# Visualization (Exploring Co-variation)

Peter Ganong and Maggie Shi

January 27, 2026

# Introduction

# Skills hopefully acquired at the end of lecture

Take a two variables in a dataset. Visualize to learn more about how they co-vary.

Key cases of interest:

- Categorical variable and a continuous variable

- Two categorical variables

- Two continuous variables

# Categorical variable and continuous variable

# Categorical vs. continuous: roadmap

- `penguins` dataset

- Boxplots

- Densities

- Small multiples

# **penguins** dataset

```
1 url = ("https://raw.githubusercontent.com/mcnakhaee/palmerpenguins/master/p
2 penguins = pd.read_csv(url)
3 penguins.head()
```

|   | species | island | bill_length_mm | bill_depth_mm | flipper_length_mm | boc |
|---|---------|--------|----------------|---------------|-------------------|-----|
| 0 | Adelie | Torgersen | 39.1 | 18.7 | 181.0 | 375 |
| 1 | Adelie | Torgersen | 39.5 | 17.4 | 186.0 | 380 |
| 2 | Adelie | Torgersen | 40.3 | 18.0 | 195.0 | 325 |
| 3 | Adelie | Torgersen | NaN | NaN | NaN | NaN |
| 4 | Adelie | Torgersen | 36.7 | 19.3 | 193.0 | 345 |

# **penguins** dataset

species appears to be a categorical variable

```
1  penguins['species'].value_counts()
```

```
species
Adelie       152
Gentoo       124
Chinstrap     68
Name: count, dtype: int64
```
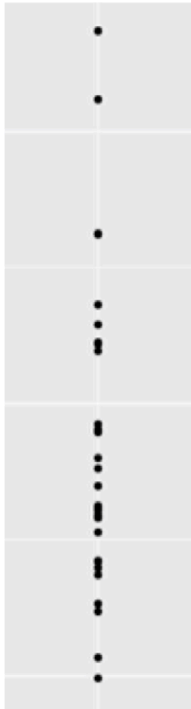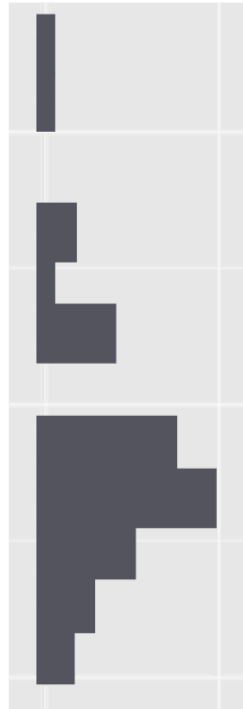
Discussion question: is it a Nominal or Ordinal variable?

# Categorical & continuous: box plot

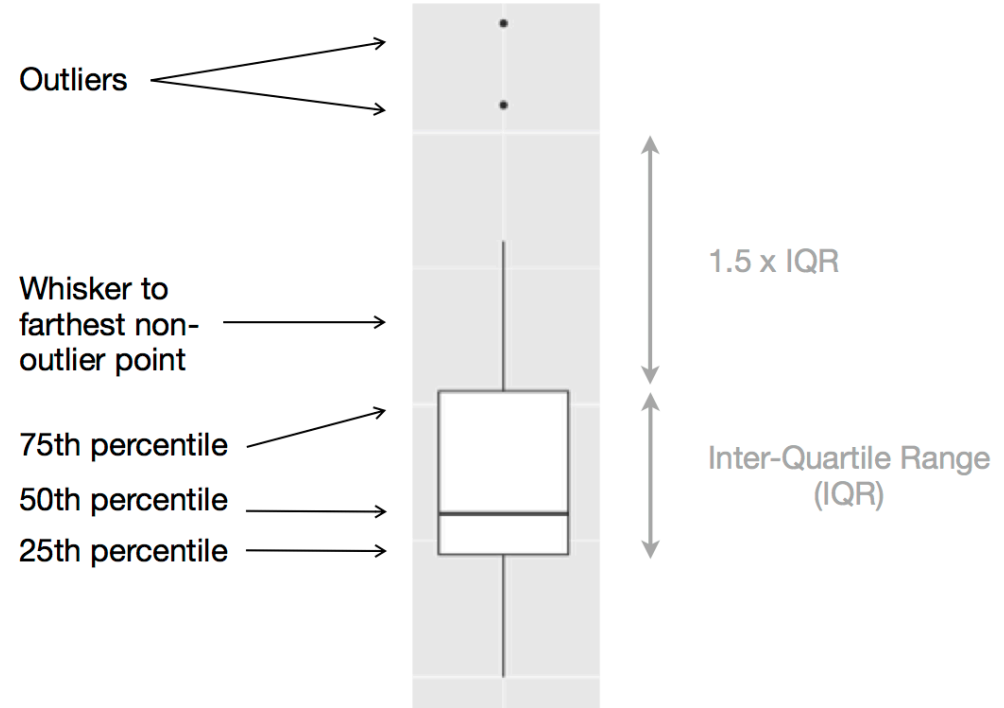The actual values in a distribution
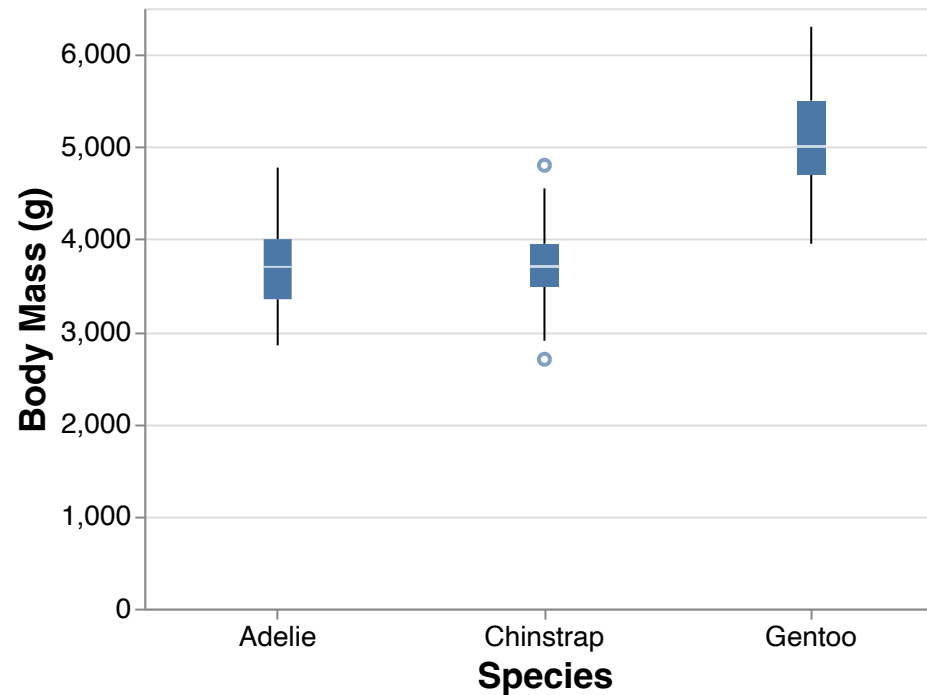
How a histogram would display the values (rotated)

