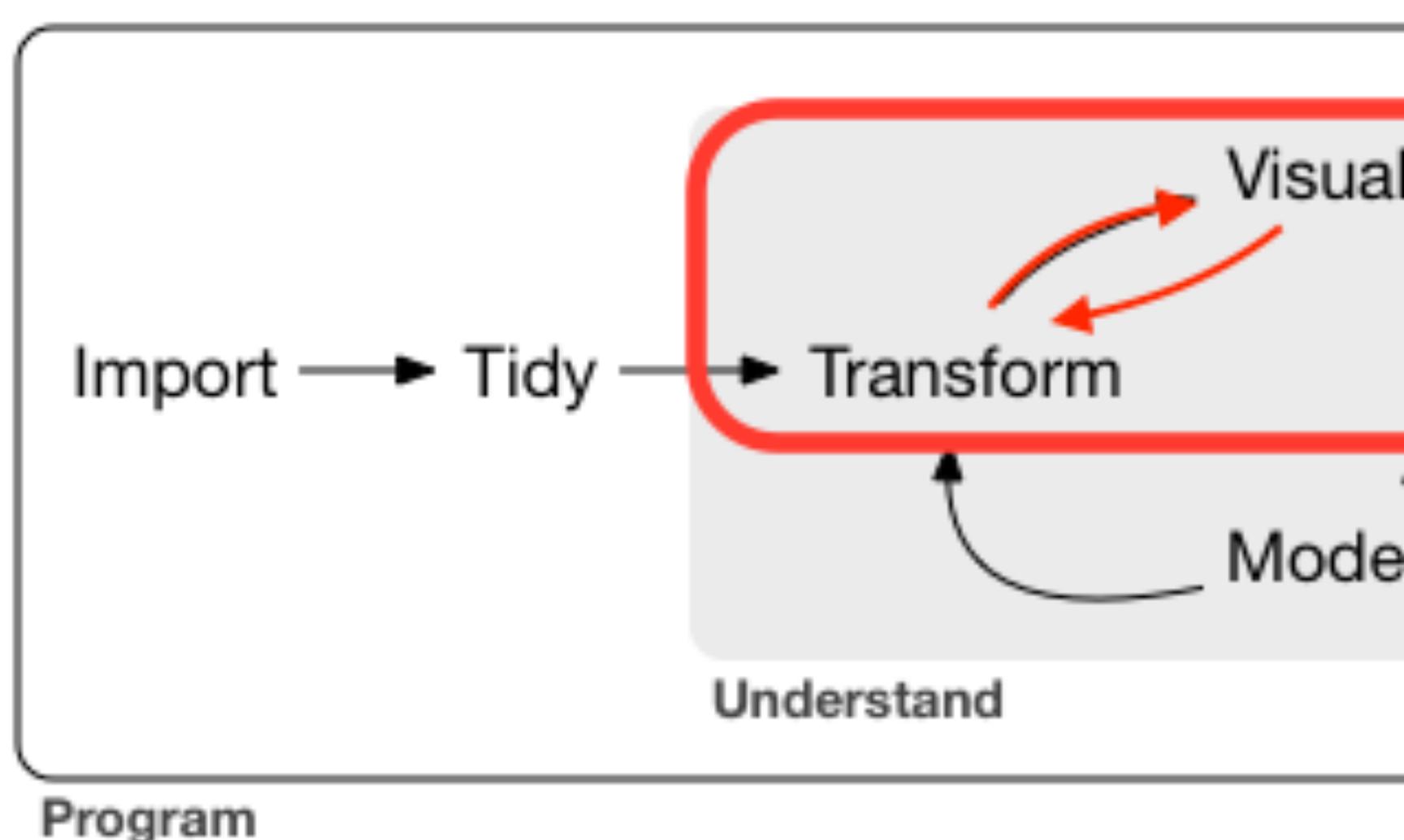
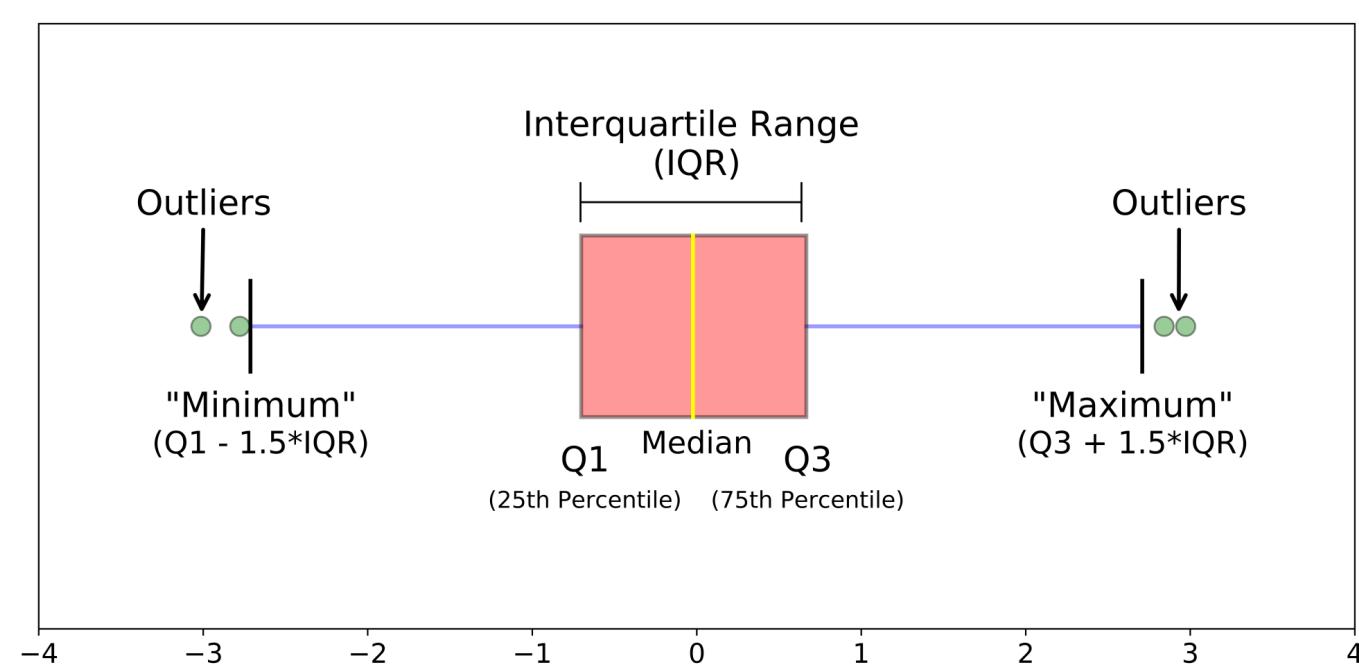
# EDA: **EXPLORATORY DATA**

## Exploratory data

EDA is an iterative, cyclical process of inform  
and relationships between th

- Generate questions about your data
- Search for answers by trying different ways of visualizing your data
- Use this to refine your questions and generate new questions

# Process of E



R for Data Science, Wickham an

# Example data: Diamonds

Prices and measurements of diamonds

```
## #     tibble: 53,940 x 10
##       carat   cut      color clarity depth table
##       <dbl> <ord>    <ord> <ord>   <dbl> <
## 1 0.23   Ideal     E      SI2      61.5
## 2 0.21   Premium   E      SI1      59.8
## 3 0.23   Good      E      VS1      56.9
## 4 0.290  Premium   I      VS2      62.4
## 5 0.31   Good      J      SI2      63.3
## 6 0.24   Very Good J      VVS2     62.8
## 7 0.24   Very Good I      VVS1     62.3
## 8 0.26   Very Good H      SI1      61.9
## 9 0.22   Fair       E      VS2      65.1
## 10 0.23  Very Good H      VS1      59.4
## # ... with 53,930 more rows
```

Available in ggplot2 package

## Question

- What type of variation occurs with covariation?
- What type of covariation occurs with regression?

## Useful definitions

- A **variable** is a quantity, quality or characteristic that is measured
- A **value** is the state of a variable at a particular time or place when it is measured
- An **observation** is a set of measurements made under similar conditions

## Types of variables

- Categorical — qualitative
  - ◆ Nominal — unordered categories
  - ◆ Ordinal — ordered or ranked categories
- Numeric — quantitative and continuous
  - ◆ Continuous — infinite possible values
  - ◆ Discrete — countable possible values
- Time — dates and times

