



# THINKING ABOUT DATA VISUALIZATION

DANIEL KERCHNER, SENIOR SOFTWARE DEVELOPER  
EMILY BLUMENTHAL, DATA SERVICES LIBRARIAN  
GW LIBRARIES AND ACADEMIC INNOVATION

# YOUR EXPERIENCES

---

Share your data visualization experiences and workshop goals



# WORKSHOP OUTLINE

.....

Our content today is divided into four parts. Each part will be described with examples.

## 01 | Data Visualization Process

Walking through the visualization process, from setting goals to visualizing the data.

## 02 | Basic Design Considerations

Best practices and accessibility considerations in data visualization.

## 03 | Solutions to Common Problems

How to resolve, or even better, avoid common problems in data visualization.

## 04 | Visualization Resources

Workshops, tools, and other resources for data visualization.

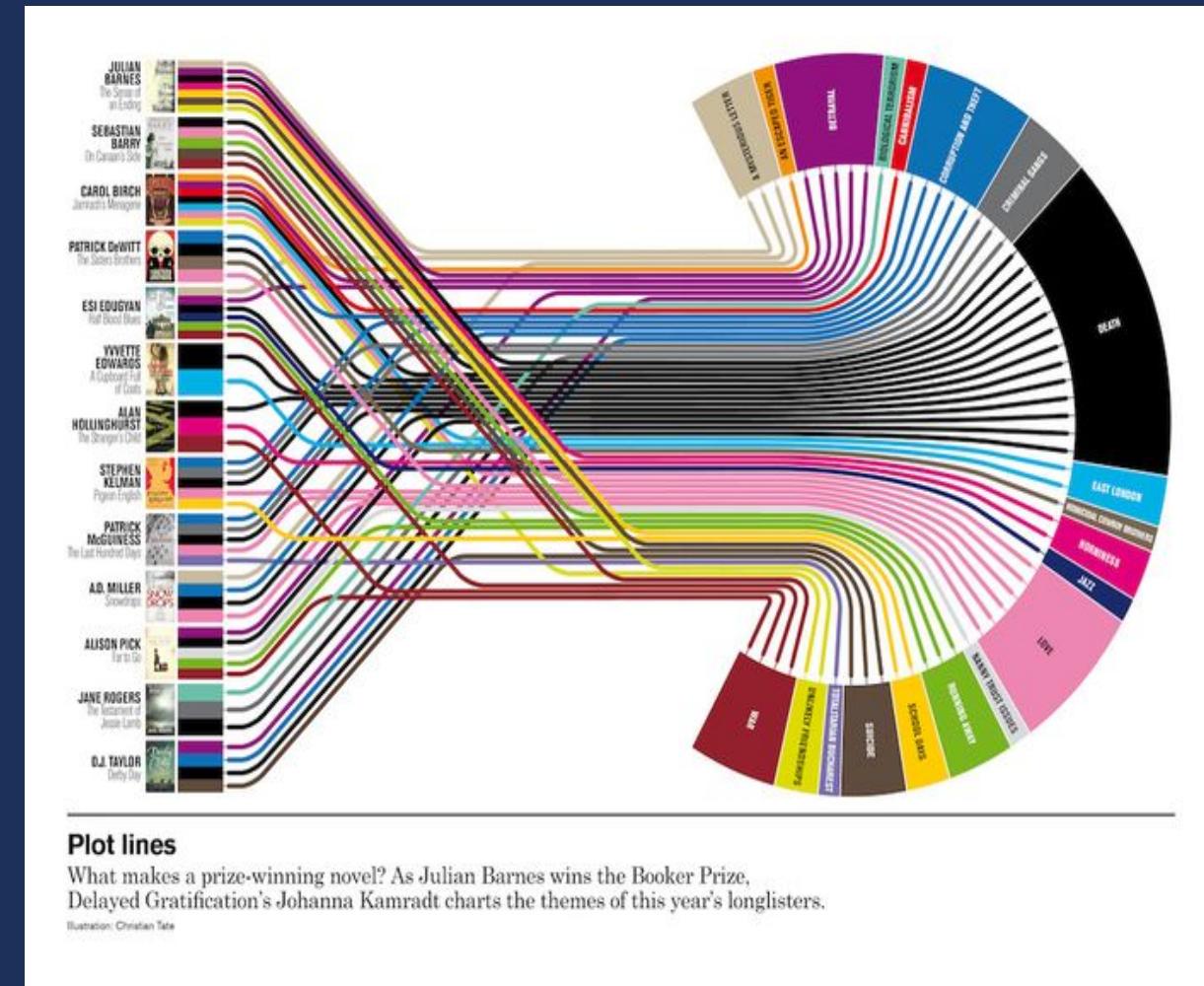
# WHAT'S WRONG WITH THESE VISUALIZATIONS?

Let's evaluate some data visualizations.



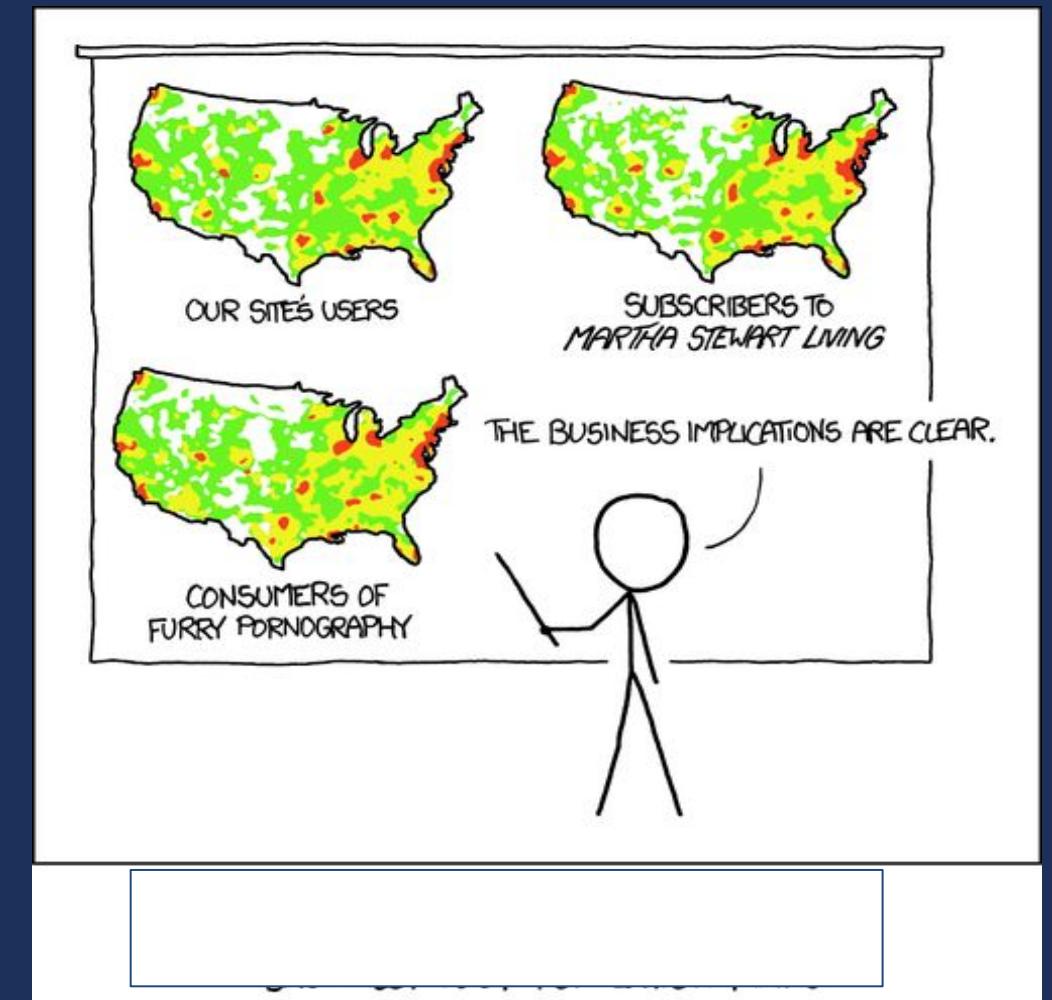
INCORRECT SCALE

<https://twitter.com/indiainpixels/stat/us/1657102345396682753>



TOO BUSY

<https://www.slow-journalism.com/infographics/booker-prize-the-infographics>

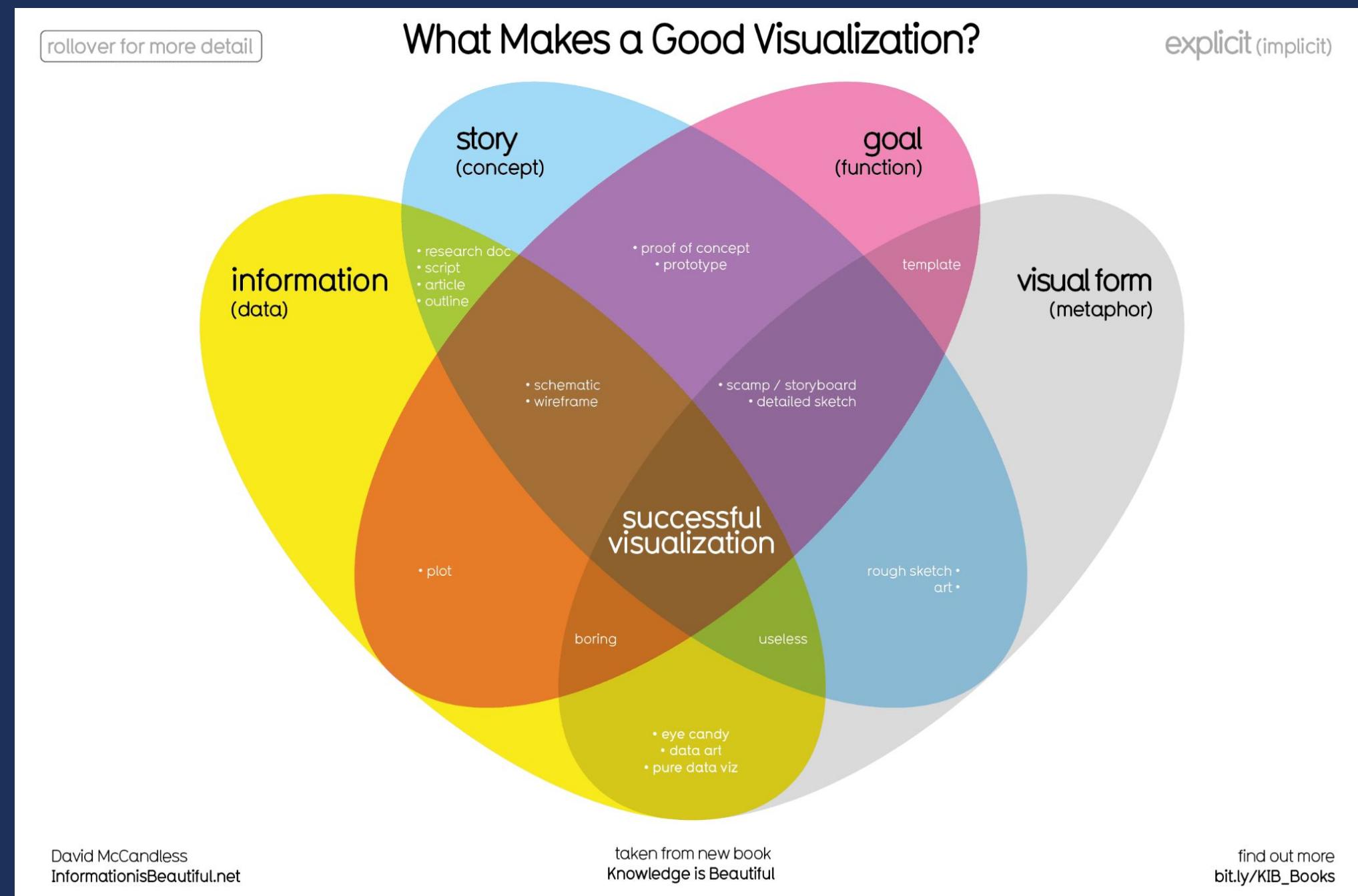


WRONG ANALYSIS

<https://xkcd.com/1138>

# WHAT MAKES A GOOD VISUALIZATION?

- A good visualization:
- has clear purpose and focus
  - contains only useful and accurate information
  - is structured correctly
  - is formatted accessibly



# WORKSHOP OUTLINE

.....

Our content today is divided into four parts. Each part will be described with examples.

01

## Data Visualization Process

Walking through the visualization process, from setting goals to visualizing the data.

02

## Basic Design Considerations

Best practices and accessibility considerations in data visualization.

03

## Solutions to Common Problems

How to resolve, or even better, avoid common problems in data visualization.

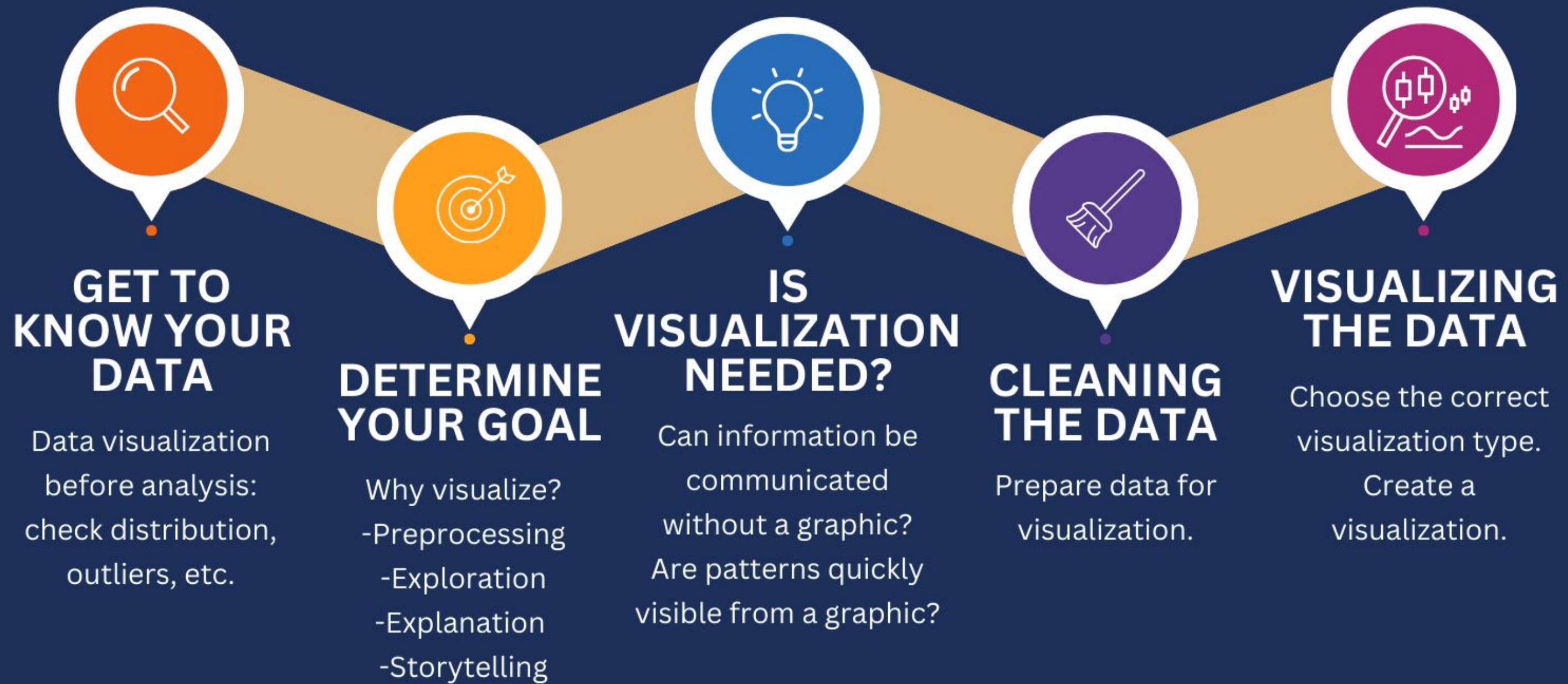
04

## Visualization Resources

Workshops, tools, and other resources for data visualization.

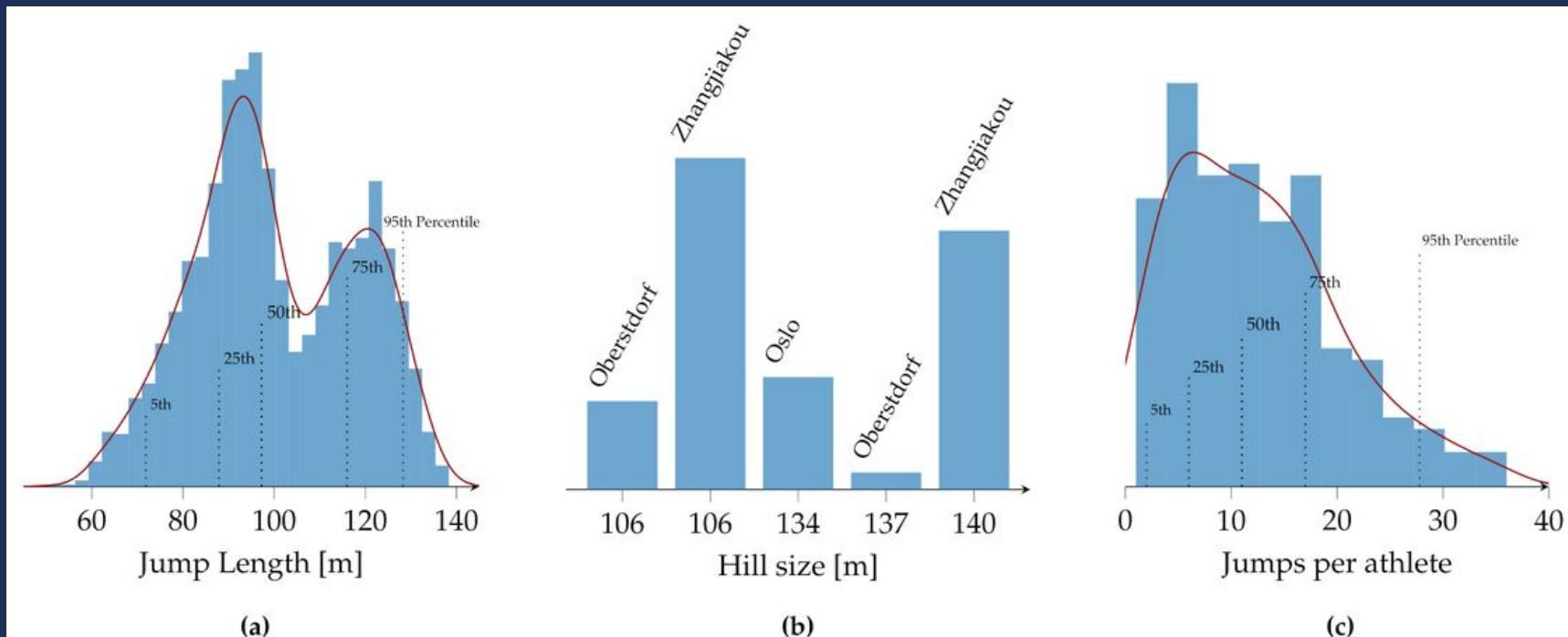
# DATA VISUALIZATION PROCESS

Data visualizations from start to finish

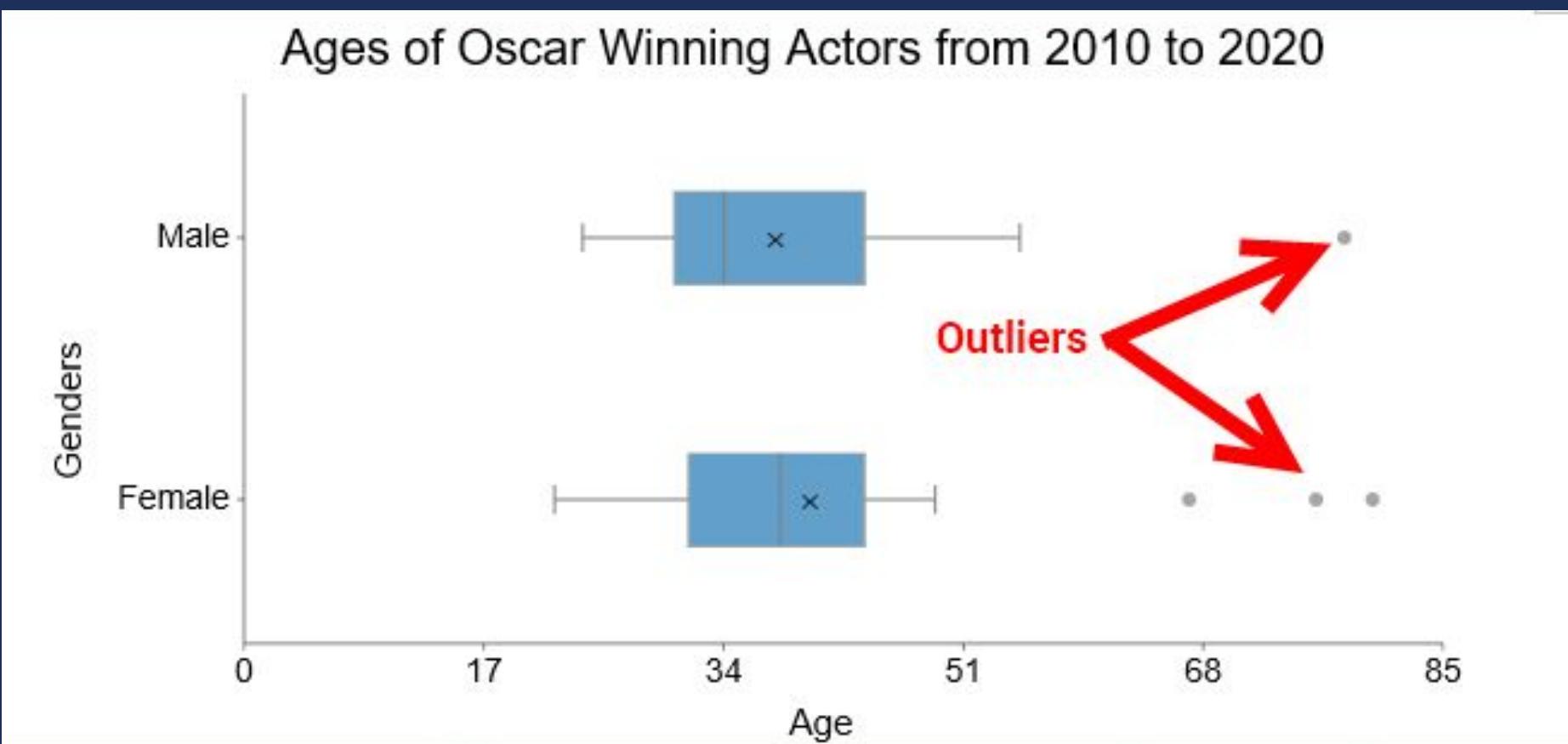


Visualization Goal

# GETTING TO KNOW YOUR DATA



Link J, Schwinn L, Pulsmeyer F, Kautz T, Eskofier BM. xLength: Predicting Expected Ski Jump Length Shortly after Take-Off Using Deep Learning. Sensors. 2022; 22(21):8474. <https://doi.org/10.3390/s22218474>



<https://chartexpo.com/blog/box-plot-outliers#>

## Identifying Outliers

Checking the  
Distribution

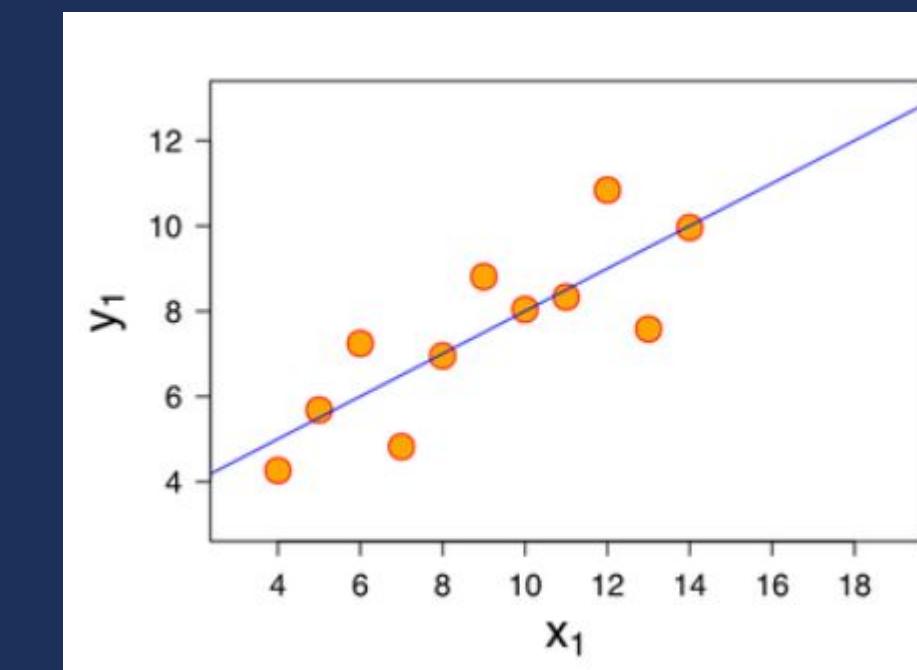
## Visualization Goal

# GETTING TO KNOW YOUR DATA

I	
x	y
10.0	8.04
8.0	6.95
13.0	7.58
9.0	8.81
11.0	8.33
14.0	9.96
6.0	7.24
4.0	4.26
12.0	10.84
7.0	4.82
5.0	5.68

For all four datasets:

Property	Value	Accuracy
Mean of x	9	exact
Sample variance of x: $s_x^2$	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of y: $s_y^2$	4.125	$\pm 0.003$
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression: $R^2$	0.67	to 2 decimal places



- Explore relationships between variables
  - Visualize trends or patterns
  - **Anscombe's Quartet**

## Visualization Goal

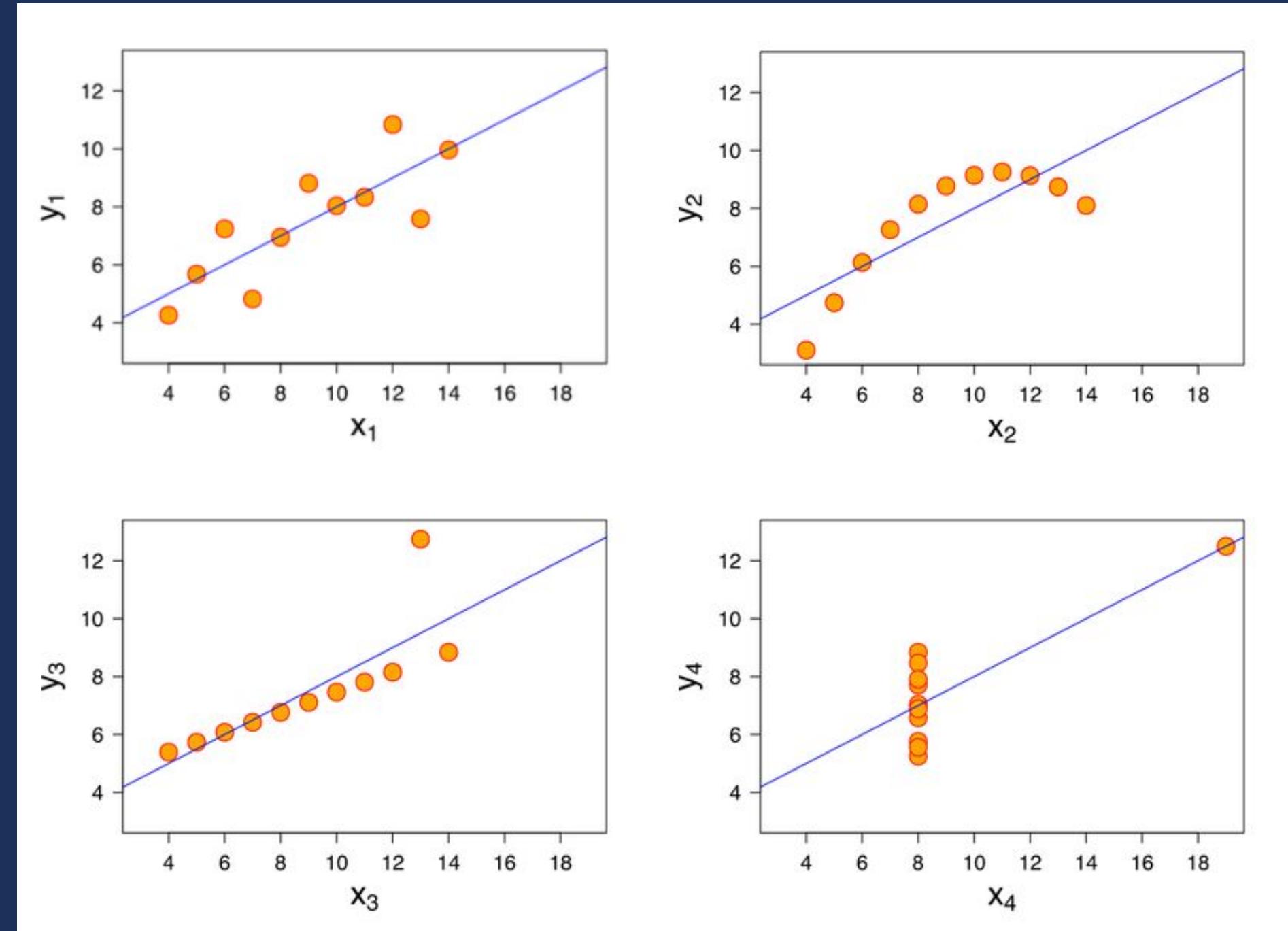
# GETTING TO KNOW YOUR DATA

Anscombe's quartet							
I		II		III		IV	
x	y	x	y	x	y	x	y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

[https://en.wikipedia.org/wiki/Anscombe%27s\\_quartet](https://en.wikipedia.org/wiki/Anscombe%27s_quartet)

For all four datasets:

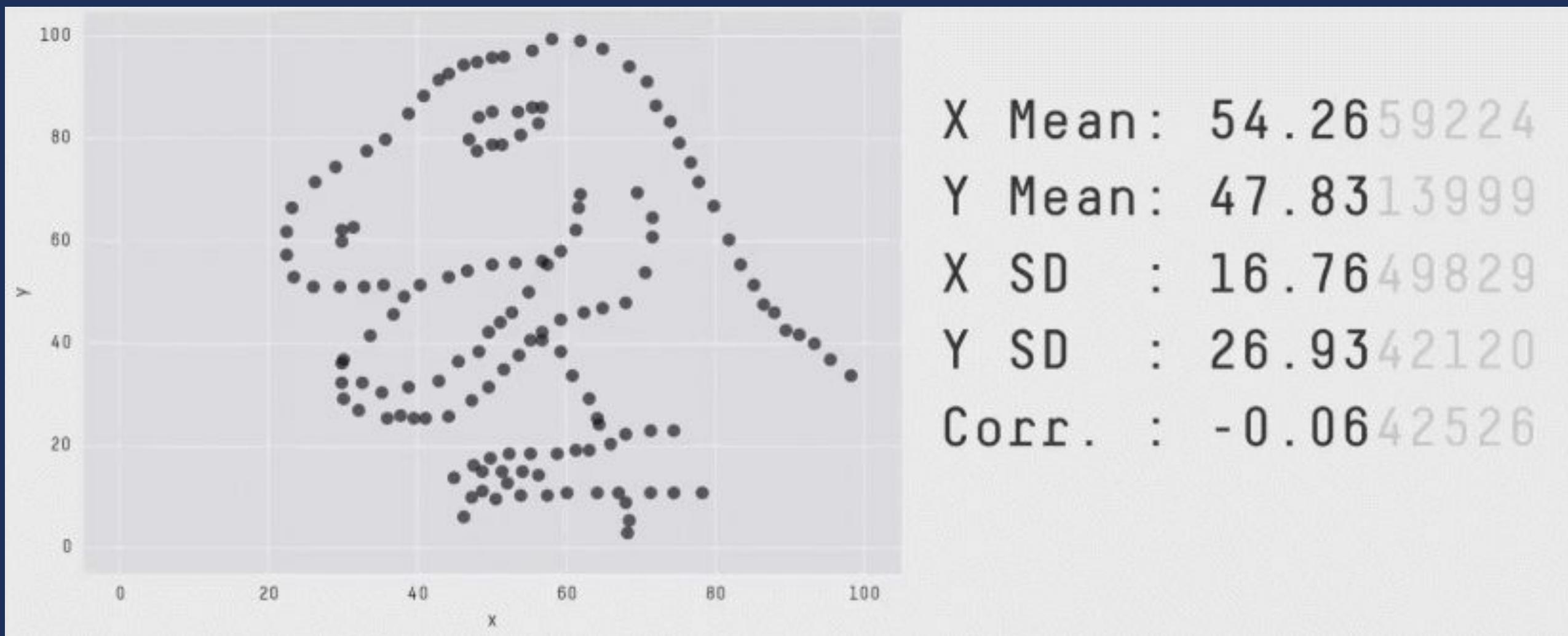
Property	Value	Accuracy
Mean of $x$	9	exact
Sample variance of $x$ : $s_x^2$	11	exact
Mean of $y$	7.50	to 2 decimal places
Sample variance of $y$ : $s_y^2$	4.125	$\pm 0.003$
Correlation between $x$ and $y$	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression: $R^2$	0.67	to 2 decimal places



## Visualization Goal

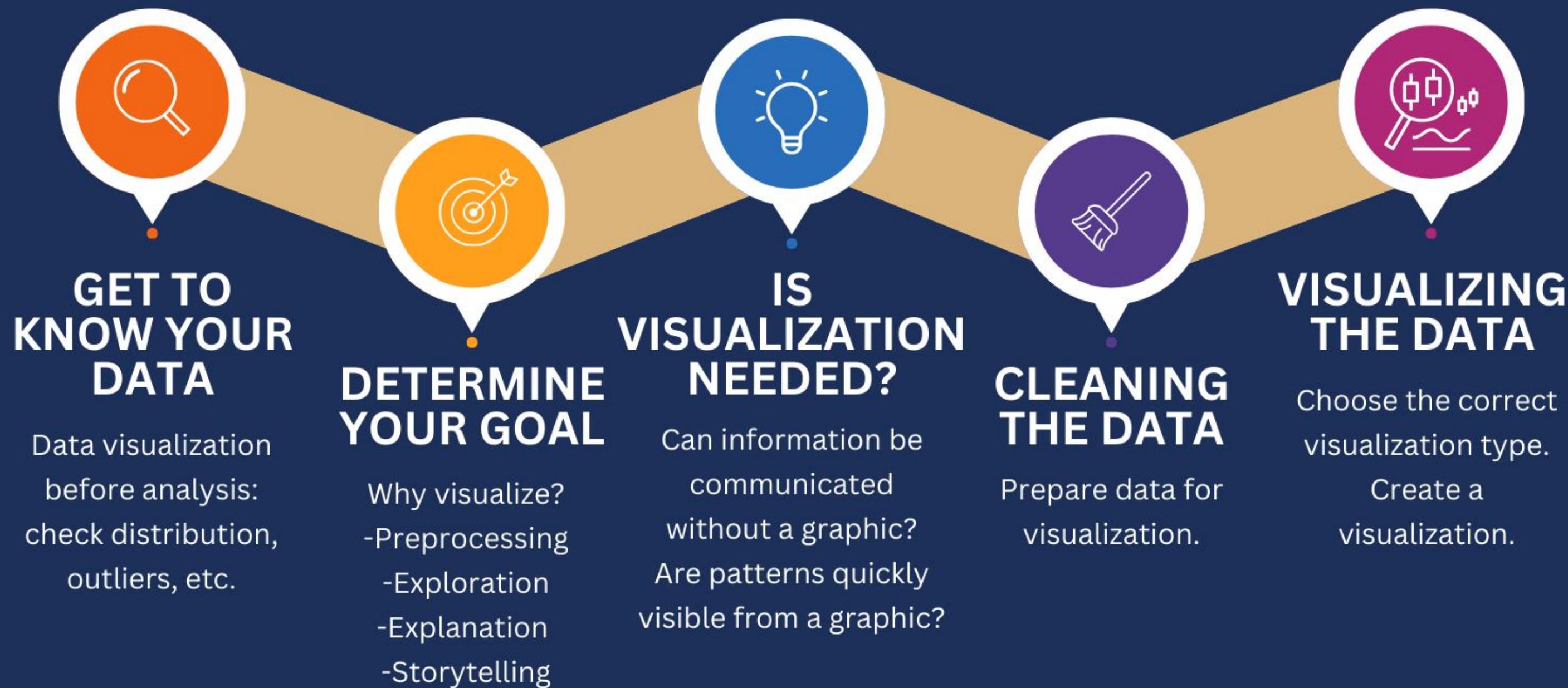
# GETTING TO KNOW YOUR DATA

- Explore relationships between variables
  - Visualize trends or patterns
  - **Anscombe's Quartet**



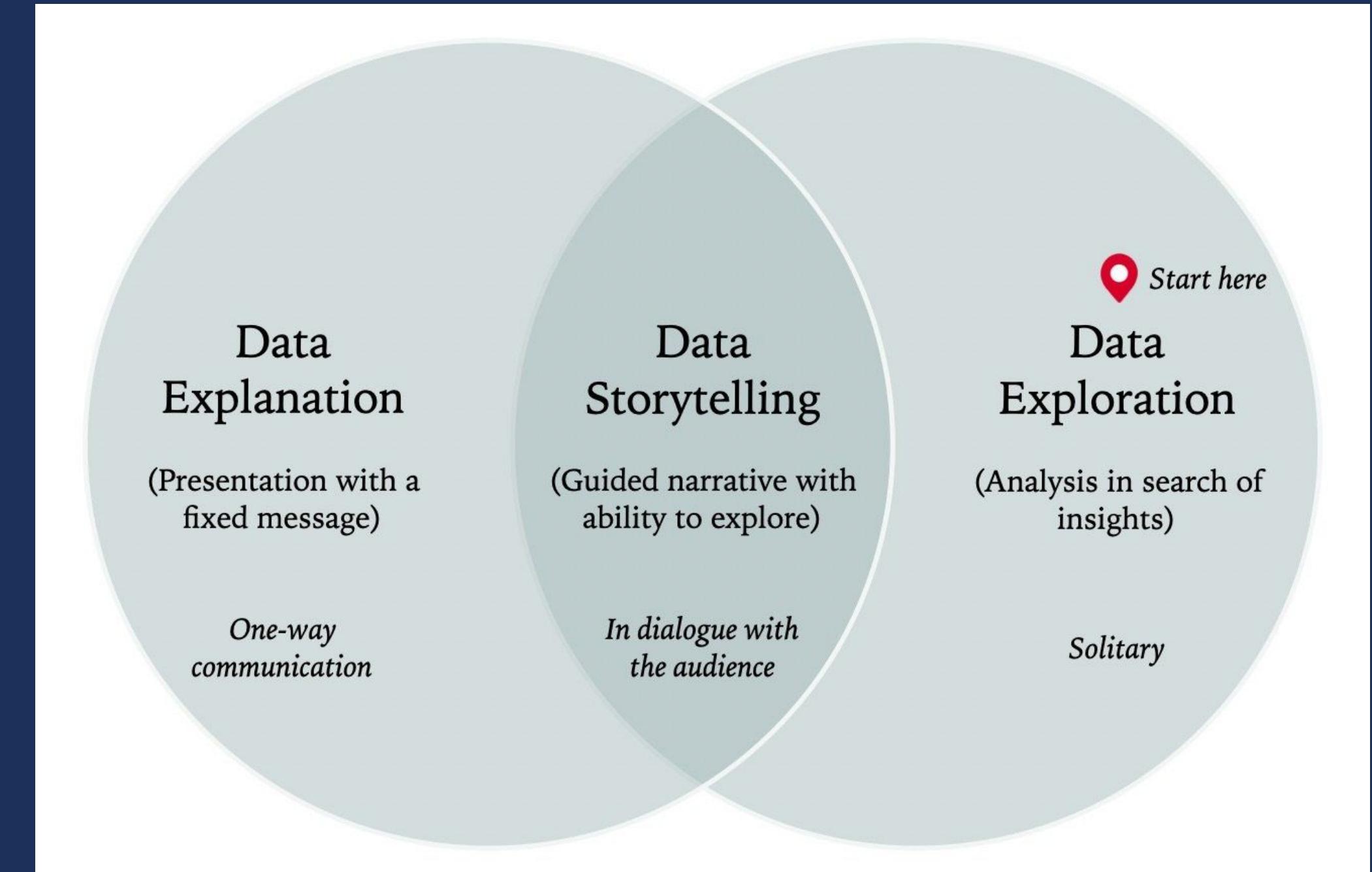
# DATA VISUALIZATION PROCESS

Data visualizations from start to finish



# DATA ANALYSIS

- Purposes of data visualizations:
  - Exploratory
  - Explanatory
  - Storytelling
- For all purposes, consider:
  - Who is the target audience?
  - Do the data support the message?