How a boxplot would display the values

Outliers

Whisker to farthest non-outlier point

75th percentile

50th percentile

25th percentile

1.5 x IQR

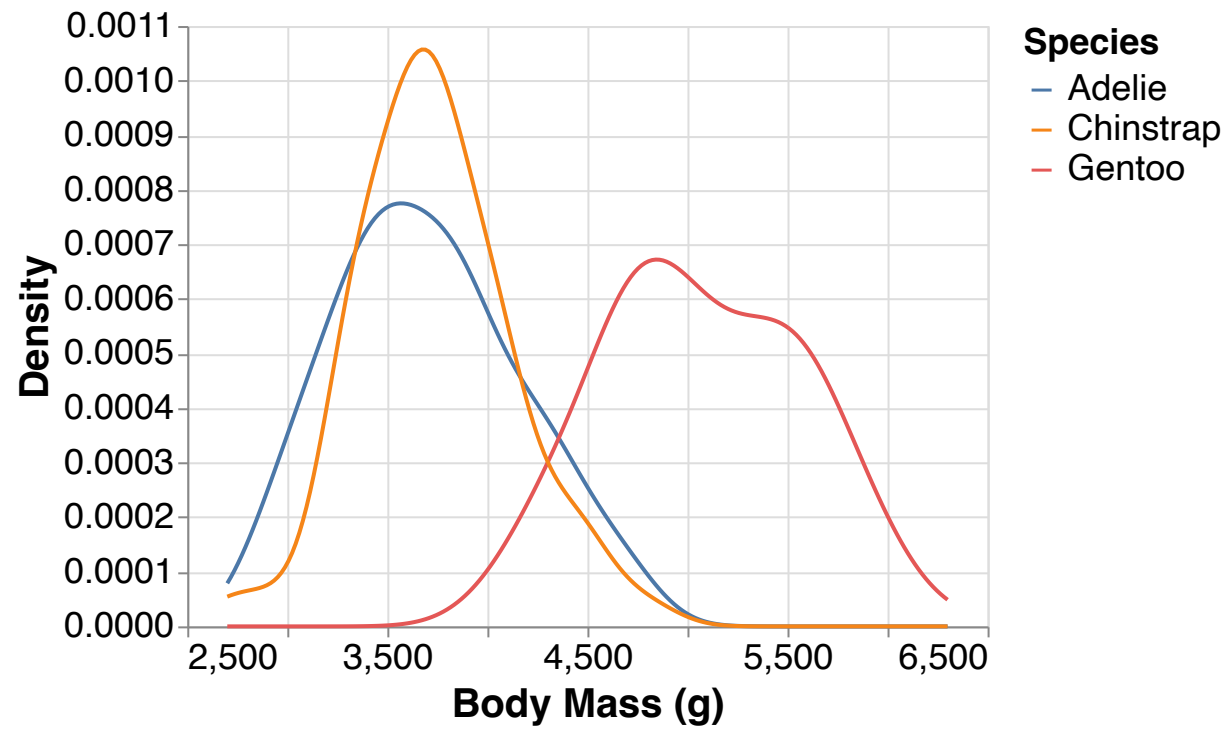Inter-Quartile Range (IQR)

# mark_boxplot()

```
1  alt.Chart(penguins).mark_boxplot().encode(
2      alt.X('species:N', title="Species"),
3      alt.Y('body_mass_g:Q', title="Body Mass (g)"),
4  )
```



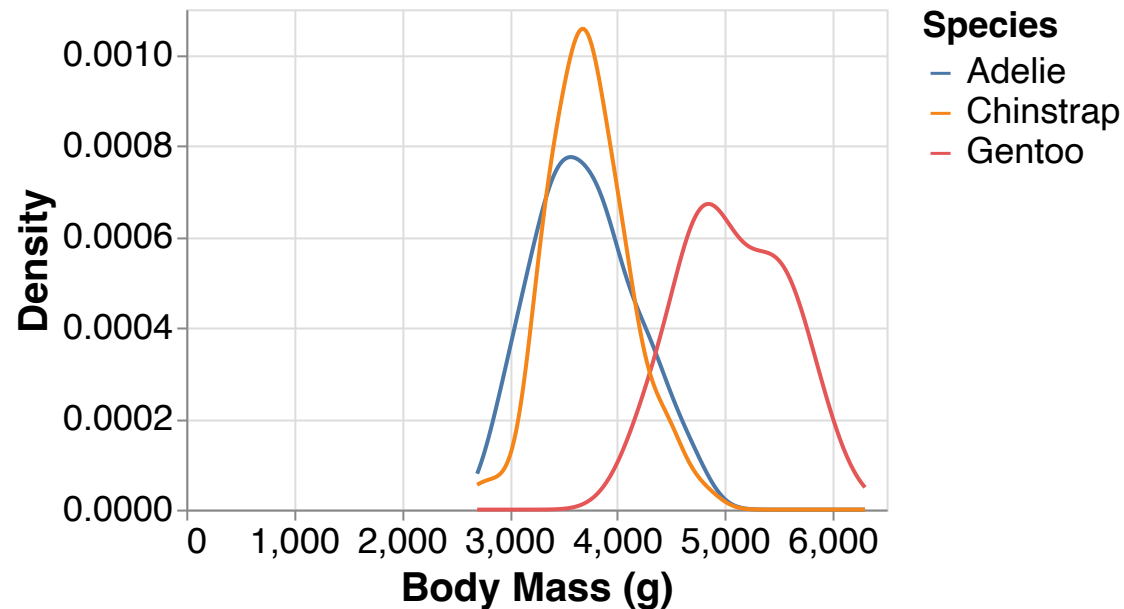Discussion question: what is the headline message from this graph?
Submessages?

# transform_density()

```python
alt.Chart(penguins).transform_density(
    'body_mass_g',
        groupby=['species'],
        as_=['body_mass_g2', 'density']
    ).mark_line().encode(
        alt.X('body_mass_g2:Q', title = "Body Mass (g)"),
        alt.Y('density:Q', title = "Density"),
        alt.Color('species:N', title = "Species")
    )
```
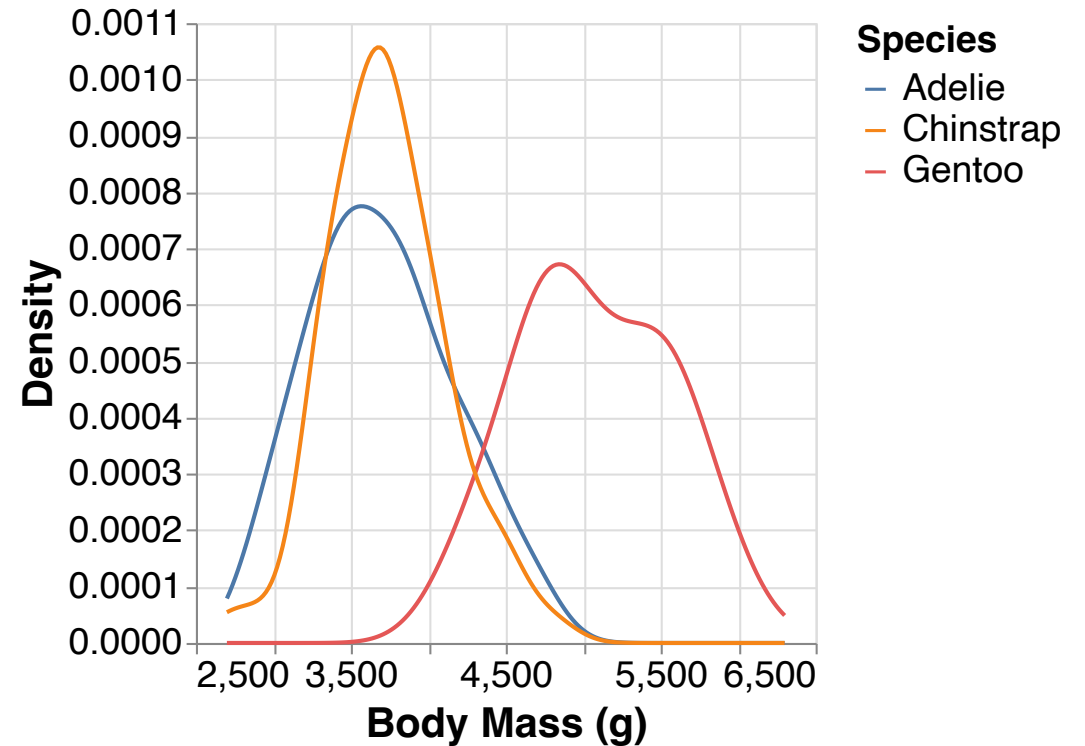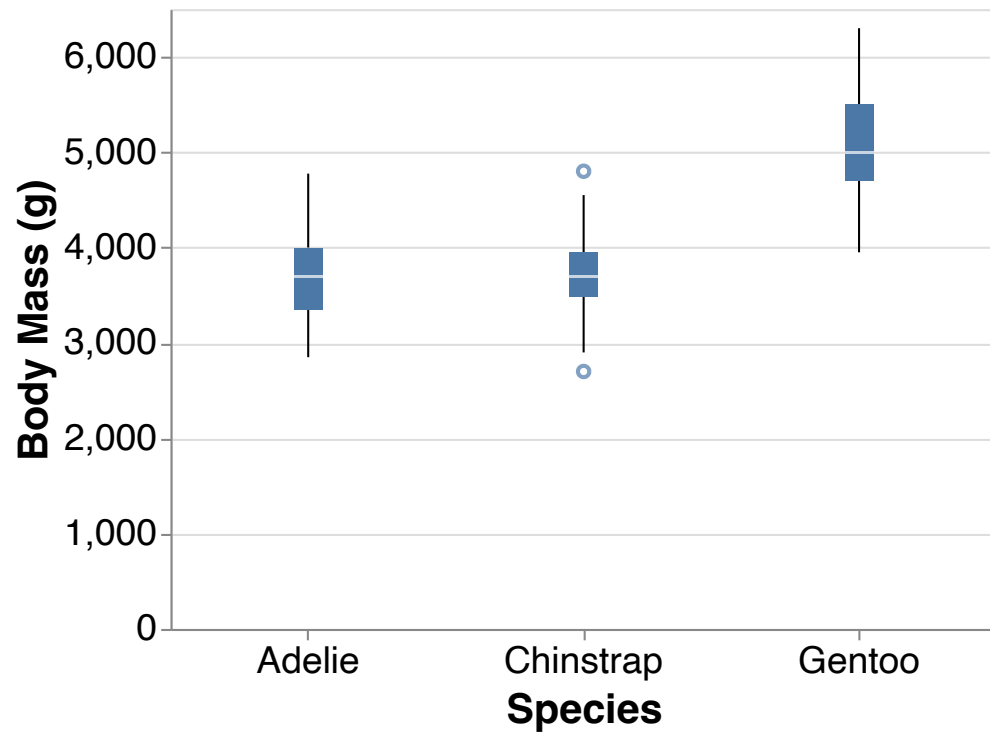
# transform_density(), scale to 0

```
1  alt.Chart(penguins).transform_density(
2      'body_mass_g',
3          groupby=['species'],
4          as_=['body_mass_g', 'density']
5      ).mark_line().encode(
6          alt.X('body_mass_g:Q', scale=alt.Scale(zero=True), title = "Body Mass (g)"),
7          alt.Y('density:Q', title = "Density"),
8          alt.Color('species:N', title = "Species")
9      )
```