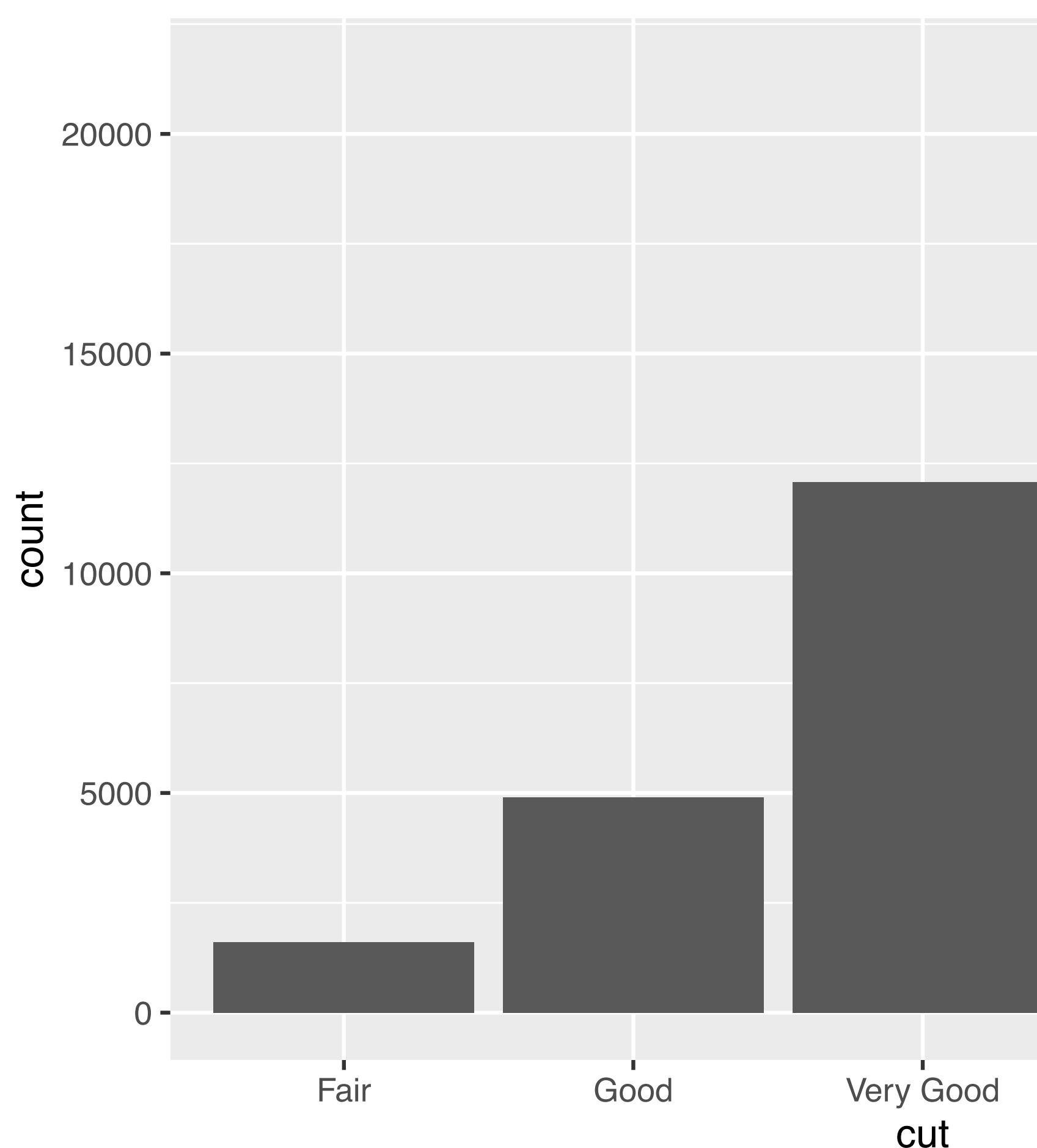
# Variation

Variation is the tendency of values of measurement to meas

- Natural variation
- Measurement error
- Time and location
- Between subjects
- Between condition

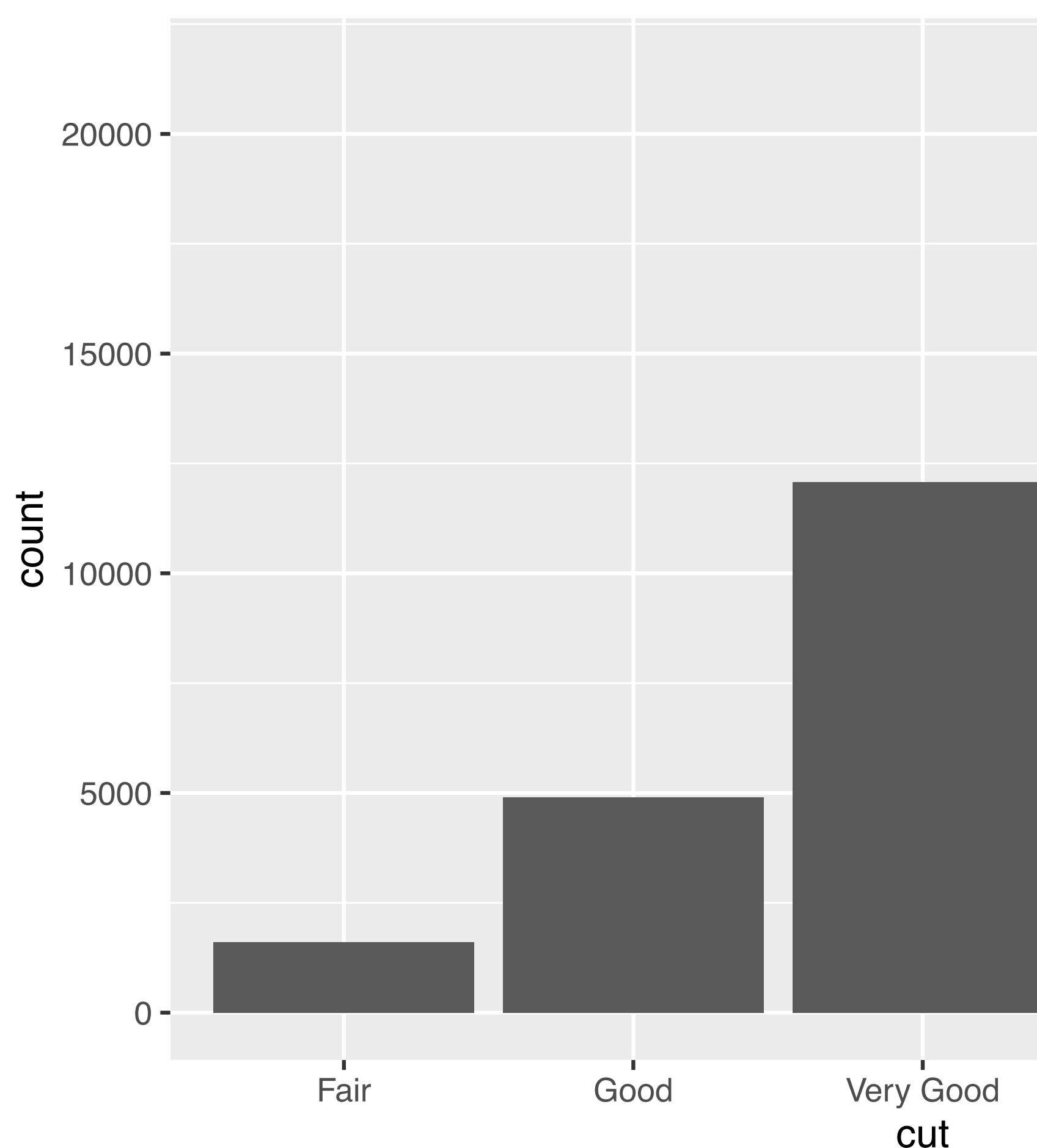
## Categorical: *diamond*

What is the distribution of cut quality?