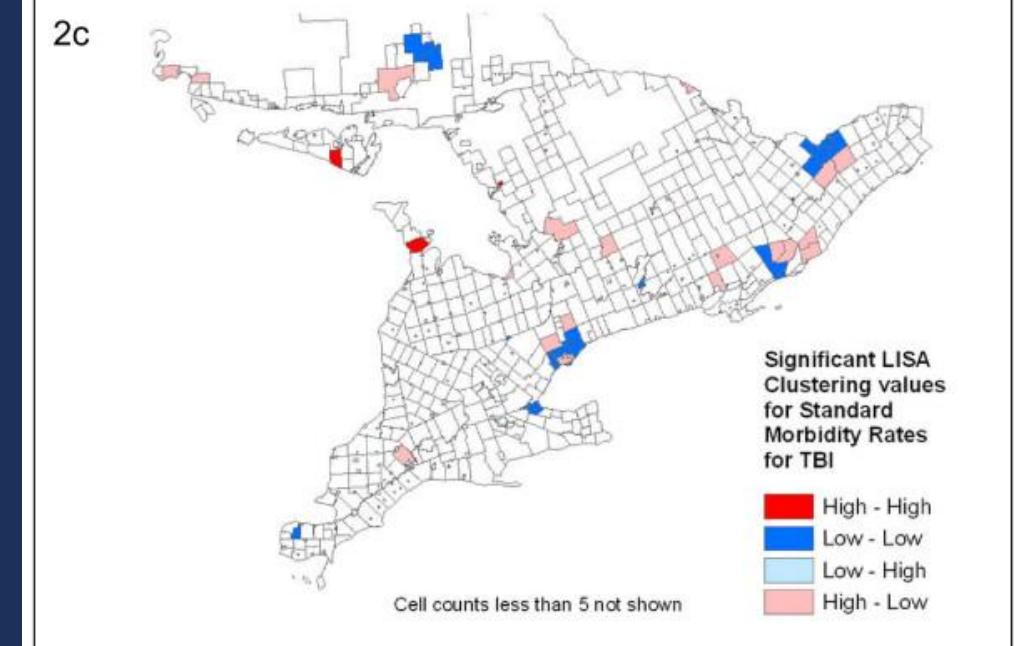
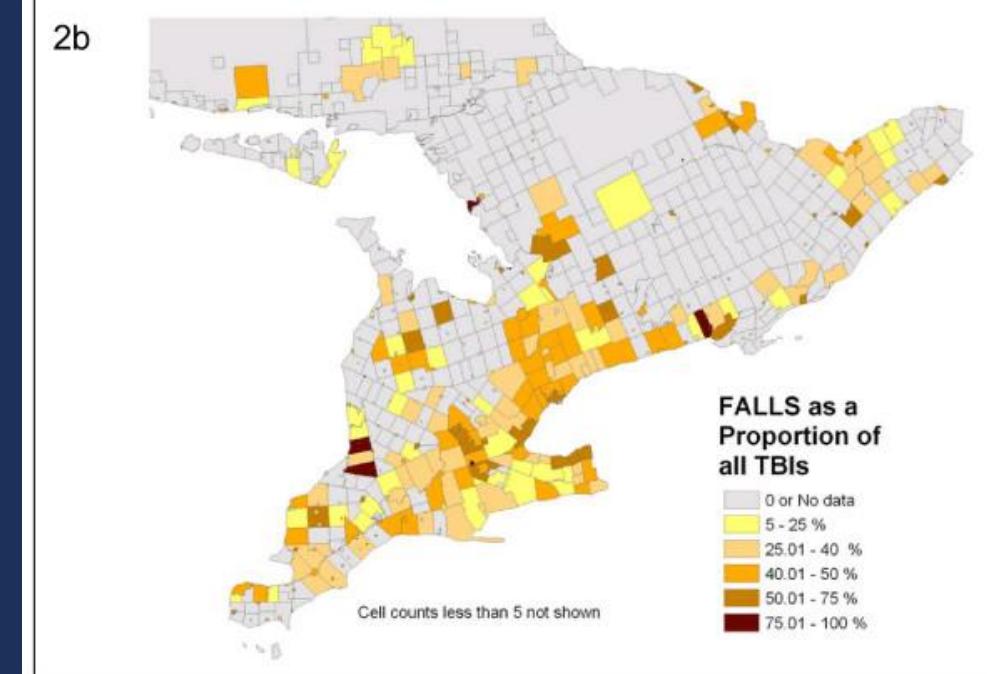
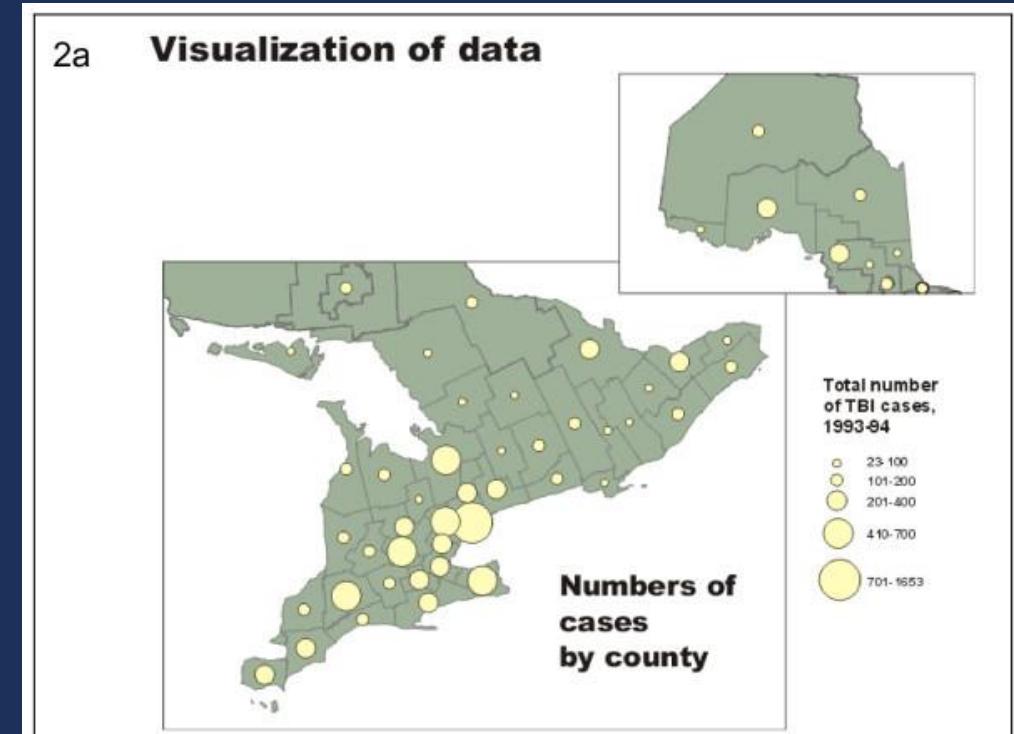


## Visualization Goal

# DATA EXPLORATION

- Consider:
  - What features of the data are you comparing and contrasting?
  - What relationships within the data are you exploring?

Colantonio, Angela & Moldofsky, Byron & Escobar, Michael & Vernich, Lee & Chipman, Mary & McLellan, Barry. (2011). Using geographical information systems mapping to identify areas presenting high risk for traumatic brain injury. Emerging themes in epidemiology. 8. 7. 10.1186/1742-7622-8-7.



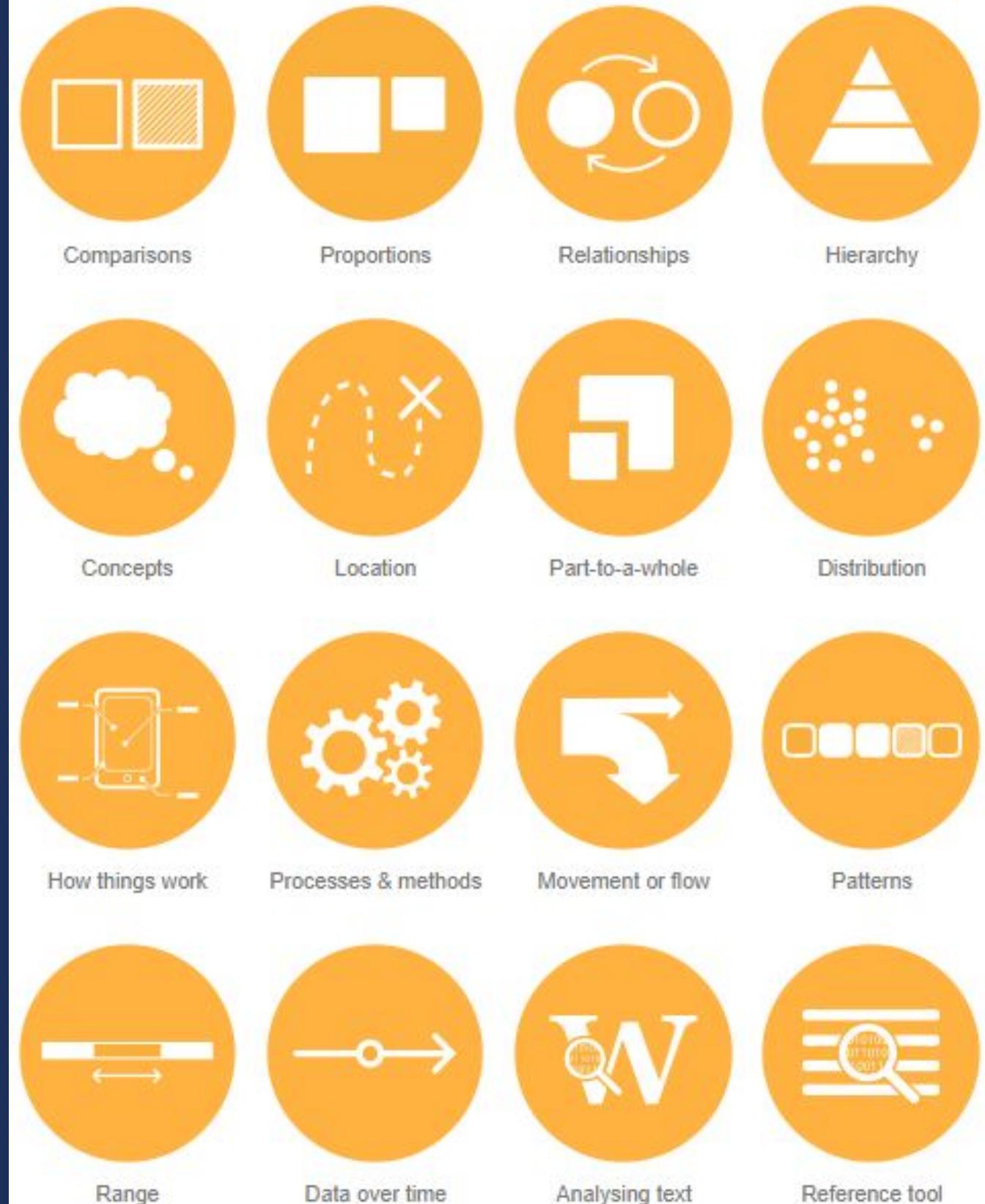
## Visualization Goal

# DATA EXPLANATION

- Different types of **explanatory** data visualizations.
- Consider:
  - What comparisons or relationships do you want to show?
  - Who is your audience?
  - What visualizations are appropriate for your data type?

## What do you want to show?

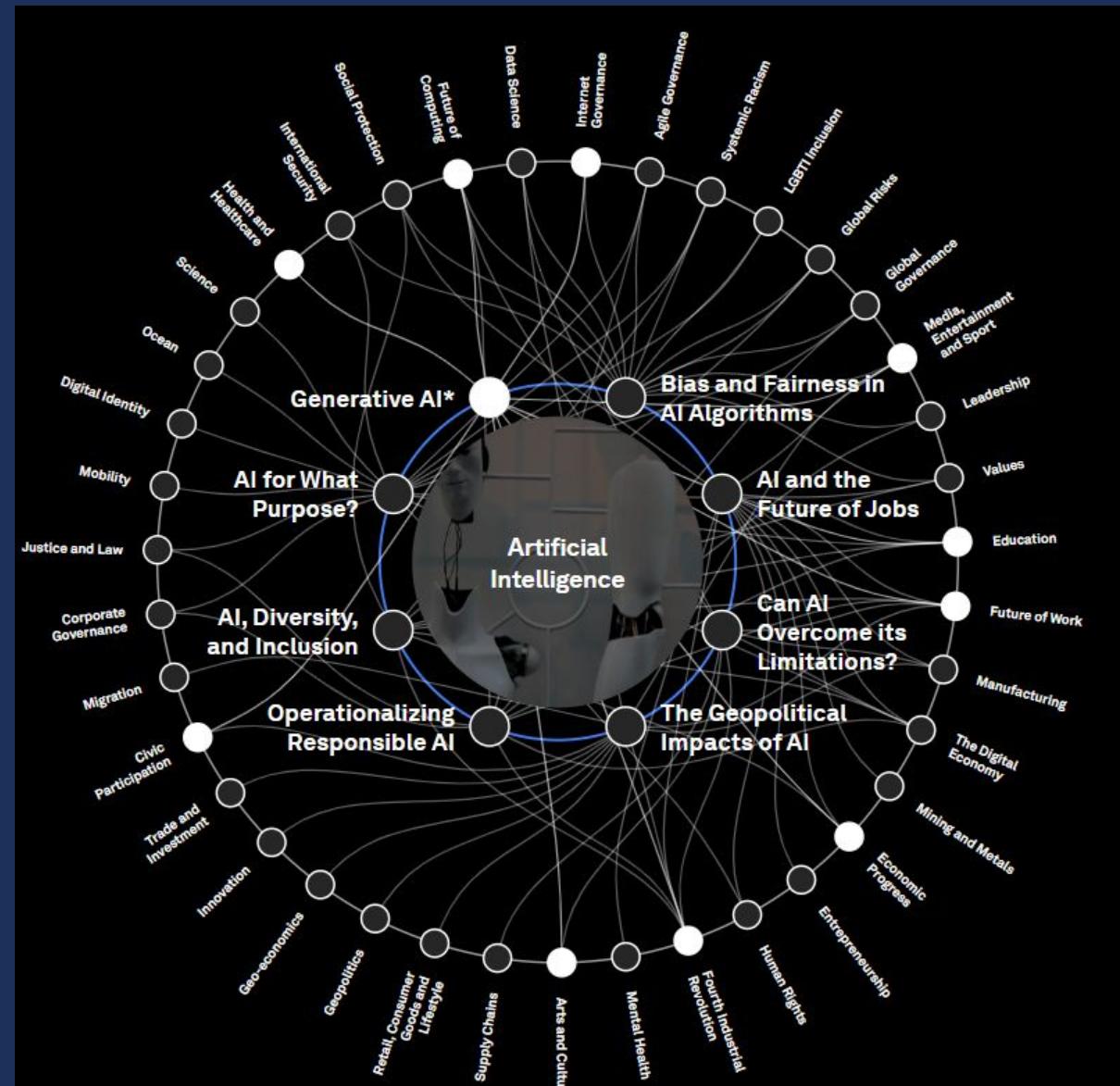
Here you can find a list of charts categorised by their data visualization functions or by what you want a chart to communicate to an audience. While the allocation of each chart into specific functions isn't a perfect system, it still works as a useful guide for selecting chart based on your analysis or communication needs.



# Visualization Goal

# DATA STORYTELLING

- Consider:
    - What story are you trying to tell?
    - Who is your audience?



<https://intelligence.weforum.org/topics/a1Gb0000000pTDREA2/key-issues/a1G680000000Ne9EAE>

[https://www.youtube.com/watch?v=hRWEteXYD\\_Y](https://www.youtube.com/watch?v=hRWEteXYD_Y)

## Visualization Goal

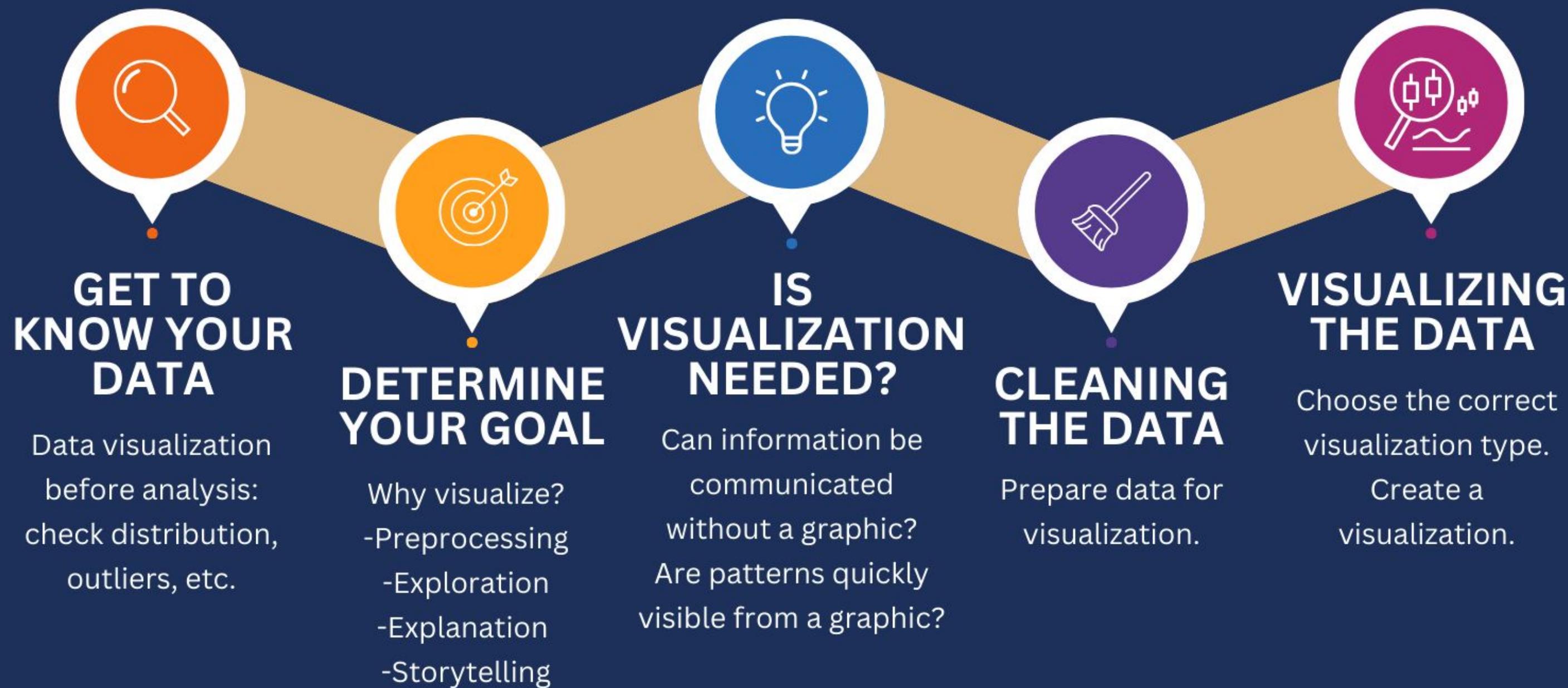
# DATA STORYTELLING

- Interactive tools may aid in data storytelling



# DATA VISUALIZATION PROCESS

Data visualizations from start to finish



# IS VISUALIZATION NEEDED?

---

- **Data visualizations** must be interpretable and relevant
  - Goal-oriented
  - Clearly display patterns or comparisons

# WORKSHOP OUTLINE

.....

Our content today is divided into four parts. Each part will be described with examples.

01

## Data Visualization Process

Walking through the visualization process, from setting goals to visualizing the data.

02

## Basic Design Considerations

Best practices and accessibility considerations in data visualization.

03

## Solutions to Common Problems

How to resolve, or even better, avoid common problems in data visualization.

04

## Visualization Resources

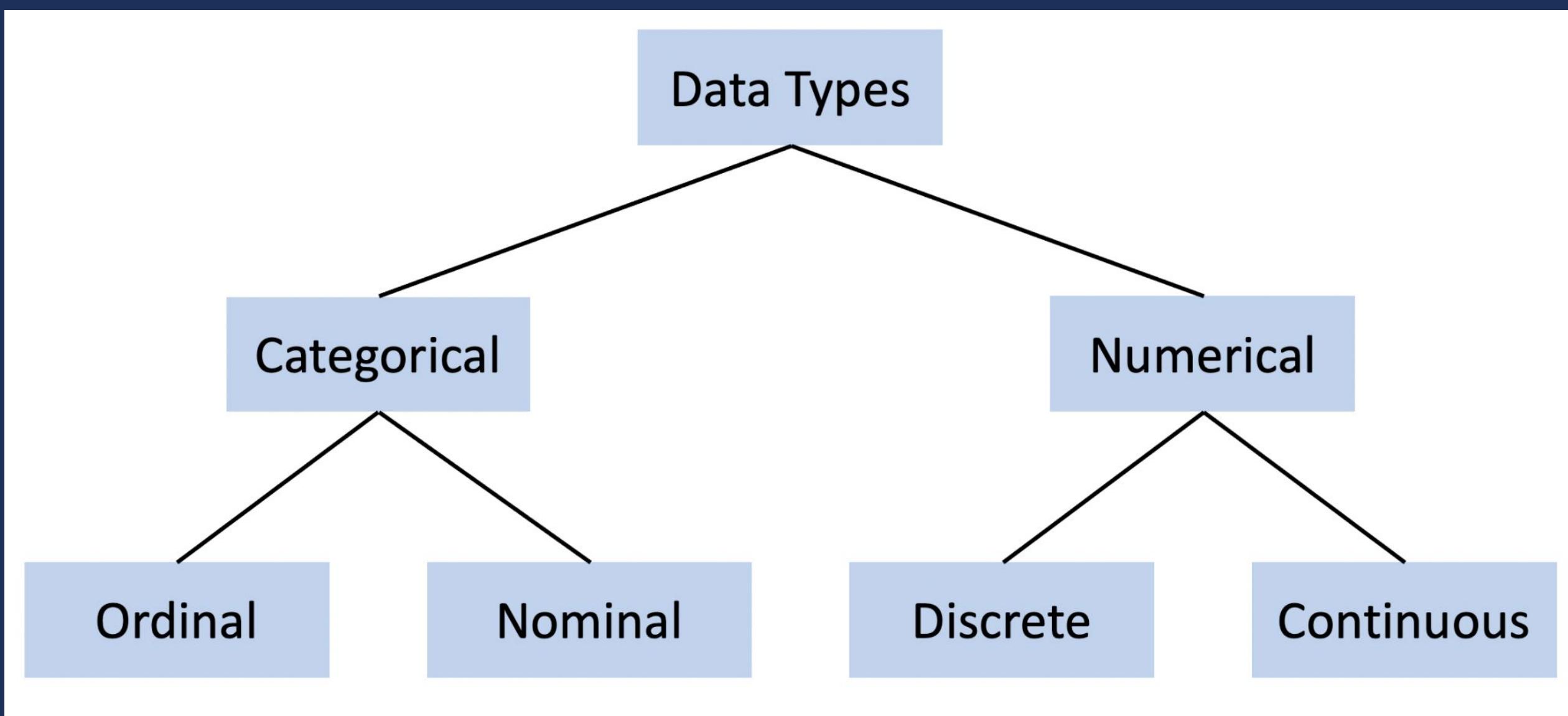
Workshops, tools, and other resources for data visualization.

# CONSIDERATIONS AND CONCEPTS

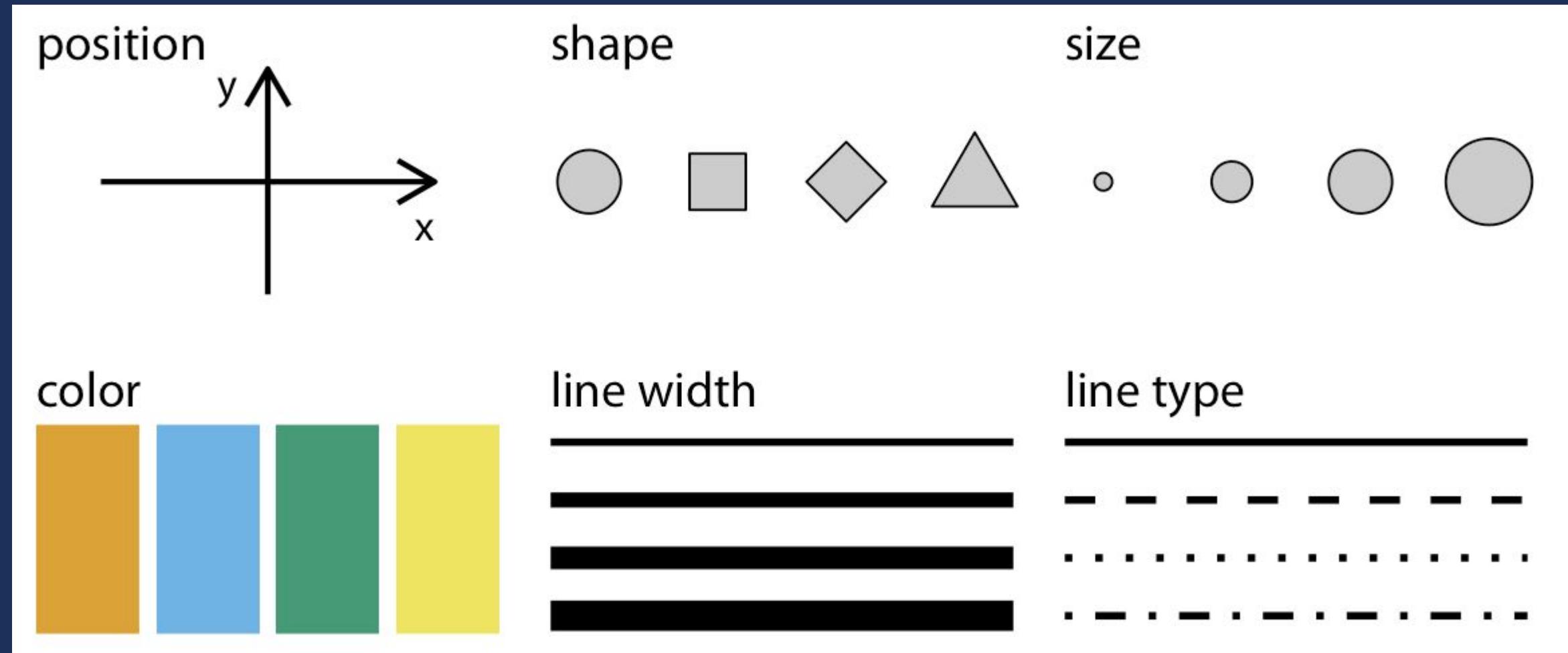
---

# A review of variable types

---

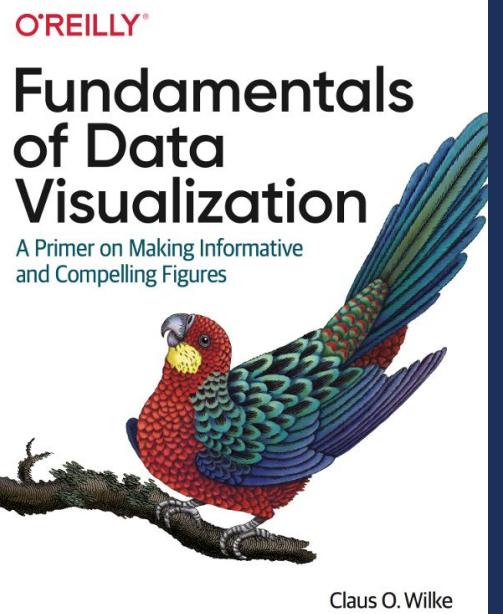


# Aesthetic properties



## Acknowledgement:

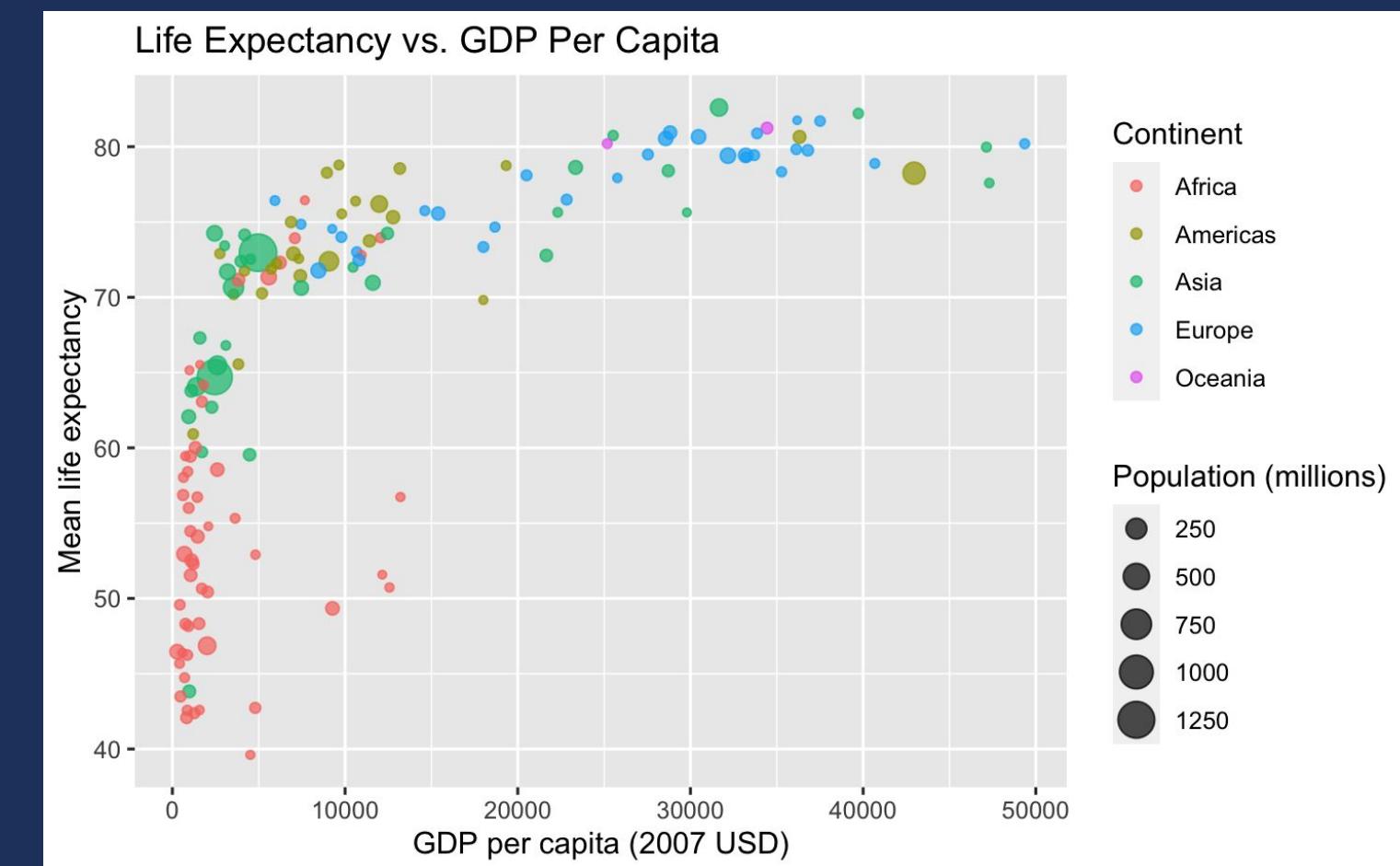
Many of the upcoming images  
are from this book, available  
online at [clauswilke.com/dataviz](http://clauswilke.com/dataviz)