Discussion question: what if we required the x-axis range to include zero?
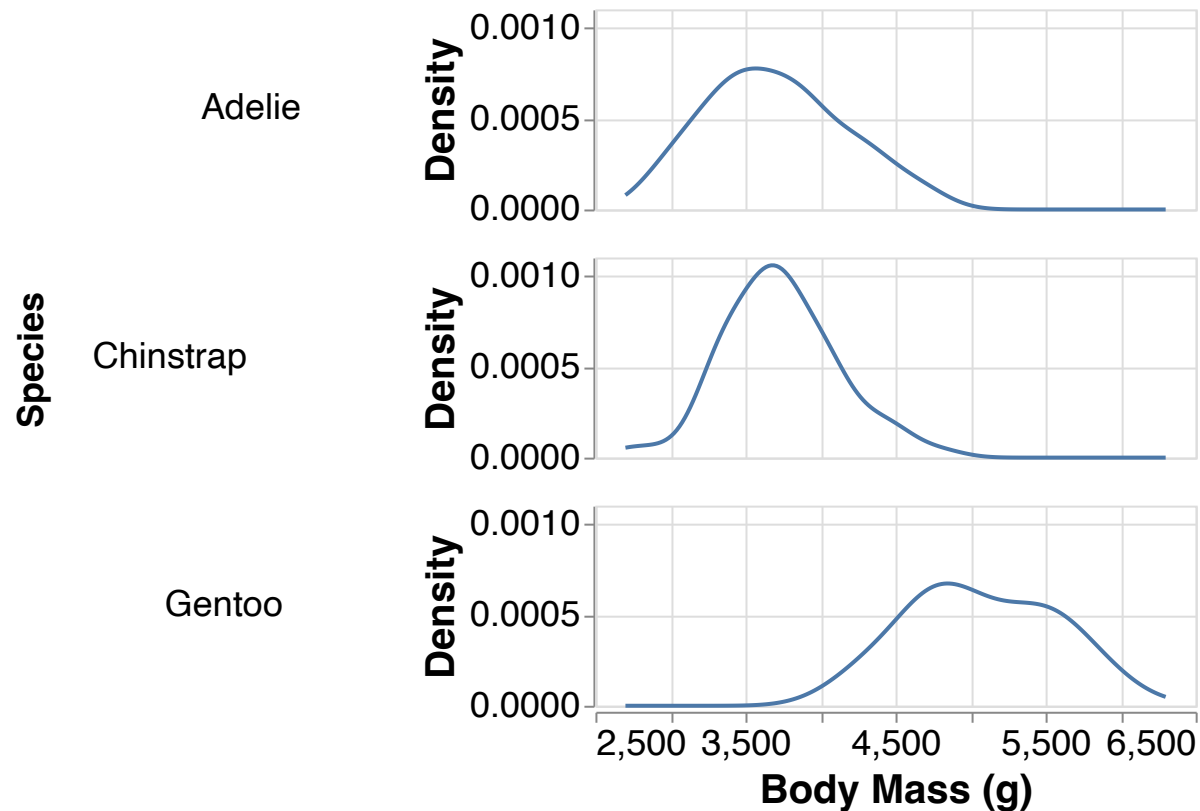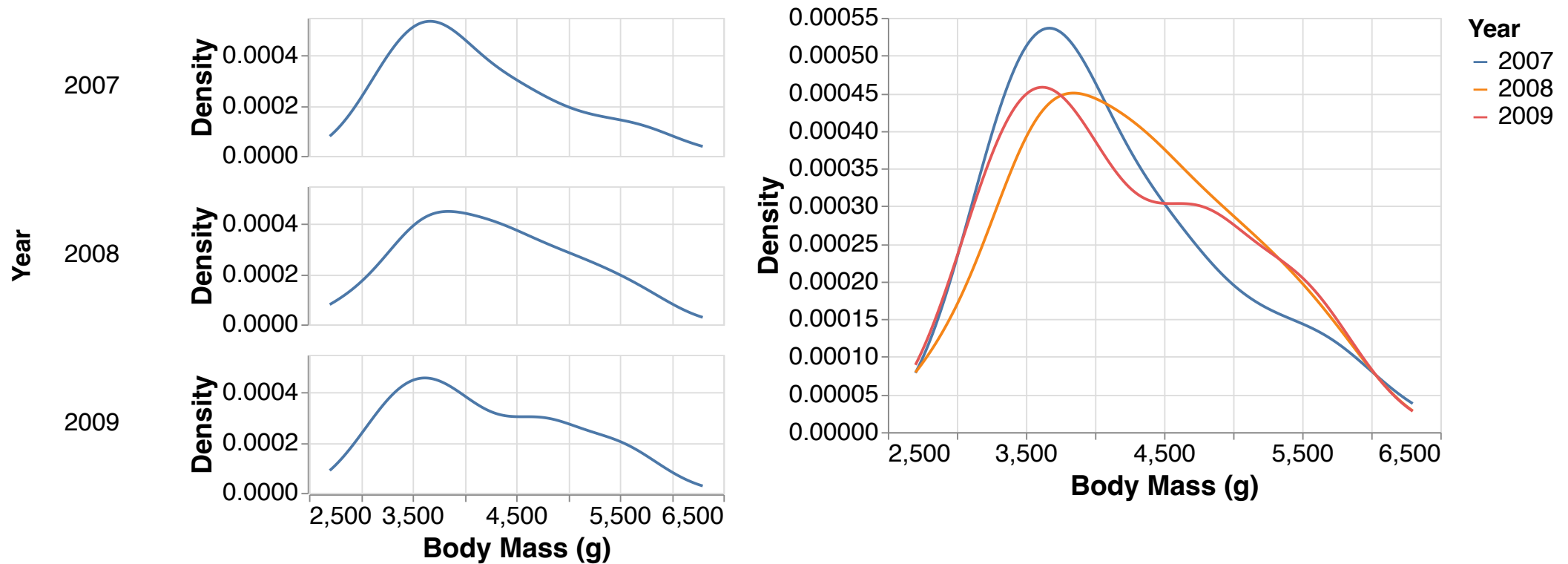Would that improve or reduce clarity? Why?

# Boxplot or density plots?



Discussion question: what messages come through more with the box plot? Through the density plot?