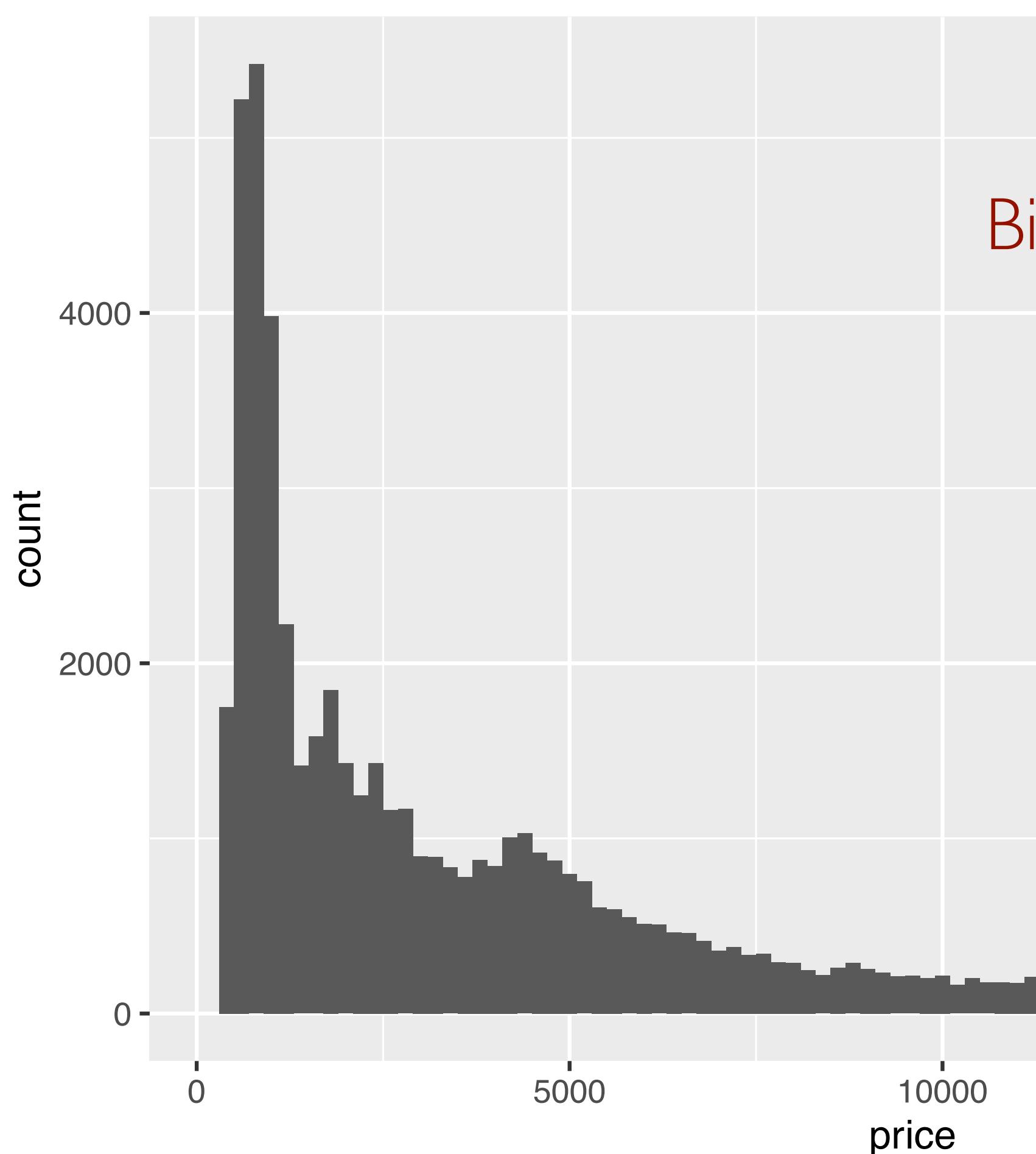
## Categorical: diamond

“Ideal” is the most common quality



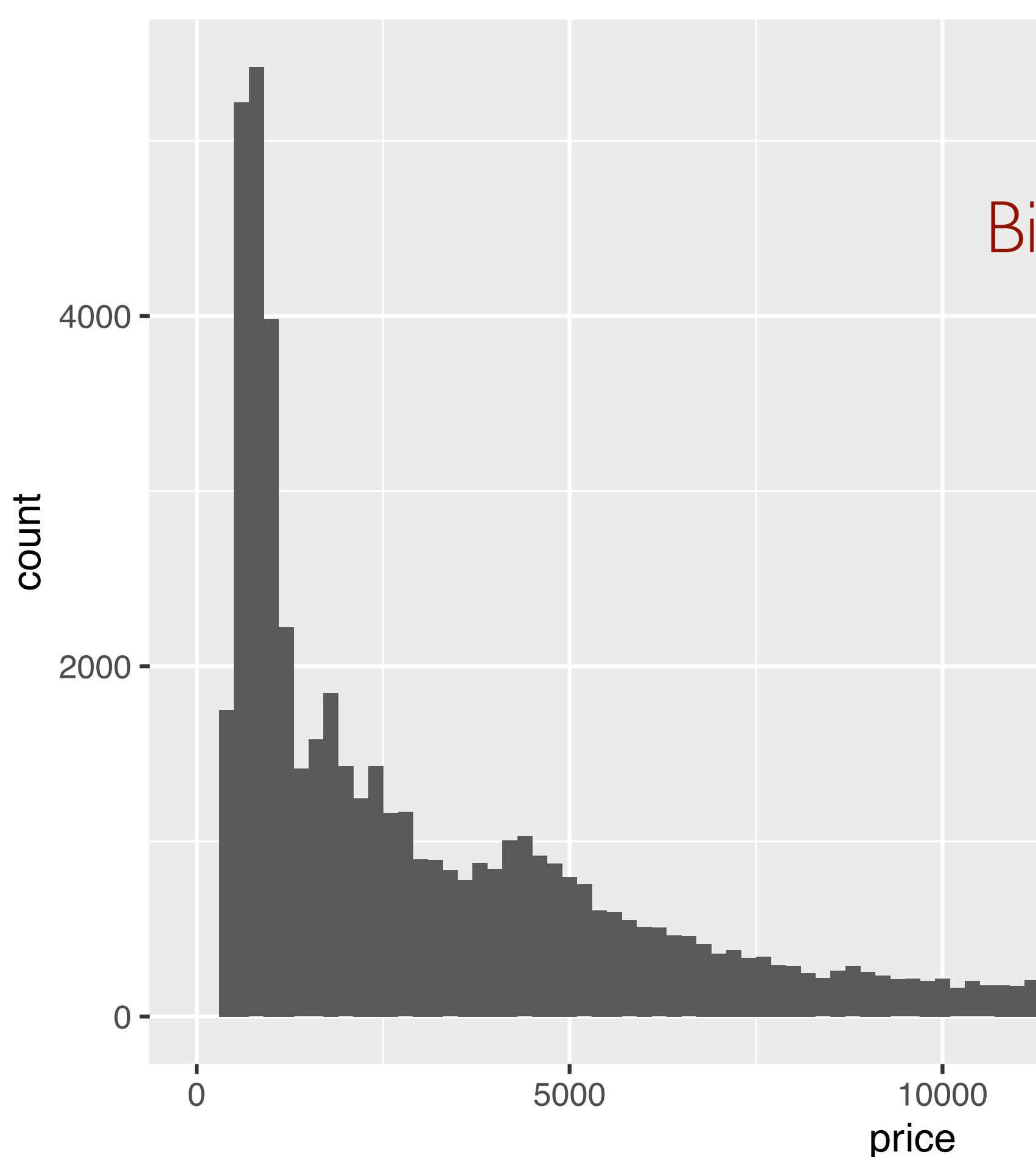
## Continuous: *diamonds*

What is the distribution of price a



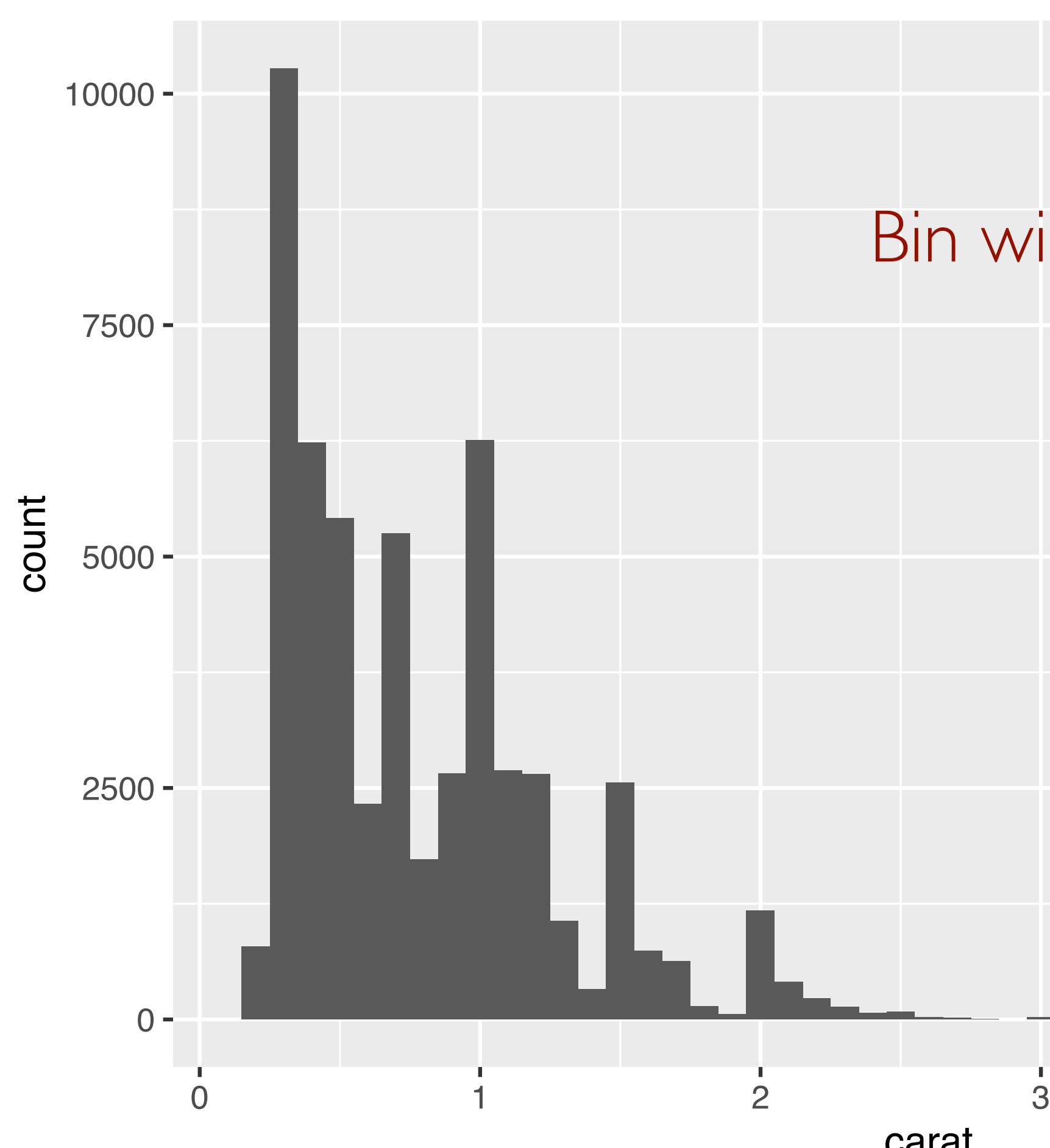
## Continuous: diam

Price is heavily right-skewed with a long tail



## Continuous: *diamond*

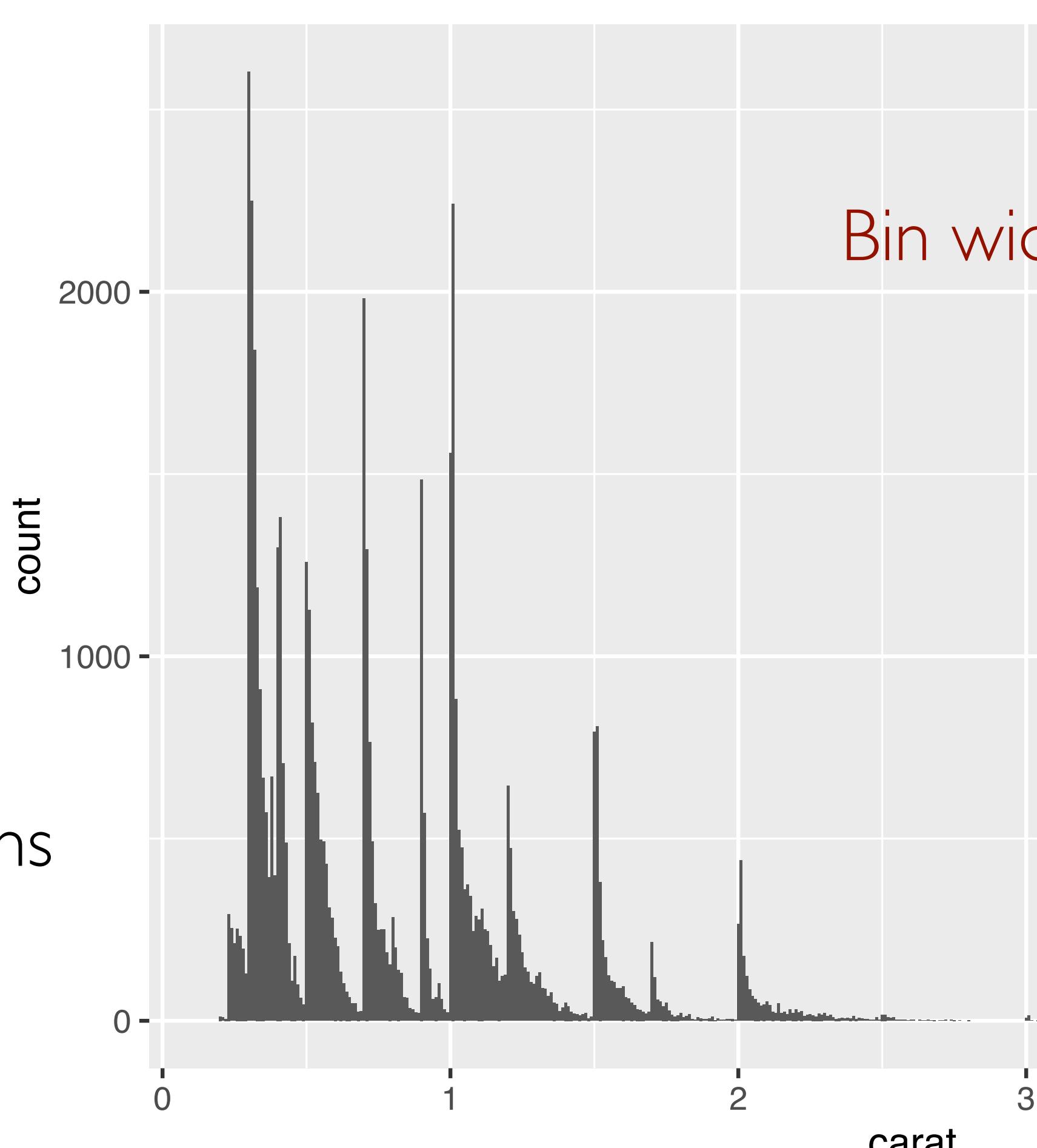
What is the distribution of carat a