# Variable → Aesthetic property

**Size**    **Color**    **Y**    **X**

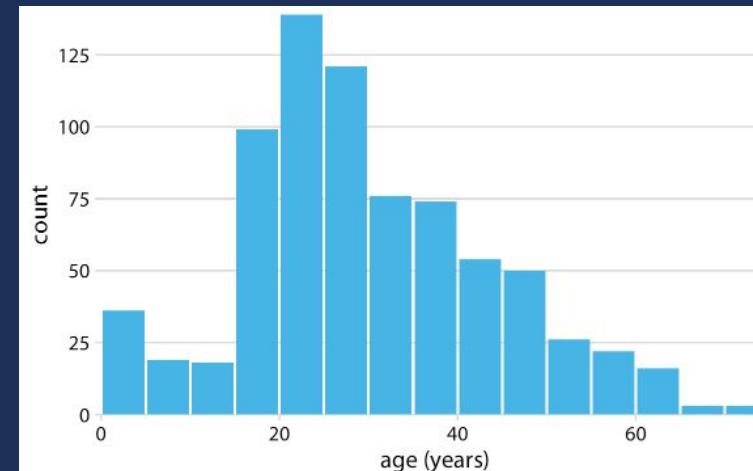
country	pop	continent	lifeExp	gdpPerCap	gdp
Afghanistan	31889923	Asia	43.828	974.5803	3.107929e+10
Albania	3600523	Europe	76.423	5937.0295	2.137641e+10
Algeria	33333216	Africa	72.301	6223.3675	2.074449e+11
Angola	12420476	Africa	42.731	4797.2313	5.958390e+10
Argentina	40301927	Americas	75.320	12779.3796	5.150336e+11
Australia	20434176	Oceania	81.235	34435.3674	7.036584e+11
Austria	8199783	Europe	79.829	36126.4927	2.962294e+11
Bahrain	708573	Asia	75.635	29796.0483	2.111268e+10
Bangladesh	150448339	Asia	64.062	1391.2538	2.093118e+11
Belgium	10392226	Europe	79.441	33692.6051	3.501412e+11



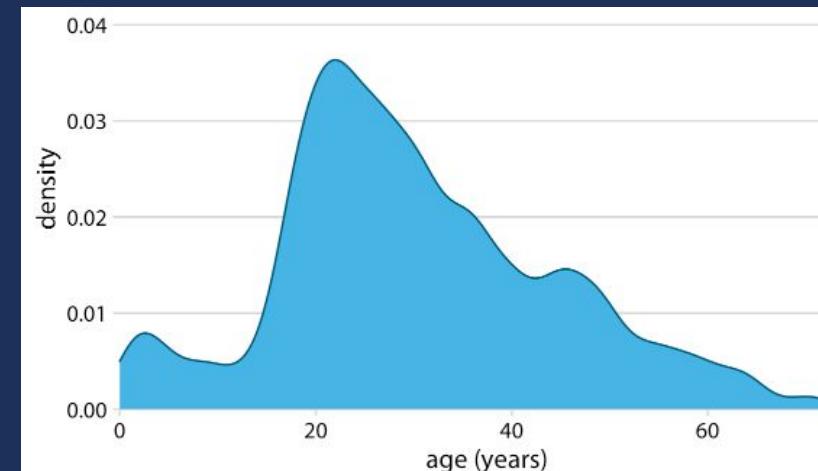
# Describing data with visualizations

## Single variable

### Numeric - Continuous

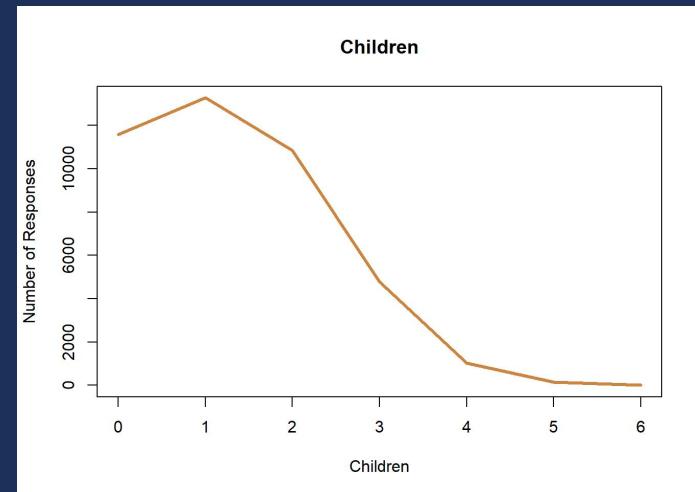
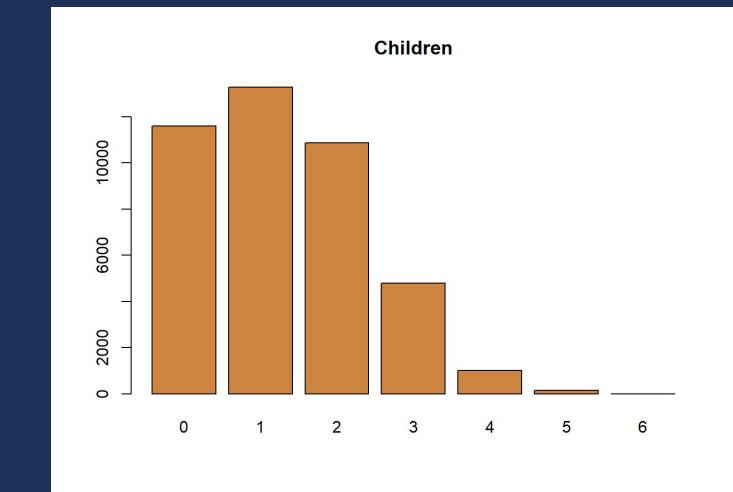


Histogram



Kernel Density Estimate

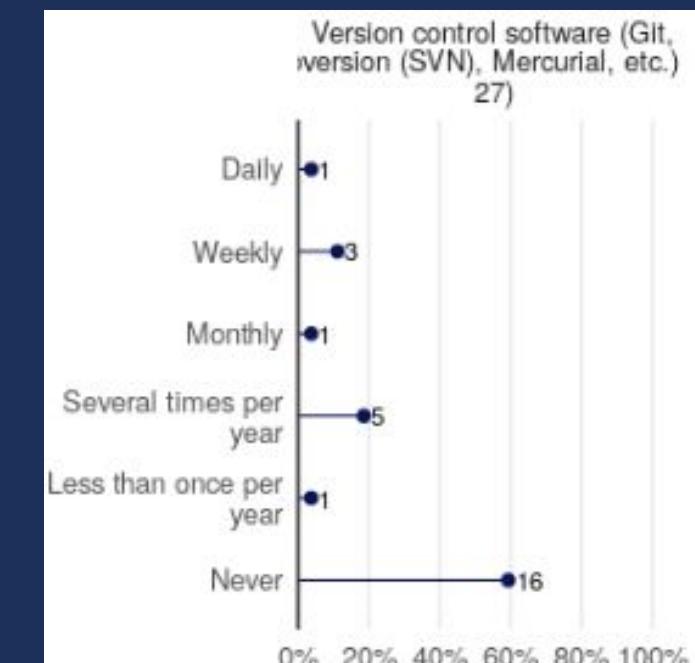
### Numeric - Discrete



### Categorical - Nominal



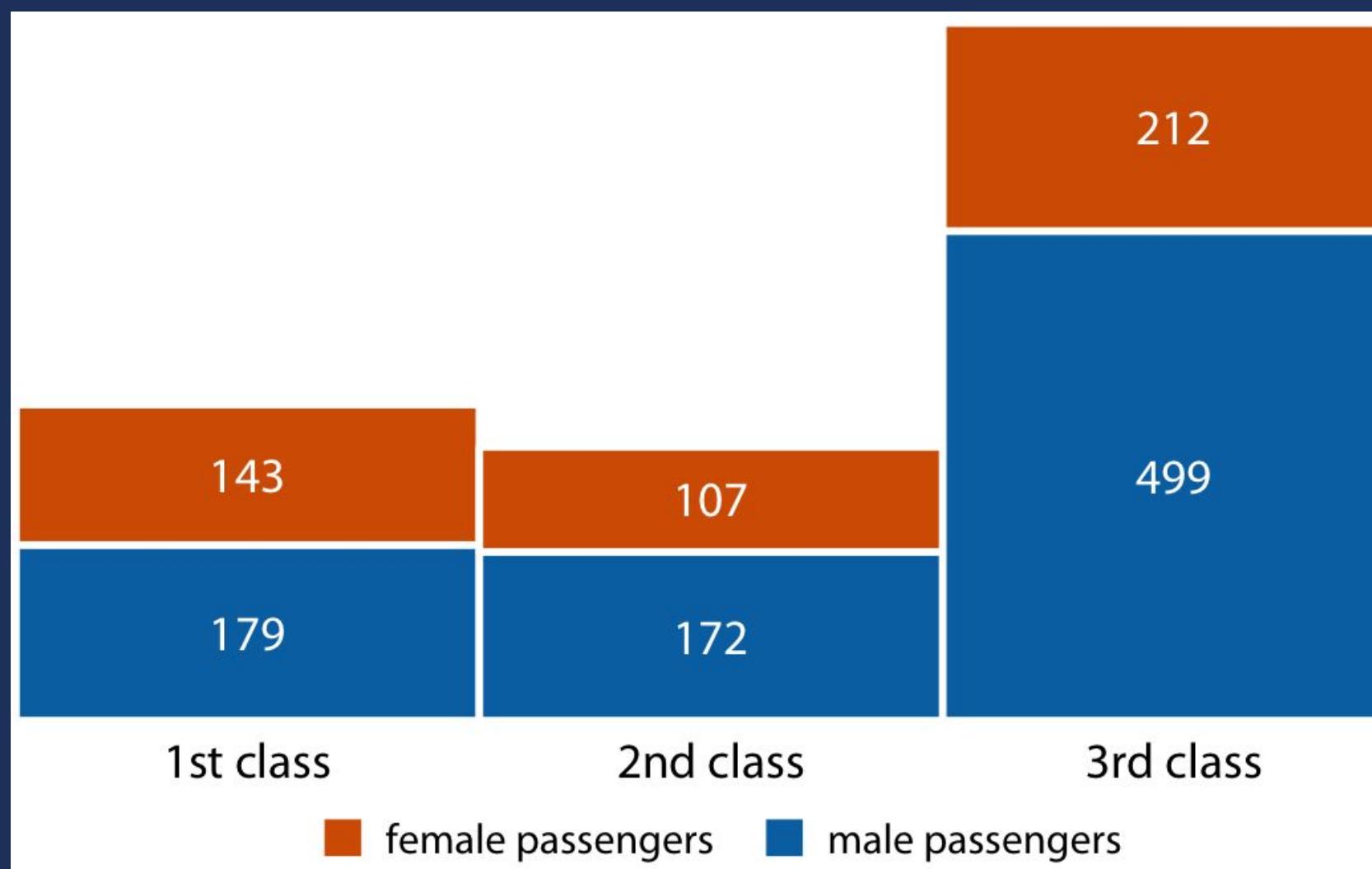
### Categorical - Ordinal



# Describing data with visualizations

## Relationship between categorical variables

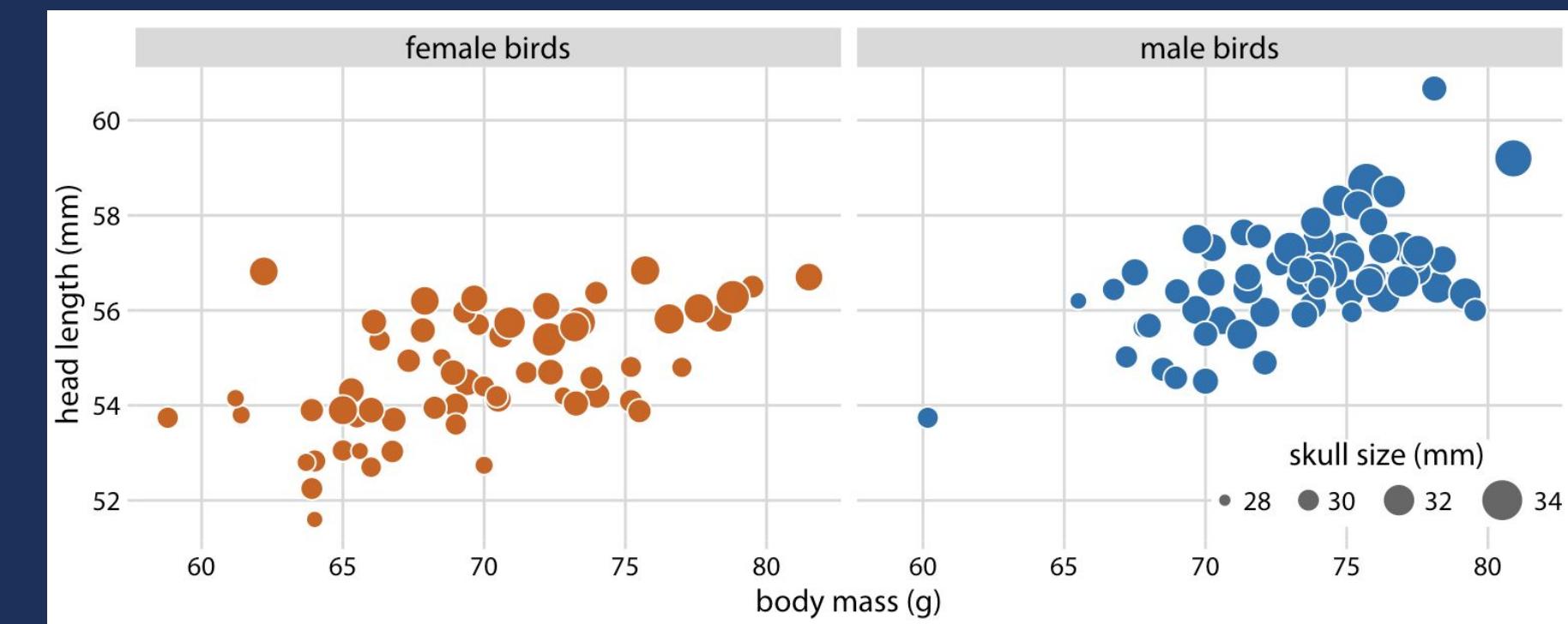
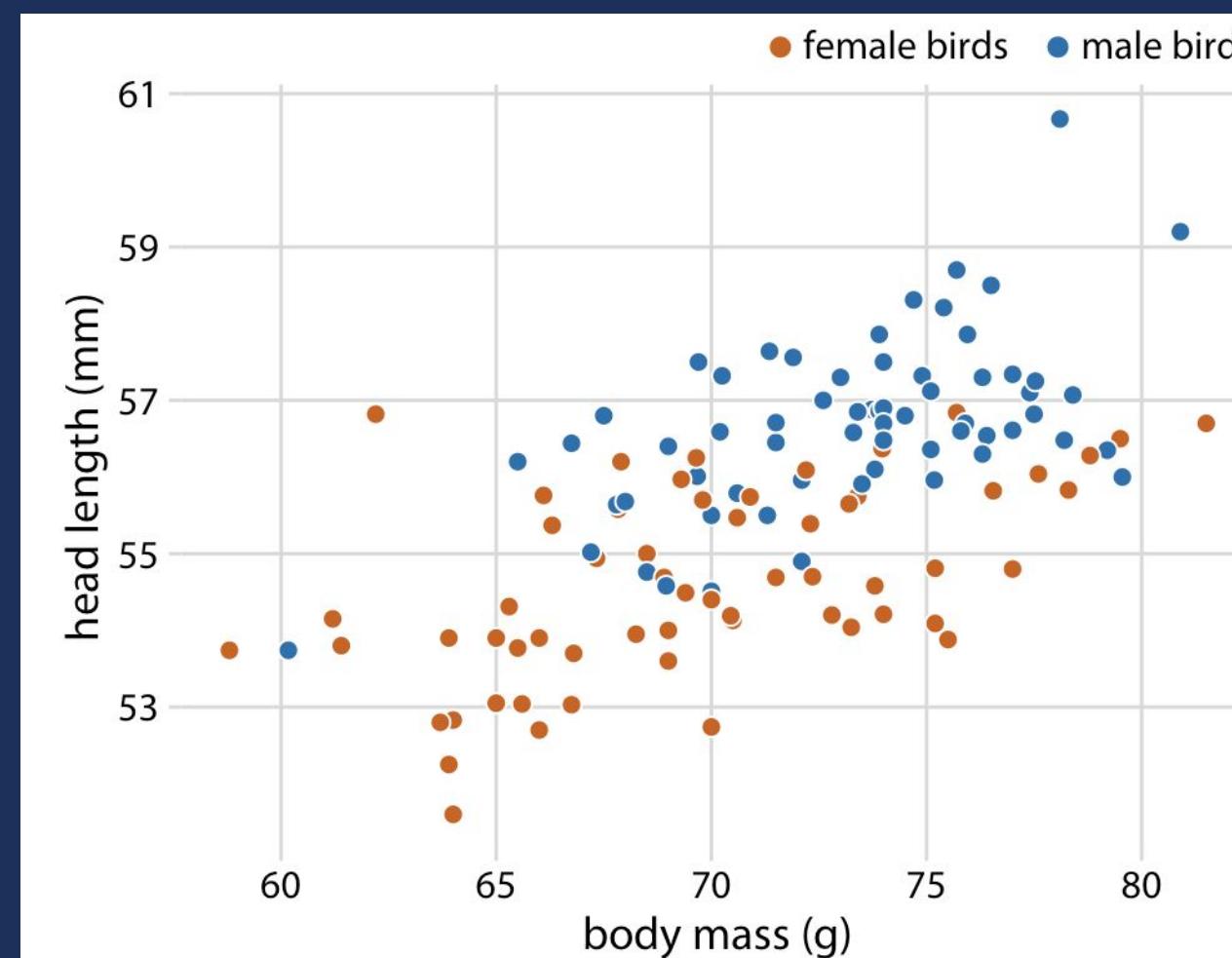
Categorical ~ Categorical



# Describing data with visualizations

## Relationship between continuous variables

Numeric/Continuous ~ Numeric/Continuous

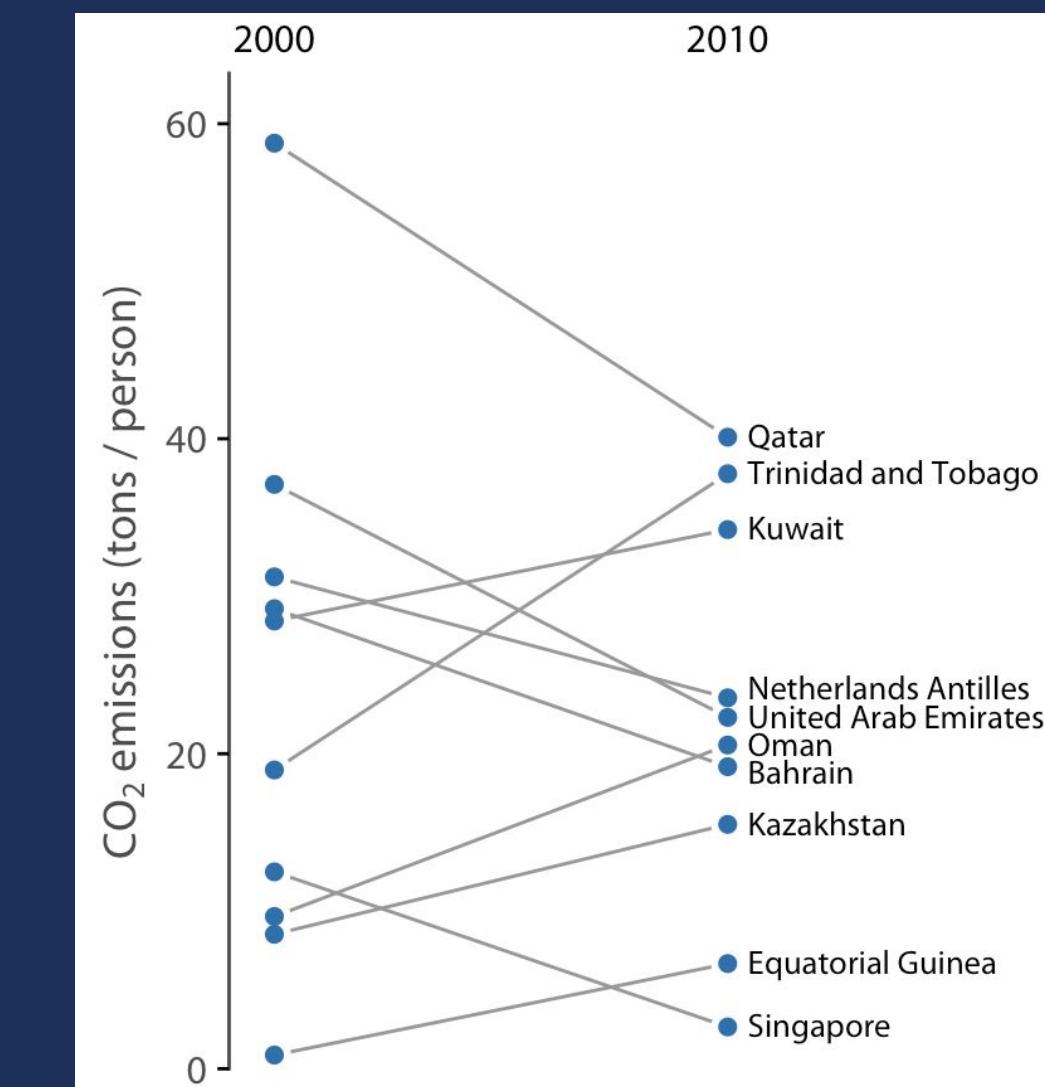
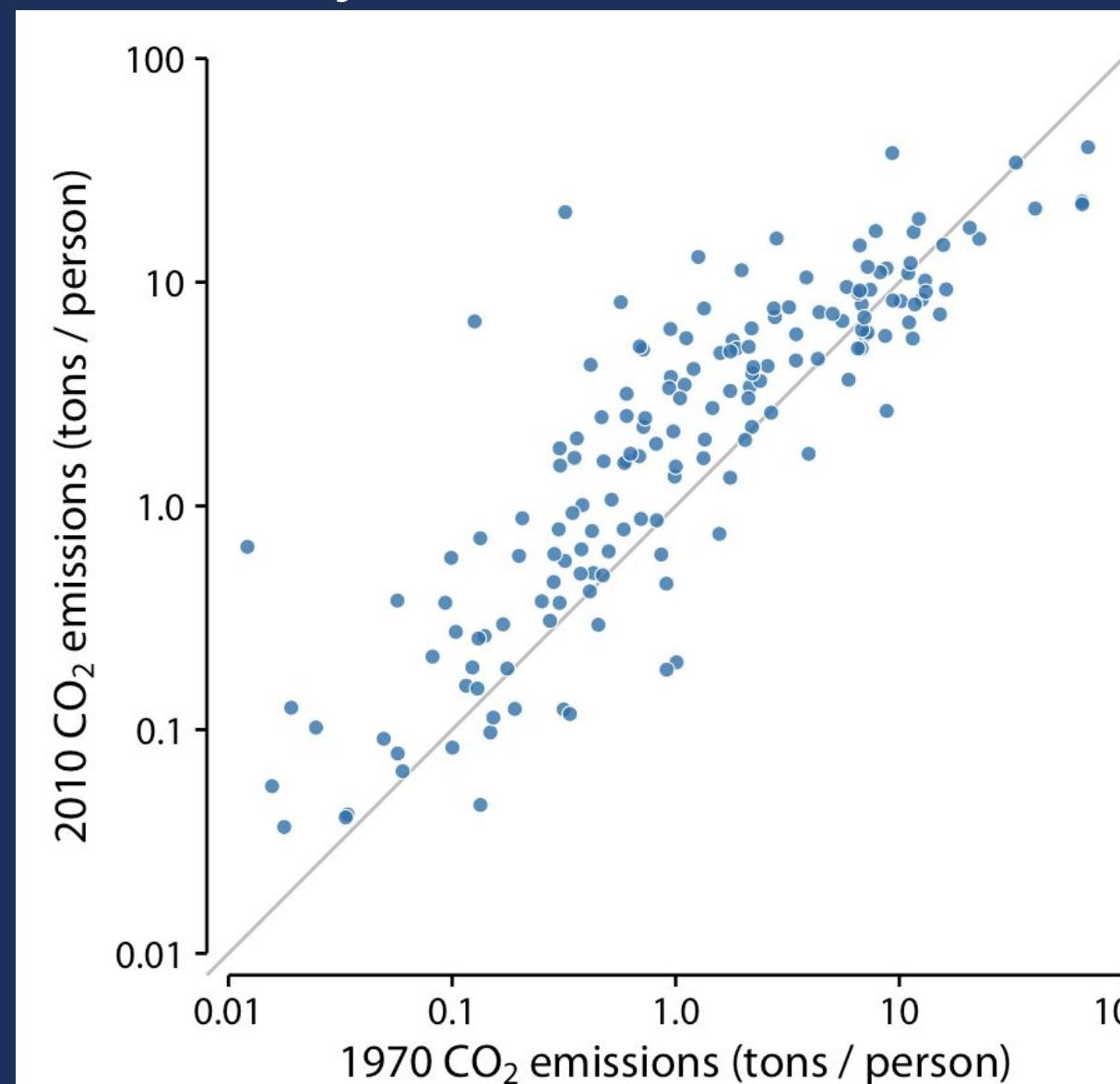


Bubble plot + faceting

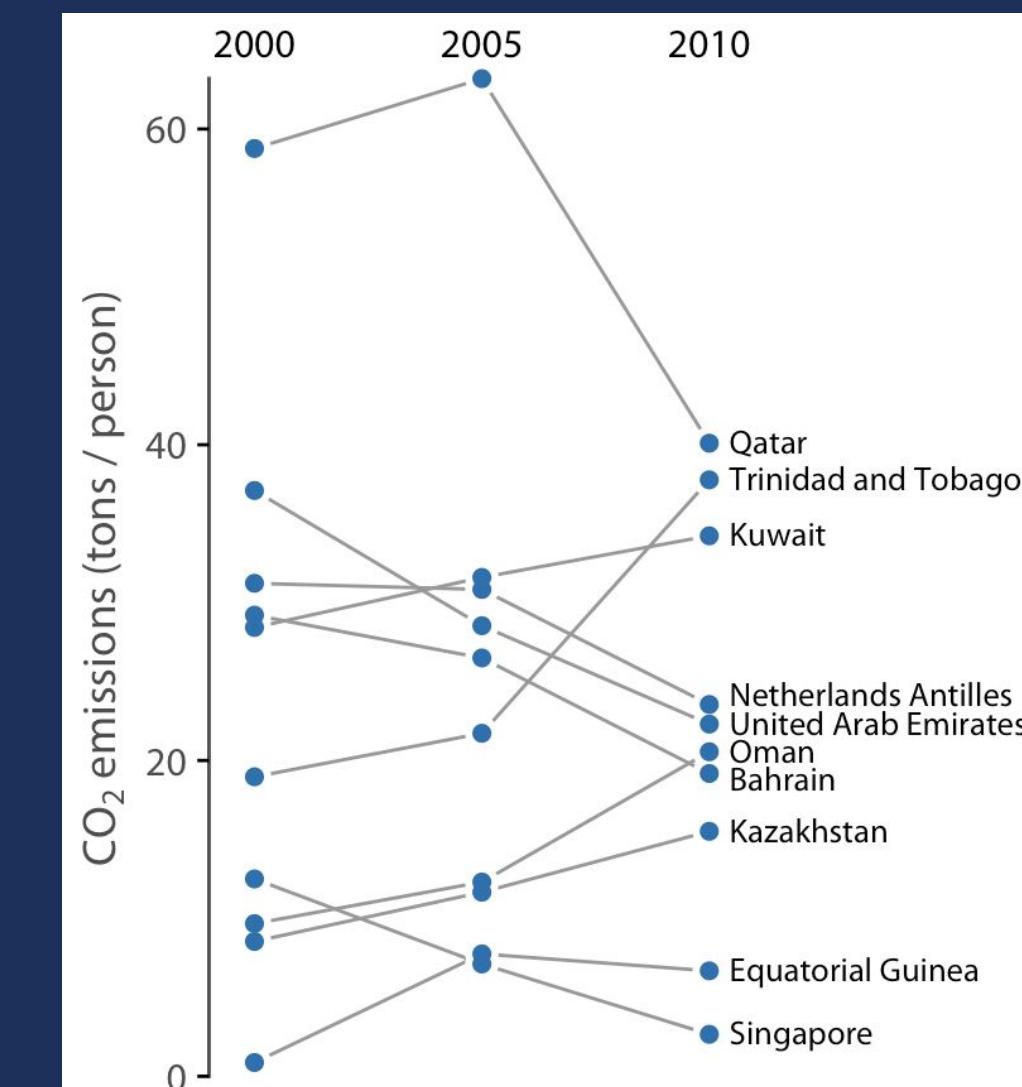
# Describing data with visualizations

## Relationship between "paired" continuous variables

x & y variables have the same scale,  
so an x=y line is useful

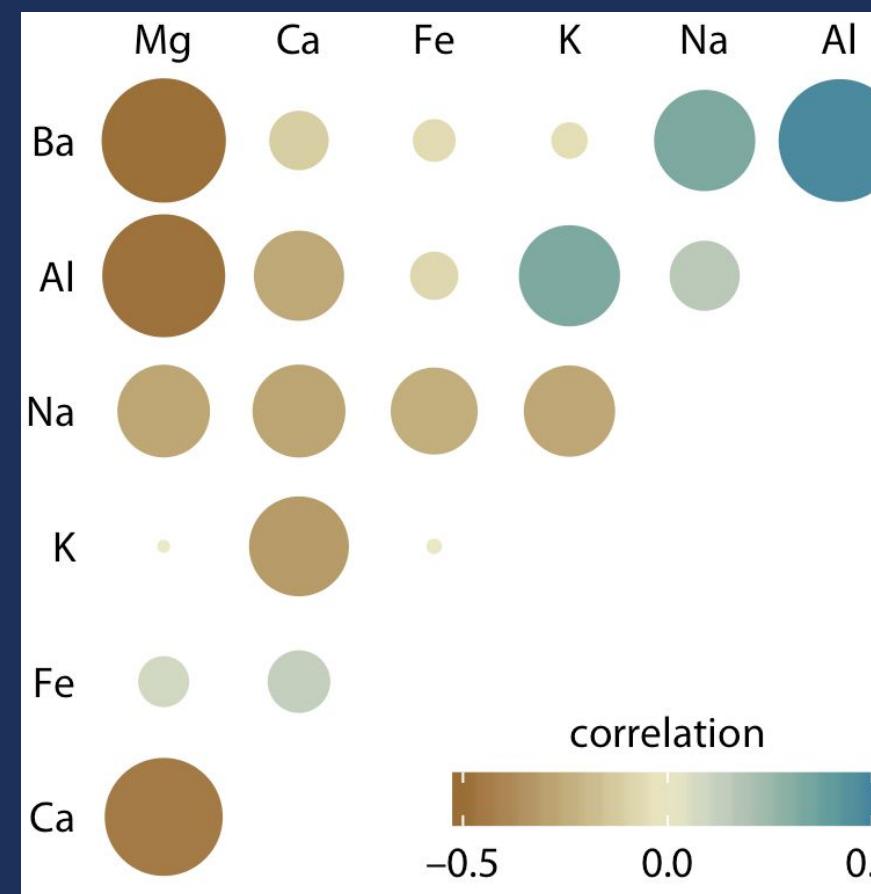
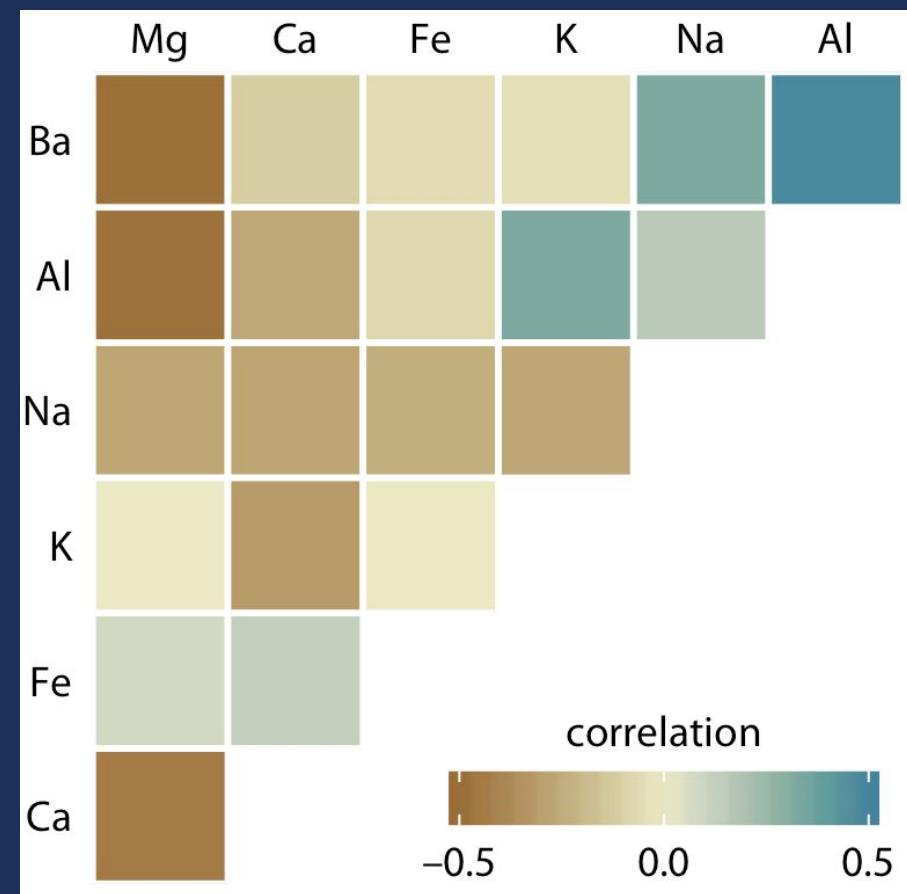