# `alt.Row`: small multiples

```python
alt.Chart(penguins).transform_density(
    'body_mass_g',
    groupby=['species'],
    as_=['body_mass_g', 'density']
).mark_line().encode(
    alt.X('body_mass_g:Q', title = "Body Mass (g)"),
    alt.Y('density:Q', title = "Density"),
    alt.Row('species:N', header=alt.Header(labelAngle=0), title = "Species")
)
```
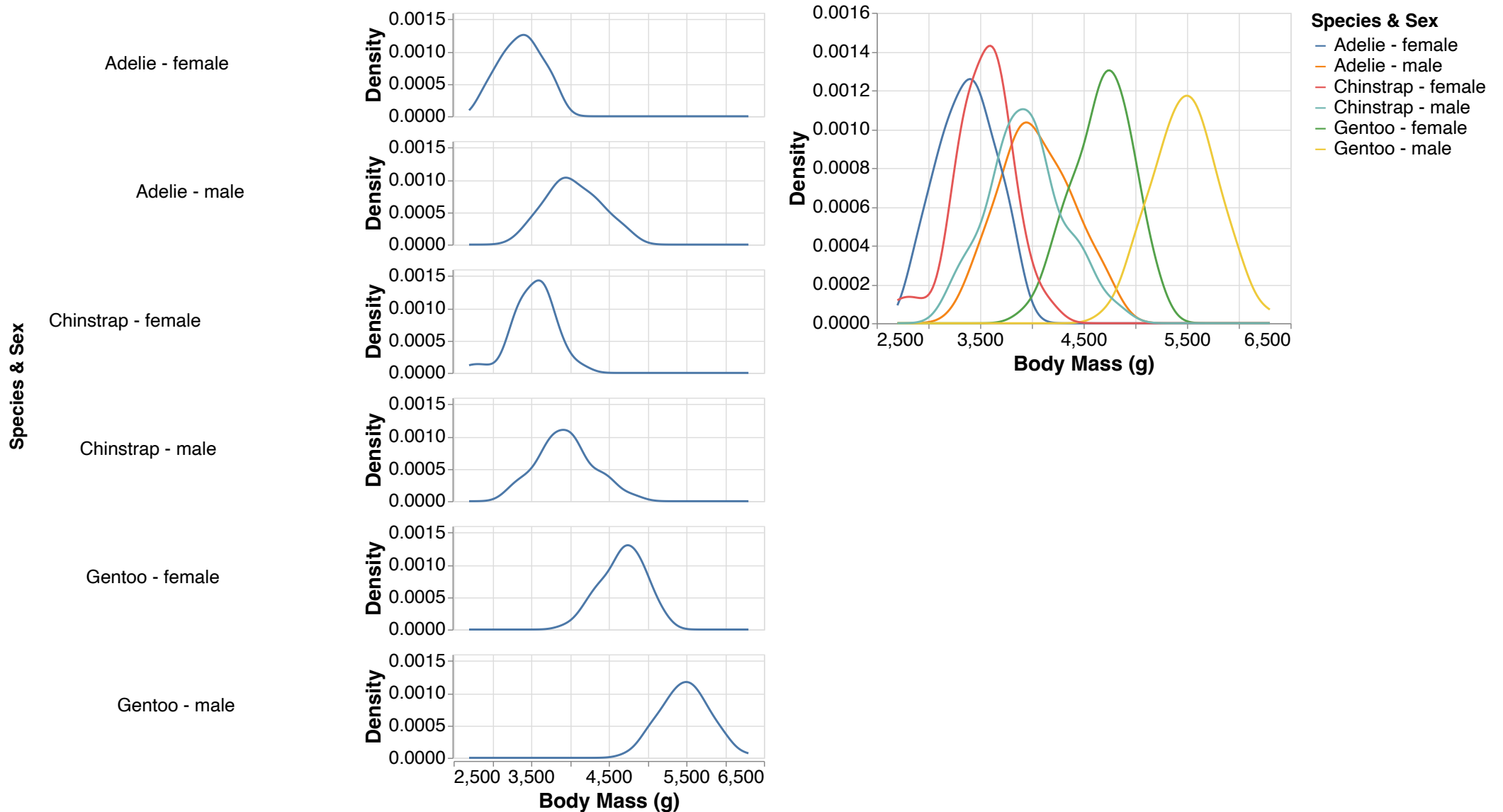
# By **year**: colors or small multiples?



Discussion question: these two graphs show identical information. Which do you prefer, and why?

# Colors or small multiples?

# Two Categorical Variables

# Two categorical variables: roadmap

- Two ways to encode frequency as a third dimension:
  `diamonds`

  - `size`

  - `color`

- A word of caution against 3D graphs

# How is cut related to color? **Size**

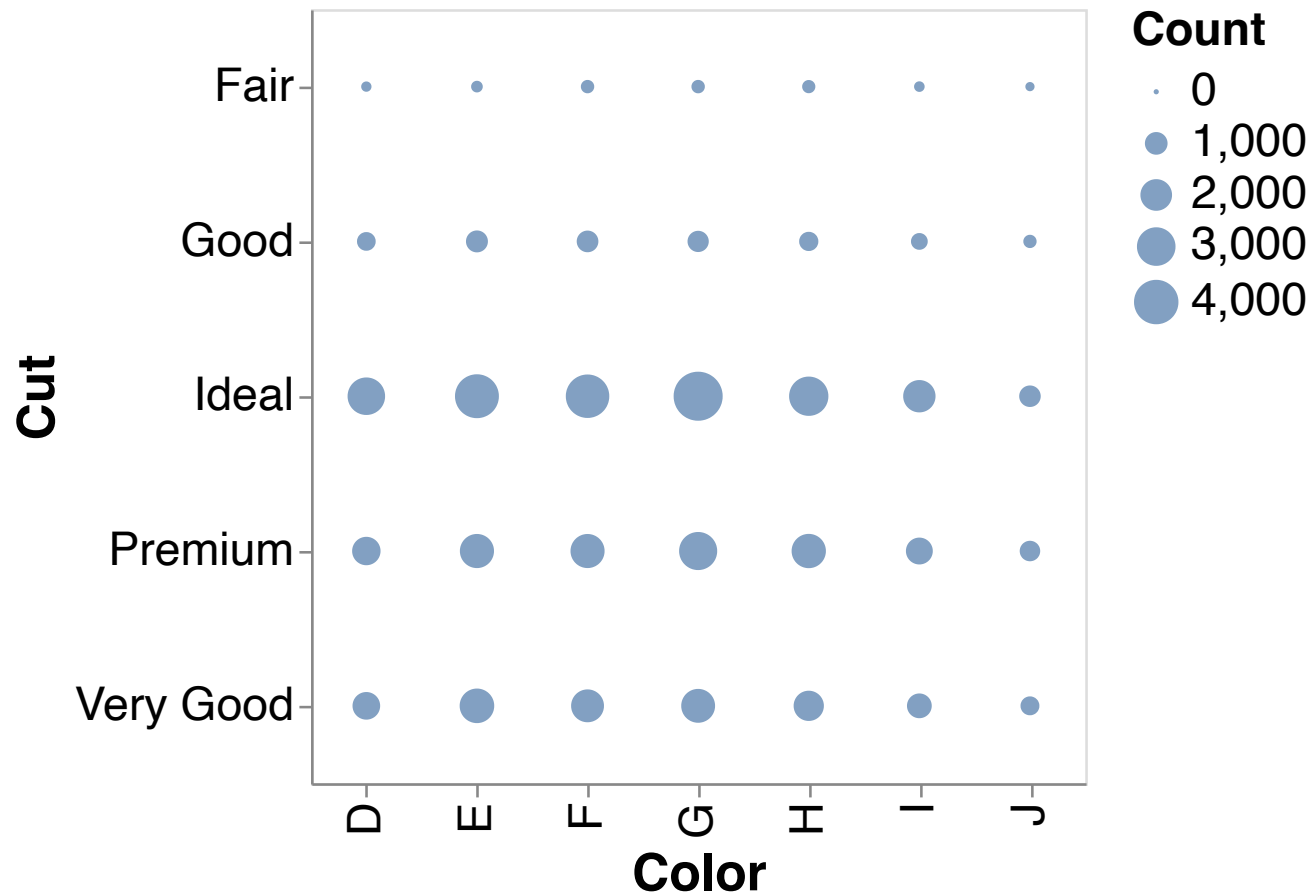In `diamonds` dataset, `color` and `cut` are both categorical

```
1  diamonds_grouped = diamonds.groupby(['color','cut']).size().reset_index().rename(columns={0:'N'})
2  diamonds_grouped
```

|    | color | cut       | N    |
|----|-------|-----------|------|
| 0  | D     | Fair      | 163  |
| 1  | D     | Good      | 662  |
| 2  | D     | Very Good | 1513 |
| 3  | D     | Premium   | 1603 |
| 4  | D     | Ideal     | 2834 |
| 5  | E     | Fair      | 224  |
| 6  | E     | Good      | 933  |
| 7  | E     | Very Good | 2400 |
| 8  | E     | Premium   | 2337 |
| 9  | E     | Ideal     | 3903 |
| 10 | F     | Fair      | 312  |
| 11 | F     | Good      | 909  |
| 12 | F     | Very Good | 2164 |
| 13 | F     | Premium   | 2331 |