## Continuous: *diamonds*

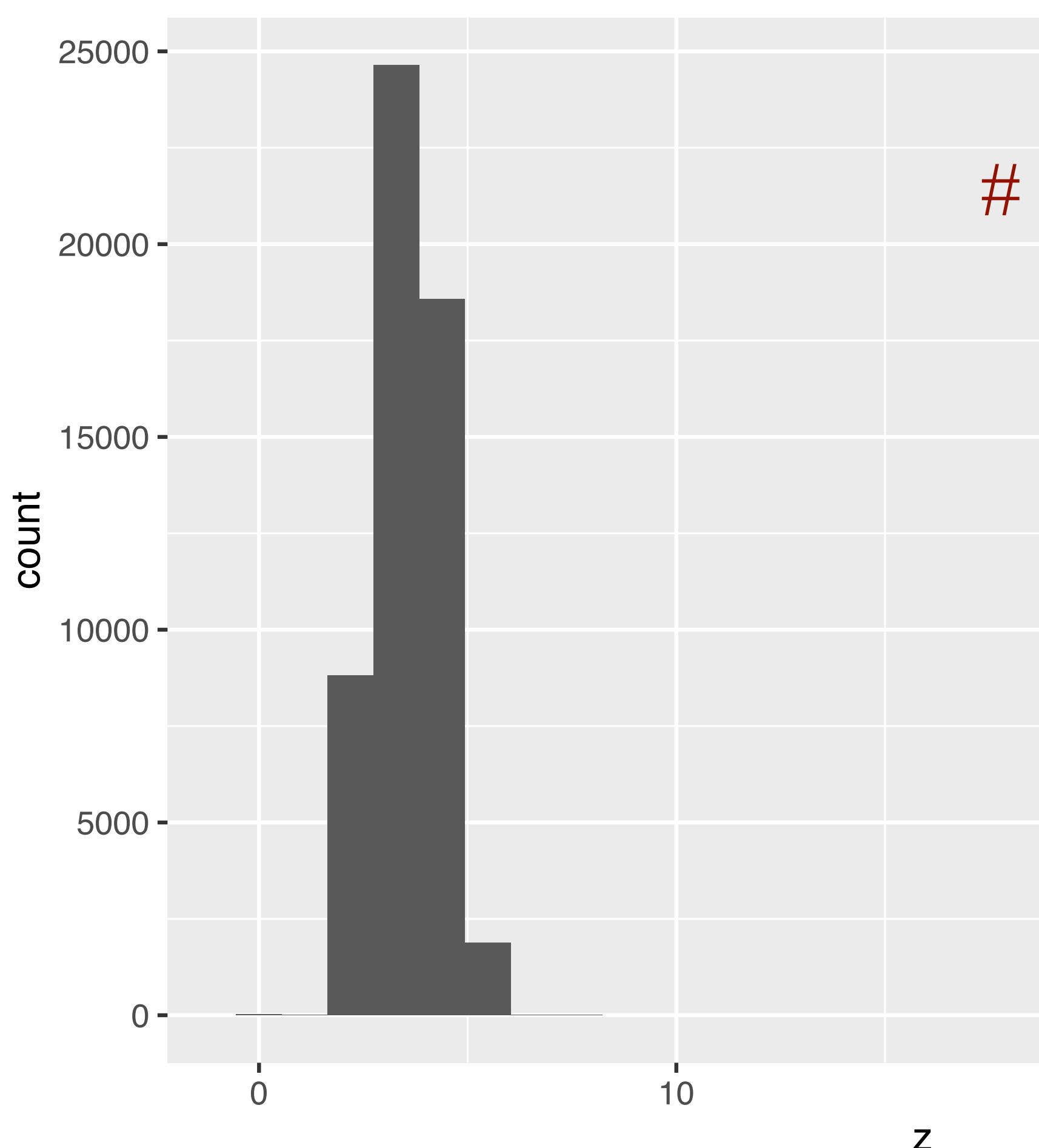
Diamonds are often cut in quadrilaterals

Try different bin widths  
or # of bins



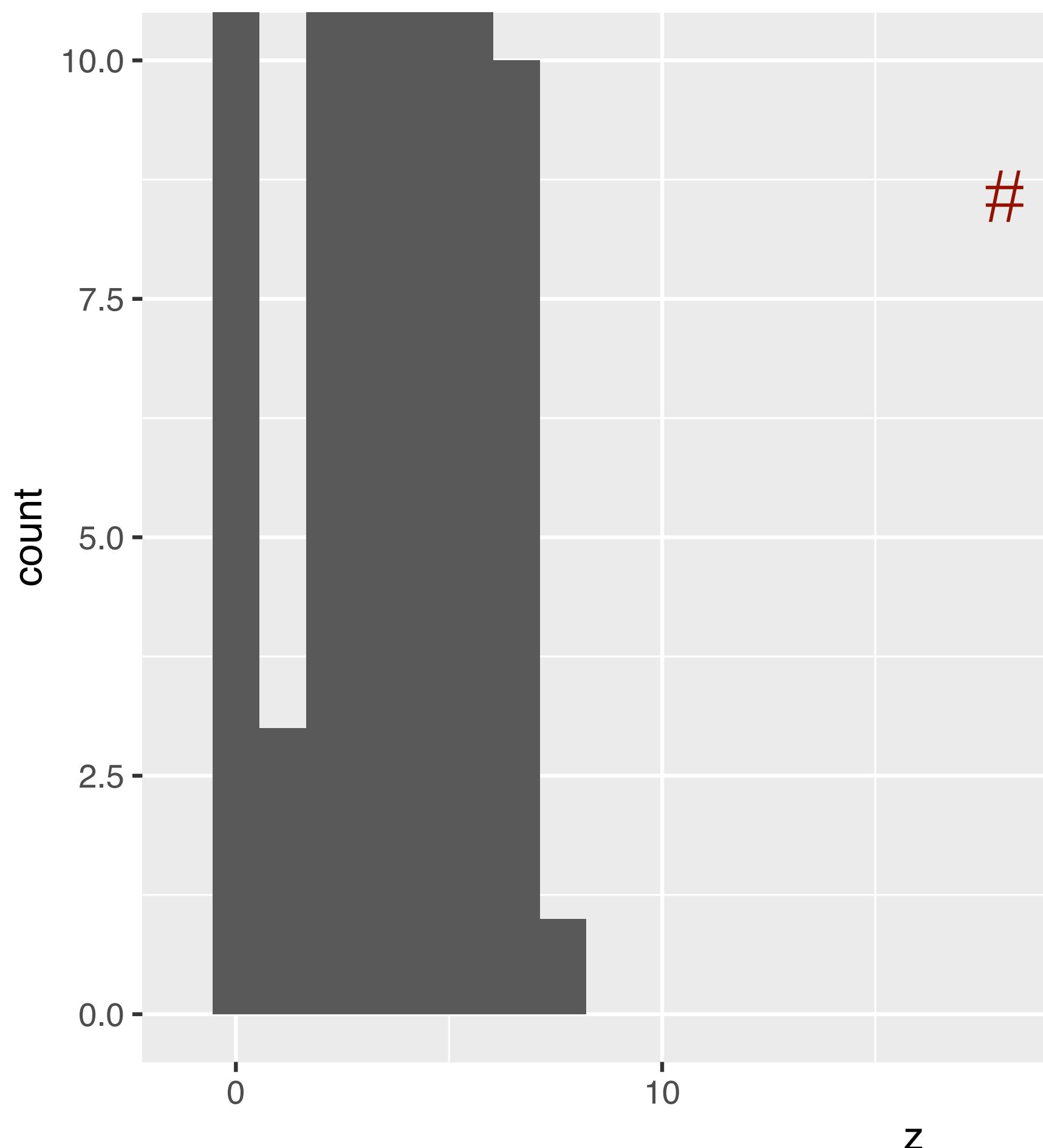
Identify outliers

Investigate z-dimension measure



# Identify outliers

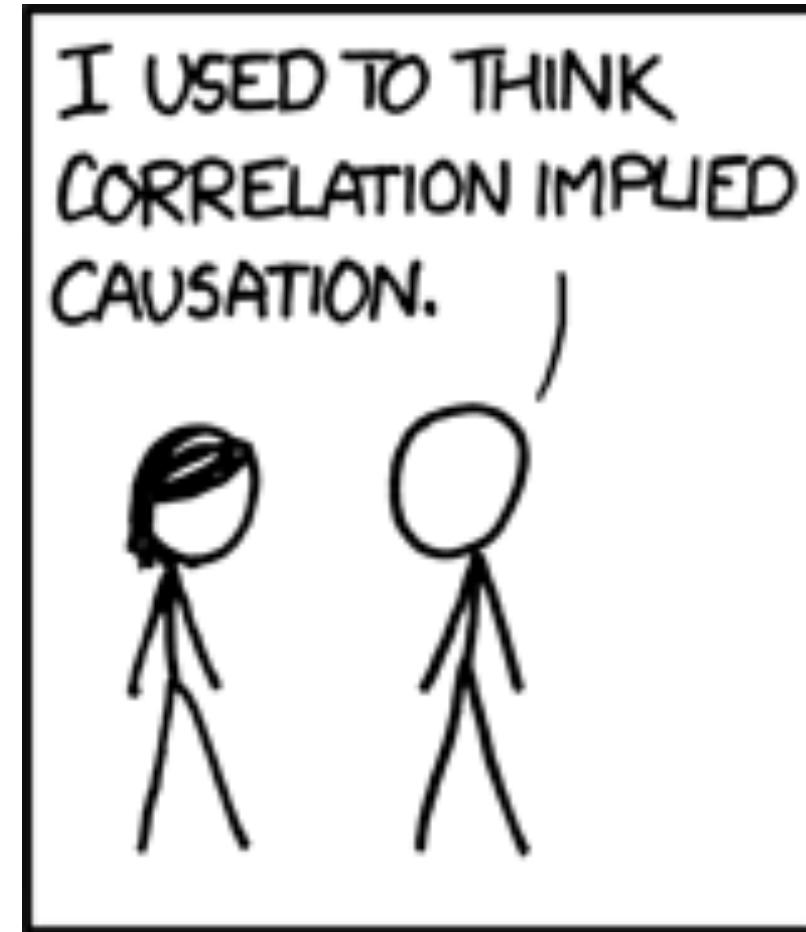
# Diamonds with $z > 30$ or $z < -30$



# Covariation

Covariation is the tendency of values of two or more variables to vary together in a regular way.