Slopegraph



# Describing data with visualizations

## *Correlation* between continuous variables

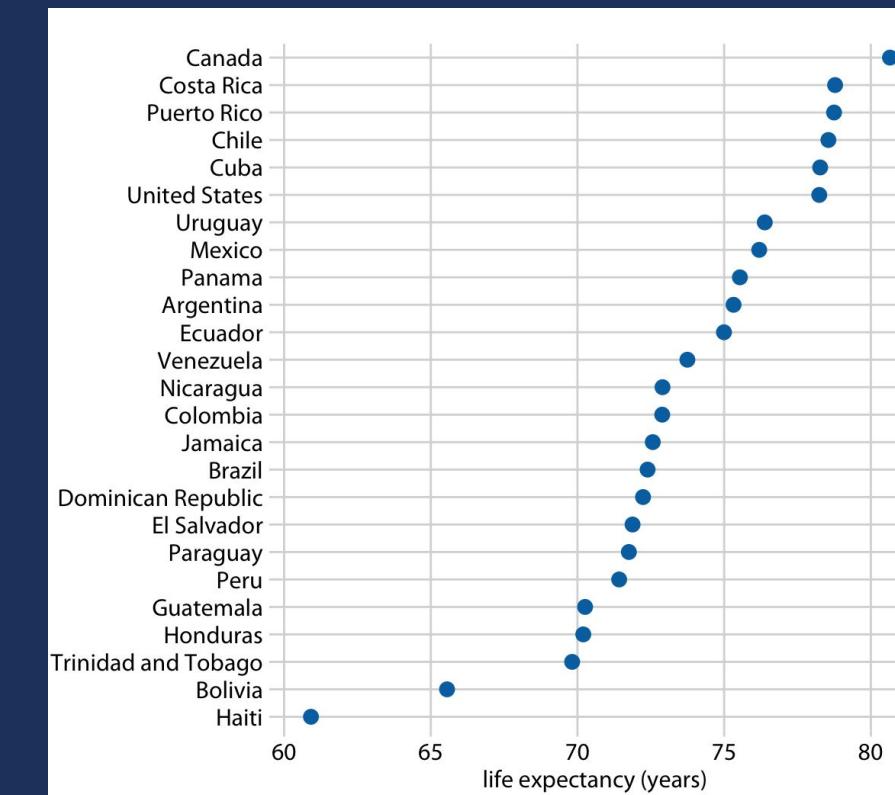
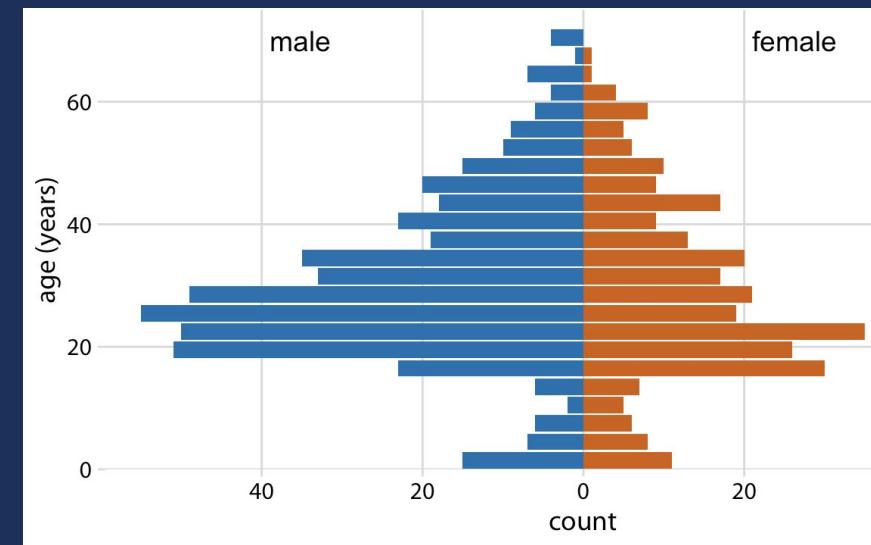
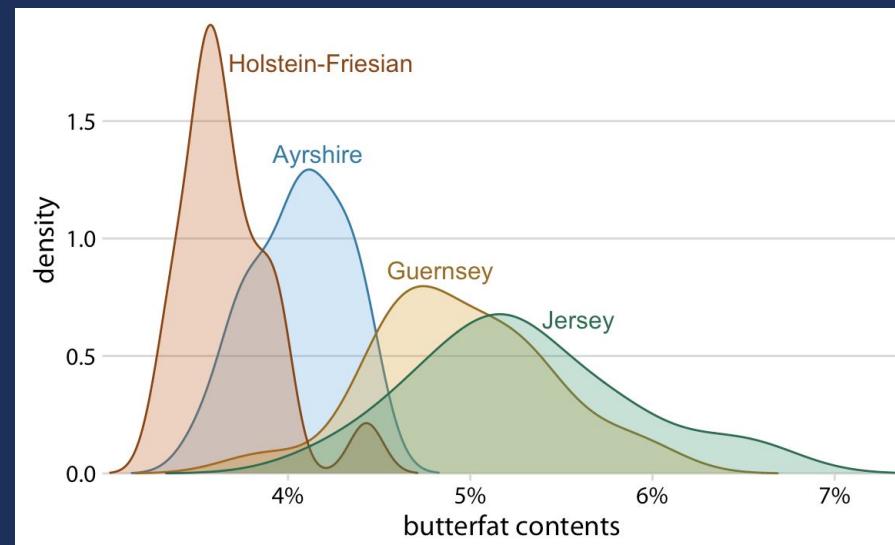


Correlogram

# Describing data with visualizations

## Relationship between continuous & categorical

### Numeric ~ Categorical

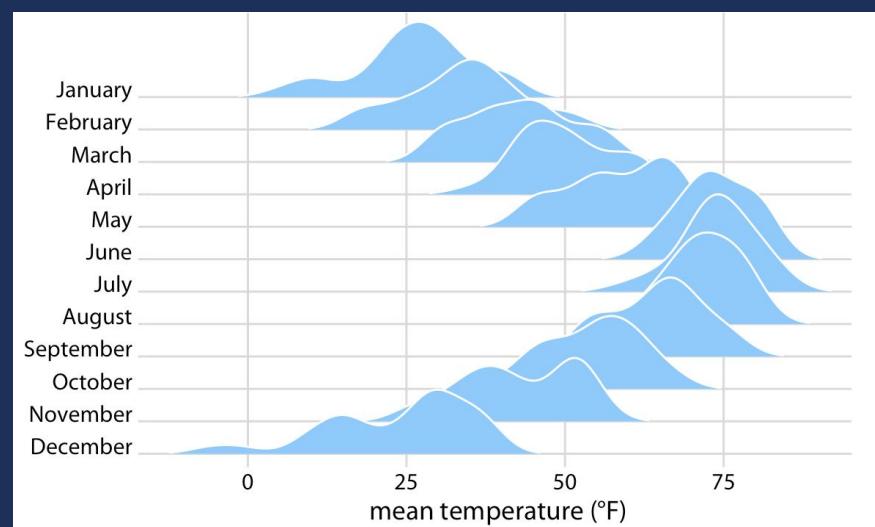
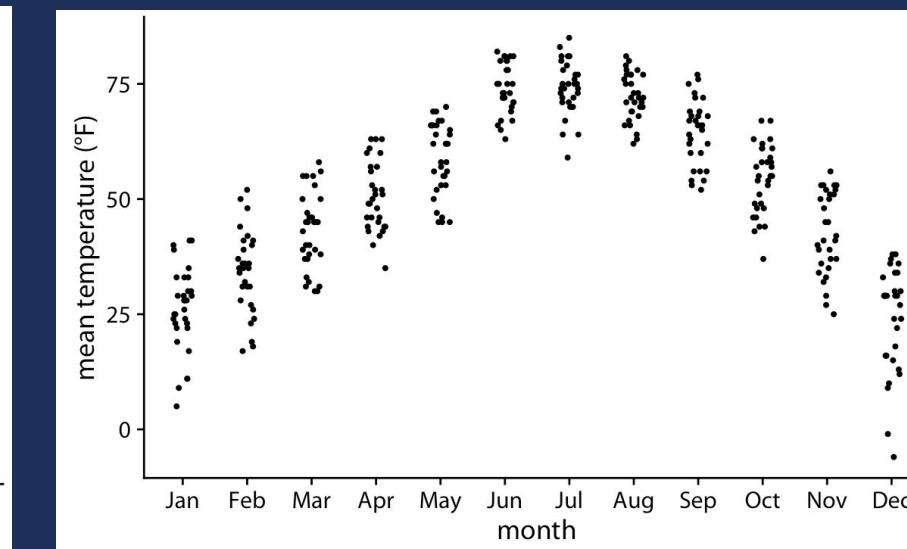
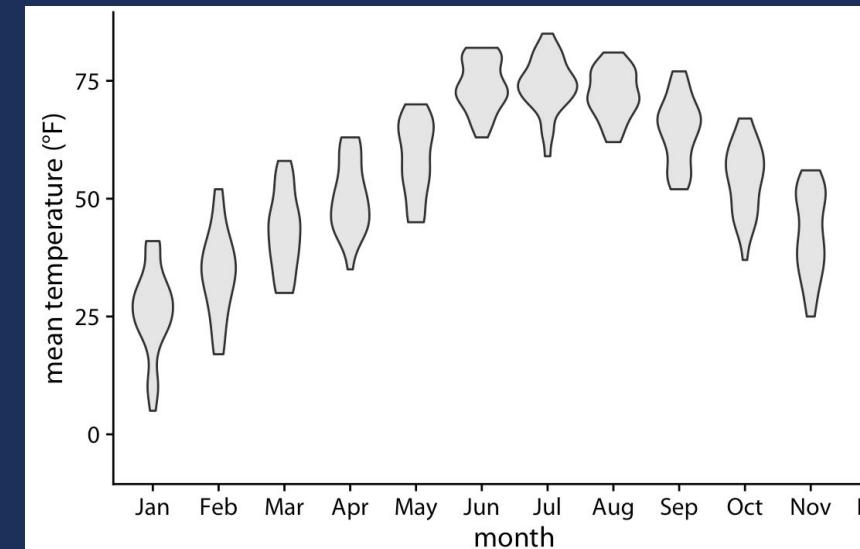
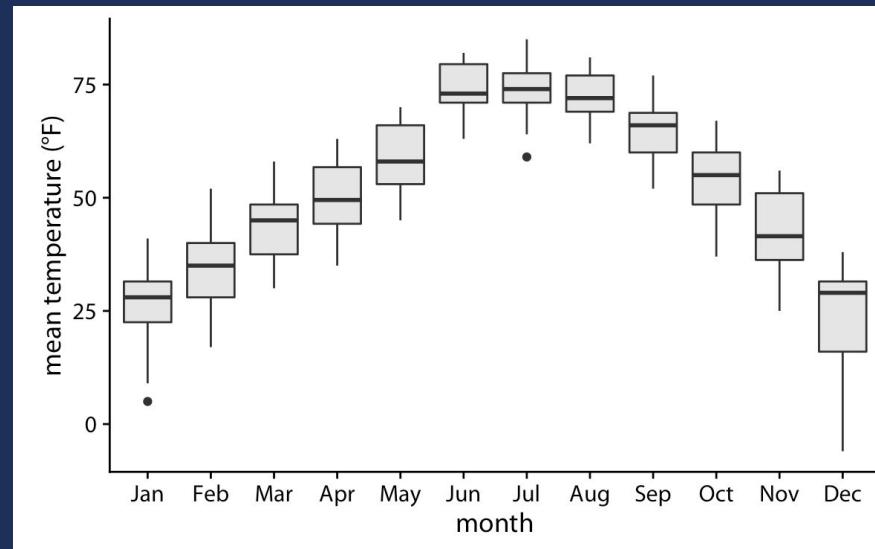


Dot plot

# Describing data with visualizations

## Relationship between two variables

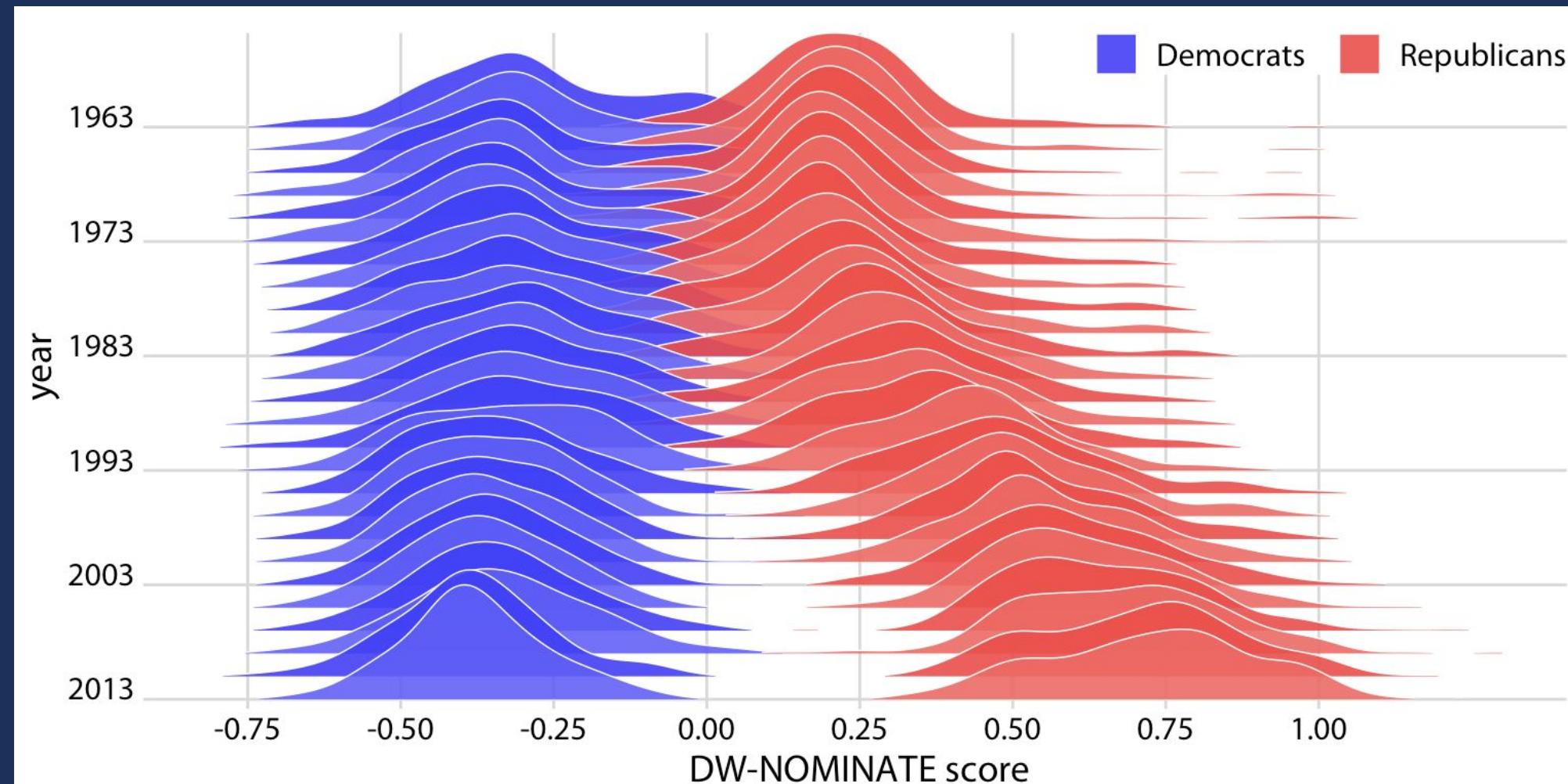
Numeric ~ Categorical



# Describing data with visualizations

---

Which types of relationships are these visualizations showing?

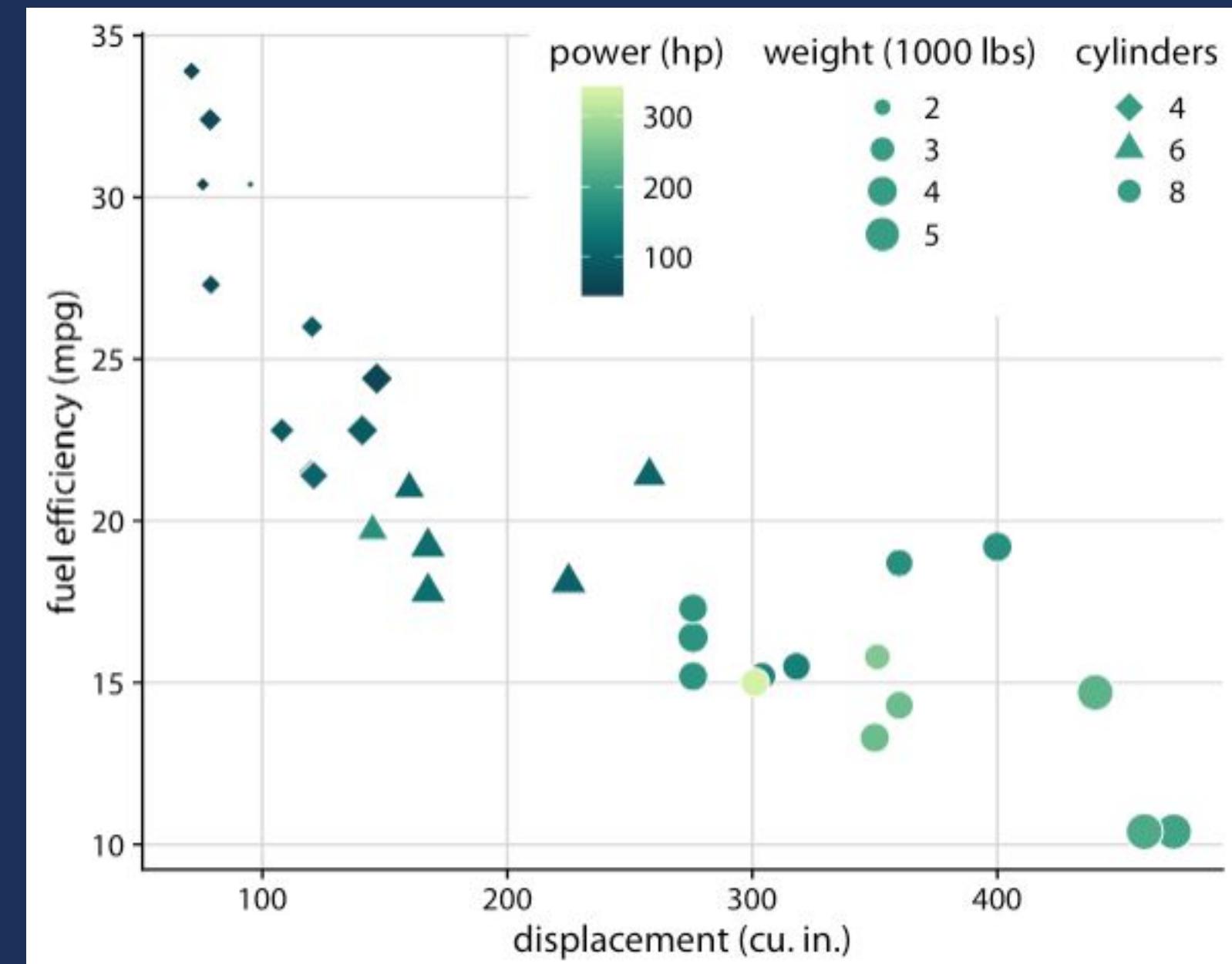


Ridgeline Plot

# Describing data with visualizations

## Relationship between three or more variables

How many variables do we have here, which types, and which aesthetics are mapped to them?

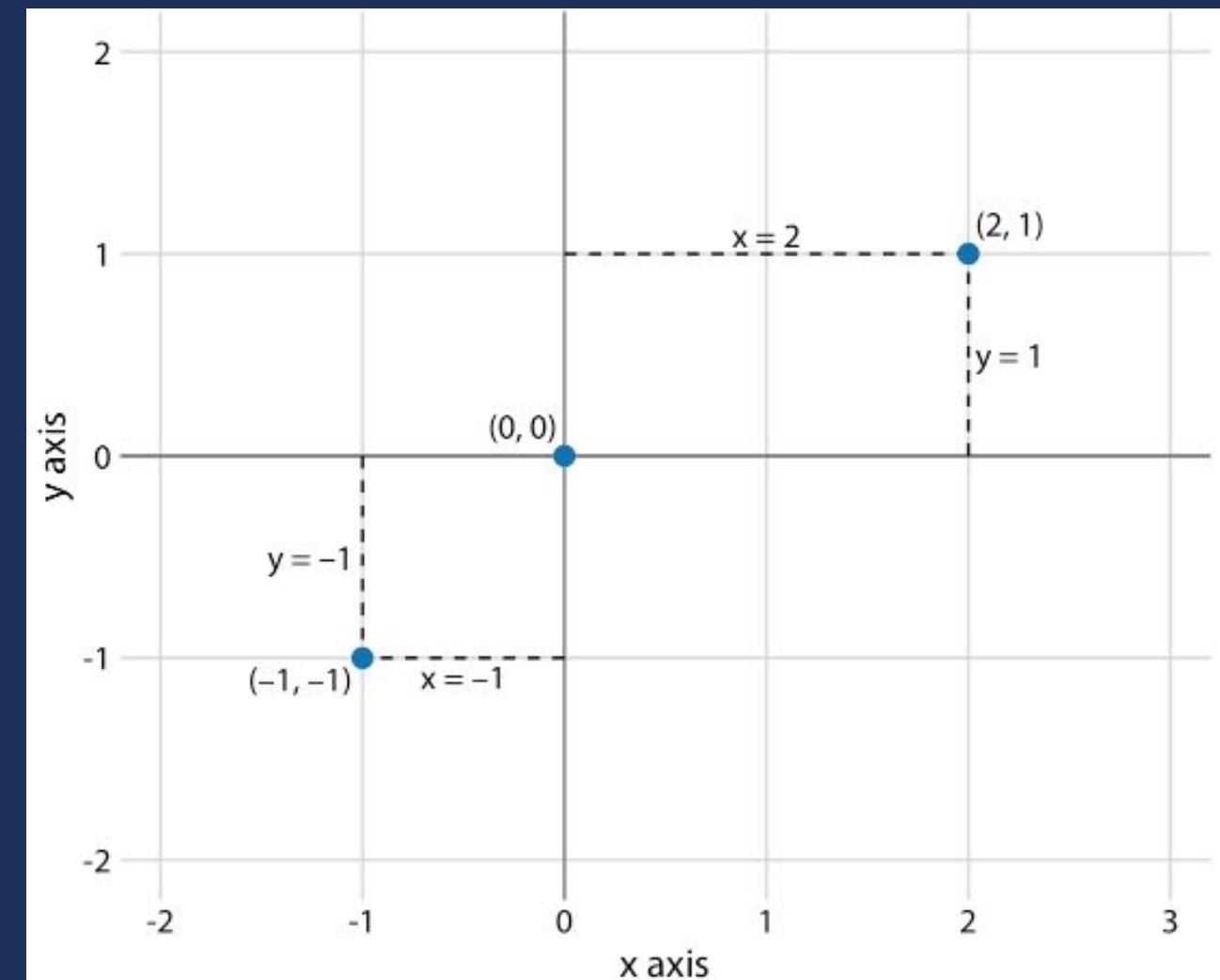


# Position

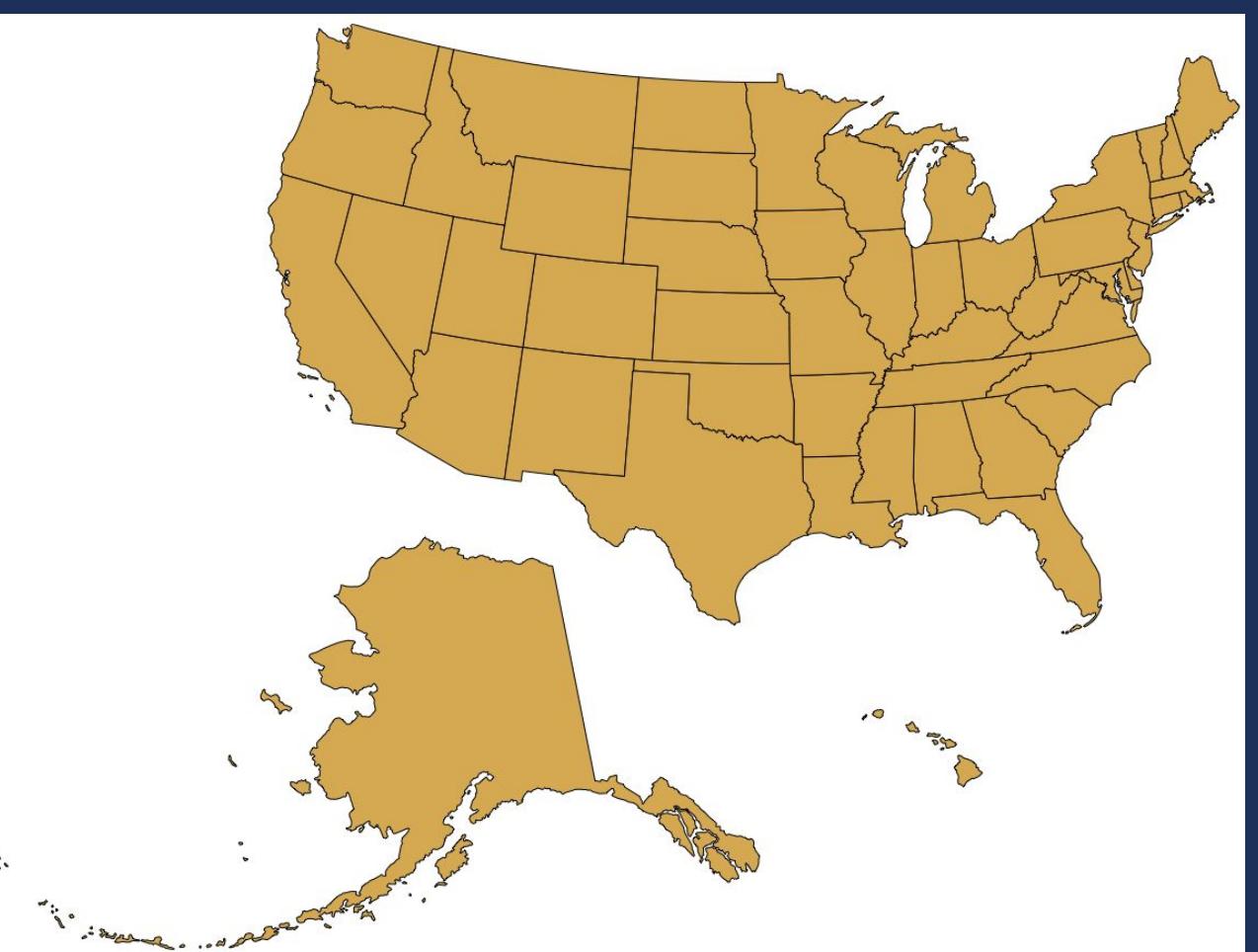
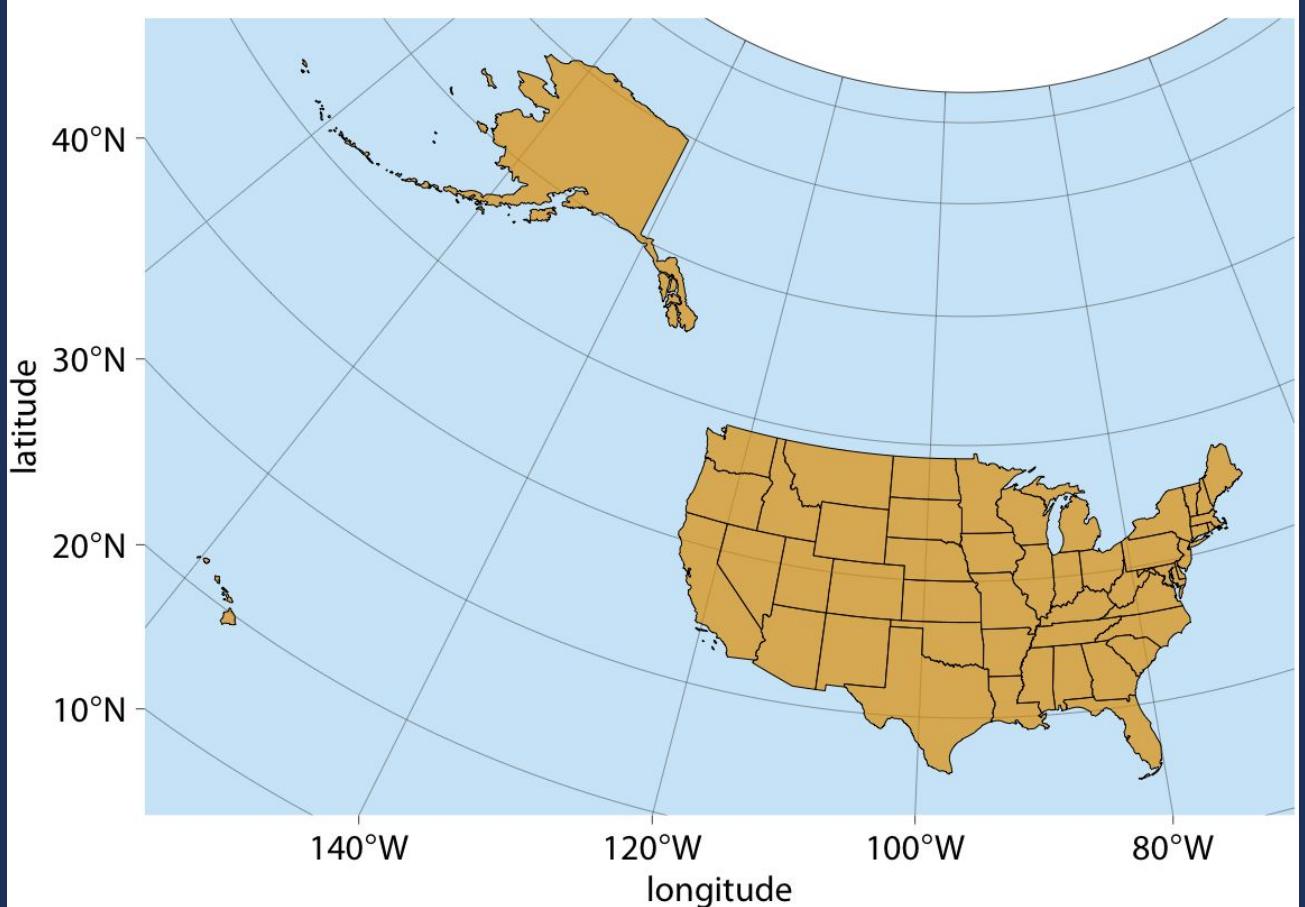
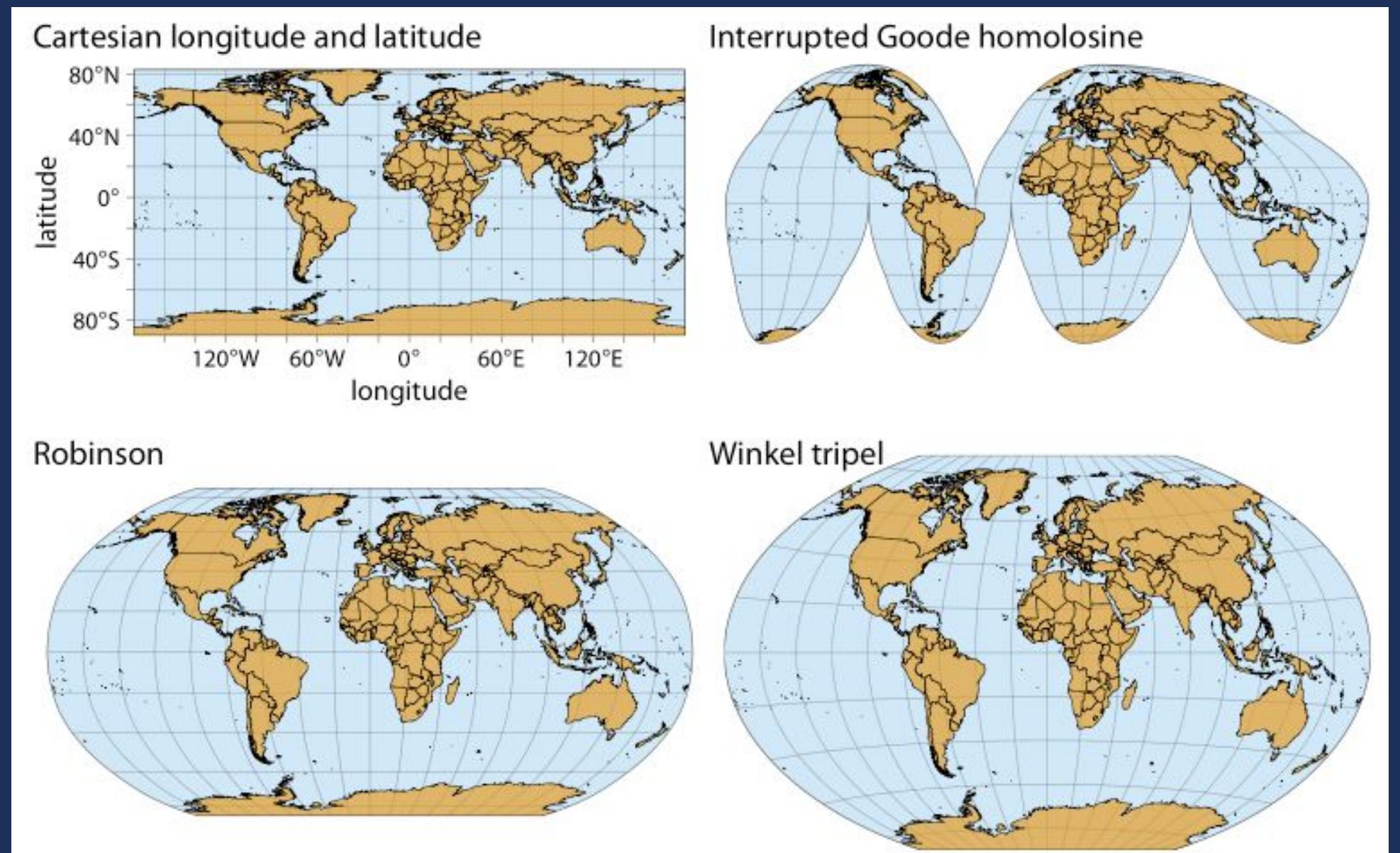
---

## Types of position:

- Cartesian (x/y)
- Radial/Polar (r, theta)
- Geographic (lat/long)