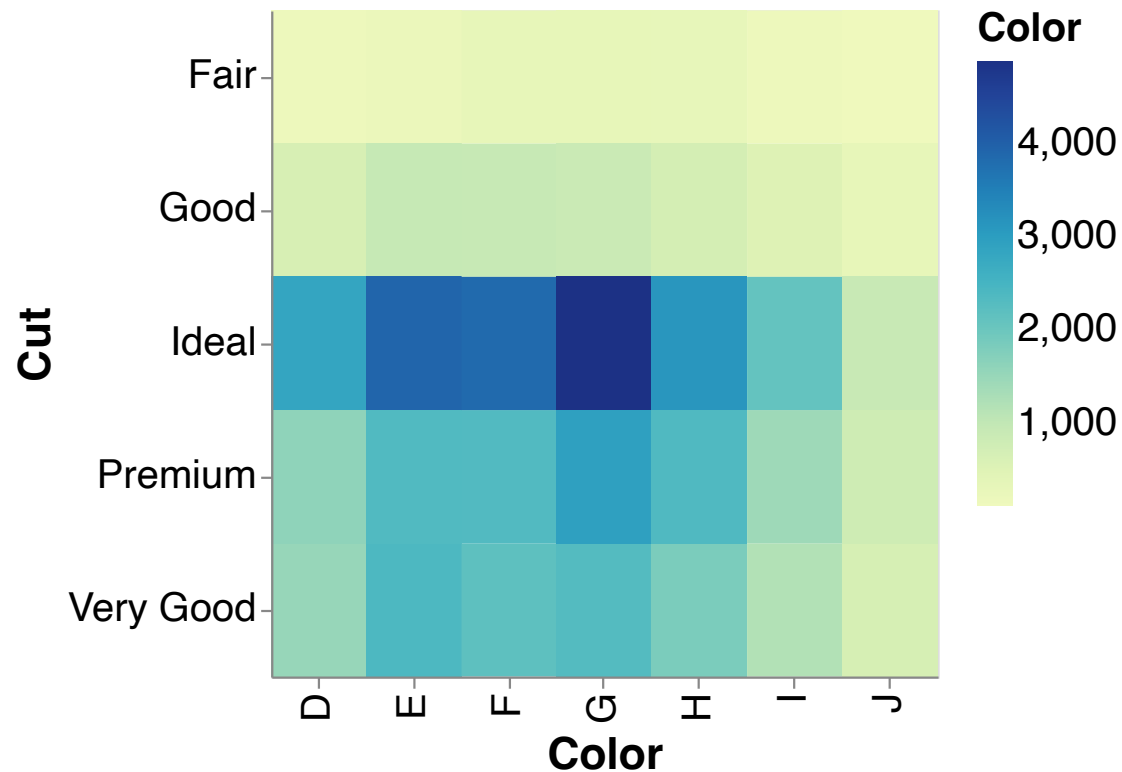| | color | cut | N |
|---|---|---|---|
| 14 | F | Ideal | 3826 |
| 15 | G | Fair | 314 |
| 16 | G | Good | 871 |
| 17 | G | Very Good | 2299 |
| 18 | G | Premium | 2924 |
| 19 | G | Ideal | 4884 |
| 20 | H | Fair | 303 |
| 21 | H | Good | 702 |
| 22 | H | Very Good | 1824 |
| 23 | H | Premium | 2360 |
| 24 | H | Ideal | 3115 |
| 25 | I | Fair | 175 |
| 26 | I | Good | 522 |
| 27 | I | Very Good | 1204 |
| 28 | I | Premium | 1428 |
| 29 | I | Ideal | 2093 |
| 30 | J | Fair | 119 |
| 31 | J | Good | 307 |
| 32 | J | Very Good | 678 |
| 33 | J | Premium | 808 |
| 34 | J | Ideal | 896 |

# How is cut related to color? Color

```
1  alt.Chart(diamonds_grouped).mark_circle().encode(
2      alt.X('color:N', title = "Color"),
3      alt.Y('cut:N', title = "Cut"),
4      alt.Size('N:Q', title = "Count"))
```
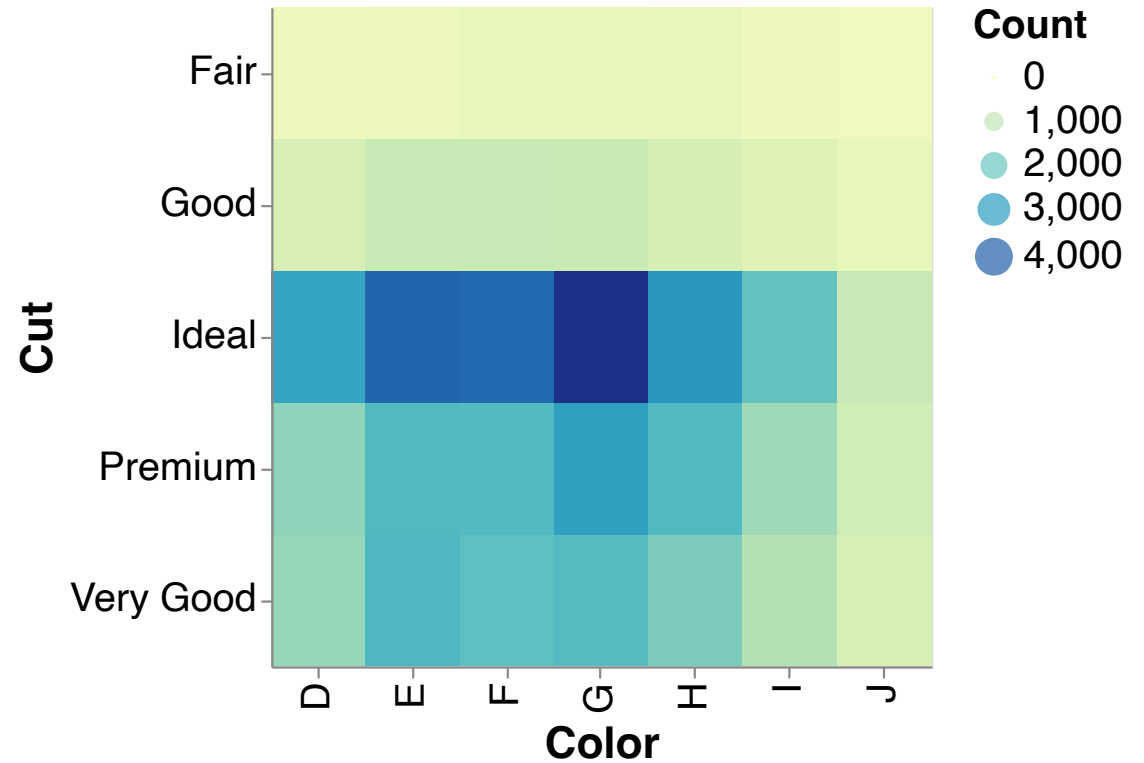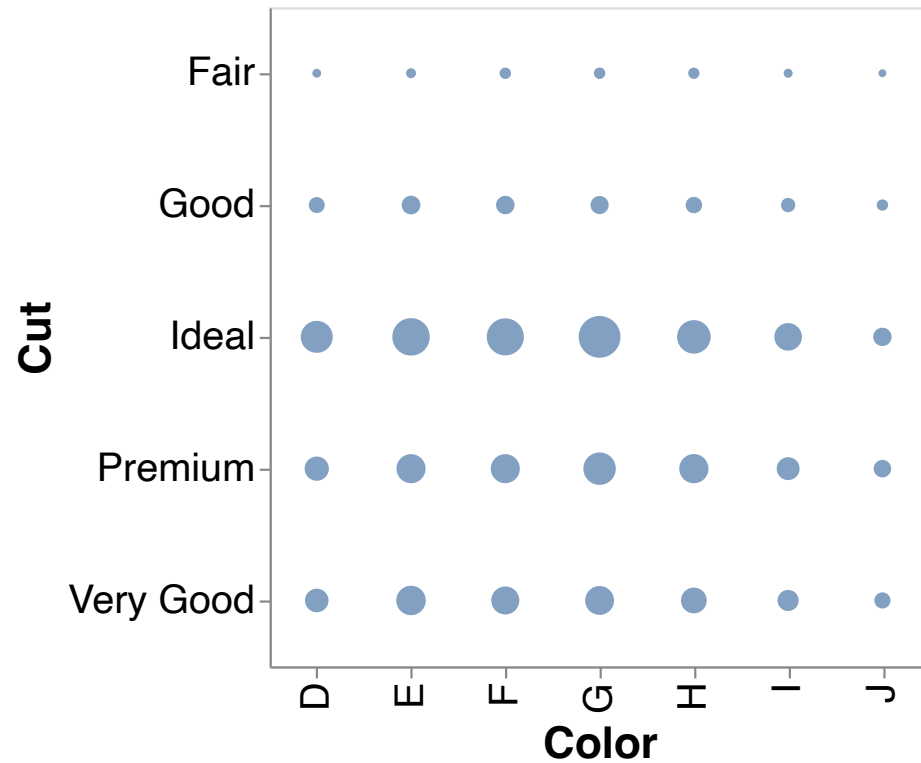
# How is cut related to color?

```
1  alt.Chart(diamonds_grouped).mark_rect().encode(
2      alt.X('color:N', title = "Color"),
3      alt.Y('cut:N', title = "Cut"),
4      alt.Color('N:Q', title = "Color"))
```



Discussion question: what diamond types are most common?