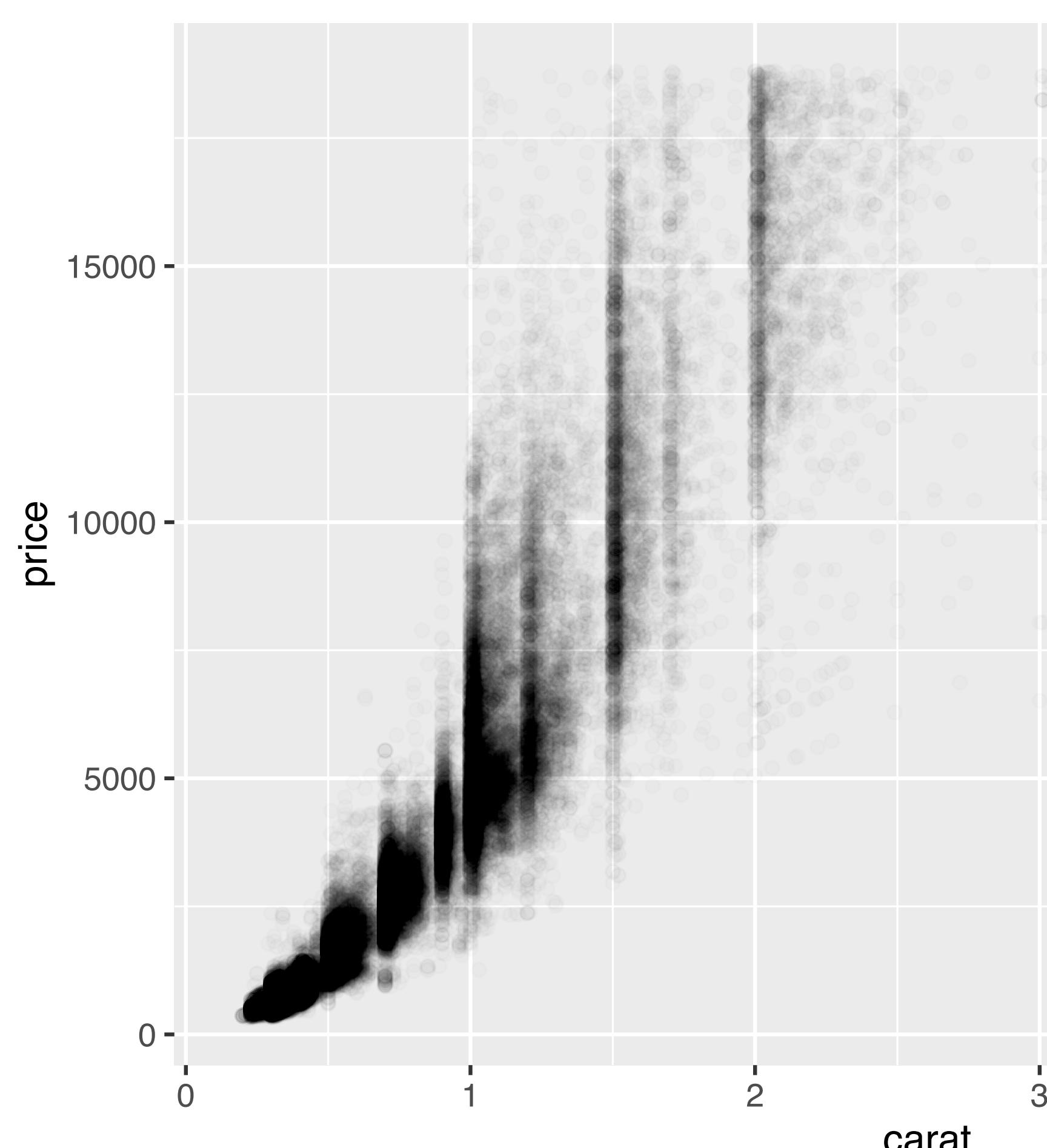
- Correlation and dependence
  - ◆ Variables share a statistical relationship
- Causation
  - ◆ One variable directly influences another
- Confounding
  - ◆ A third variable influences both variables in the study



<https://xkcd.com/55>

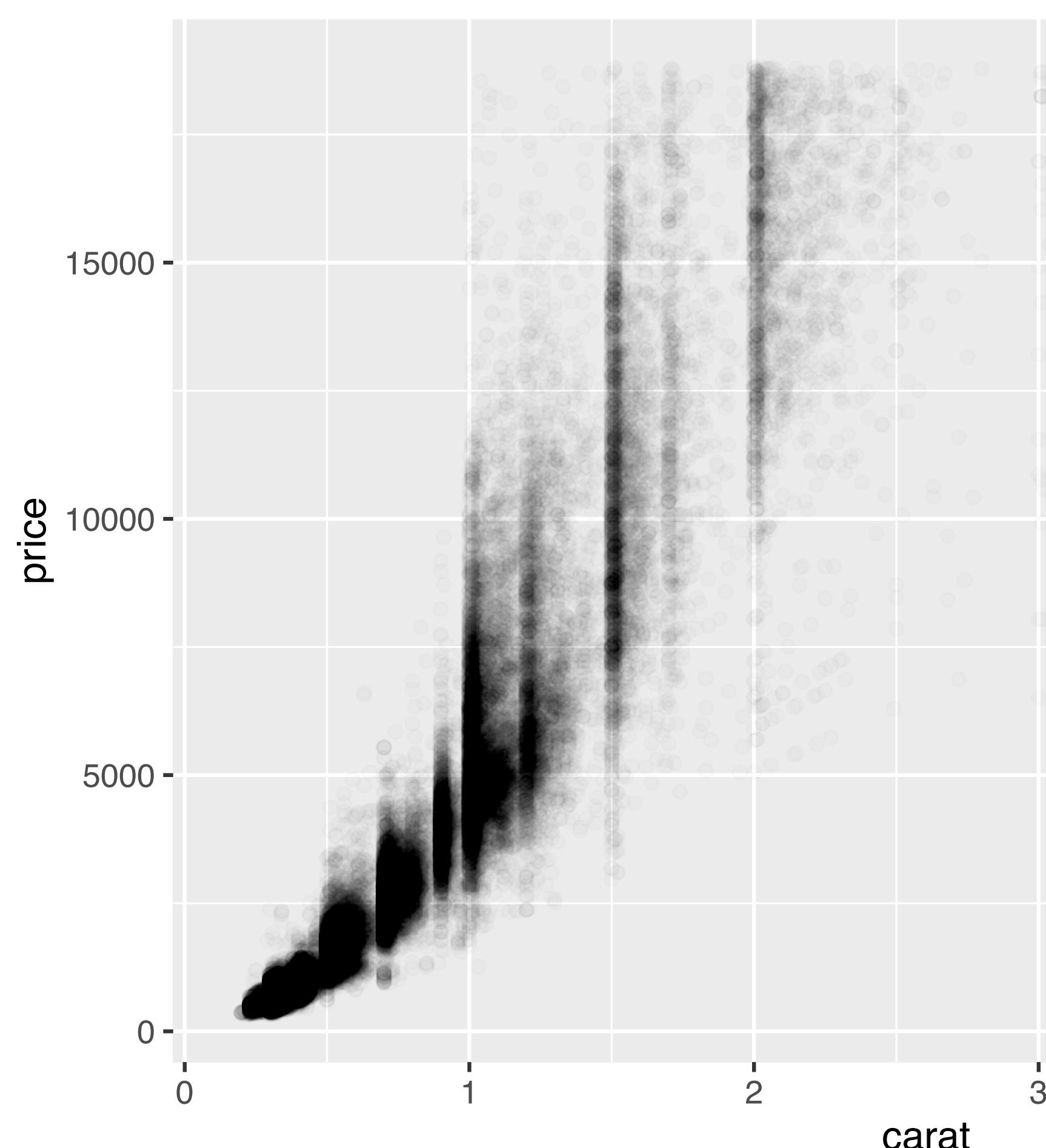
## Two continuous: $price$ vs $carat$

What is the relationship between price and carat?