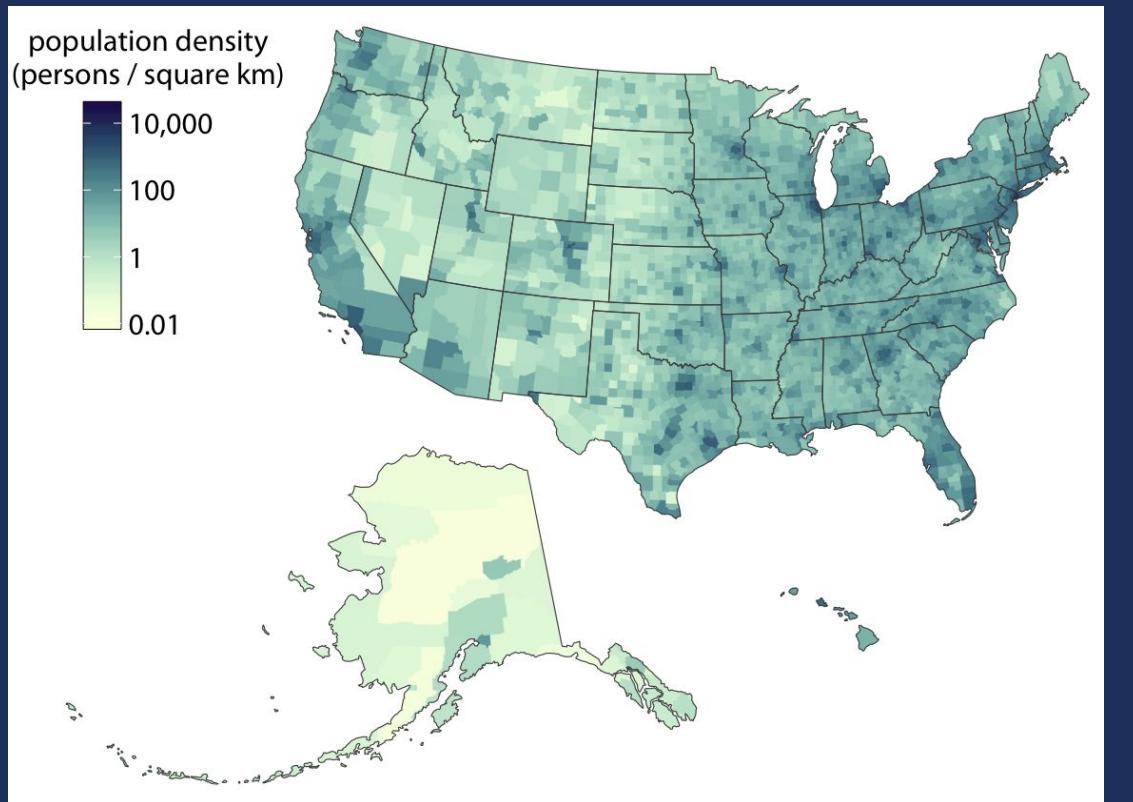


# Geospatial projections



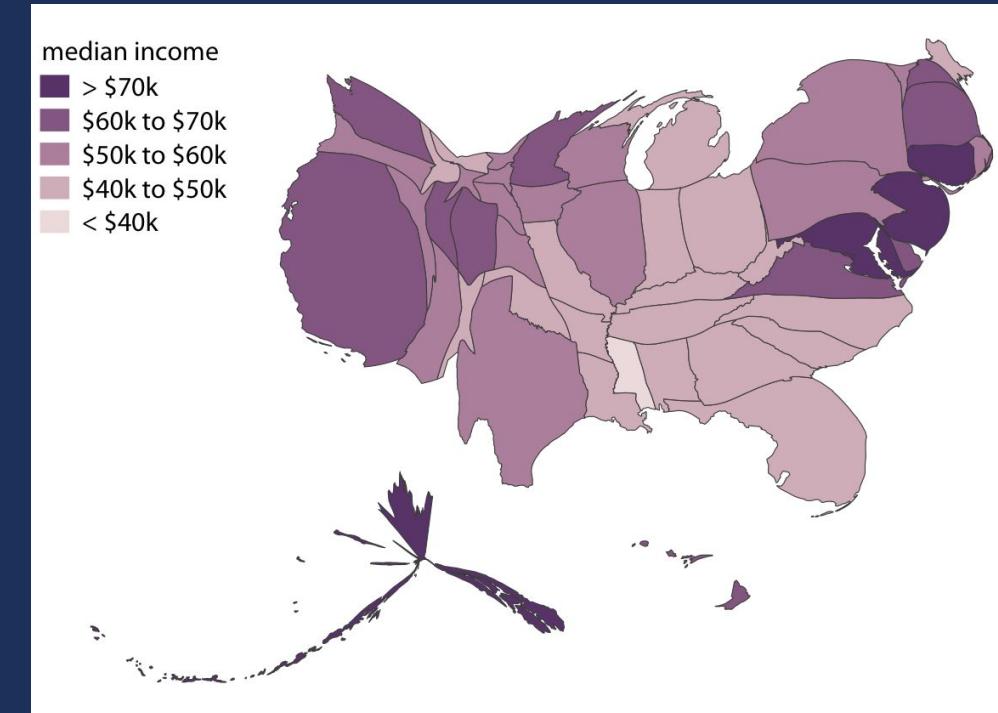
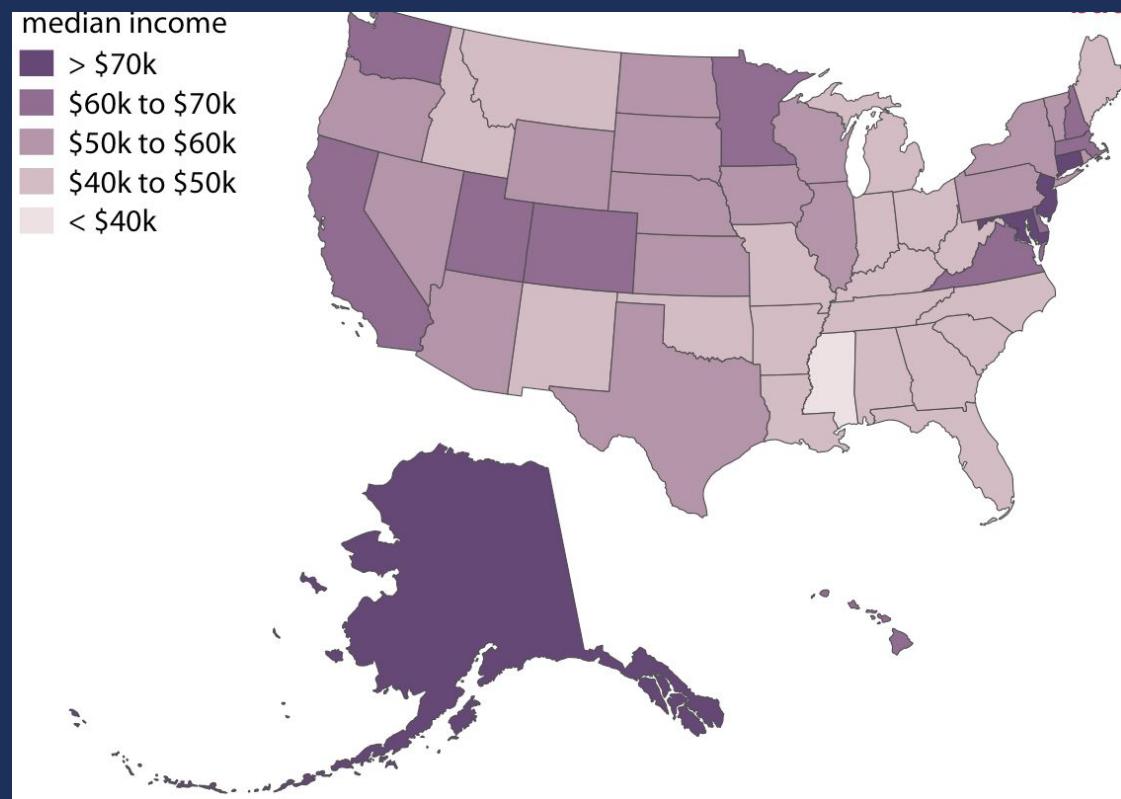
# Geospatial visualizations

---

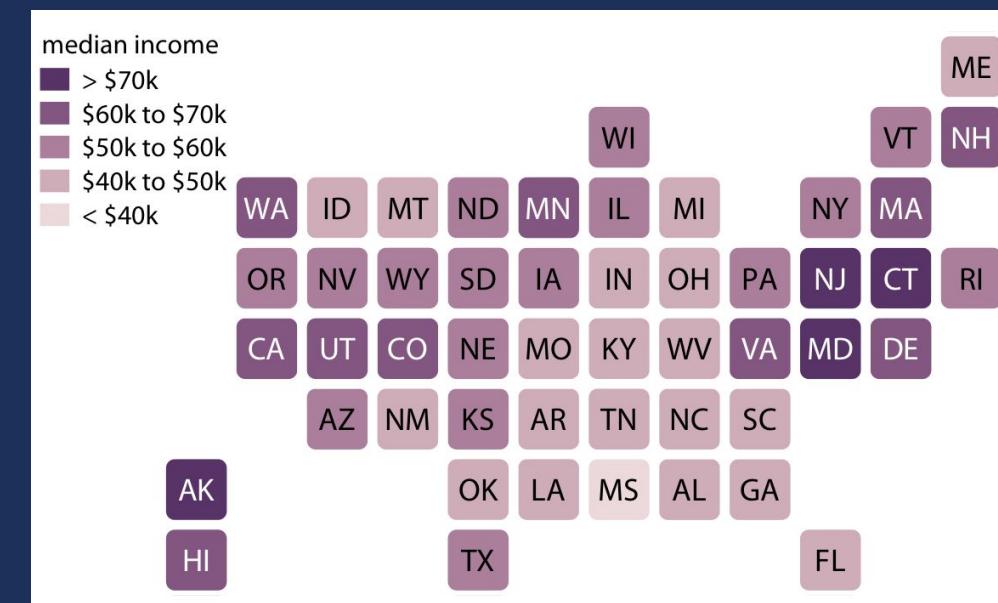


Choropleth

# Geospatial visualizations



cartogram



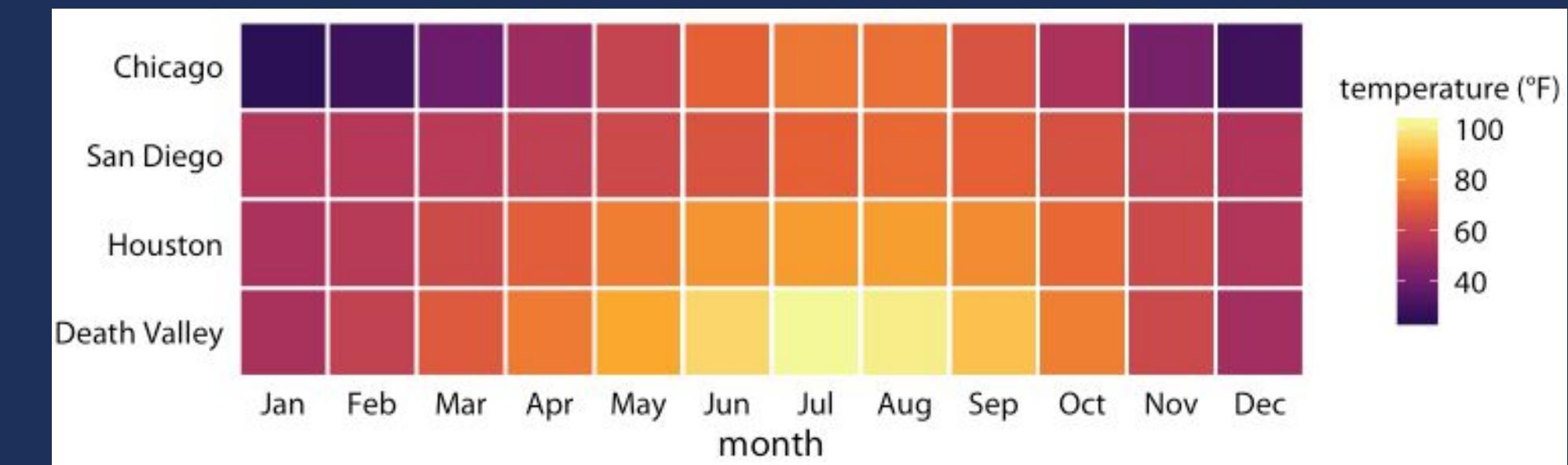
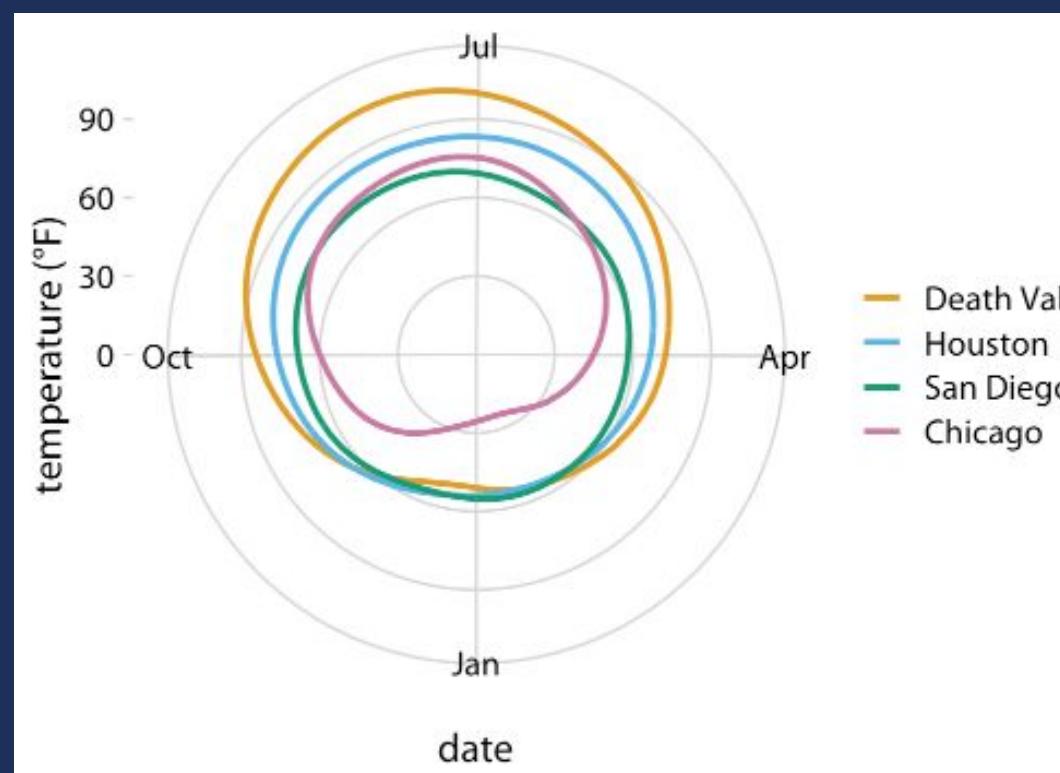
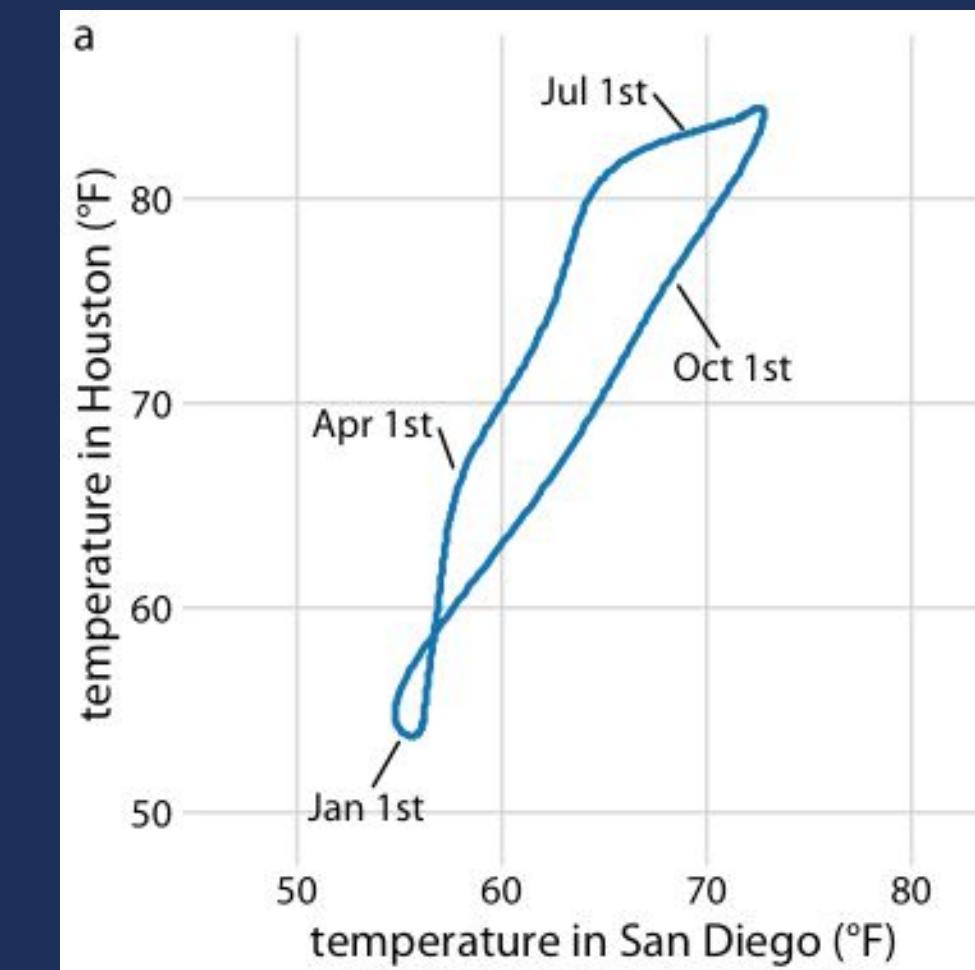
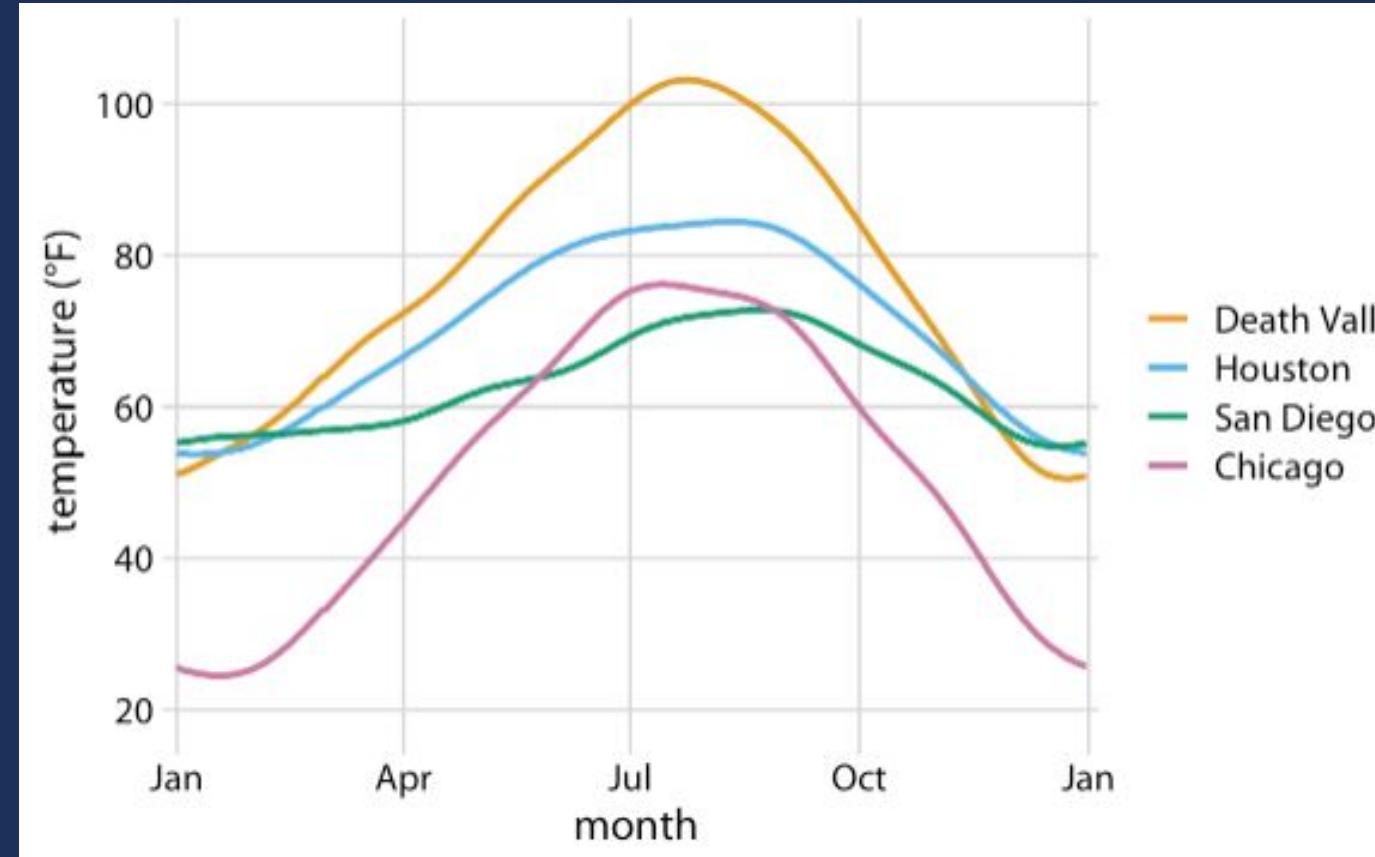
cartogram heatmap

# How might you represent this data?

---

Table 2.2: First 12 rows of a dataset listing daily temperature normals for four weather stations. Data source: NOAA.

Month	Day	Location	Station ID	Temperature
Jan	1	Chicago	USW00014819	25.6
Jan	1	San Diego	USW00093107	55.2
Jan	1	Houston	USW00012918	53.9
Jan	1	Death Valley	USC00042319	51.0
Jan	2	Chicago	USW00014819	25.5
Jan	2	San Diego	USW00093107	55.3
Jan	2	Houston	USW00012918	53.8
Jan	2	Death Valley	USC00042319	51.2
Jan	3	Chicago	USW00014819	25.3
Jan	3	San Diego	USW00093107	55.3
Jan	3	Death Valley	USC00042319	51.3
Jan	3	Houston	USW00012918	53.8



# Color

---

Think about what colors you might use to represent:

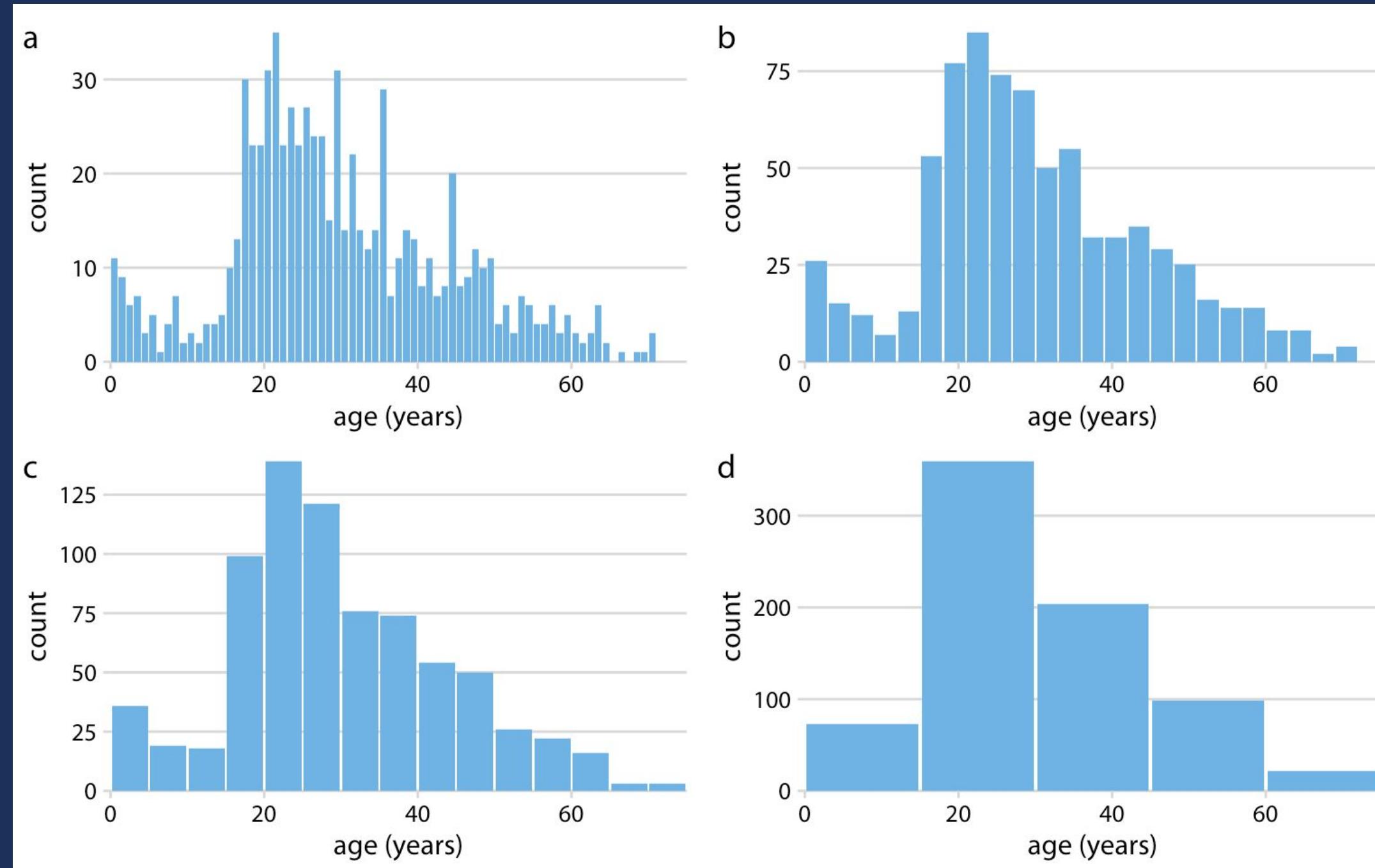
- A numeric variable (whether continuous or discrete)
- A nominal (unordered) categorical variable
- An ordinal categorical variable

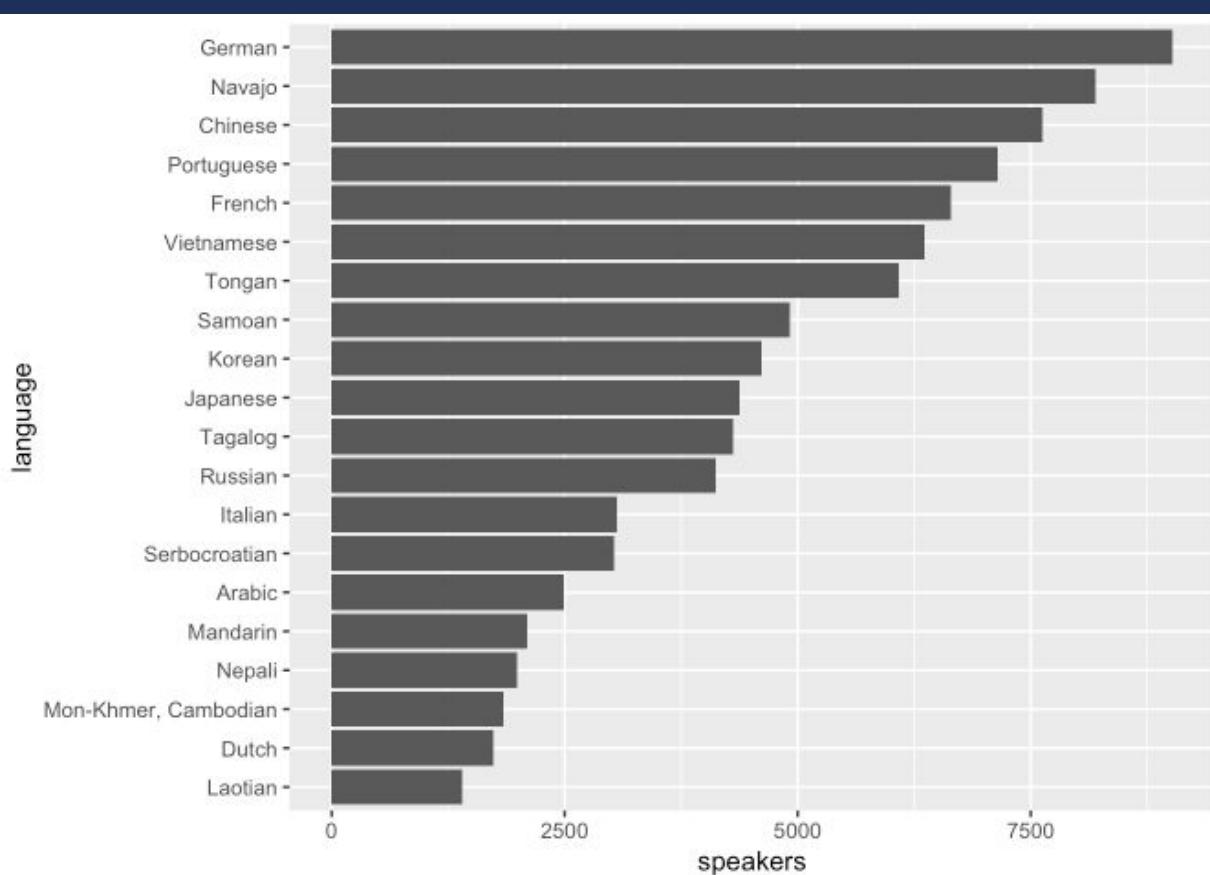
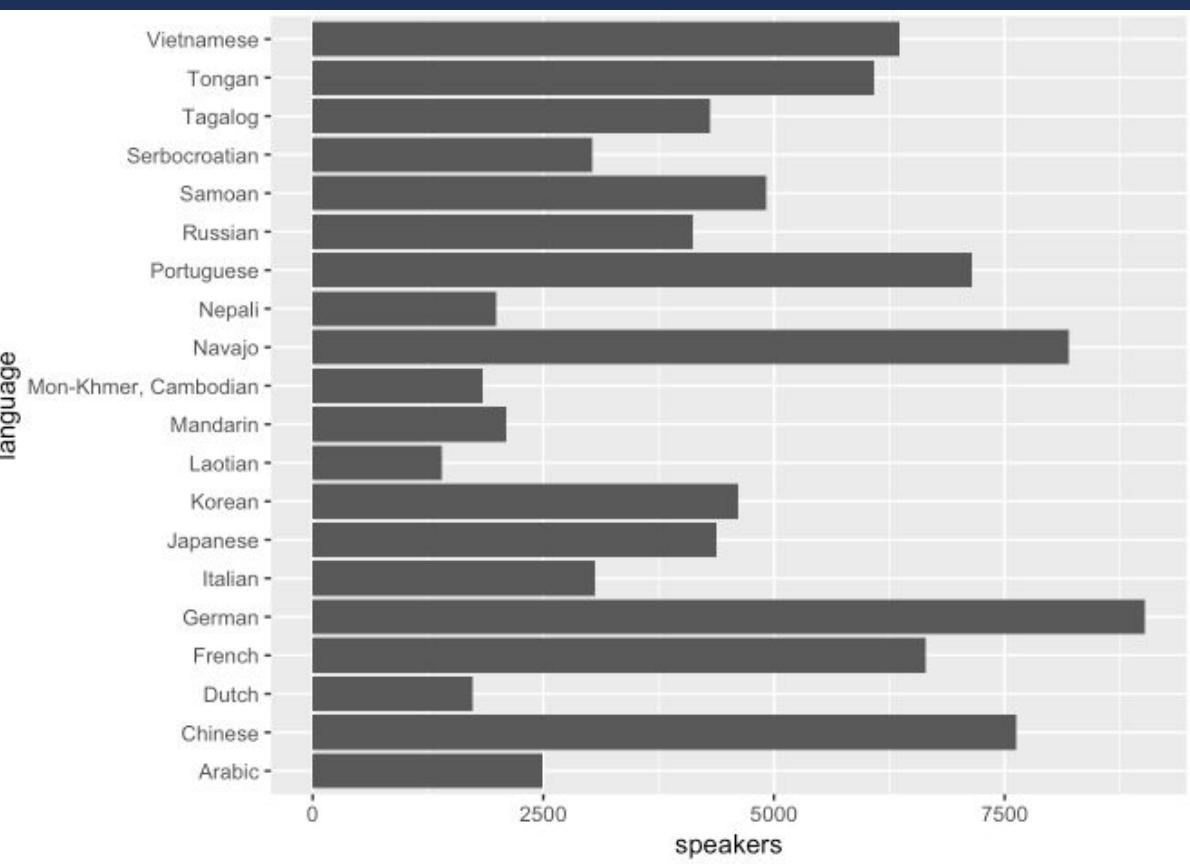
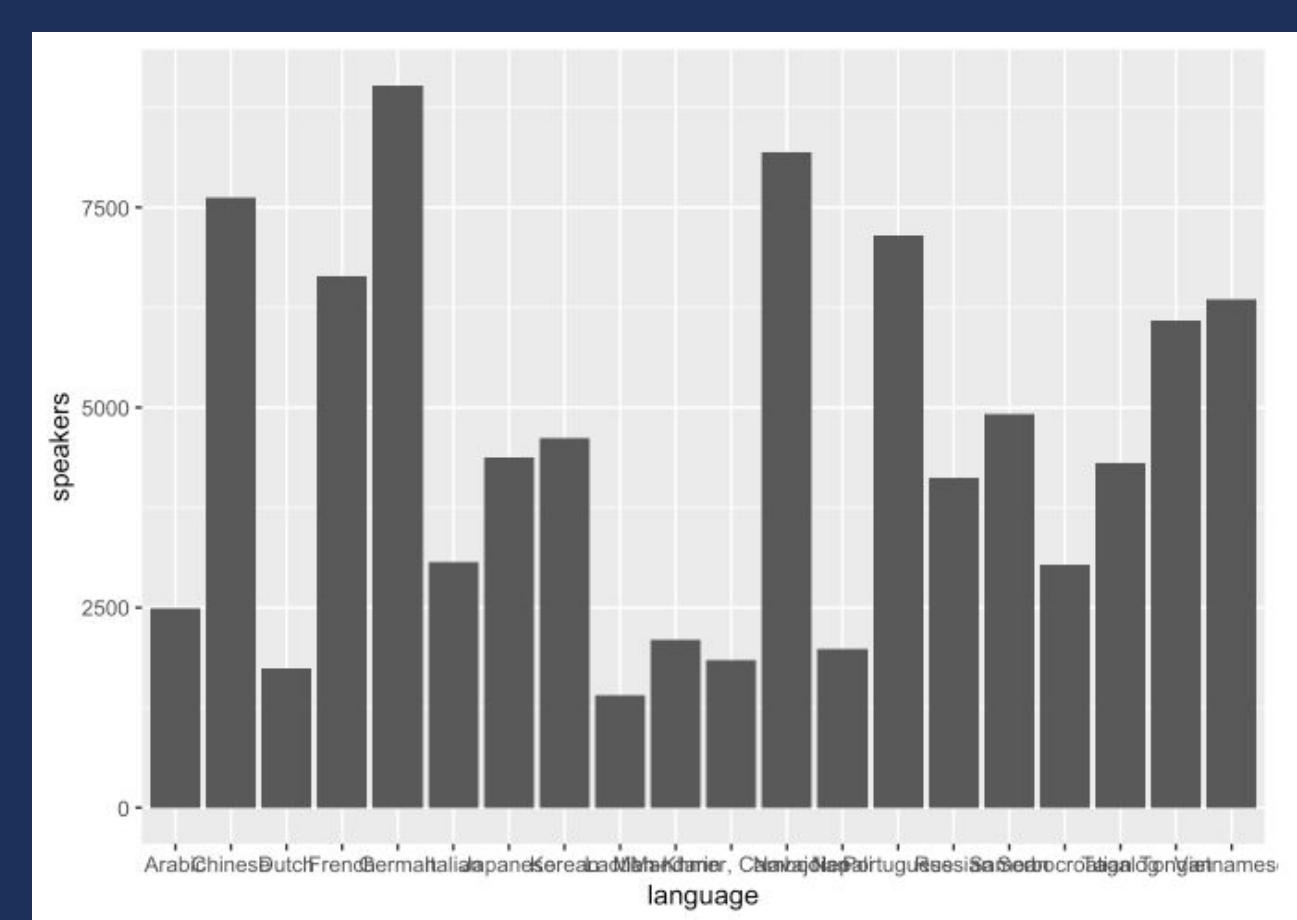