# How is cut related to color?



Discussion question: these two plots display the same information, but encoded differently. Which do you prefer?

# A word of caution: 3D graphs

You may have seen covariation between two variables depicted as a 3D plot before

# Two Categorical Variables: summary

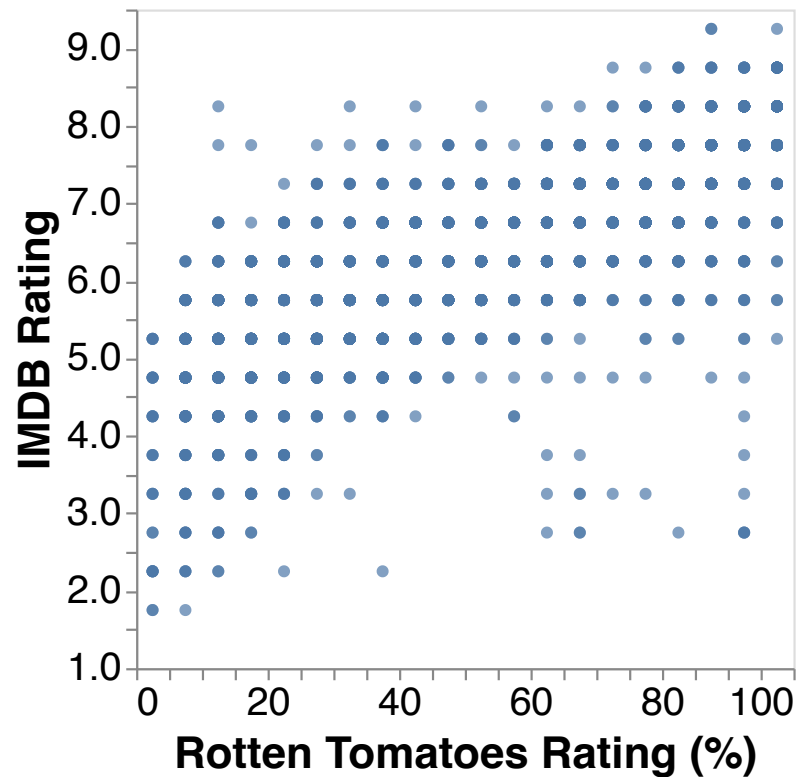- Encode frequency as `color` or `size`

- Avoid 3D representations!

# Two Continuous Variables

# Two continuous variables: roadmap

- `movies` ratings from Rotten Tomatoes and IMDB

- `diamonds`: `carat` vs `price`

# How are RT and IMDB ratings related?
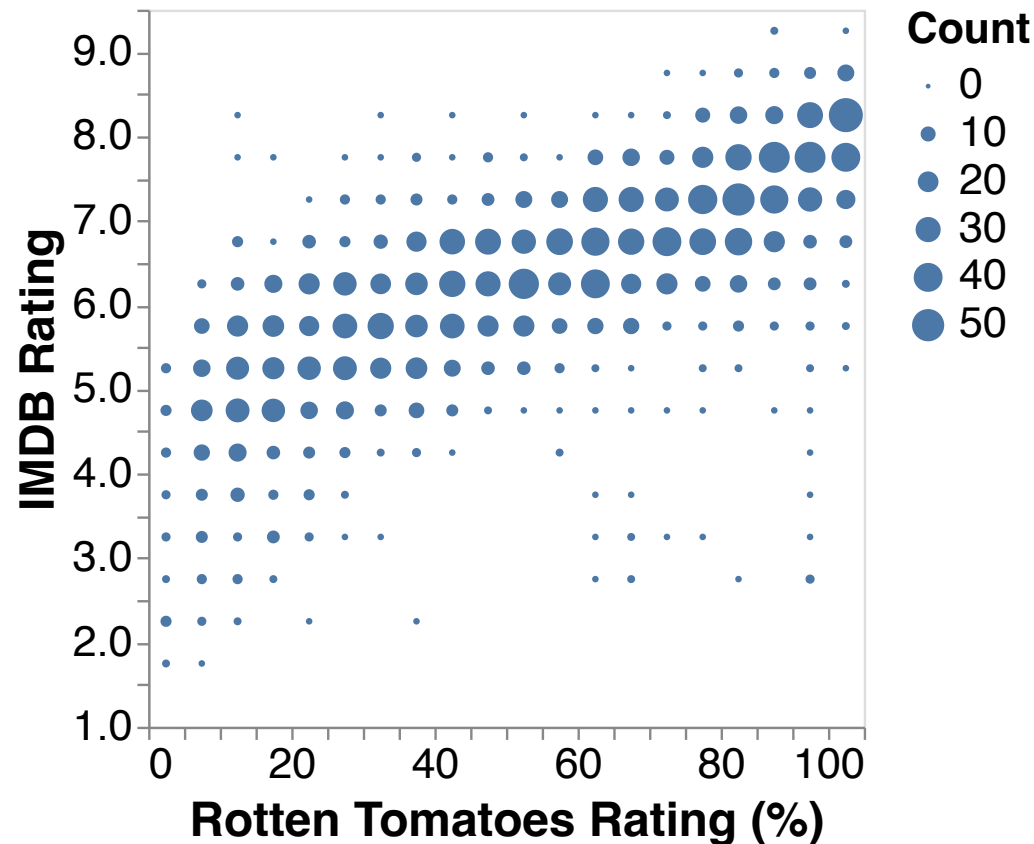
```
1  alt.Chart(movies).mark_circle().encode(
2      alt.X('Rotten_Tomatoes_Rating:Q', bin=alt.BinParams(maxbins=20), title
3      alt.Y('IMDB_Rating:Q', bin=alt.BinParams(maxbins=20), title = "IMDB Rat
4  )
```



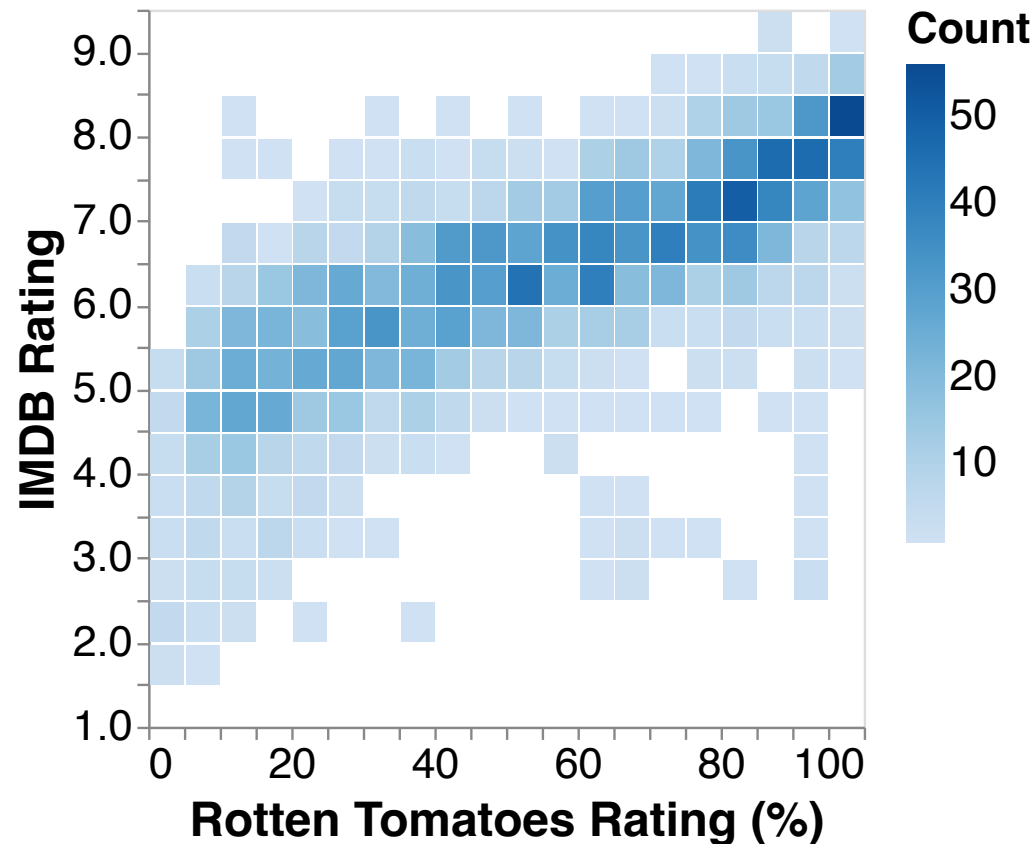Suffers from overplotting!