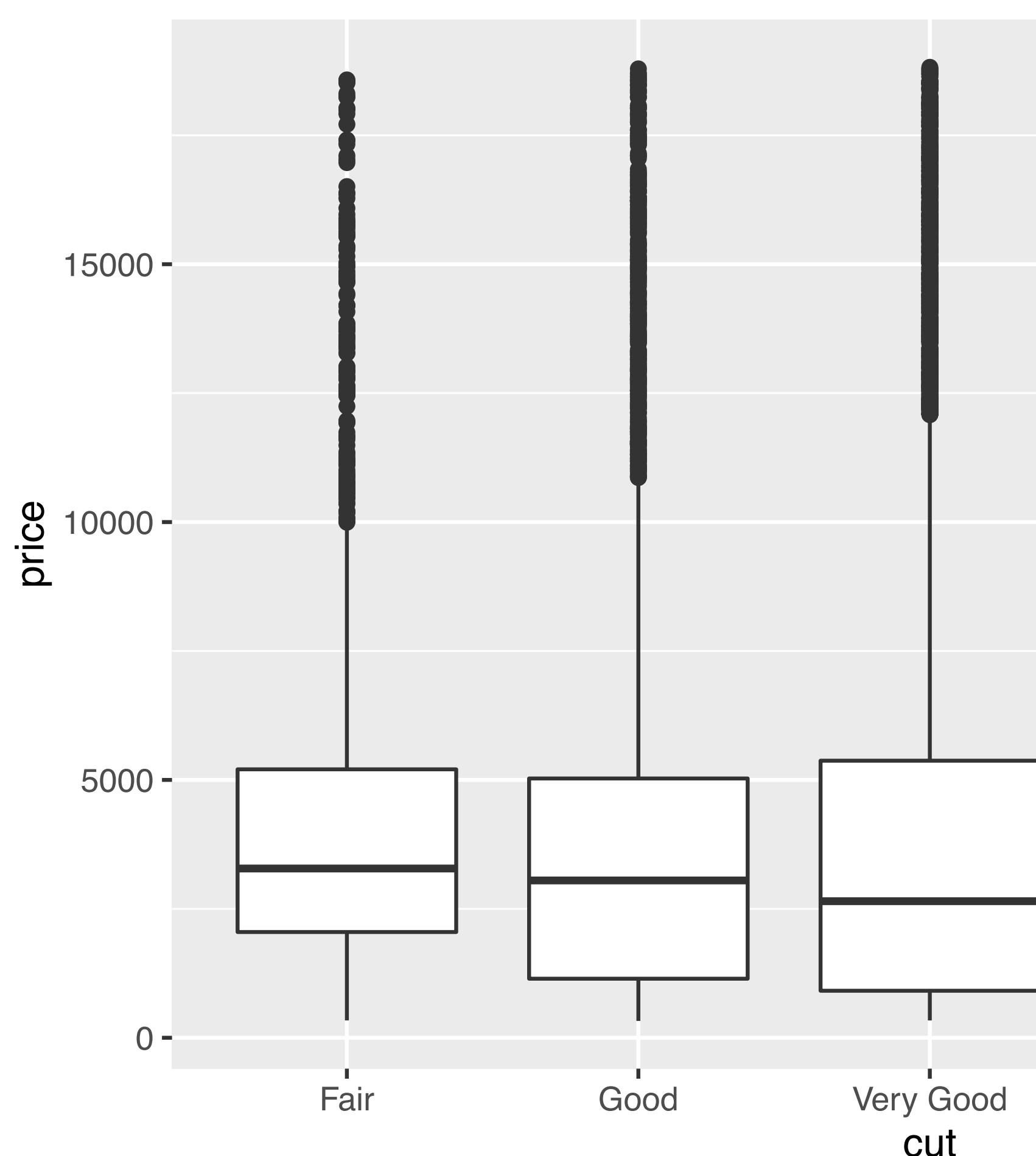
Two continuous:  $price$

Larger diamonds tend to be more expensive



# Categorical and continuous

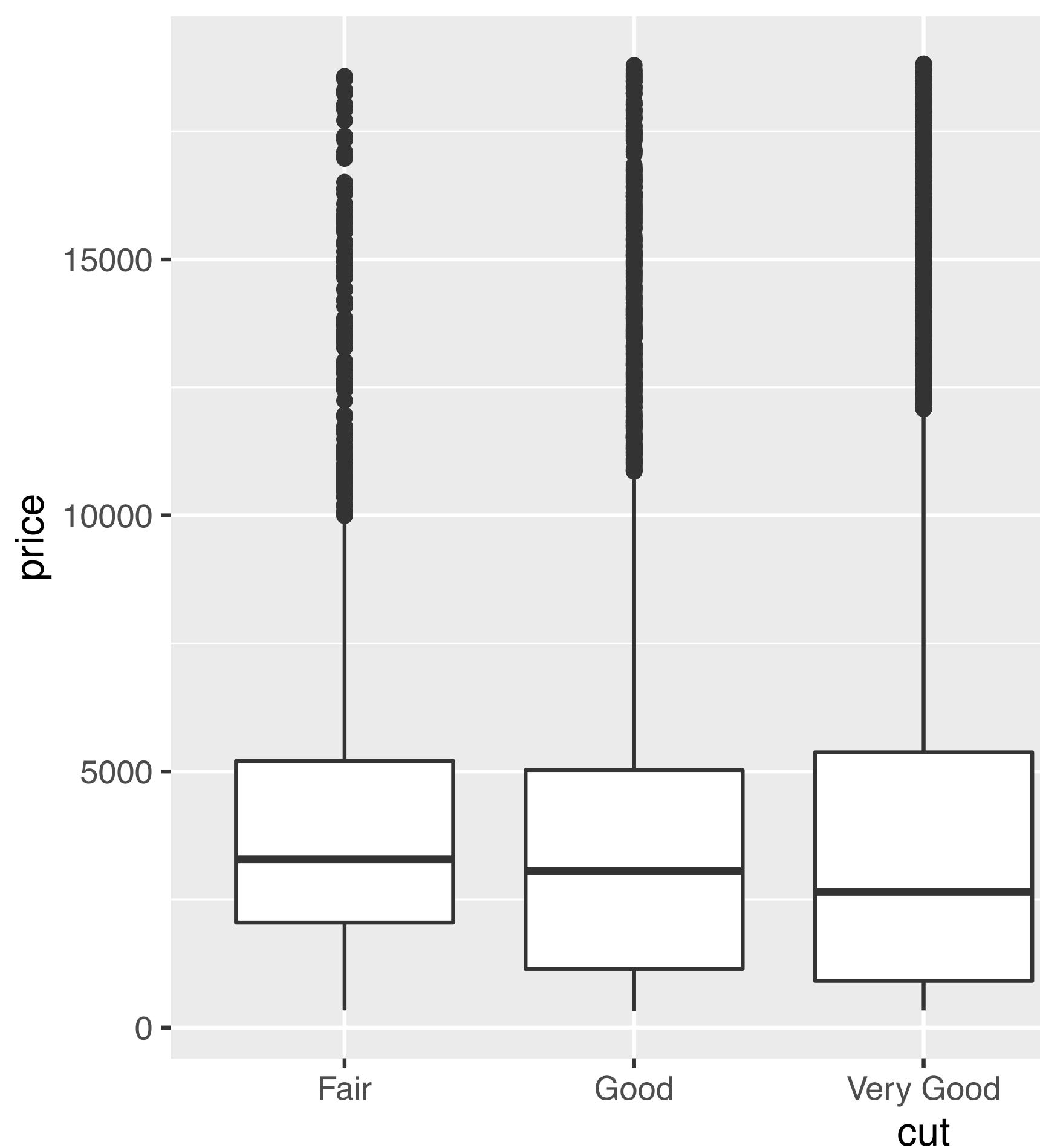
What is the relationship between



## Categorical and continuous

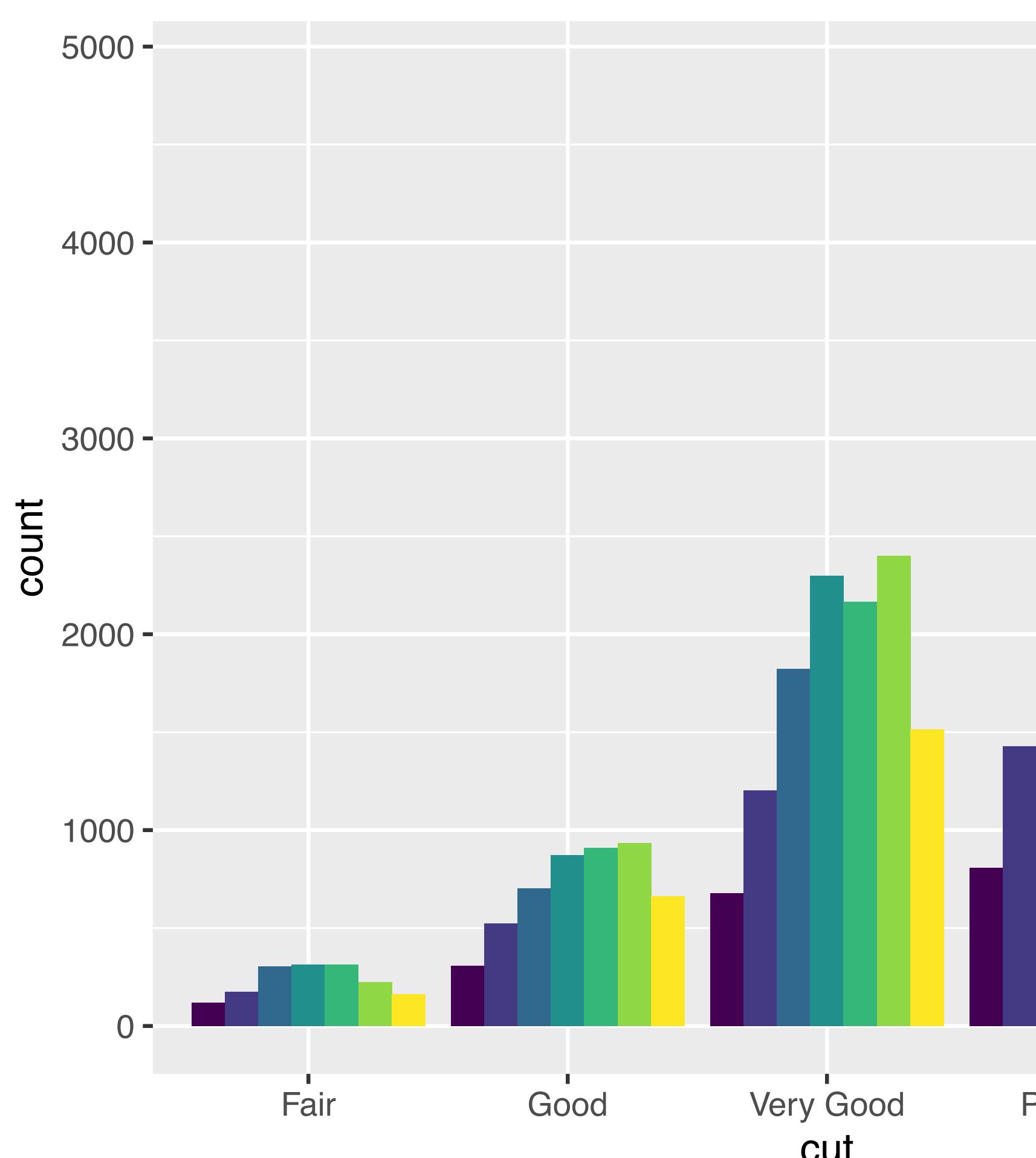
Higher quality cuts tend to be

Weird.