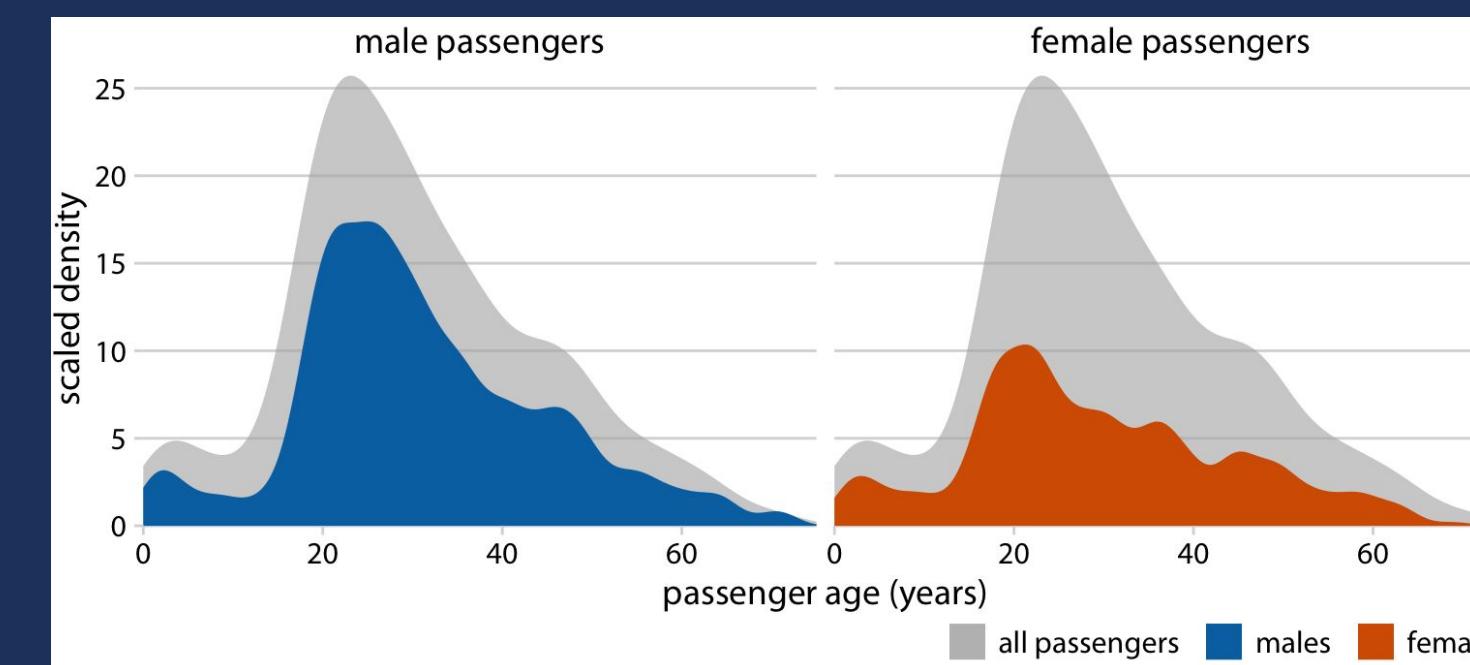
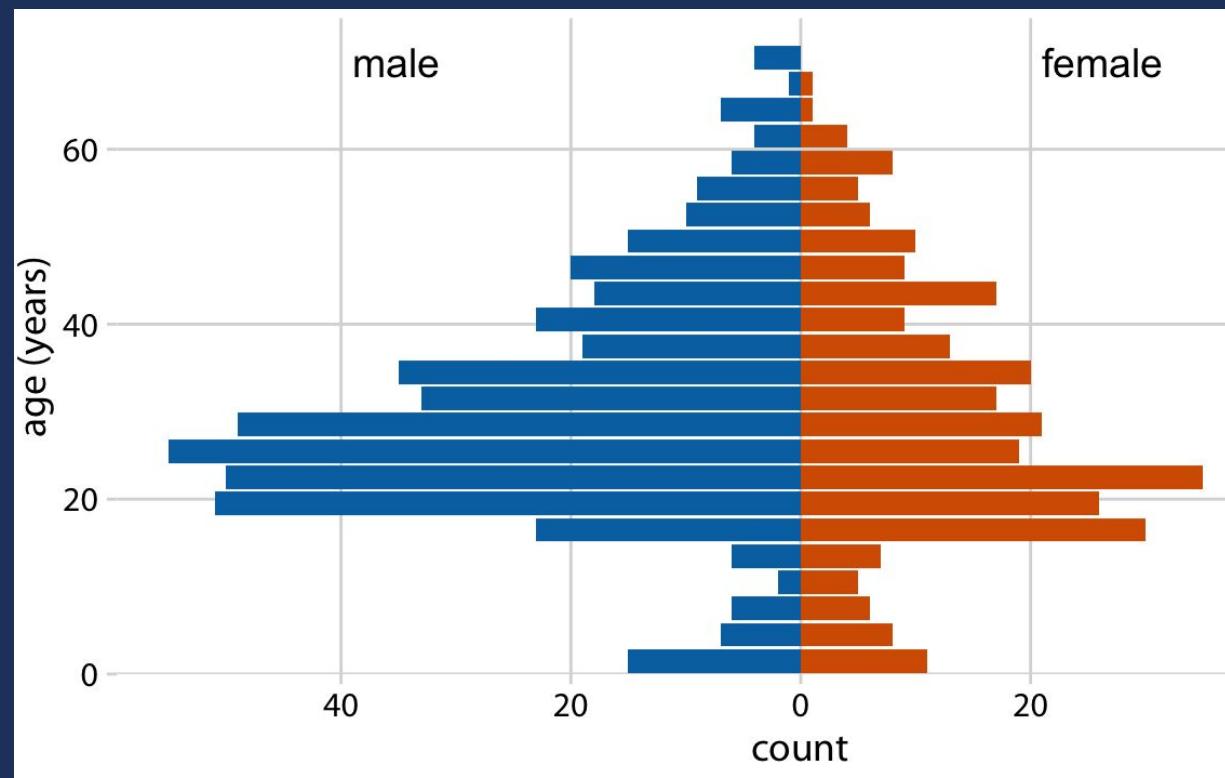
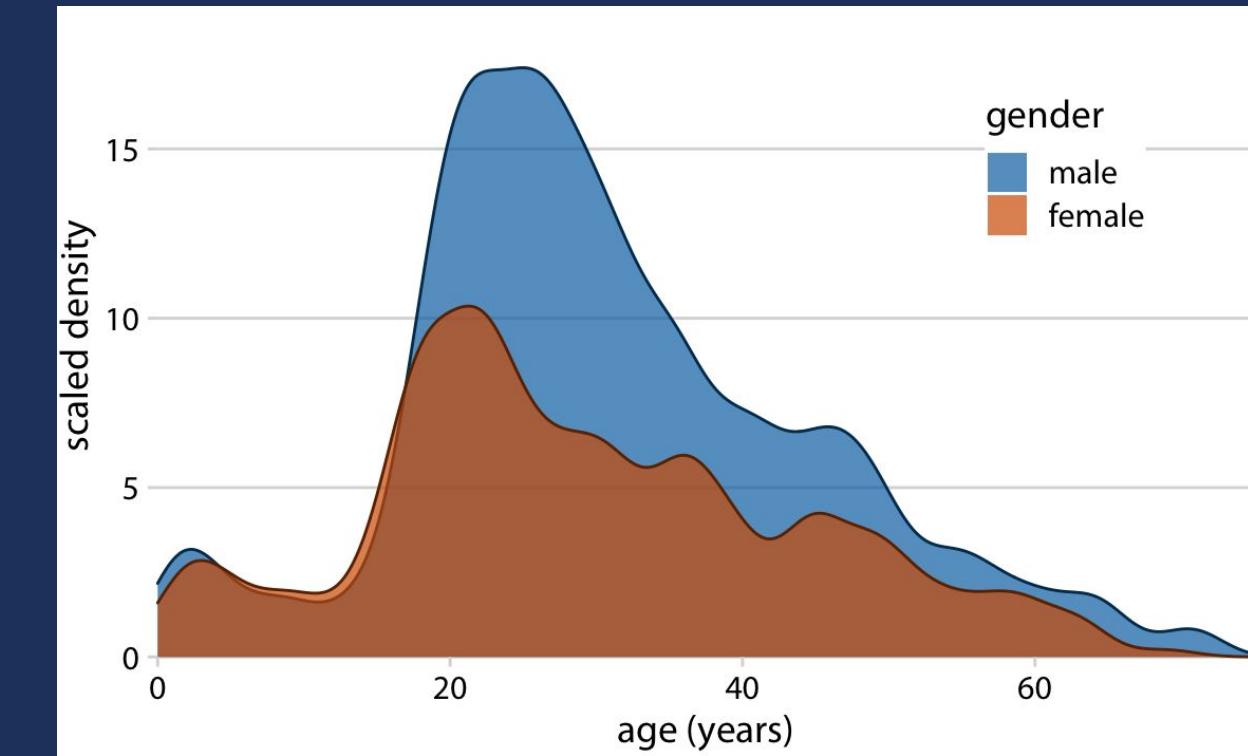
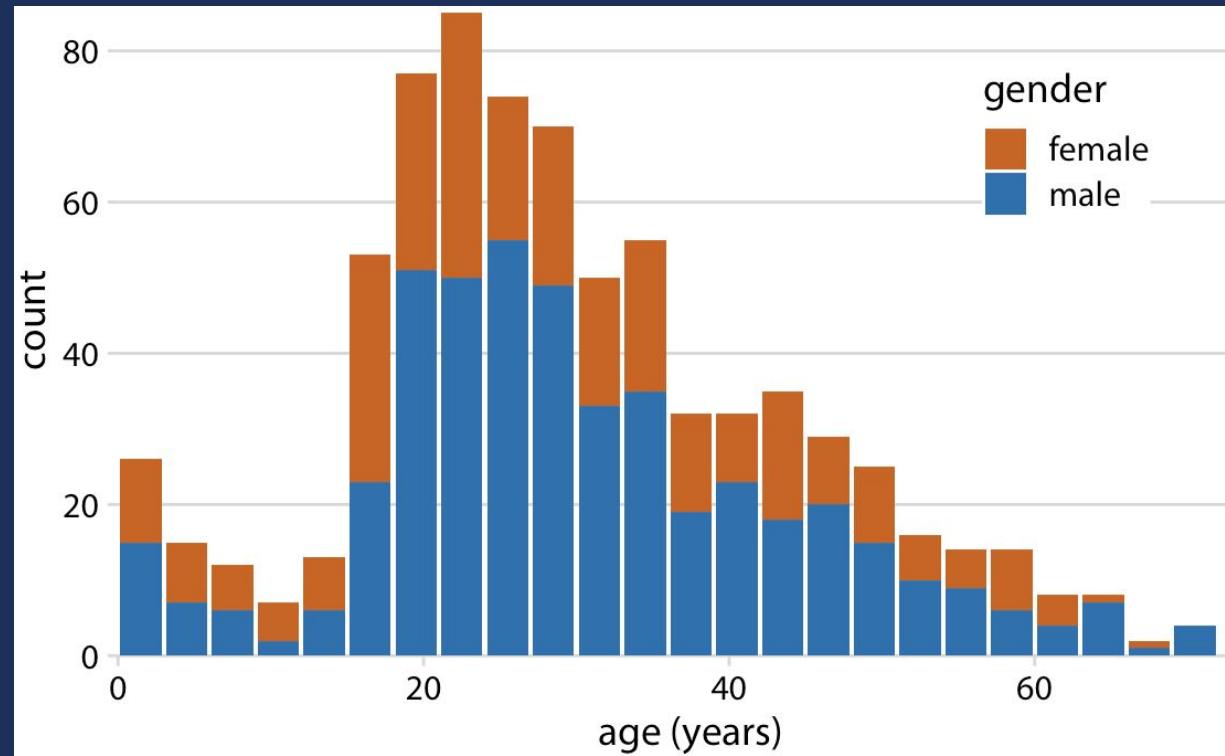
# WORKING THROUGH EXAMPLES

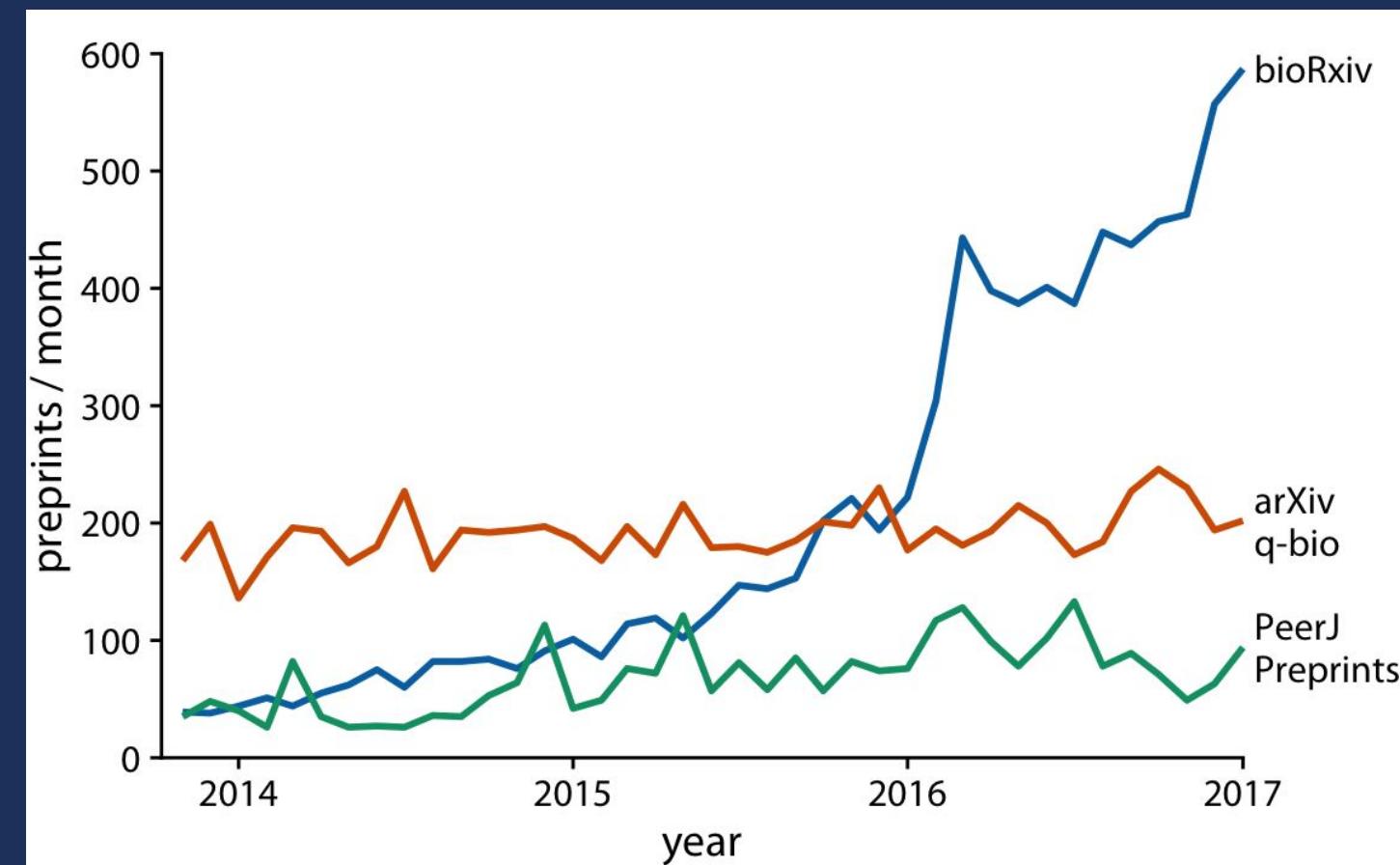
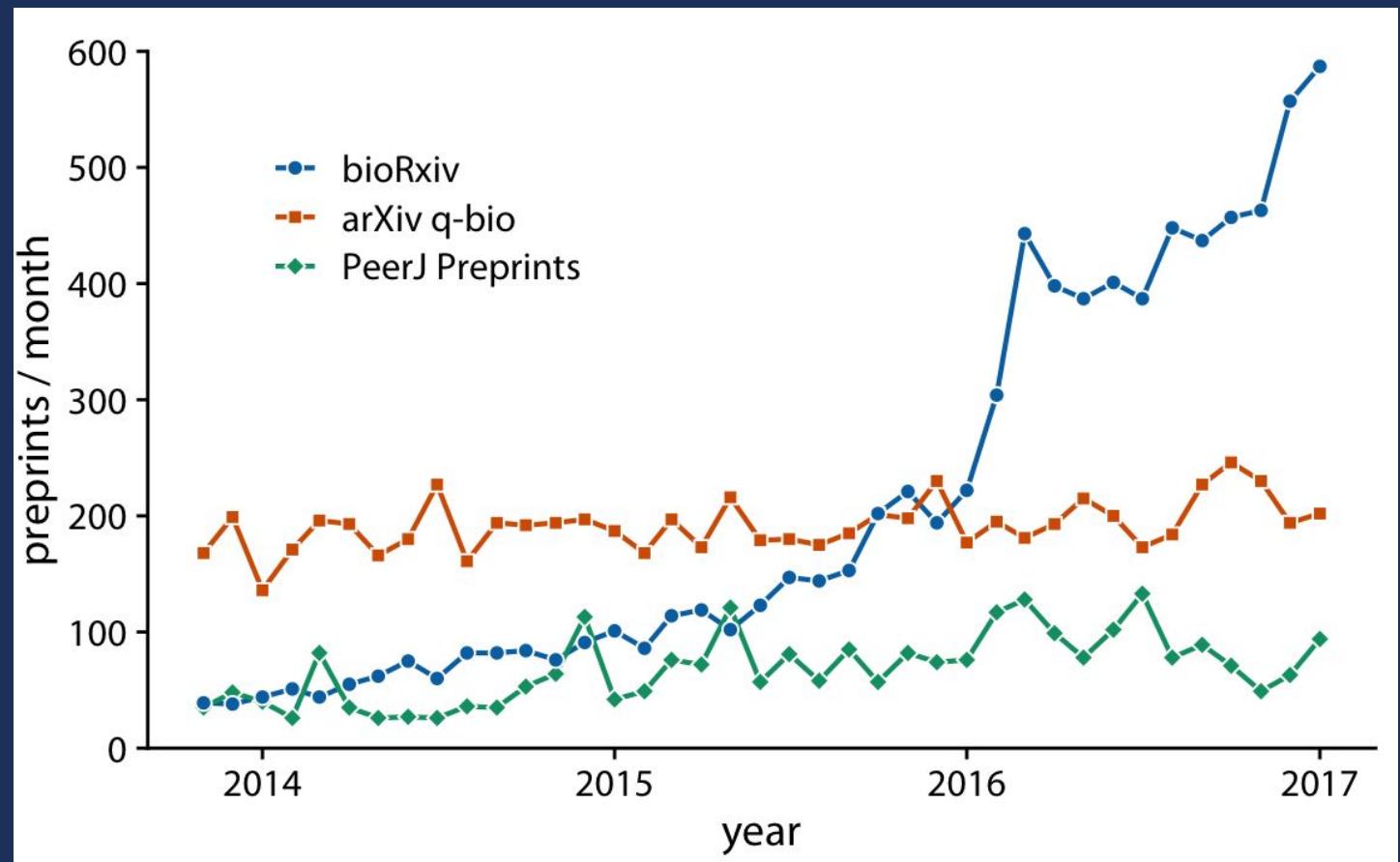
---

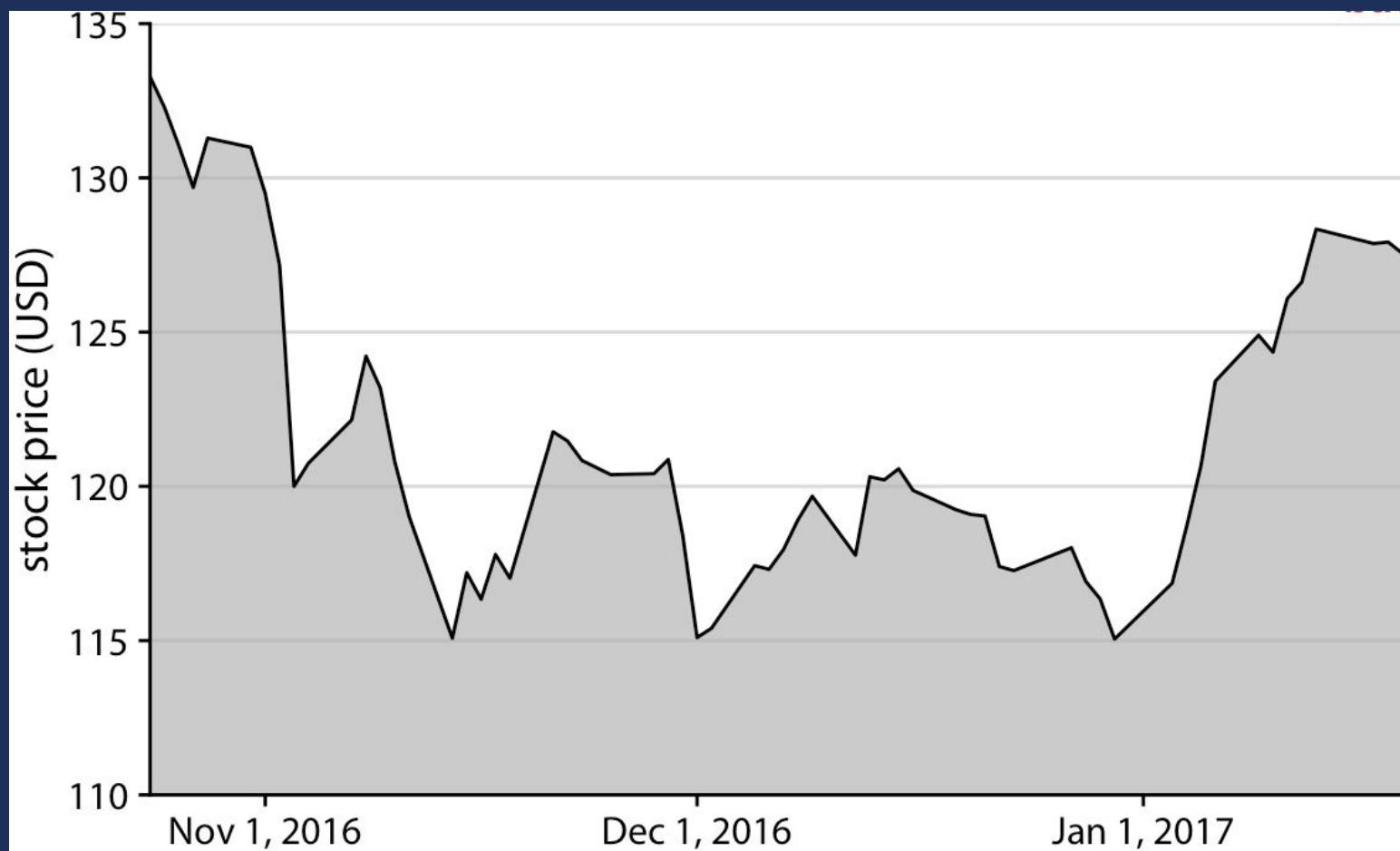
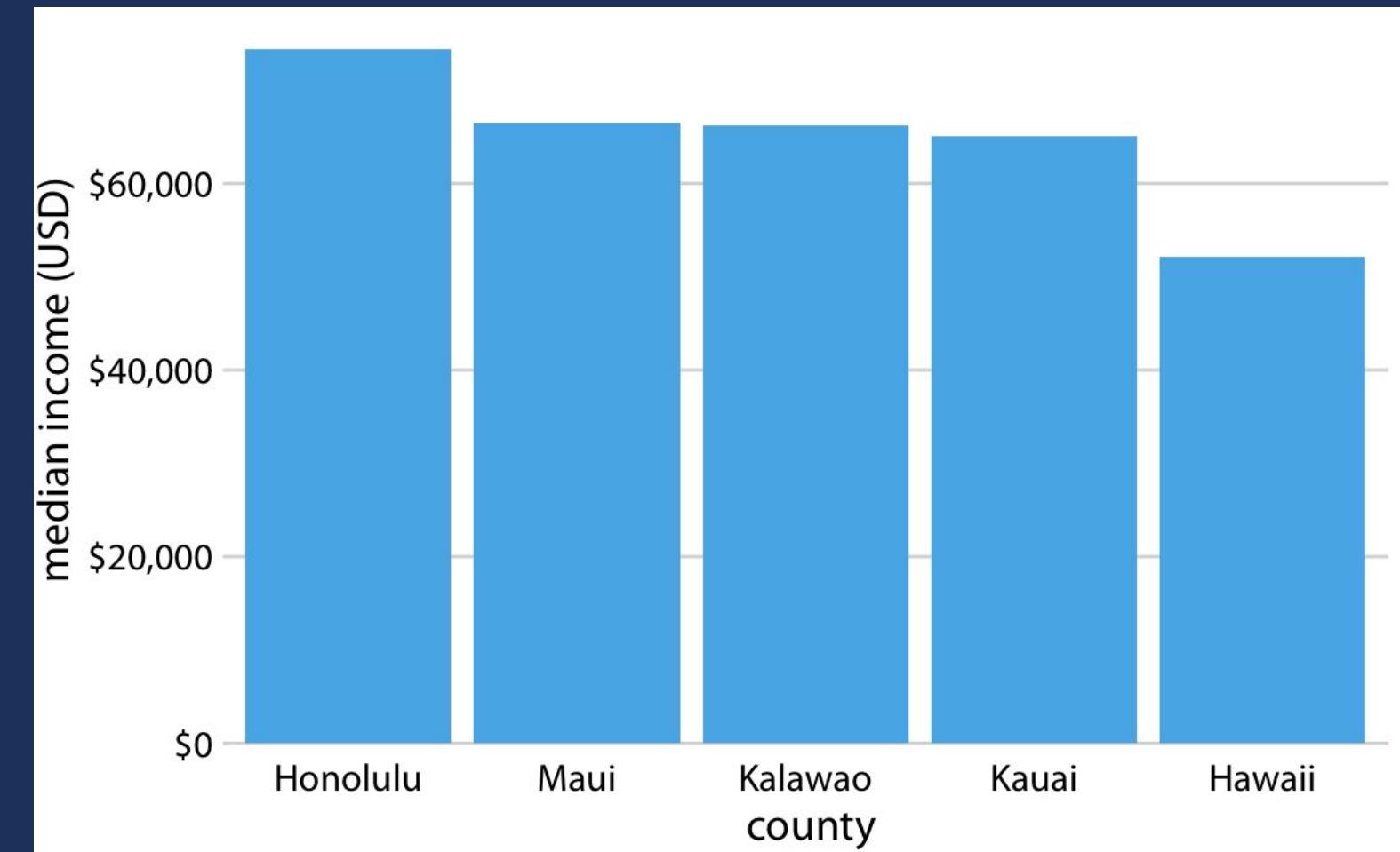
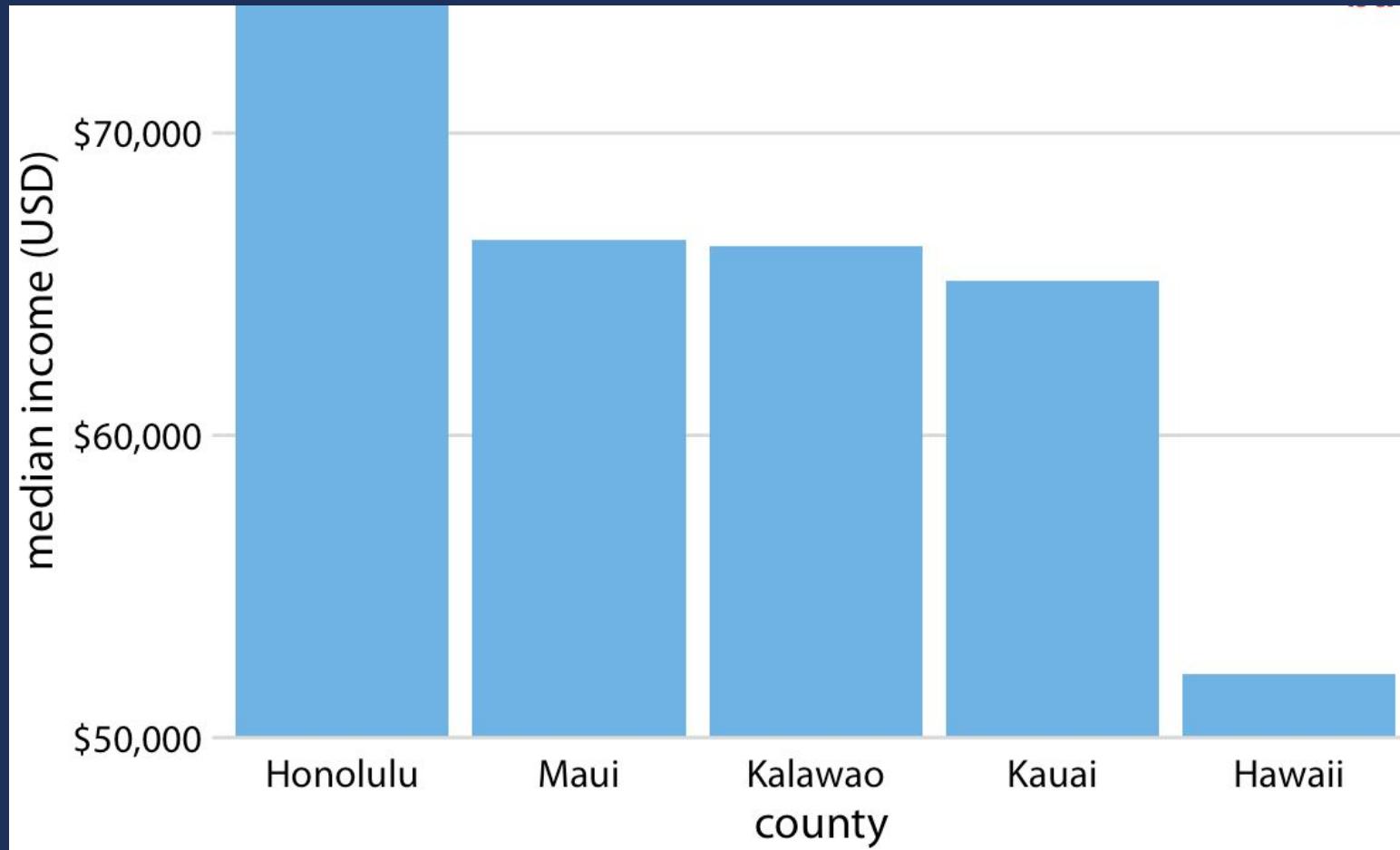
# Good, Bad, or Ugly?











# A WORD ABOUT FILE FORMATS

---

Bitmap (raster) formats:

- JPEG, GIF, TIFF, BMP, ...

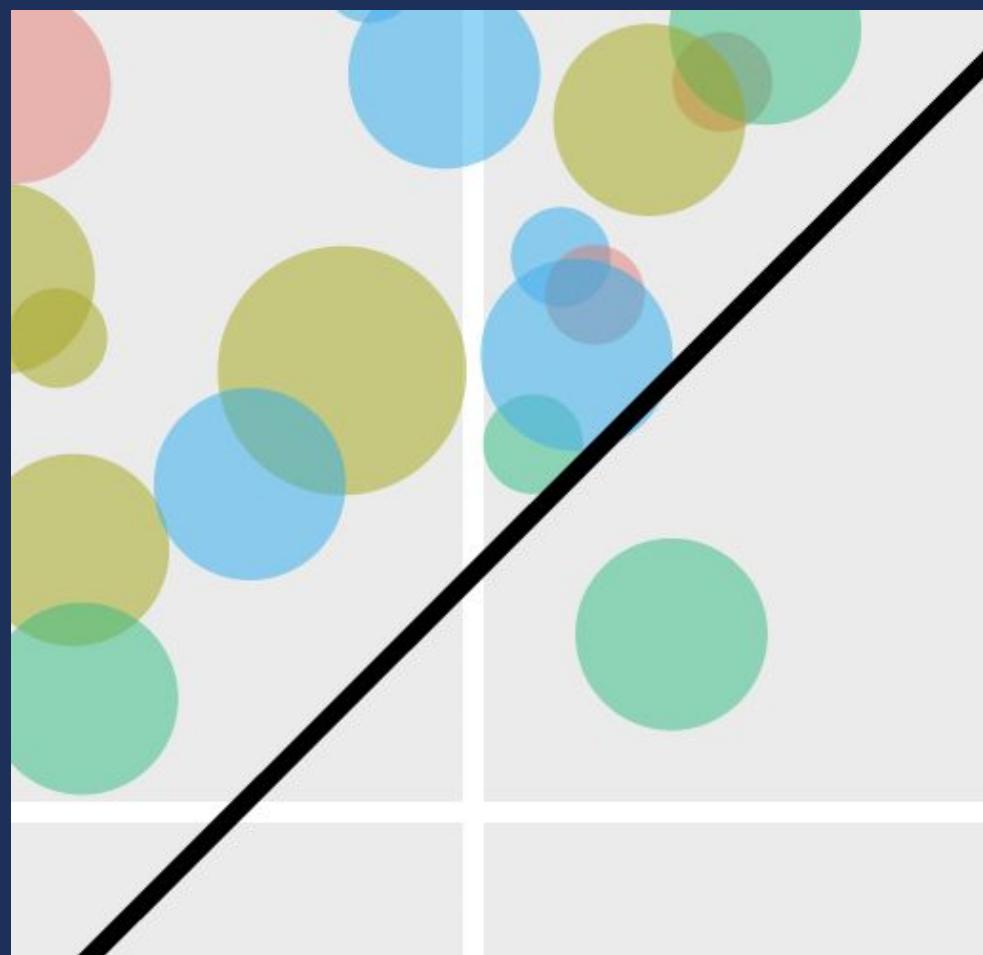
Vector formats:

- PDF, SVG, EPS

Be deliberate about height/width (aspect ratio)!

# JPEG vs. PDF

---



# WHICH VISUALIZATION SOFTWARE TO USE?

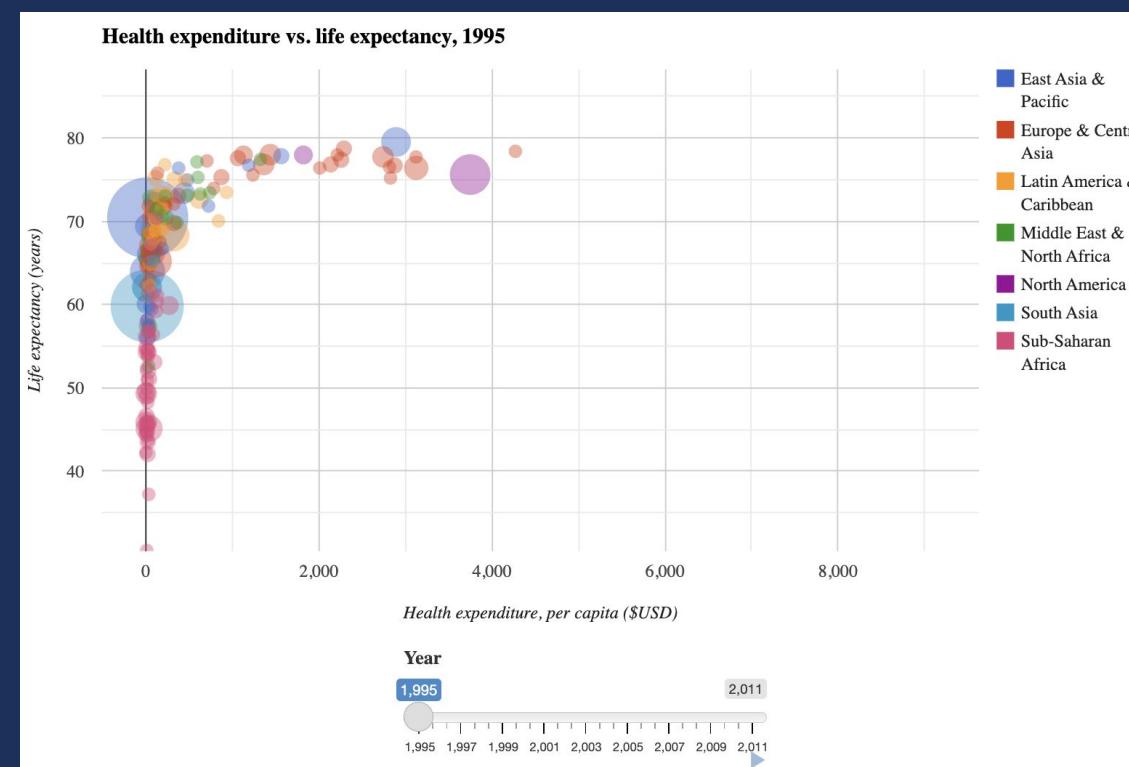
---

Consider:

- Reproducibility
  - Same data + same transformations → Substantially similar visualization  
(conveys the same message)
- Repeatability
  - Same data + same transformations, EXACTLY the same visualization

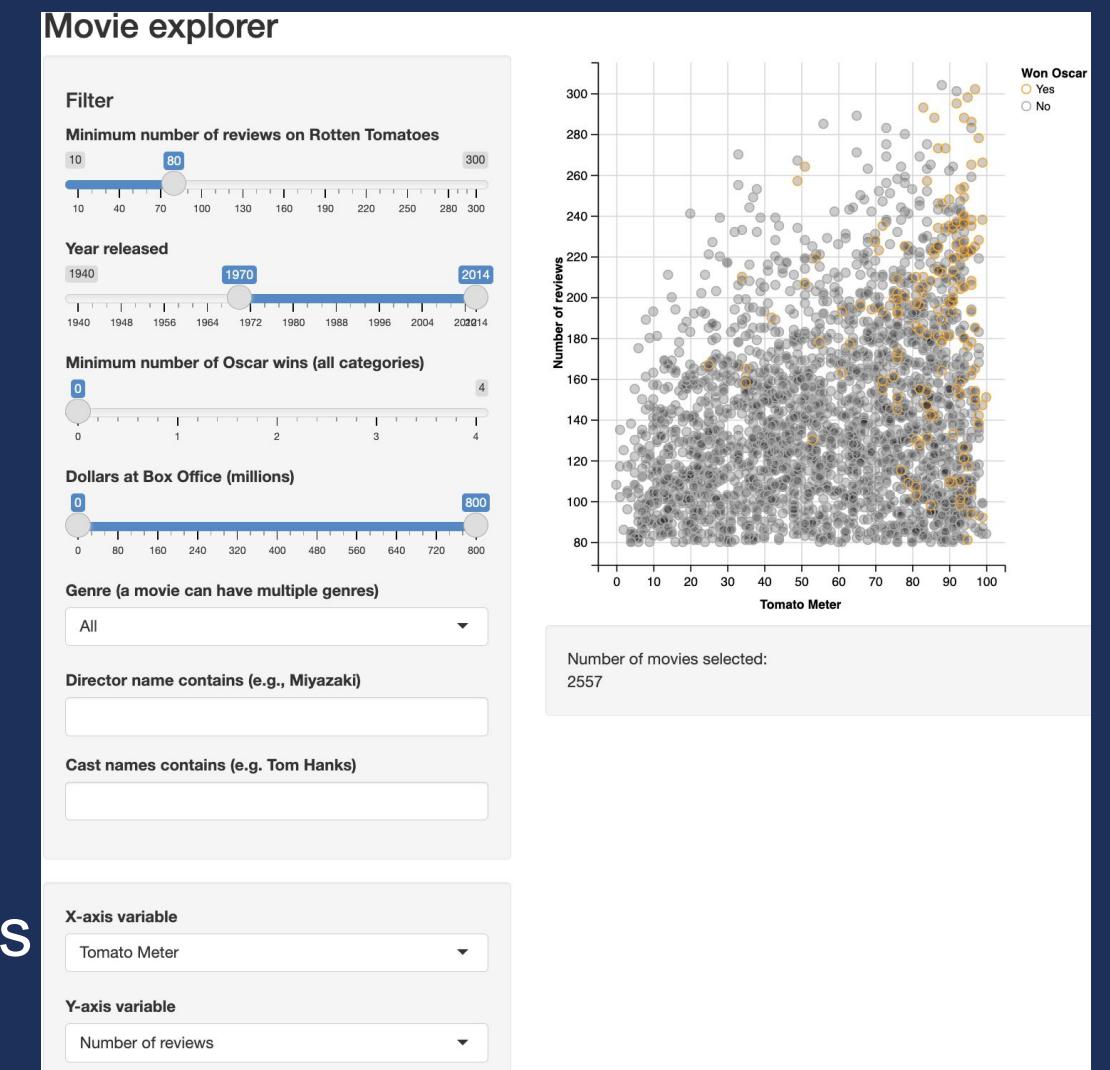
# INTERACTIVE VISUALIZATIONS

Users can engage with data to manipulate & explore data



Slide through time

Explore different relationships



# WORKSHOP OUTLINE

.....