# use `alt.Size('count()')`

```
1  alt.Chart(movies_url).mark_circle().encode(
2      alt.X('Rotten_Tomatoes_Rating:Q', bin=alt.BinParams(maxbins=20)),
3      alt.Y('IMDB_Rating:Q', bin=alt.BinParams(maxbins=20)),
4      alt.Size('count()')
5  )
```
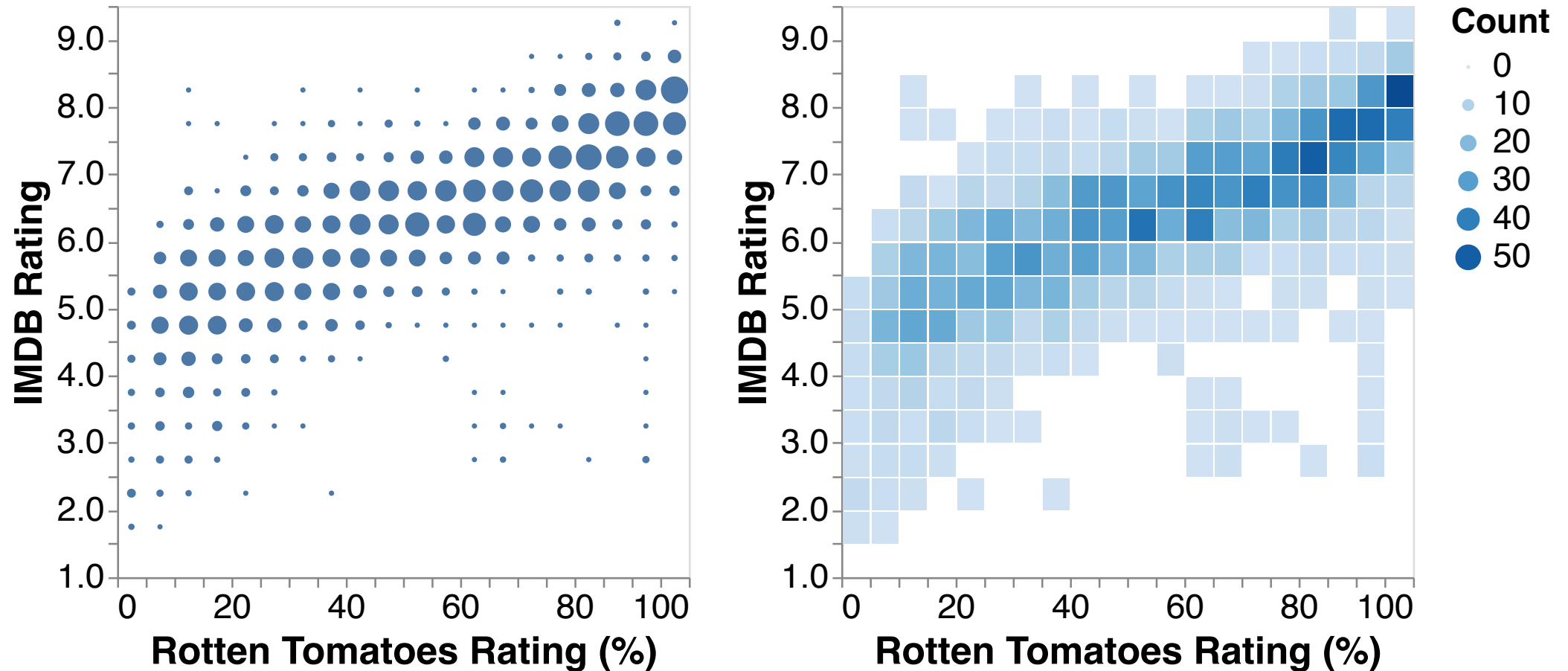
# use `alt.Color('count()')`

```
1  alt.Chart(movies_url).mark_bar().encode(
2      alt.X('Rotten_Tomatoes_Rating:Q', bin=alt.BinParams(maxbins=20), title
3      alt.Y('IMDB_Rating:Q', bin=alt.BinParams(maxbins=20), title = "IMDB Rat
4      alt.Color('count()', title = "Count")
5  )
```
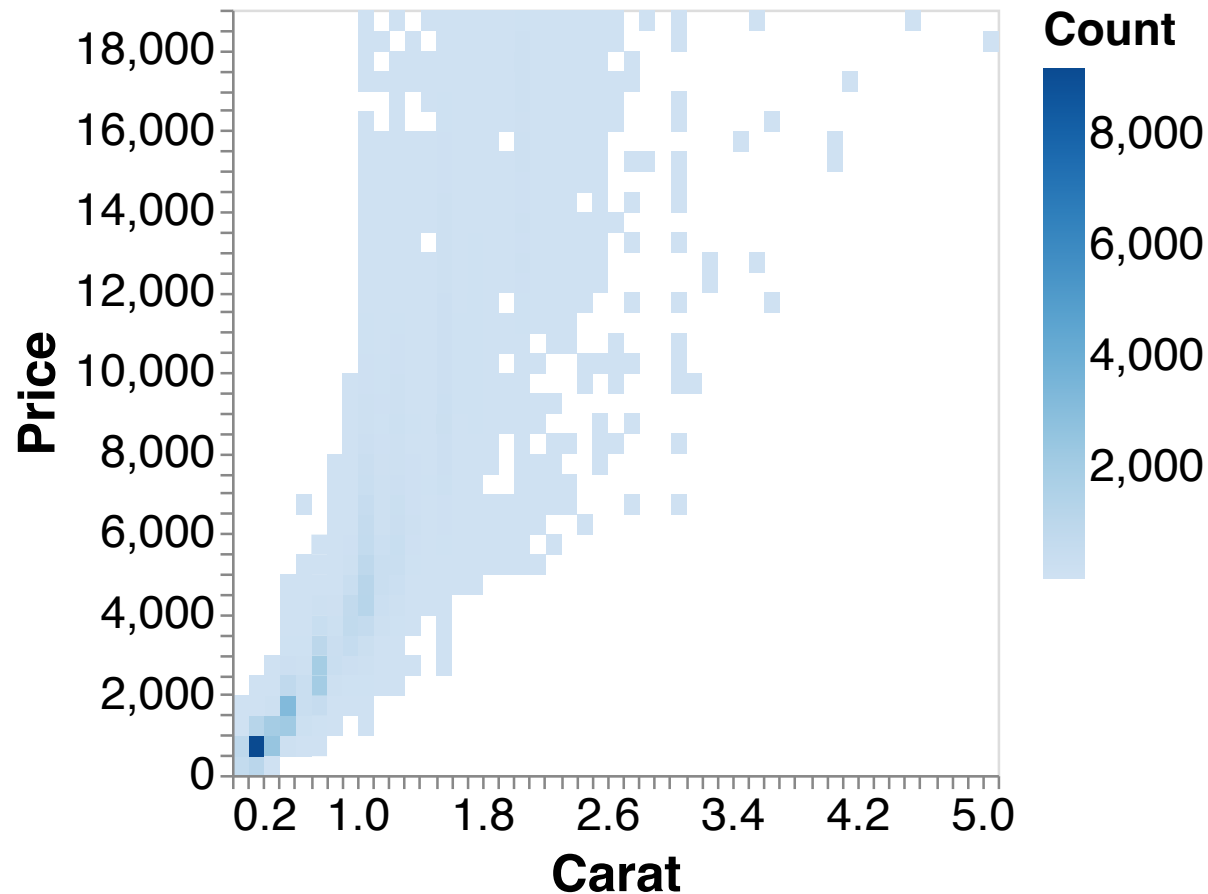
# Discussion question



Compare the *size* and *color*-based 2D histograms above. Which encoding do you prefer? Why?