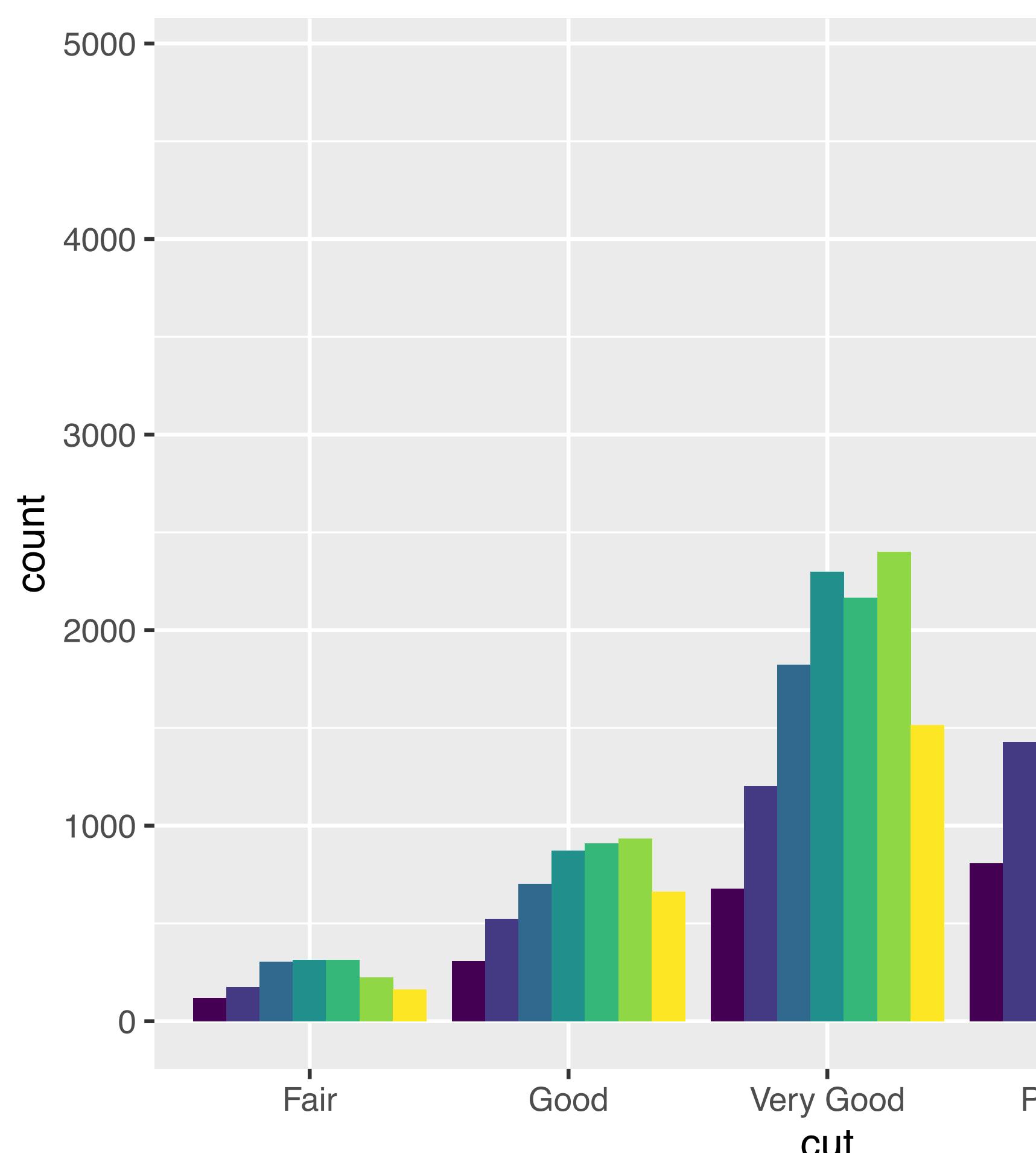
## Two categorical: count

What is the relationship between



## Two categorical: count

Higher quality cuts tend to have “G” color

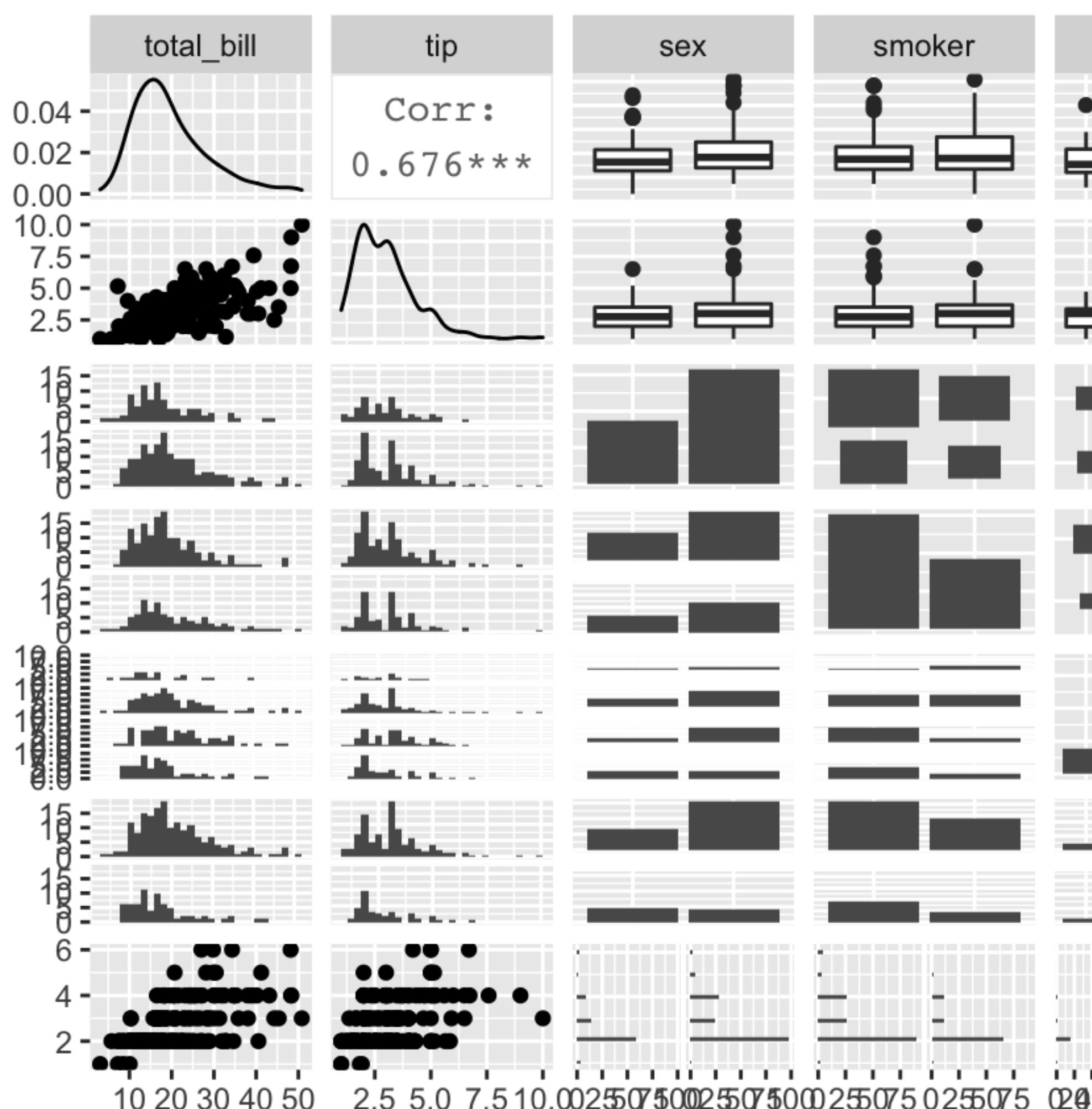


## Approach to

- State research question
- Identify target variables  
answer questions
- Visualize target variable  
answers the questions
- Tell “the story” of the d