Our content today is divided into four parts. Each part will be described with examples.

01

## Data Visualization Process

Walking through the visualization process, from setting goals to visualizing the data.

02

## Basic Design Considerations

Best practices and accessibility considerations in data visualization.

03

## Solutions to Common Problems

How to resolve, or even better, avoid common problems in data visualization.

04

## Visualization Resources

Workshops, tools, and other resources for data visualization.

# SOLUTIONS TO SOME PROBLEMS

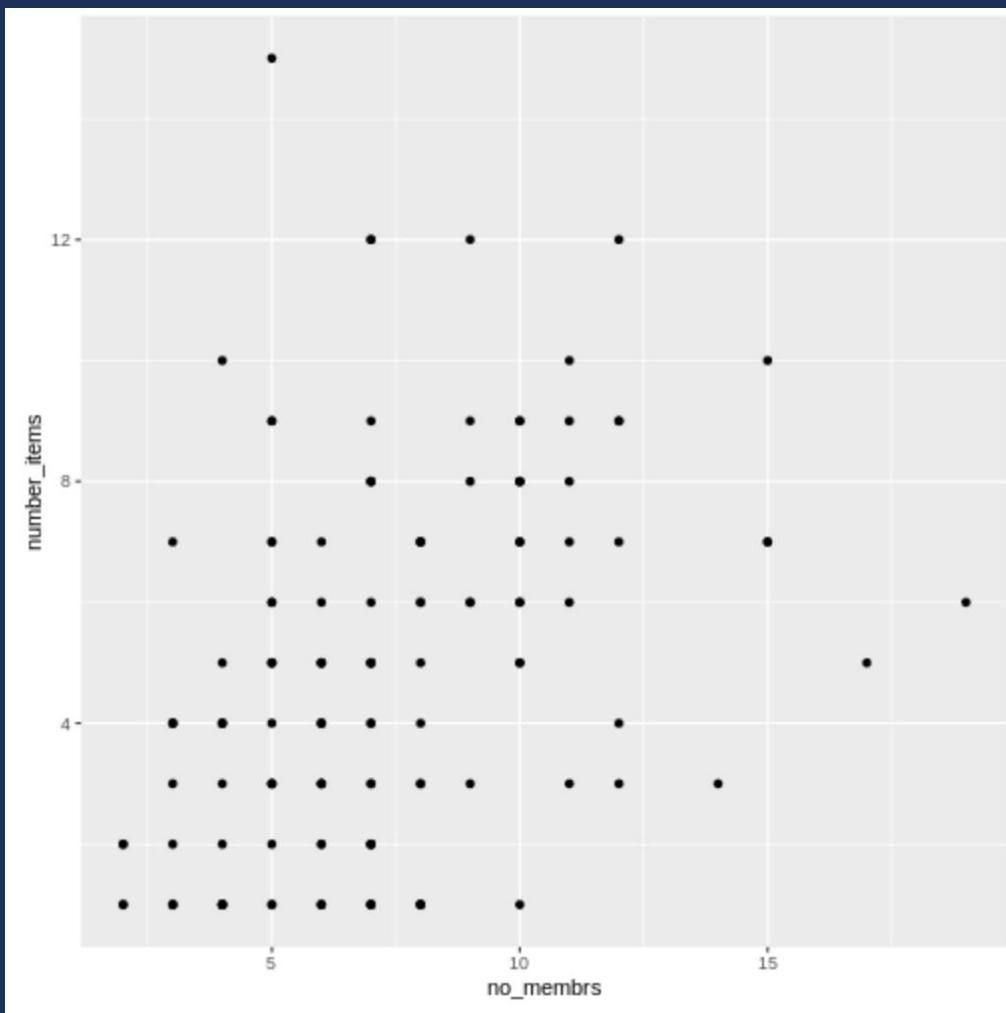
---

Problem #1: Overlapping data points

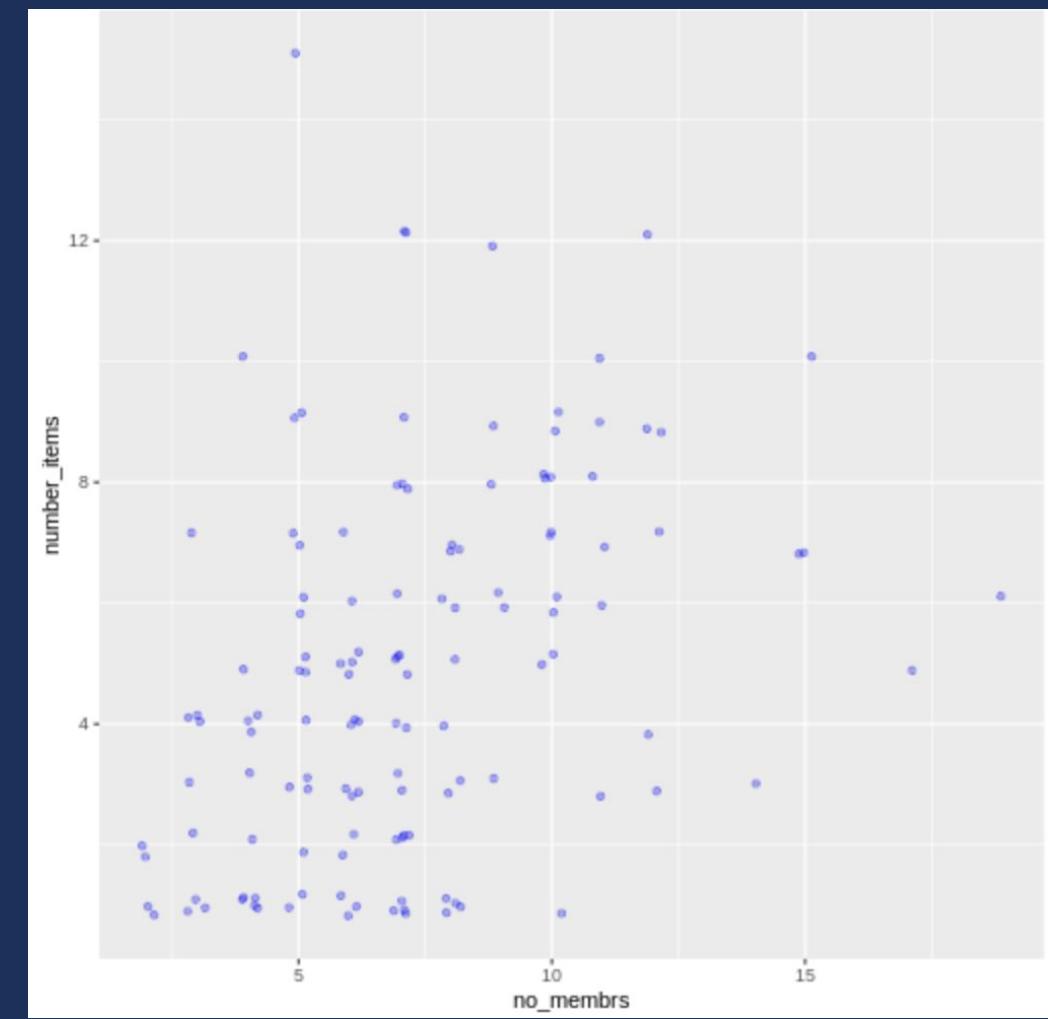
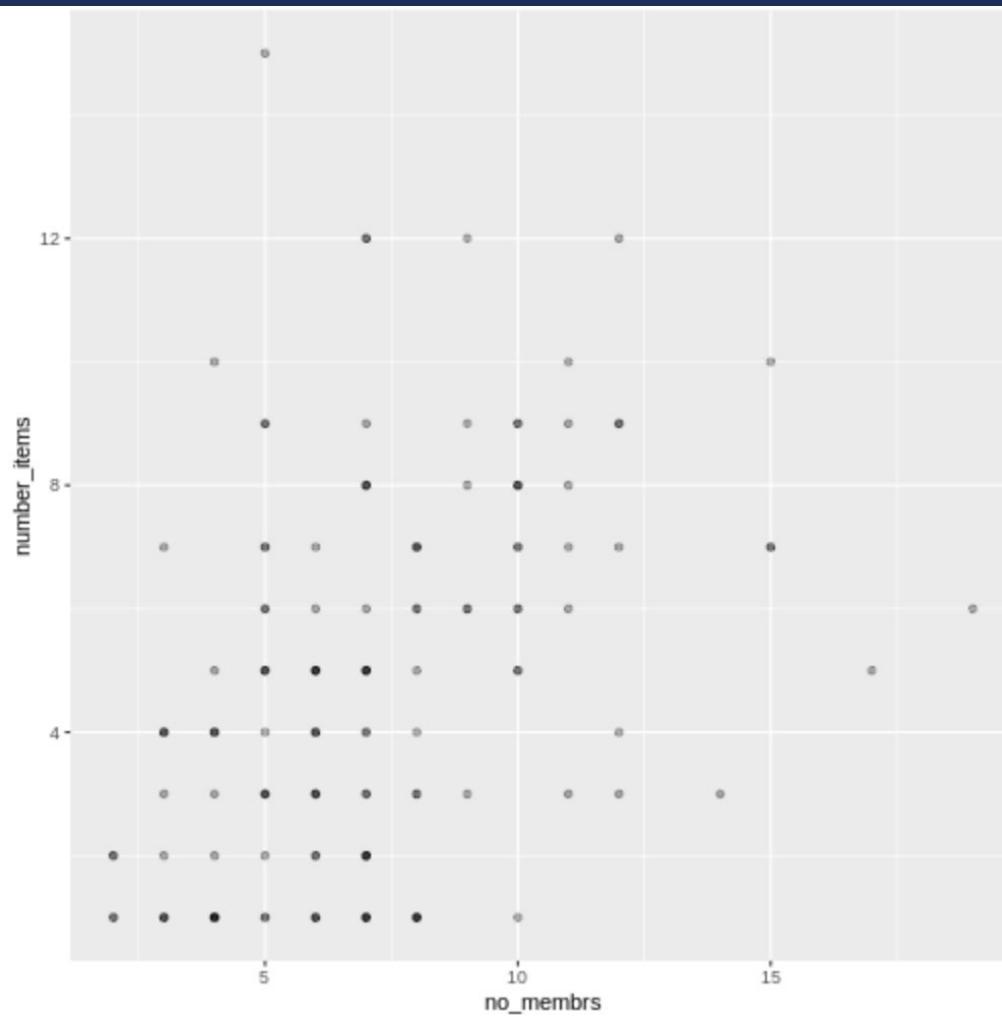
# Problem: Overlapping Data Points

---

Partial Transparency  
(a little better)



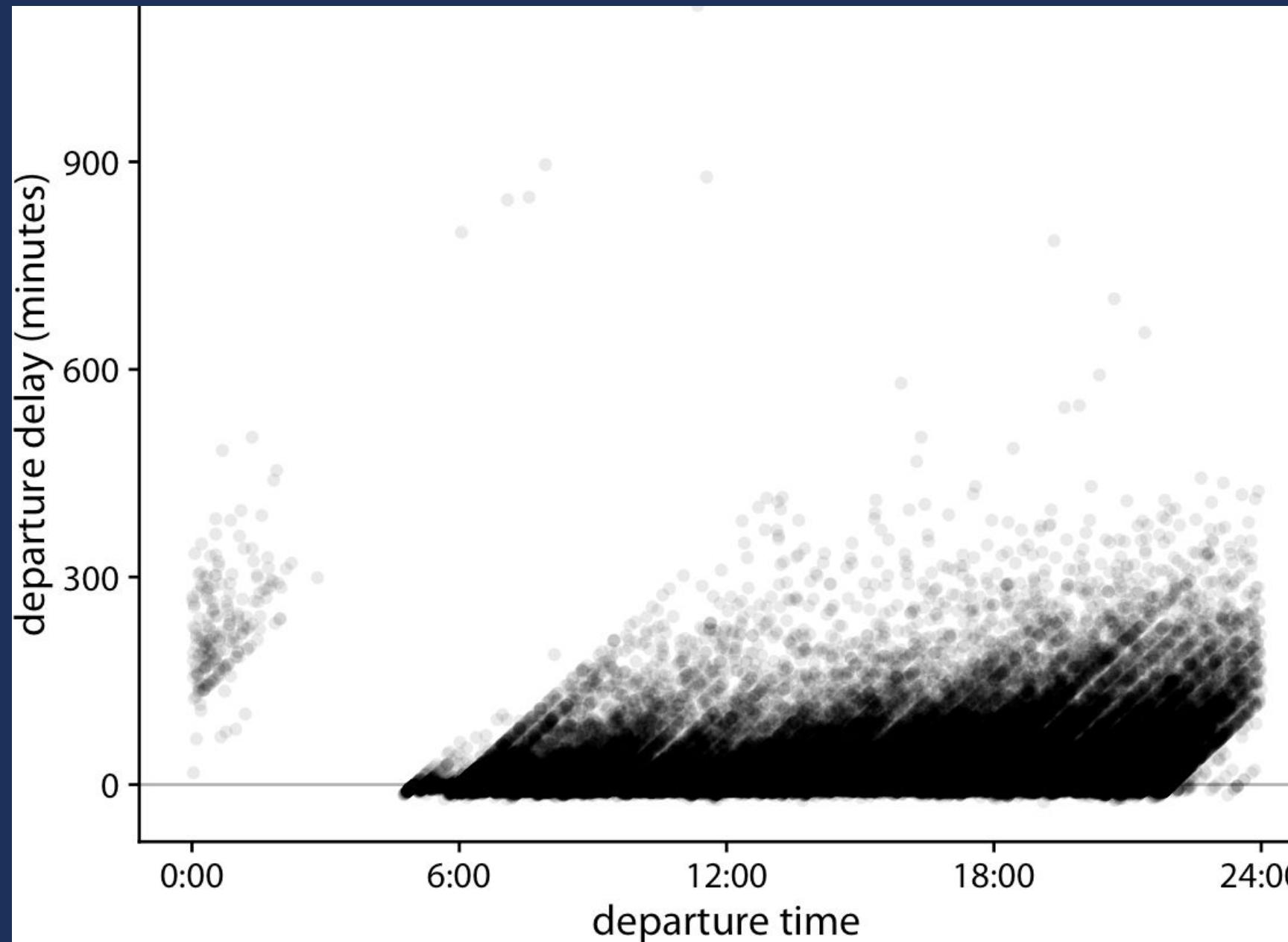
Partial Transparency  
Plus Jitter (much better)



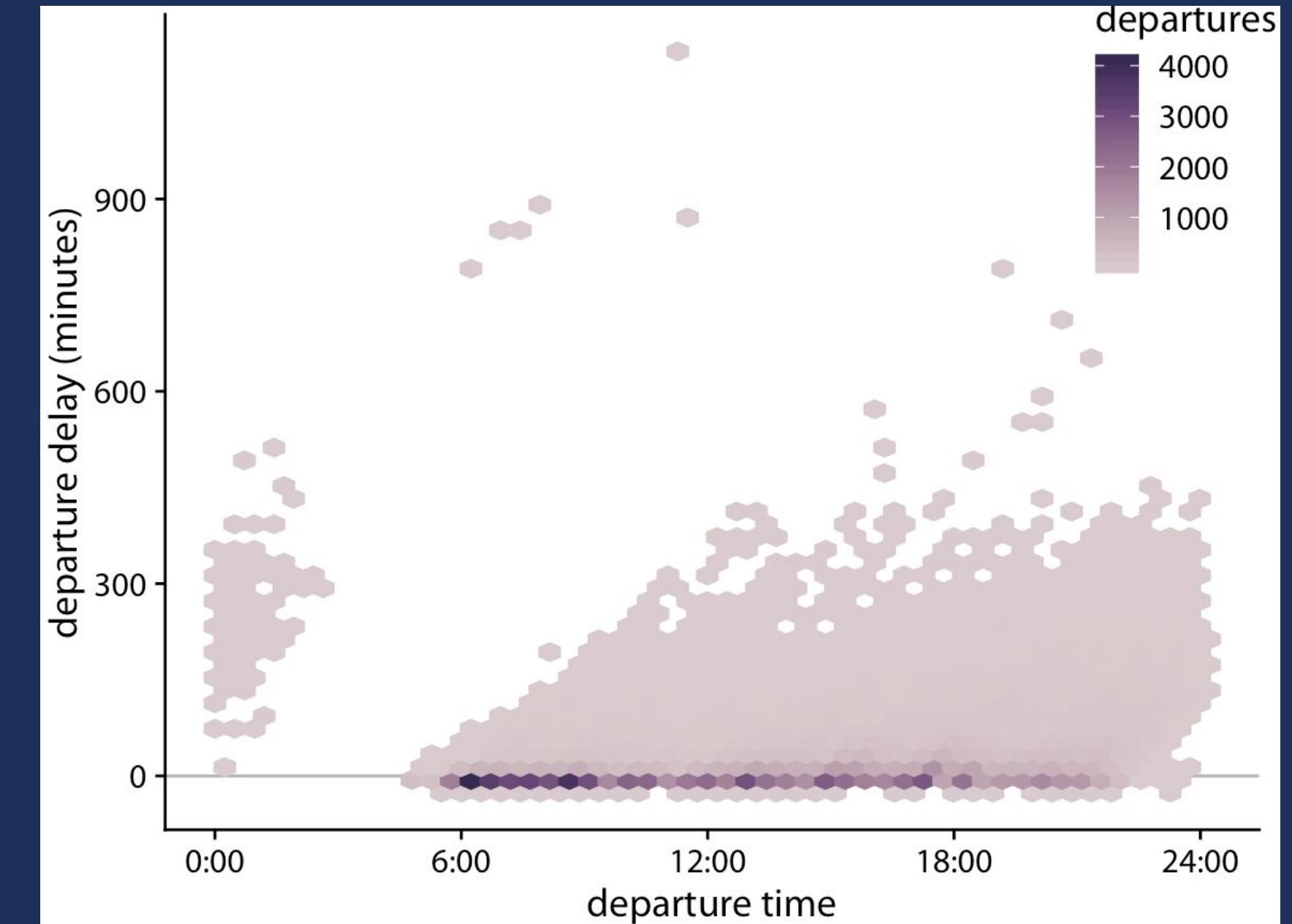
# Problem: Overlapping Data Points

---

but with a LOT of data

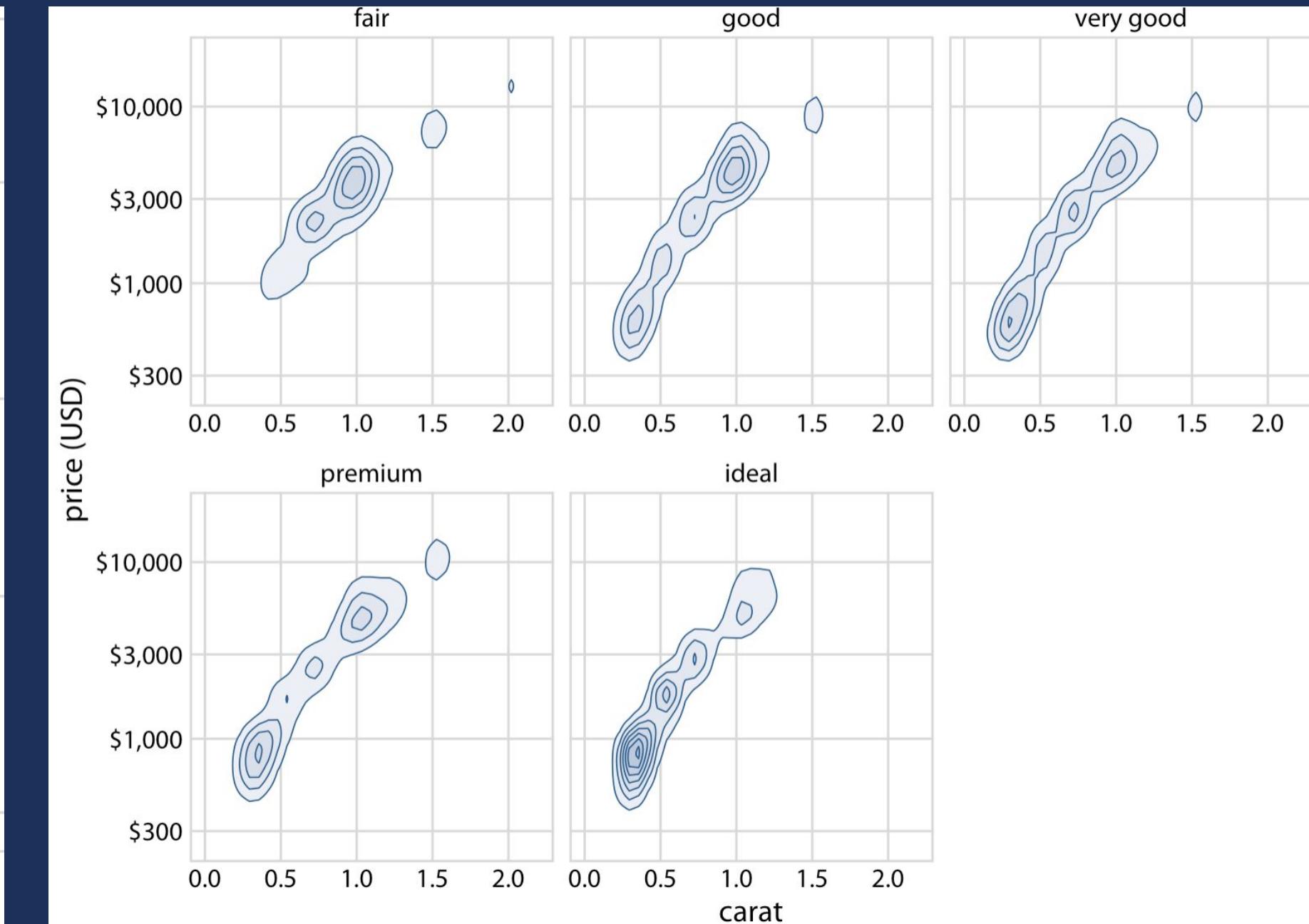
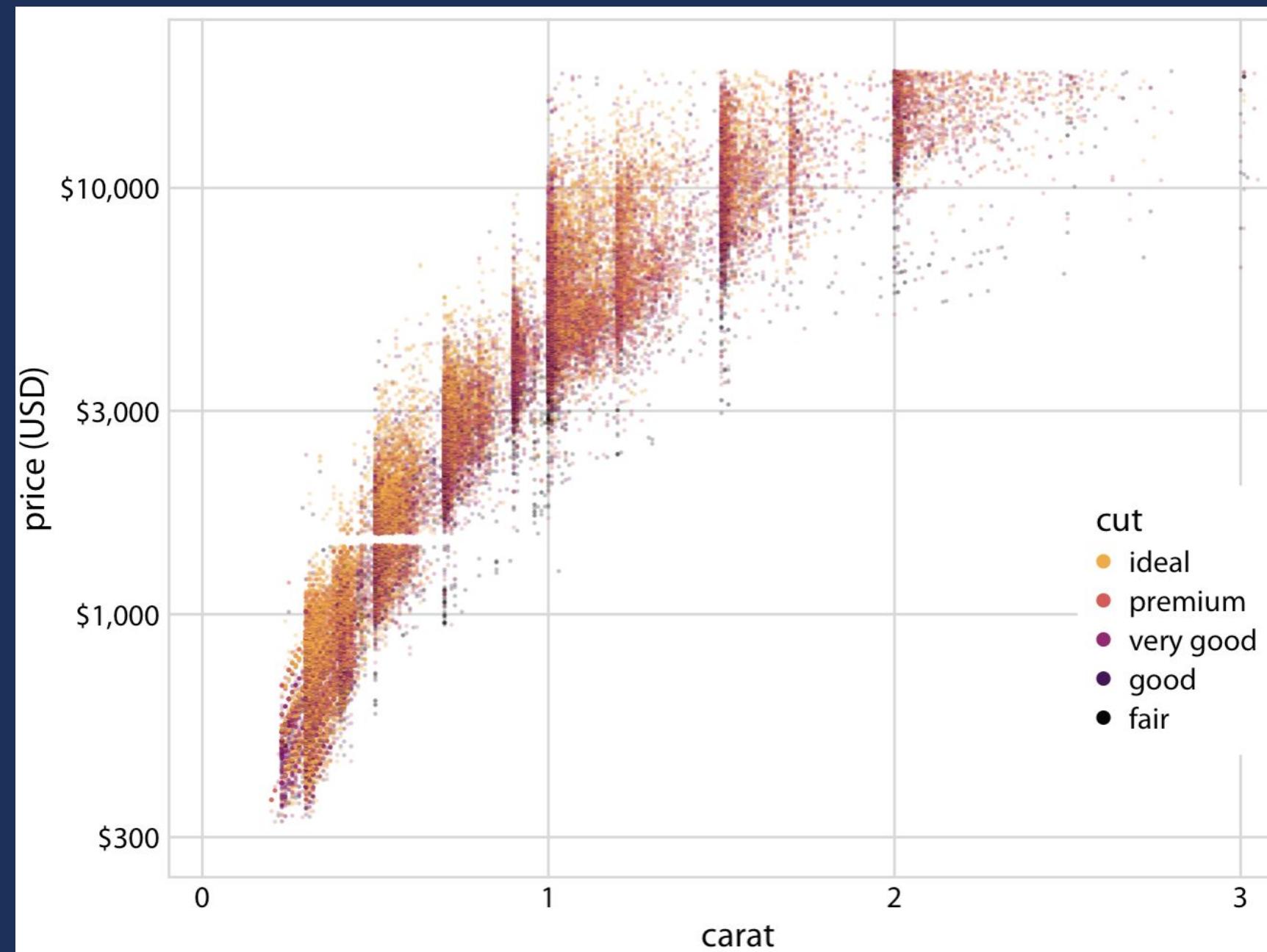


with 2-D binning



# Problem: Overlapping Data Points

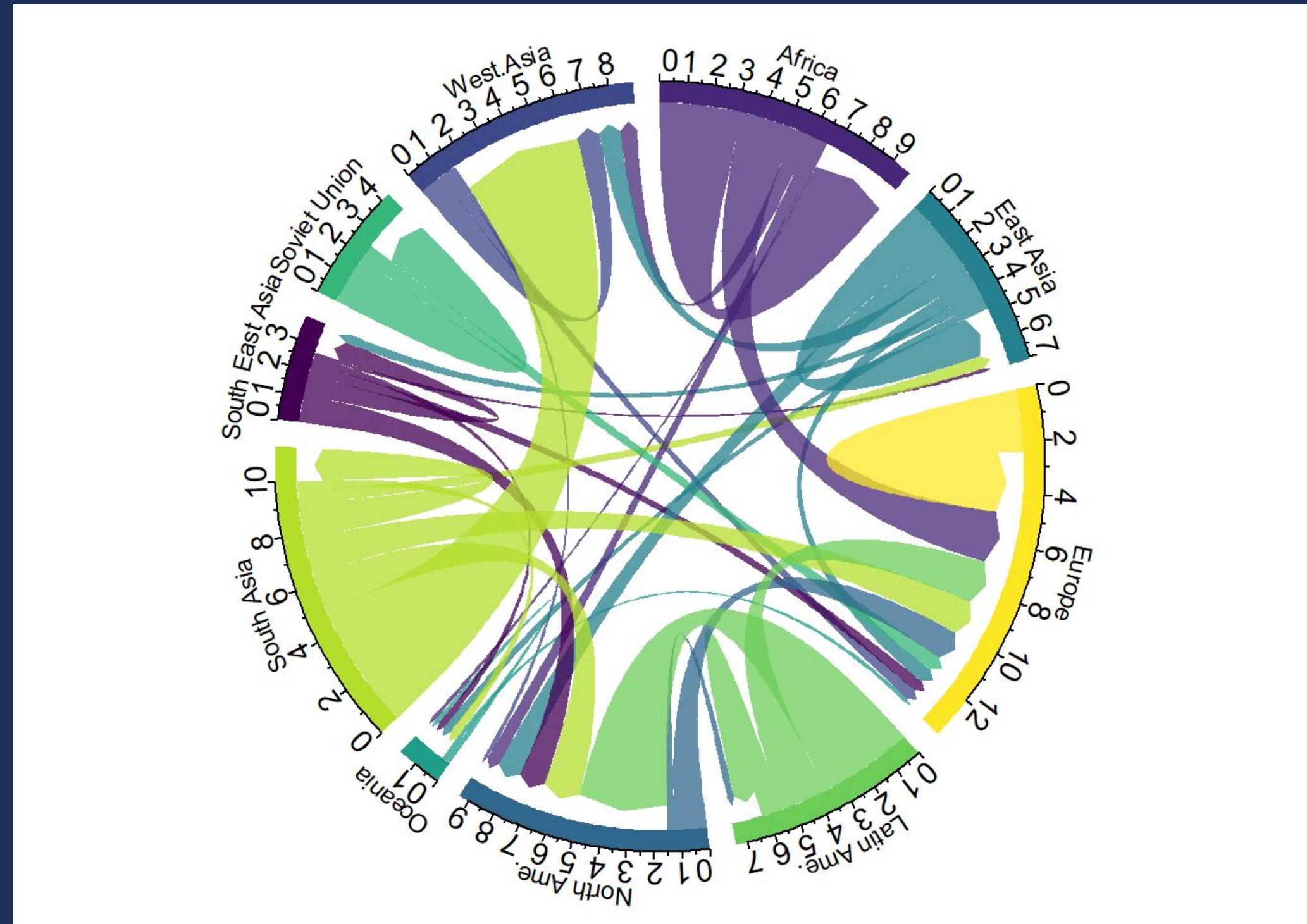
but with a LOT of data



Using density contour lines + faceting

	Africa	East Asia	Europe	Latin Ame.	North Ame.	Oceania	South Asia
Africa	3.142471	0.000000	2.107883	0.000000	0.540887	0.155988	0.000000
East Asia	0.000000	1.630997	0.601265	0.000000	0.973060	0.333608	0.000000
Europe	0.000000	0.000000	2.401476	0.000000	0.000000	0.000000	0.000000
Latin Ame.	0.000000	0.000000	1.762587	0.879198	3.627847	0.000000	0.000000
North Ame.	0.000000	0.000000	1.215929	0.276908	0.000000	0.000000	0.000000
Oceania	0.000000	0.000000	0.170370	0.000000	0.000000	0.190706	0.000000
South Asia	0.000000	0.525881	1.390272	0.000000	1.508008	0.347420	1.307950
South East Asia	0.000000	0.145264	0.468762	0.000000	1.057904	0.278746	0.000000
Soviet Union	0.000000	0.000000	0.609230	0.000000	0.000000	0.000000	0.000000
West.Aisia	0.000000	0.000000	0.449623	0.000000	0.169274	0.000000	0.000000

rowname	key	value
Oceania	Latin Ame.	0.000000
South Asia	Latin Ame.	0.000000
South East Asia	Latin Ame.	0.000000
Soviet Union	Latin Ame.	0.000000
West.Aisia	Latin Ame.	0.000000
Africa	North Ame.	0.540887
East Asia	North Ame.	0.973060
Europe	North Ame.	0.000000
Latin Ame.	North Ame.	3.627847
North Ame.	North Ame.	0.000000
Oceania	North Ame.	0.000000
South Asia	North Ame.	1.508008
South East Asia	North Ame.	1.057904
Soviet Union	North Ame.	0.000000
West.Aisia	North Ame.	0.169274
Africa	Oceania	0.155988
East Asia	Oceania	0.333608