# How is carat related to price? Raw data

```python
alt.Chart(diamonds).mark_point().encode(
    alt.X('carat:Q', title = "Carat"),
    alt.Y('price:Q', title = "Price")
)
```

# How is carat related to price? Color
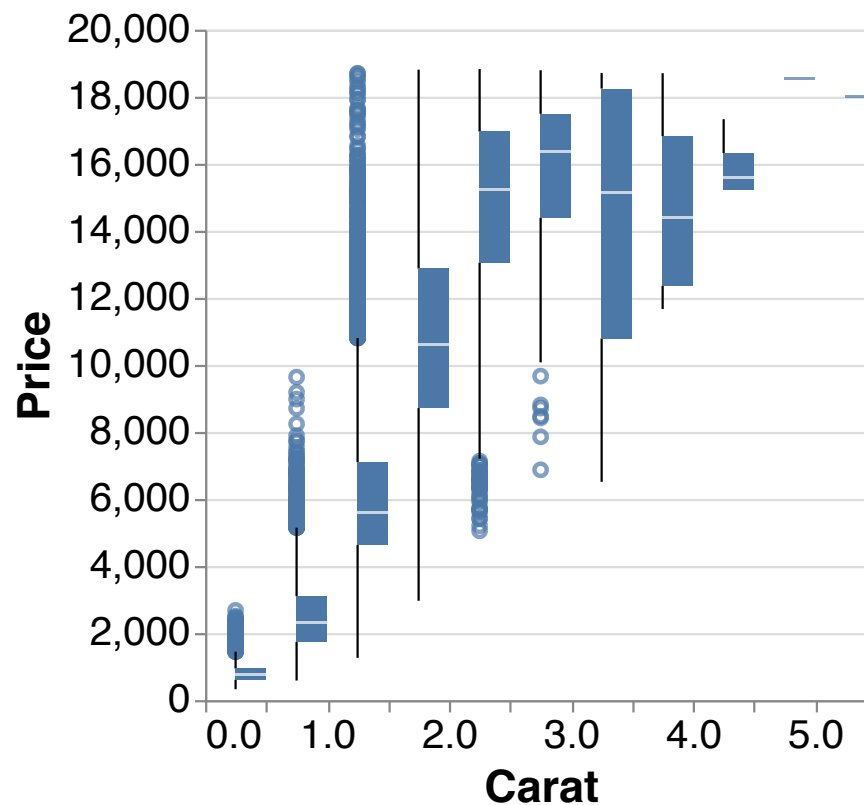
```
1  alt.Chart(diamonds).mark_rect().encode(
2      alt.X('carat:Q', bin=alt.Bin(maxbins=70), title = "Carat"),
3      alt.Y('price:Q', bin=alt.Bin(maxbins=70), title = "Price"),
4      alt.Color('count()', scale=alt.Scale(scheme='blues'), title = "Count"))
```

# How is carat related to price?

## mark_boxplot()

```python
1  alt.Chart(diamonds).mark_boxplot().encode(
2      alt.X('carat:Q', bin=alt.Bin(maxbins=10), title = "Carat"),
3      alt.Y('price:Q', title = "Price"))
```

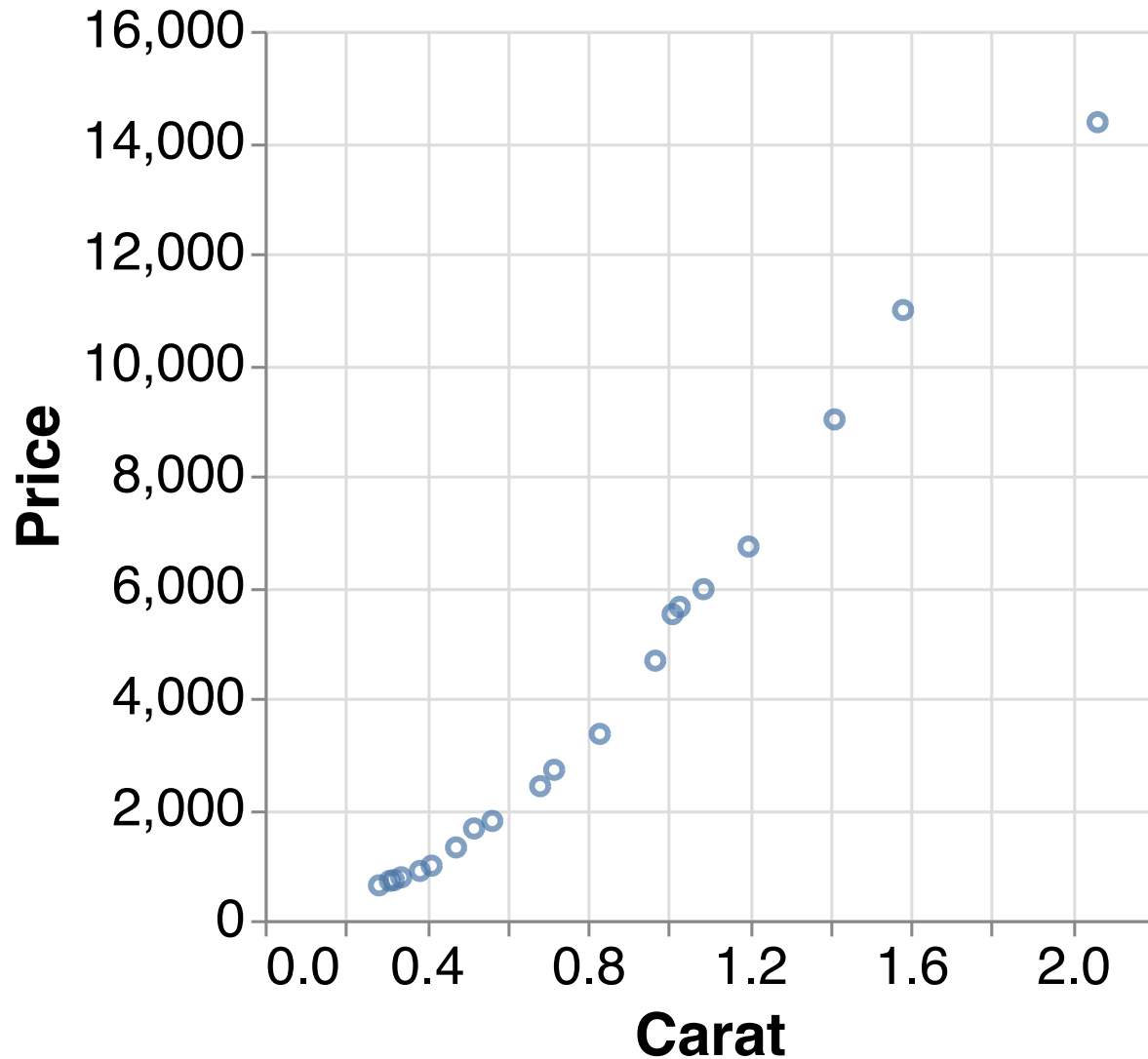# How is carat related to price? binscatter (code)

```python
1  df = diamonds
2  df['carat_bin'] = pd.qcut(df['carat'], q=20, labels=(np.arange(1, 21, 1)))
3
4  df = df.groupby('carat_bin').agg(
5      carat = ('carat', 'mean'),
6      price = ('price', 'mean')).reset_index()
7
8  alt.Chart(df).mark_point().encode(
9      alt.X('carat:Q', title = "Carat"),
10     alt.Y('price:Q', title = "Price")
11 )
```

# How is carat related to price? binscatter

- Can also create a binscatter: `binscatter` in stata and `binsreg` in R.

- Doesn't exist yet for `altair`, but easy to code up yourself

- What it does:

  1. Computes bins using quantiles of x

  2. Computes means of y within each bin

# How is carat related to price? (plot)

# Discussion question – "How is carat related to price?"

Review the `mark_rect()`, `mark_boxplot()`, and `binscatter` plots

- Headline?

- Sub-messages?

# Exploring covariation: summary

| Scenario | Functions |
| --- | --- |
| Categorical and continuous variable | `mark_boxplot()` |
| | `transform_density()` |
| | `alt.Row()` |
| Two categorical variables | `size` |
| | `color` |
| Two continuous variables | `alt.Size('count()')` |
| | `alt.Color('count()')` |
| | `mark_boxplot()` |
| | binscatter |