“INTEREST  
VISUALIZAT

# Too many plots



GGally: <https://ggobi.github.io/ggally/>

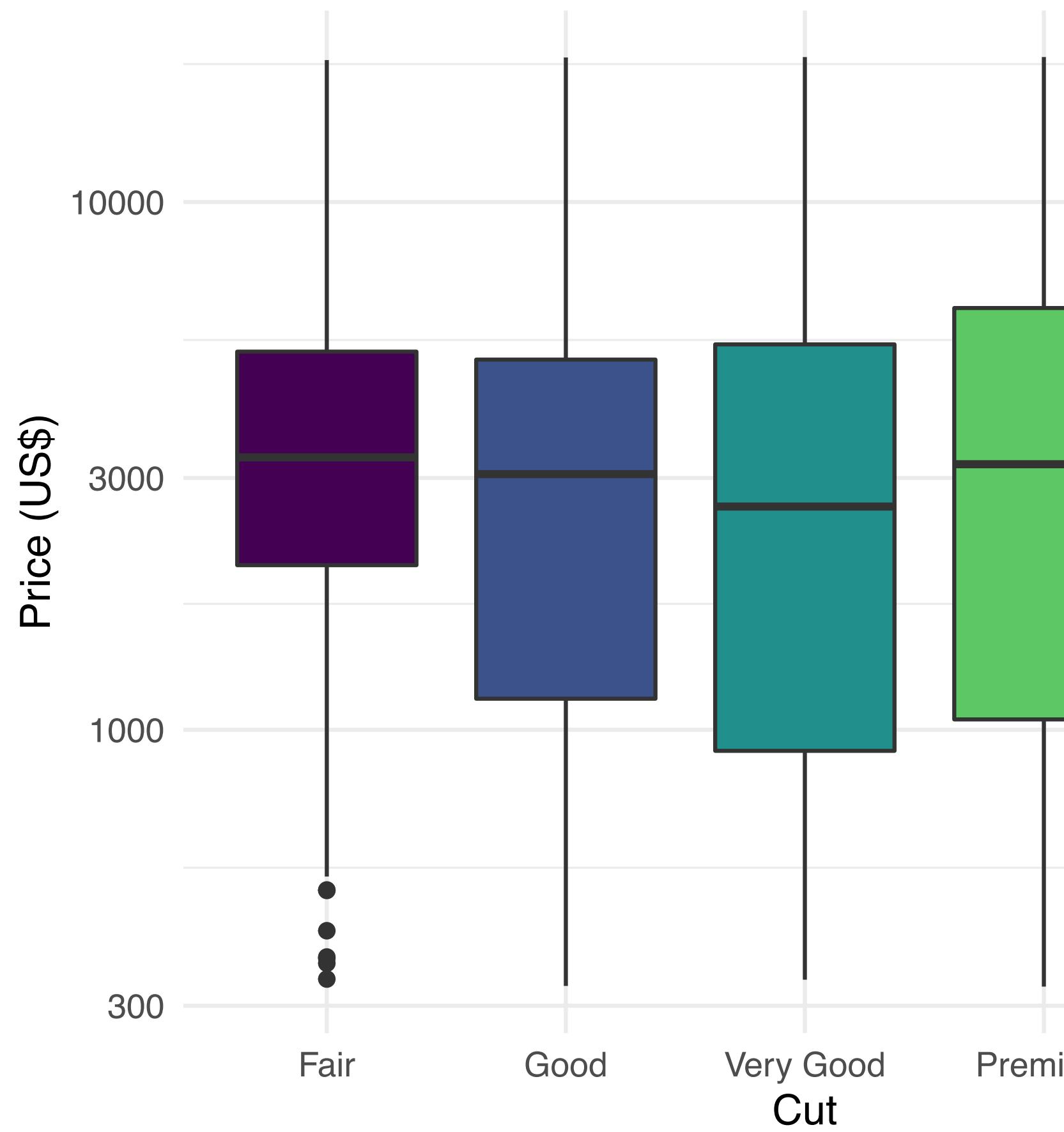
What makes a visualization?

## An “interesting” vis

- establishes a relationship
- explains something
- deviates from expectation

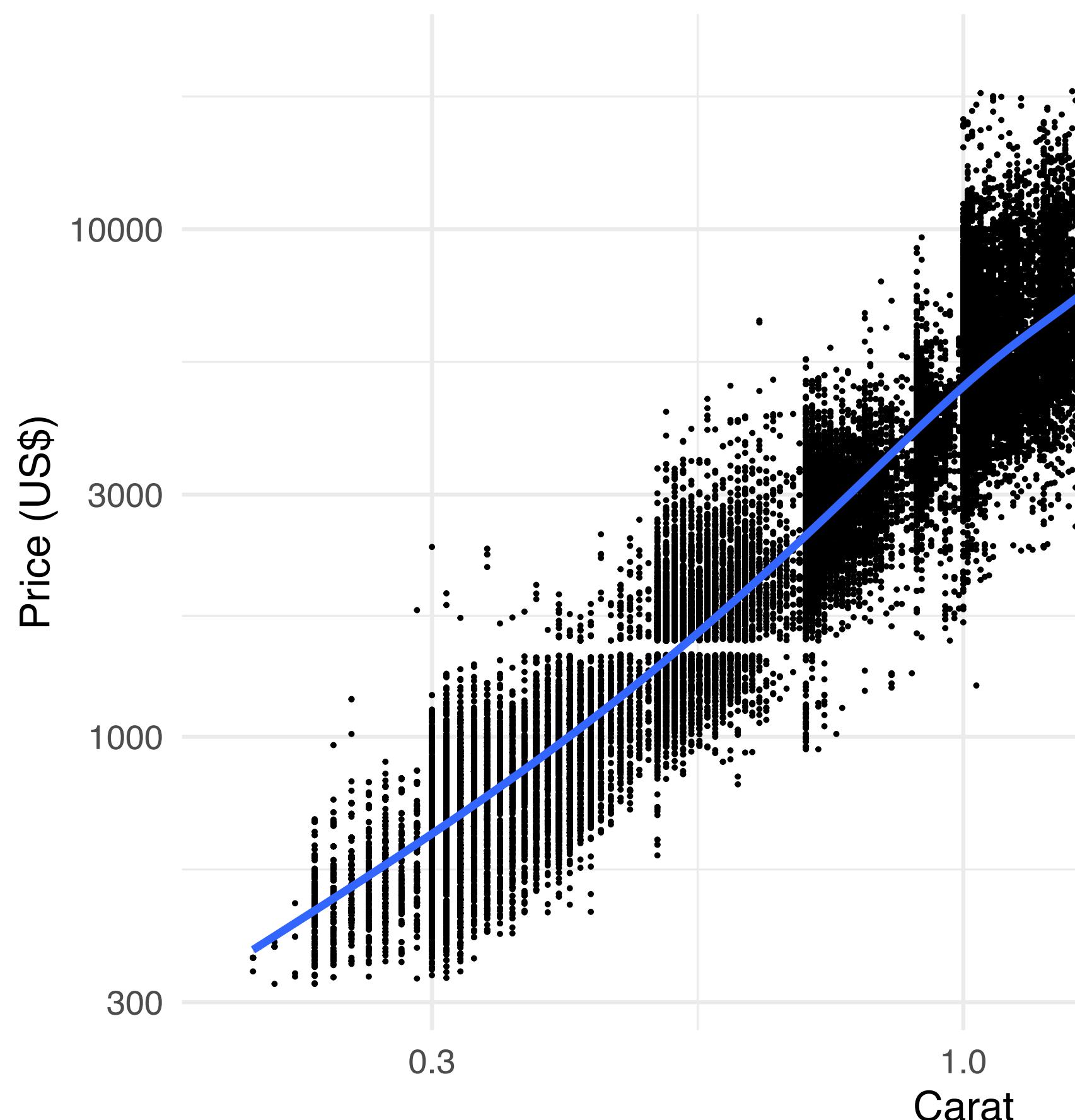
## Deviates from expectation

Higher quality cuts of diamonds tend to be more expensive



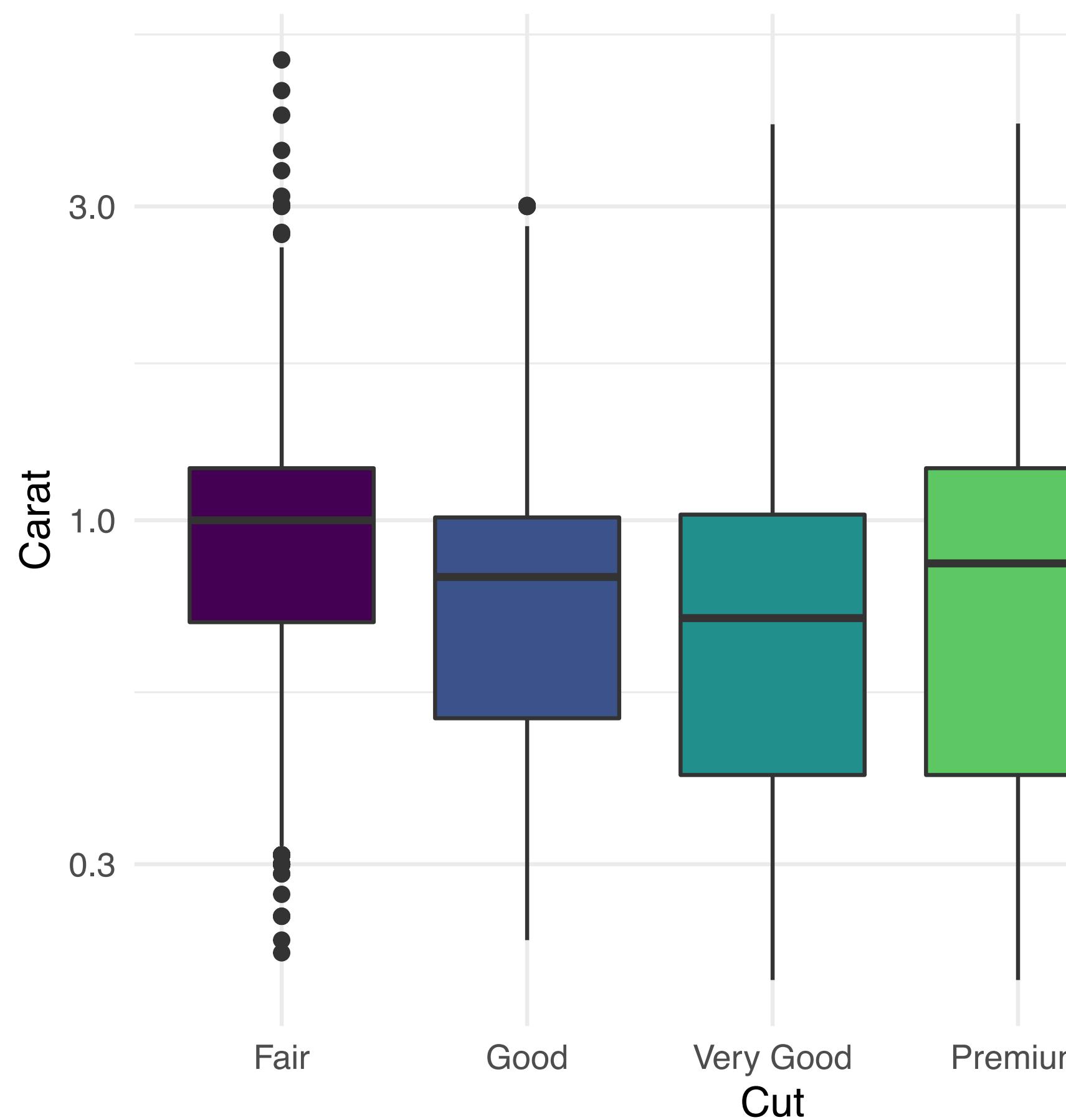
Establishes a relationship

Larger diamonds tend to be more expensive



Explains something

Higher quality cuts of diamonds tend to be larger



Explains something

For same-sized diamonds, higher quality