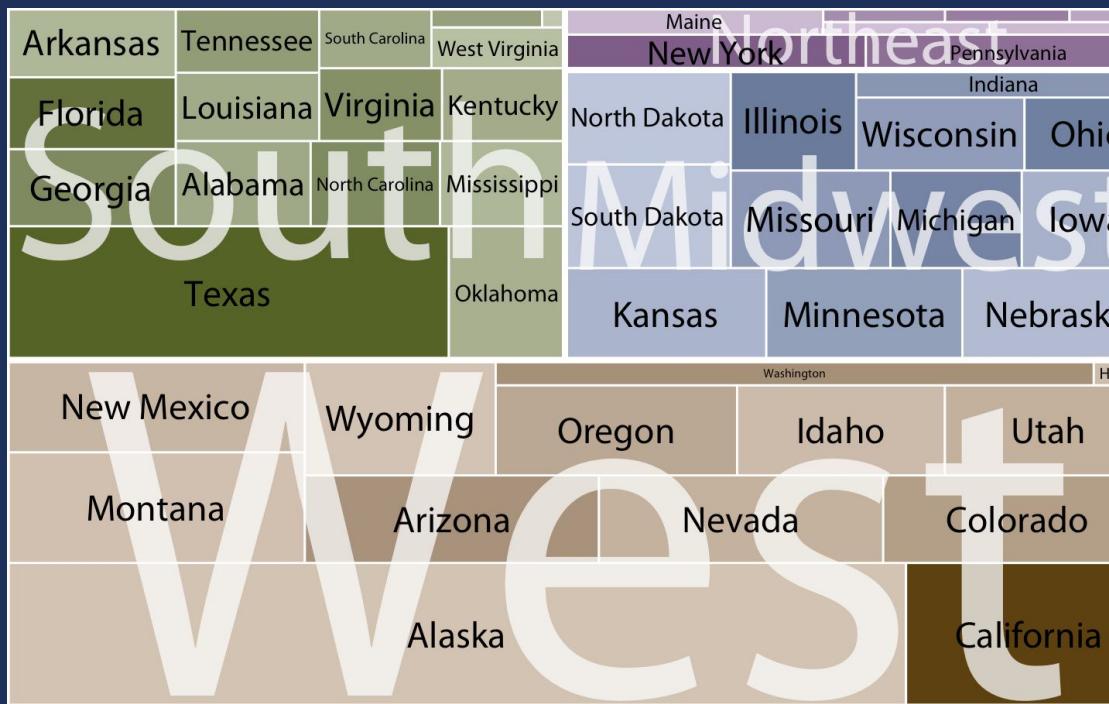


# SOLUTIONS TO SOME PROBLEMS

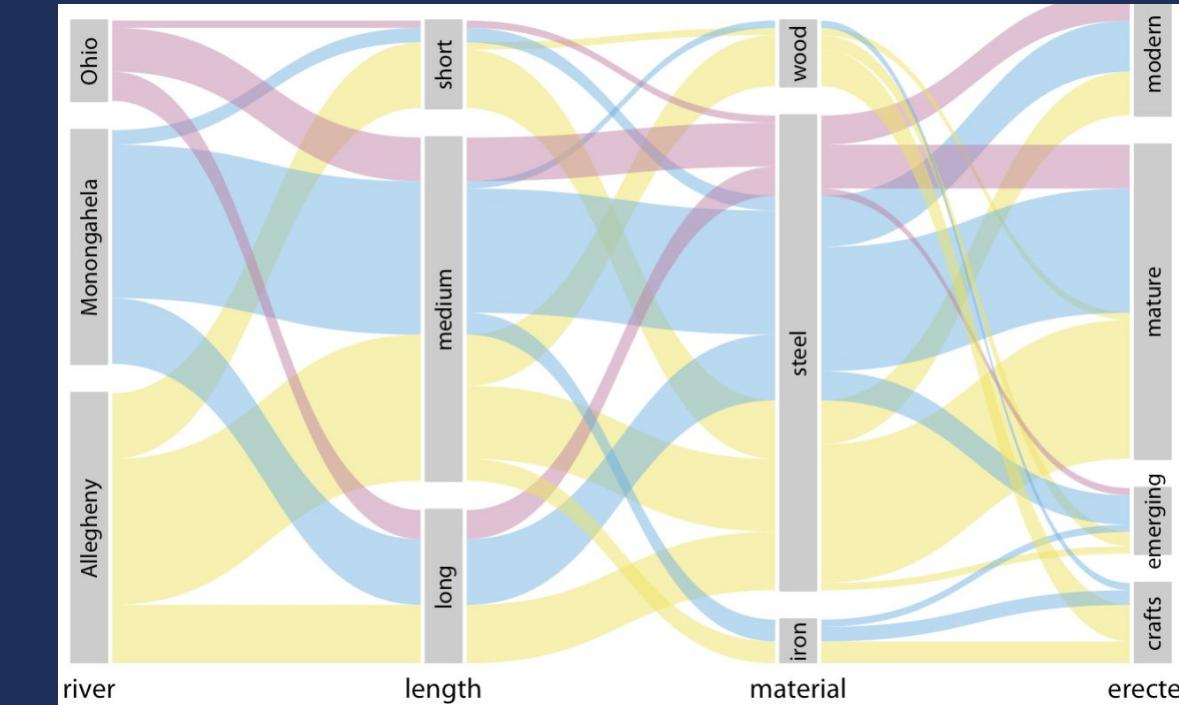
---

Problem #2: Categories with subcategories

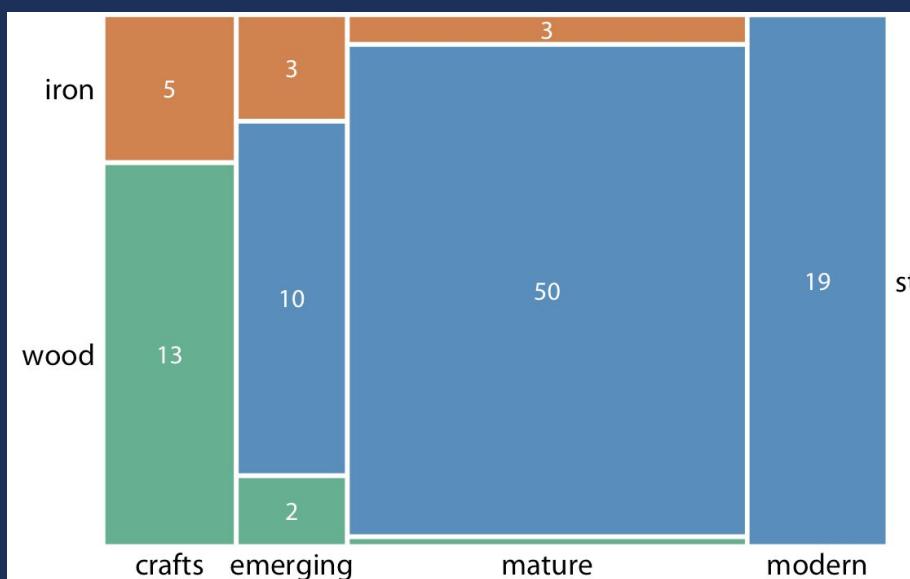
# Problem: Categories with subcategories



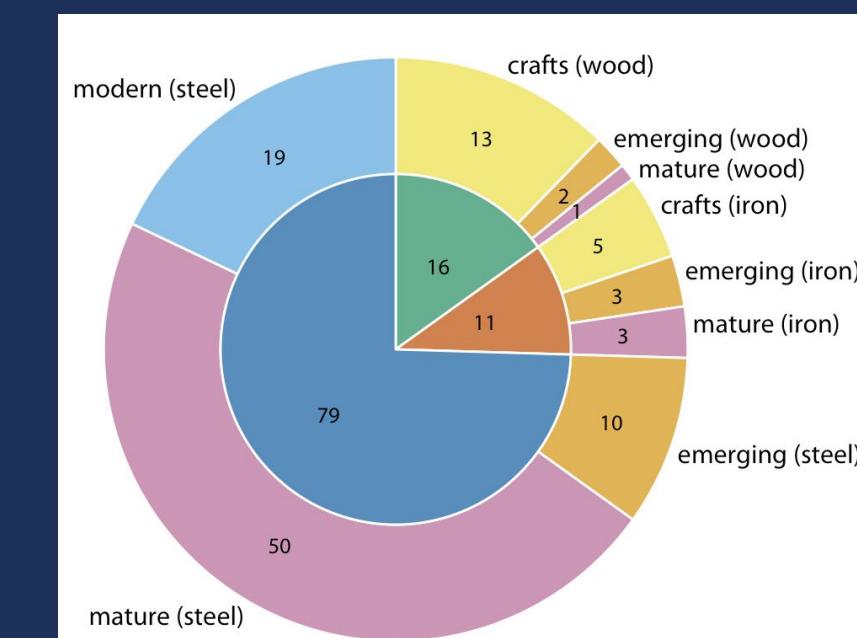
Treemap



Parallel Sets



Mosaic Plot

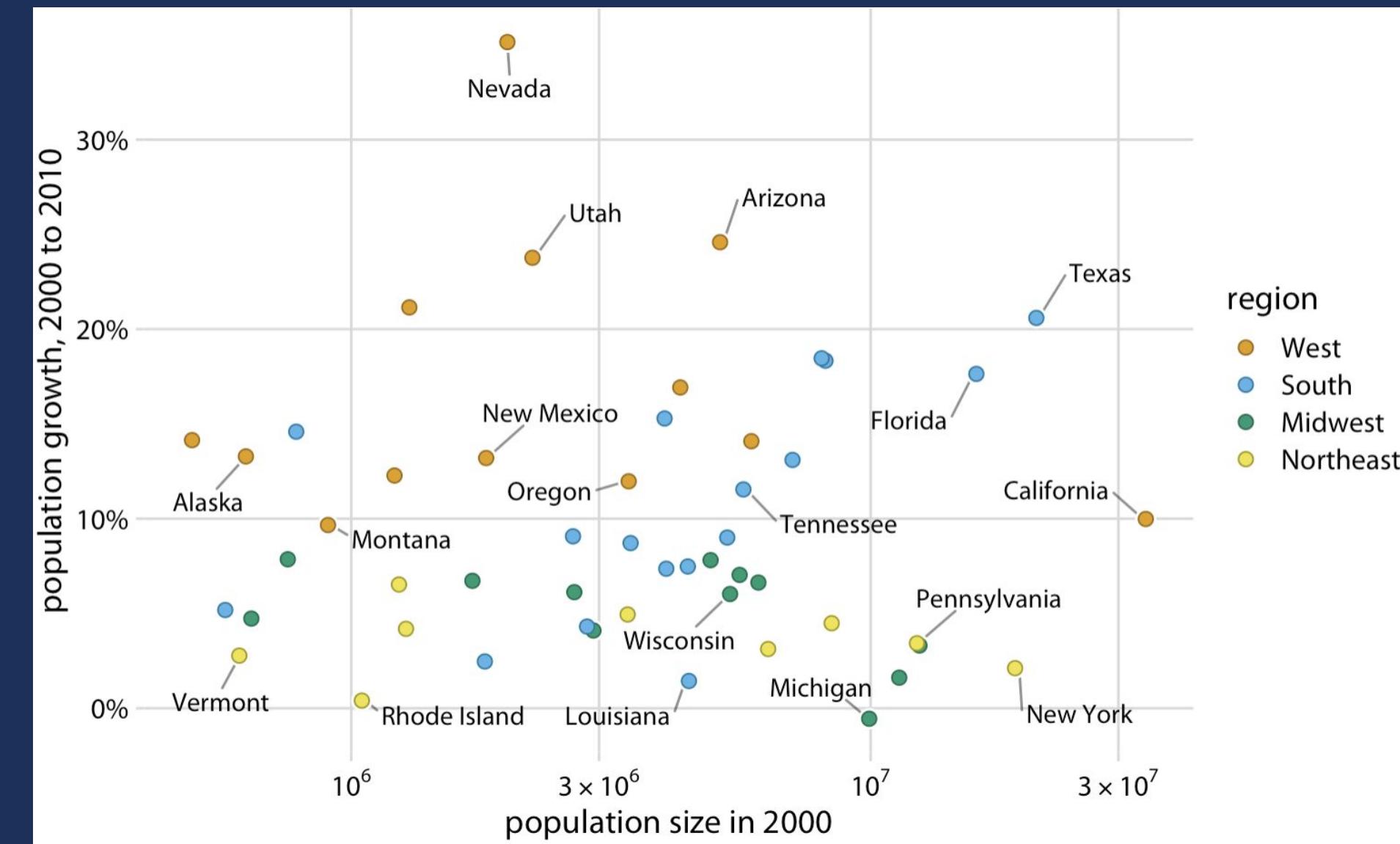
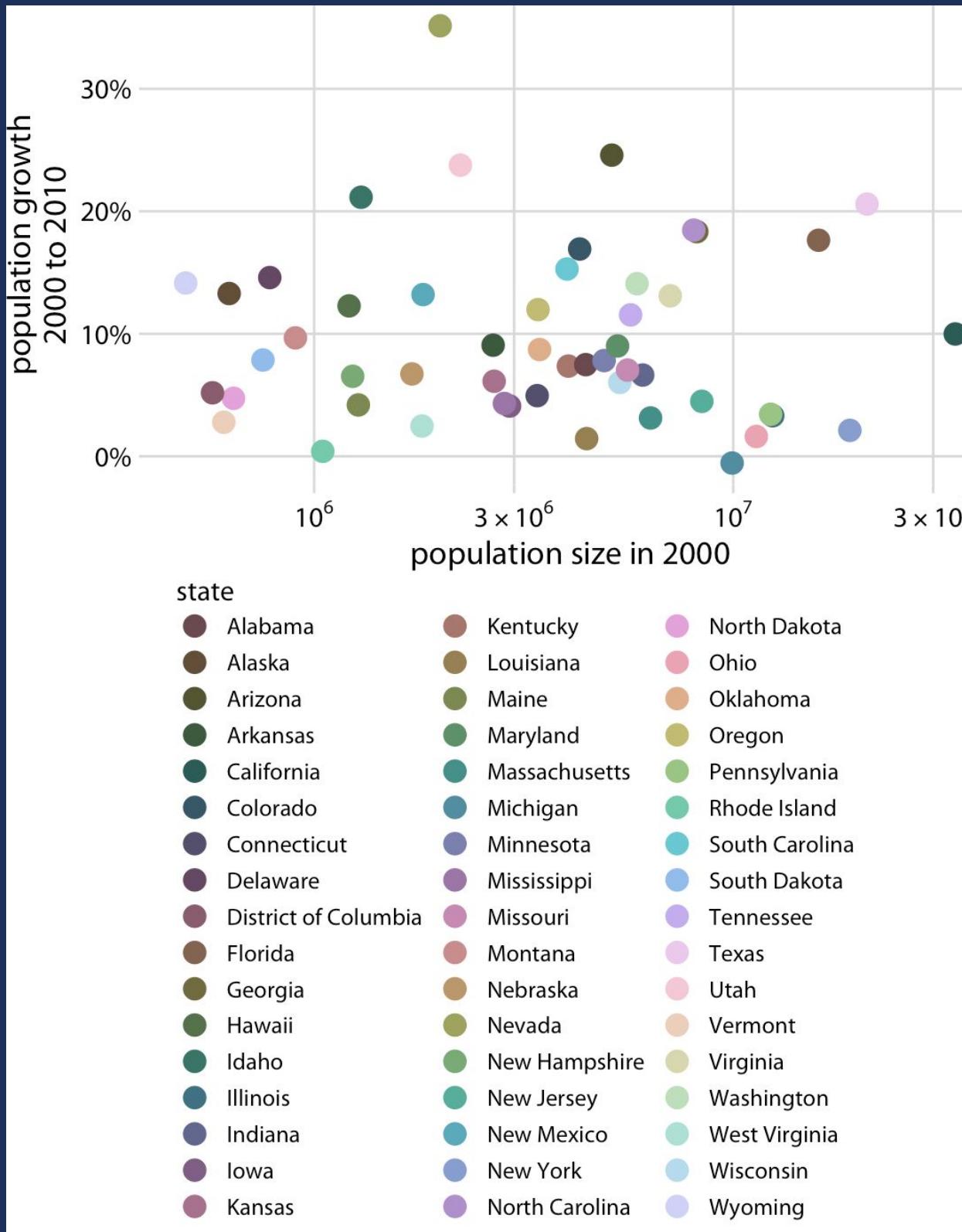


Nested Pie Chart

# SOLUTIONS TO SOME PROBLEMS

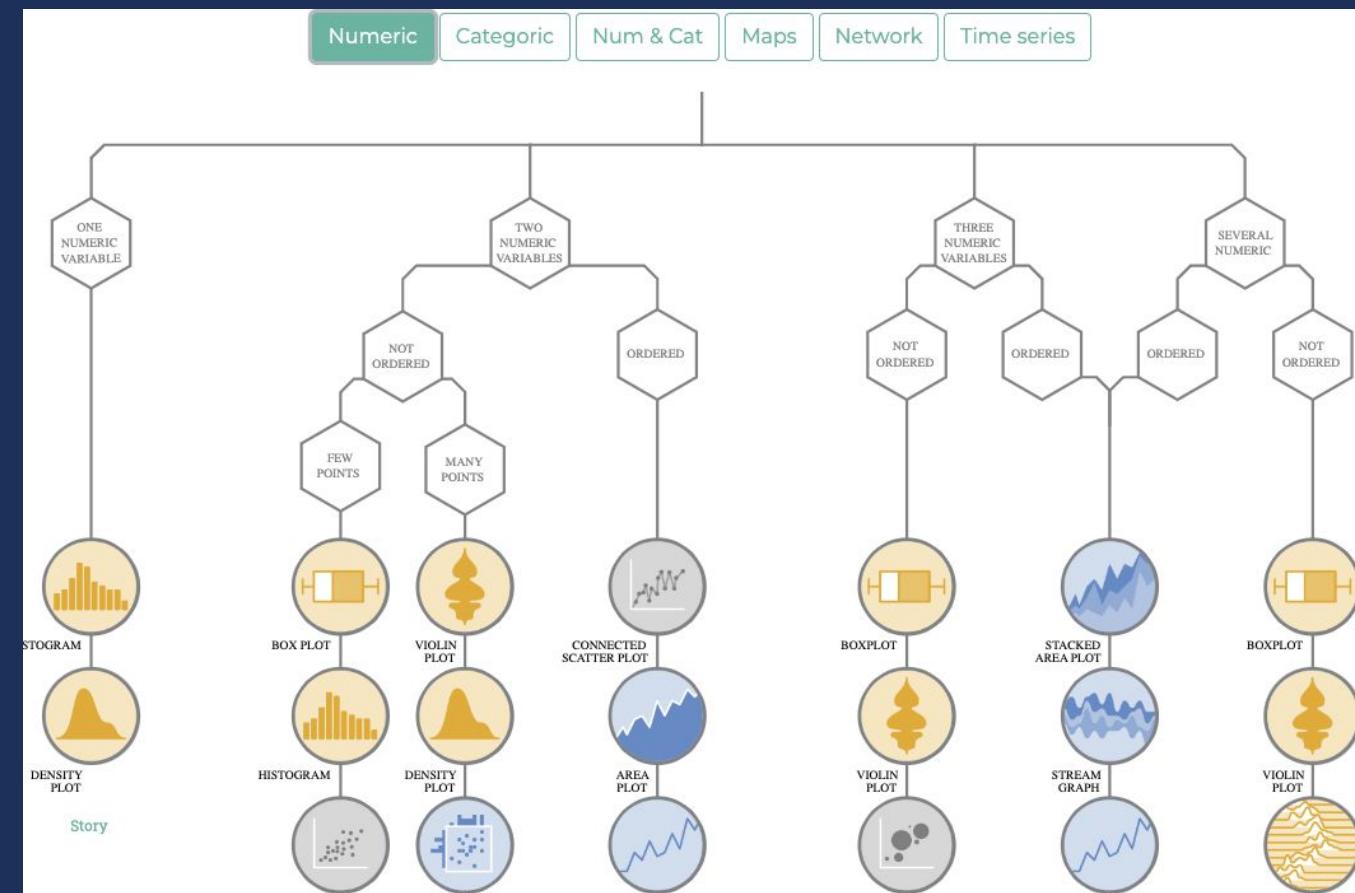
---

Problem #3: Too many categories/colors



from Data to Viz ~ [data-to-viz.com](http://data-to-viz.com)

# Decision Tree - suggested viz types



# library of viz types



# caveats

Show all	Top 10	Improvement	Misleading	Map	Bar
	<a href="#">Order your data</a>		<a href="#">To cut or not to cut?</a>		<a href="#">The spaghetti chart</a>
When displaying the value of several entities, ordering them makes the graph much more insightful.	Cutting the Y-axis is one of the most controversial practice in data viz. See why.	A line graph with too many lines becomes unreadable: it is called a spaghetti graph.		<a href="#">Pie chart</a>	The human eye is bad at reading angles. See how to replace the most criticized chart ever.
	<a href="#">Play with histogram bin size</a>		<a href="#">Do boxplots hide information?</a>		<a href="#">The problem with error bars</a>
Always try different bin sizes when you build a histogram, it can lead to different insights.	Boxplots are a great way to summarize a distribution but hide the sample size and their distribution.	Barplots with error bars must be used with great care. See why and how to replace them.		<a href="#">Too many distributions</a>	If you need to compare the distributions of many variables, don't clutter your graphic.
	<a href="#">Overplotting</a>		<a href="#">The rainbow color palette</a>		<a href="#">Faceting: horizontal or vertical?</a>
Too many points on your scatter plot makes it unreadable? Techniques exist to avoid overplotting.	Avoid the rainbow color palette when you map a numeric variable. So many better palettes exist.	Placing the individual plot horizontally or vertically is an important choice to make.		<a href="#">Don't be counter intuitive</a>	Your audience is used to a few dataviz standards. Not respecting standards can be misleading.

181 Beautiful Data Viz'es - [dataviz-inspiration.com](http://dataviz-inspiration.com)

# WORKSHOP OUTLINE

.....

Our content today is divided into four parts. Each part will be described with examples.

01

## Data Visualization Process

Walking through the visualization process, from setting goals to visualizing the data.

02

## Basic Design Considerations

Best practices and accessibility considerations in data visualization.

03

## Solutions to Common Problems

How to resolve, or even better, avoid common problems in data visualization.

04

## Visualization Resources

Workshops, tools, and other resources for data visualization.

# RECOMMENDED RESOURCES

---

Click on titles to link to the library's holdings for each resource

## Books

**Data Points** by Nathan Yau  
ISBN: 9781118462195  
Publication Date: 2013-04-15

**Fundamentals of Data Visualization** by Claus Wilke  
ISBN: 1492031089  
Publication Date: 2019-05-14

**The Visual Display of Quantitative Information** by Edward R. Tufte  
ISBN: 0961392142  
Publication Date: 2001-07-09

**Information Visualization** by Colin Ware  
ISBN: 9780123814654  
Publication Date: 2012-05-21

## R

**Data Visualisation with R** by Thomas Rahlf  
ISBN: 9783319497518  
Publication Date: 2017-03-22

**Ggplot2** by Hadley Wickham  
ISBN: 9783319242774  
Publication Date: 2016-06-08

## Python

**Python Data Visualization Cookbook - Second Edition** by Igor Milovanovic; Dimitry Foures; Giuseppe Vettigli  
ISBN: 9781784394943  
Publication Date: 2015-11-30

# CONSULTATIONS @ GWU

---

[calendly.com/gwul-coding](https://calendly.com/gwul-coding)

- Python
- R
- HTML/CSS/JavaScript
- General coding questions

[calendly.com/data-consultation](https://calendly.com/data-consultation)

- General analysis and visualization questions
- Data management
- R, GIS, SAS, SPSS, Excel, Tableau

[go.gwu.edu/dataconsulting](http://go.gwu.edu/dataconsulting)

Graduate Student Data  
Consultants

- R, Python, SAS, SPSS, STATA, Excel

# OTHER RESOURCES @ GWU

## Library Workshops

<https://library.gwu.edu/events>

## Library Research Guide

<https://libguides.gwu.edu/c.php?g=1354705>

George Washington University / Research Guides / Data Visualization / Home

### Data Visualization

A guide to data visualization principles and techniques.

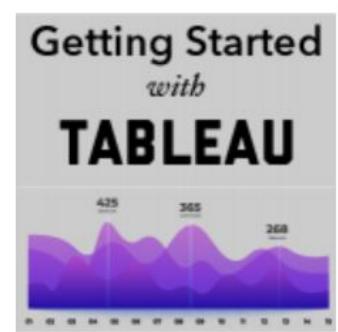
- Home**
  - Data Visualization
  - Library Resources
  - Attribution
- Best Practices**
- Visualization Tools**
- Additional Resources**
- Inspiration**

**Data Visualization**

Data visualization should be a constant consideration throughout the course of a research project. Exploring your data through visuals can help you spot patterns that might not be obvious by looking at statistics alone. Take a look at [this example](#) – even if the summary stats seem similar, the actual data can be very different!

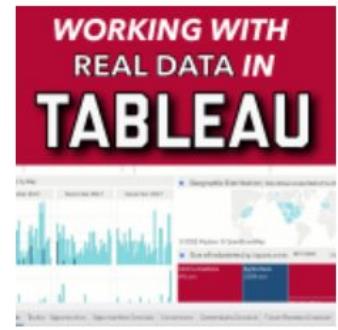
Also, remember that data visualization isn't just about understanding your data; it's often the end result of your research. So, it's a good idea to put some time and thought into your charts, graphs, and figures right from the start. Take a look at the best practices tab to consider the goal of your data visualization before you get started.

Different fields might use different software and data types, but the basic principles of good data visualization are the same. This guide can give you some tips and tools to help you create effective data visuals.



### Getting Started with Tableau

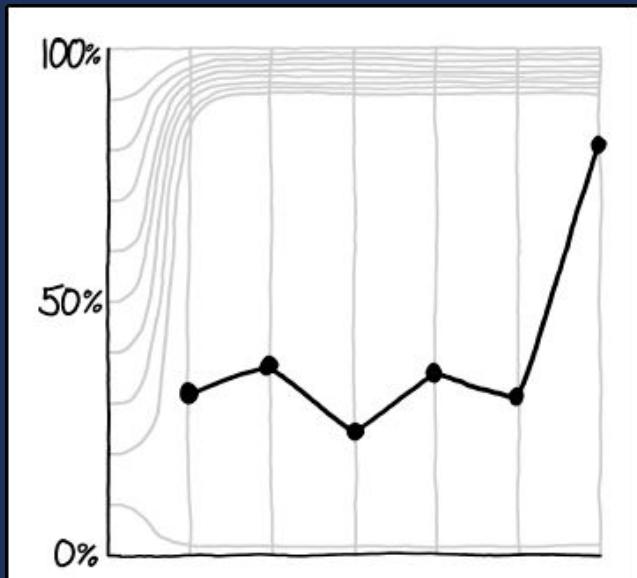
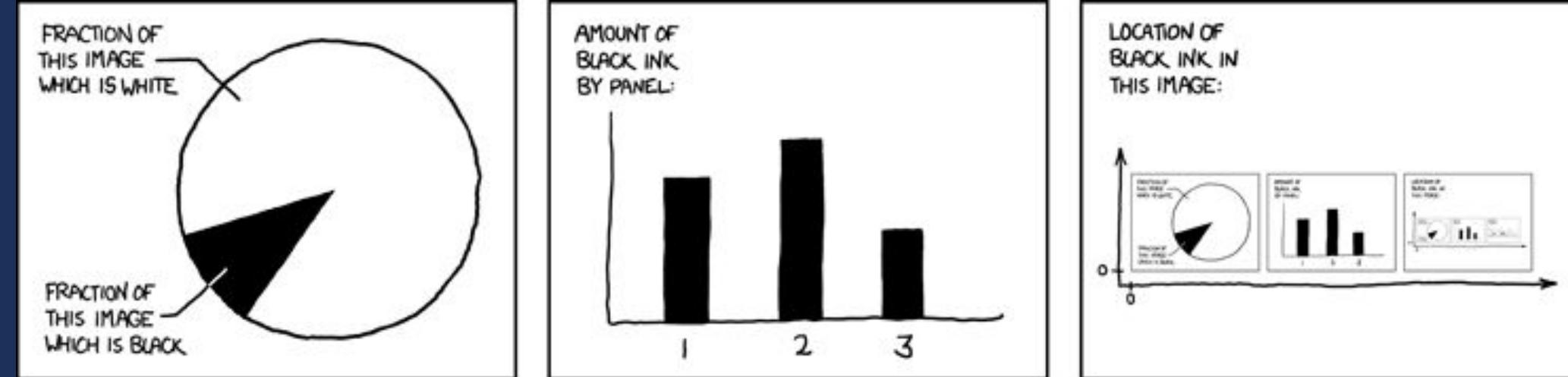
This workshop will provide a brief overview of using Tableau Public.



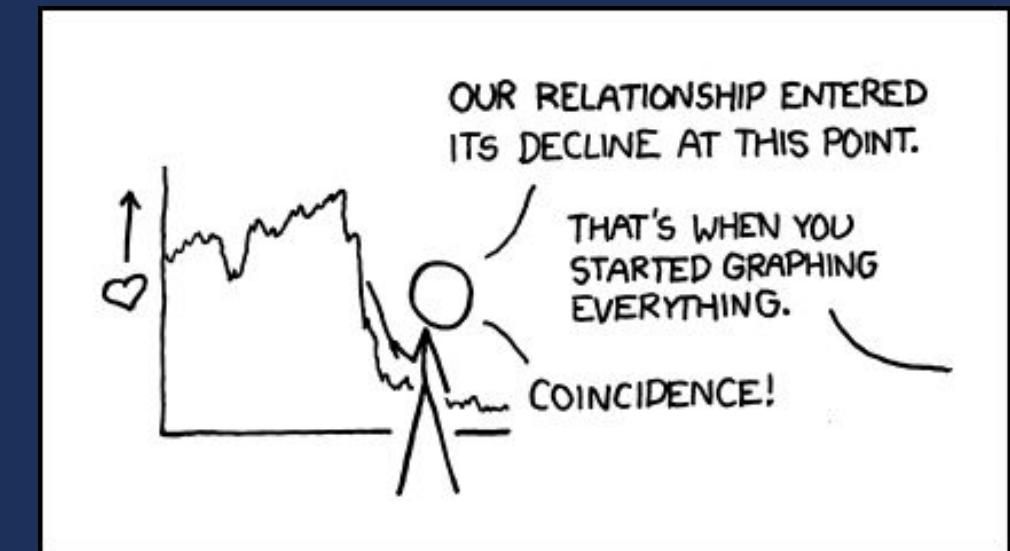
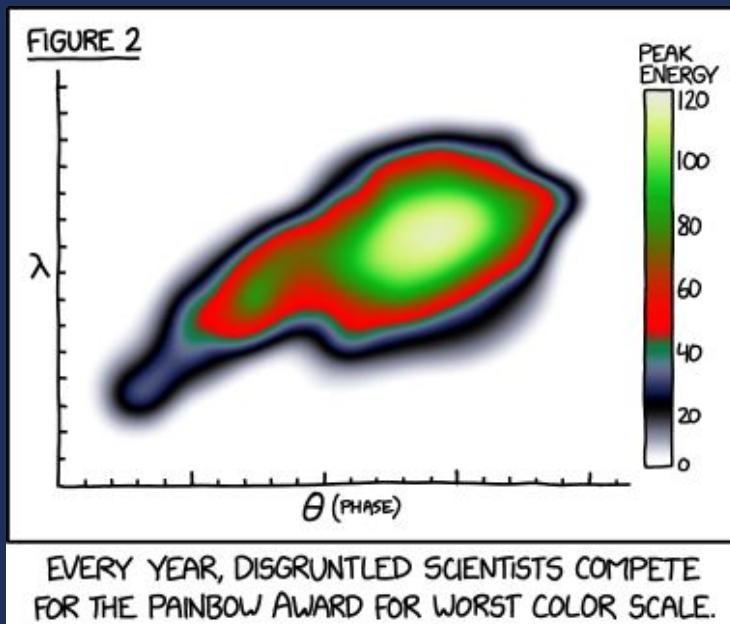
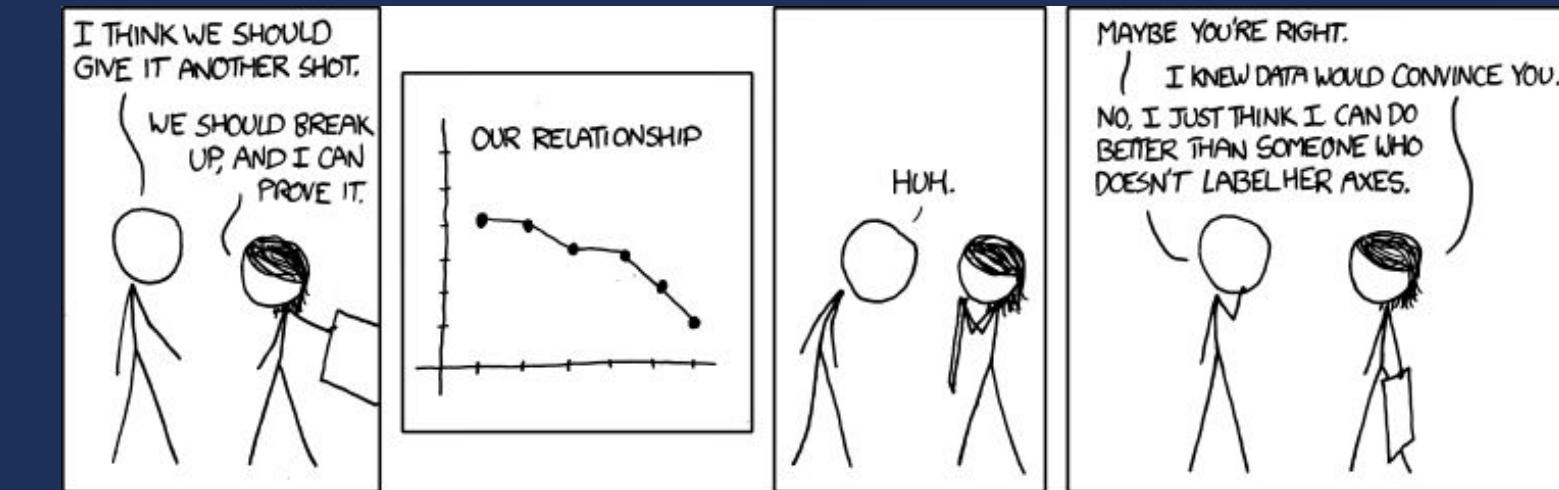
### Working with Real Data in Tableau

This workshop provides a more in-depth exploration of manipulating data in Tableau.

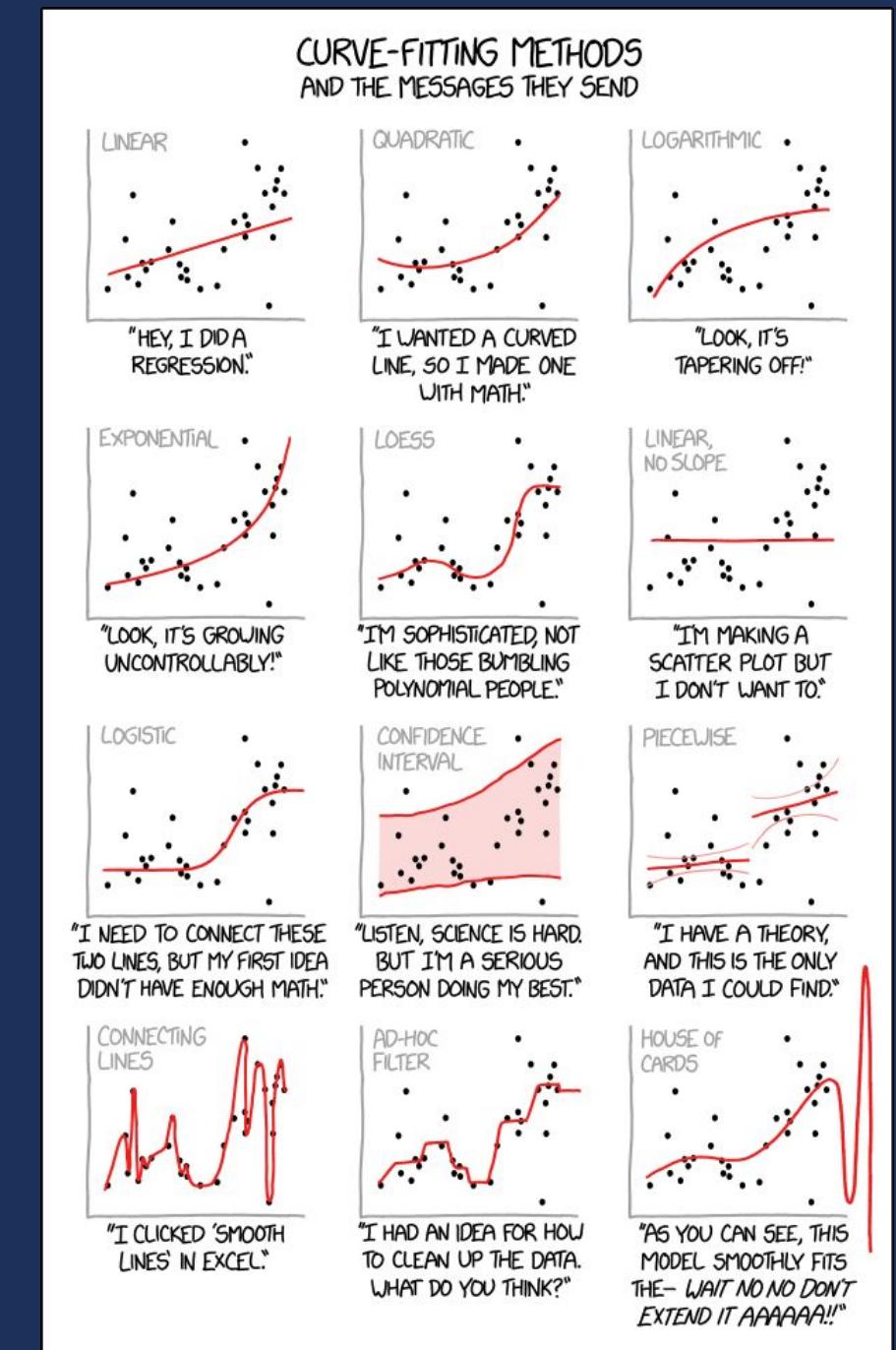
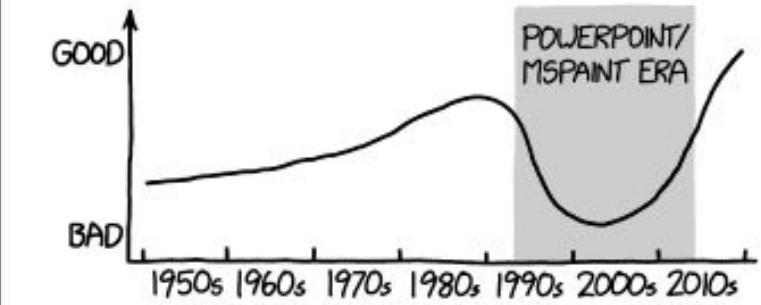
# THANK YOU!



PEOPLE HAVE WISSED UP TO THE "CAREFULLY CHOSEN Y-AXIS RANGE" TRICK, SO WE MISLEADING GRAPH MAKERS HAVE HAD TO GET CREATIVE.



## GENERAL QUALITY OF CHARTS AND GRAPHS IN SCIENTIFIC PAPERS



Images are all  
from [xkcd.com](http://xkcd.com)