

Introduction to Data Visualization

Govern for America

Jake Rozran

June 24, 2022

Who Am I?

- Please call me Jake (I'll accept Jake from State Farm)
- I am an Adjunct Professor of Data Science and Statistics @ Villanova University
- Data Science Practice Lead at a company called CivicActions
- Dad; Data Nerd; Philadelphian
- Please connect with me!
www.jakelearnsdatascience.com



Where is this "Expertise" Coming From?

Where Else is this "Expertise" Coming From?

Getting Setup for In-Class Exercise

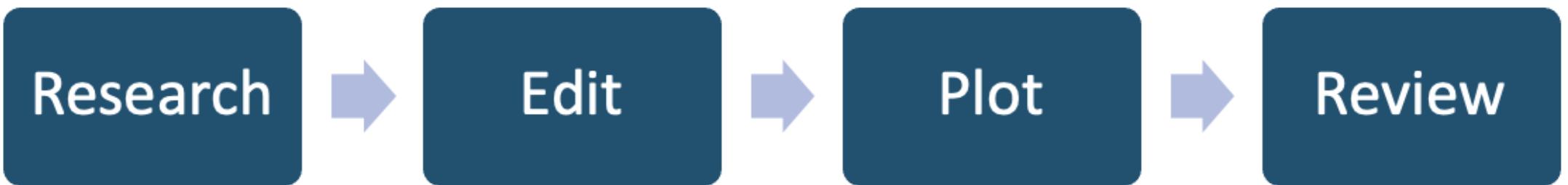
Download and open (in Excel) the Iris dataset OR

- PREFERRED - but only if you already have Microsoft Excel installed! (I know it says Google Sheets... but trust me here.)
- <https://tinyurl.com/iris-google-sheets>

Open the data in Google Sheets

- Make a copy of the spreadsheet so you can edit!
- <https://tinyurl.com/GFA-iris-google-sheets>

How to Create Effective Charts



1. Research

This is where you are finding, exploring, and understanding the limitations of your data.

| Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|--------------|-------------|--------------|-------------|---------|
| 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 5.0 | 3.6 | 1.4 | 0.2 | setosa |
| 5.4 | 3.9 | 1.7 | 0.4 | setosa |

Who can tell me something about the dataset?

2. Edit

This is where you are identifying your key message, choosing the best data to convey your message, filtering and simplifying your data, and making numerical adjustments.

What is something we can explore?

Average Sepal Length and Width by Species

Let's create this together.

| Species | Avg Sepal Length | Avg Sepal Width |
|------------|------------------|-----------------|
| setosa | 5.006 | 3.428 |
| versicolor | 5.936 | 2.770 |
| virginica | 6.588 | 2.974 |

Difference from the Average for Each Row

Let's create this together.

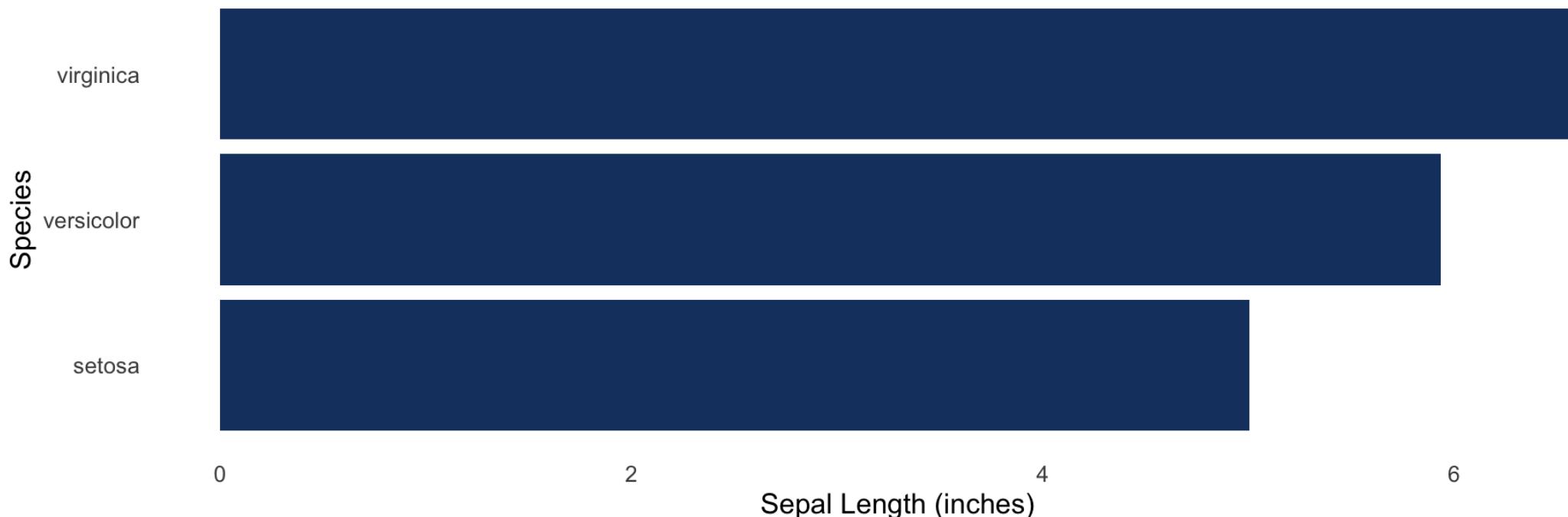
| Sepal Length | Sepal Width | Species | Avg. Length | Avg. Width | Pct. Diff. Len. | Pct. Diff. Wid. |
|--------------|-------------|---------|-------------|------------|-----------------|-----------------|
| 5.1 | 3.5 | setosa | 5.006 | 3.428 | 0.019 | 0.021 |
| 4.9 | 3.0 | setosa | 5.006 | 3.428 | -0.021 | -0.125 |
| 4.7 | 3.2 | setosa | 5.006 | 3.428 | -0.061 | -0.067 |
| 4.6 | 3.1 | setosa | 5.006 | 3.428 | -0.081 | -0.096 |
| 5.0 | 3.6 | setosa | 5.006 | 3.428 | -0.001 | 0.050 |
| 5.4 | 3.9 | setosa | 5.006 | 3.428 | 0.079 | 0.138 |

3. Plot

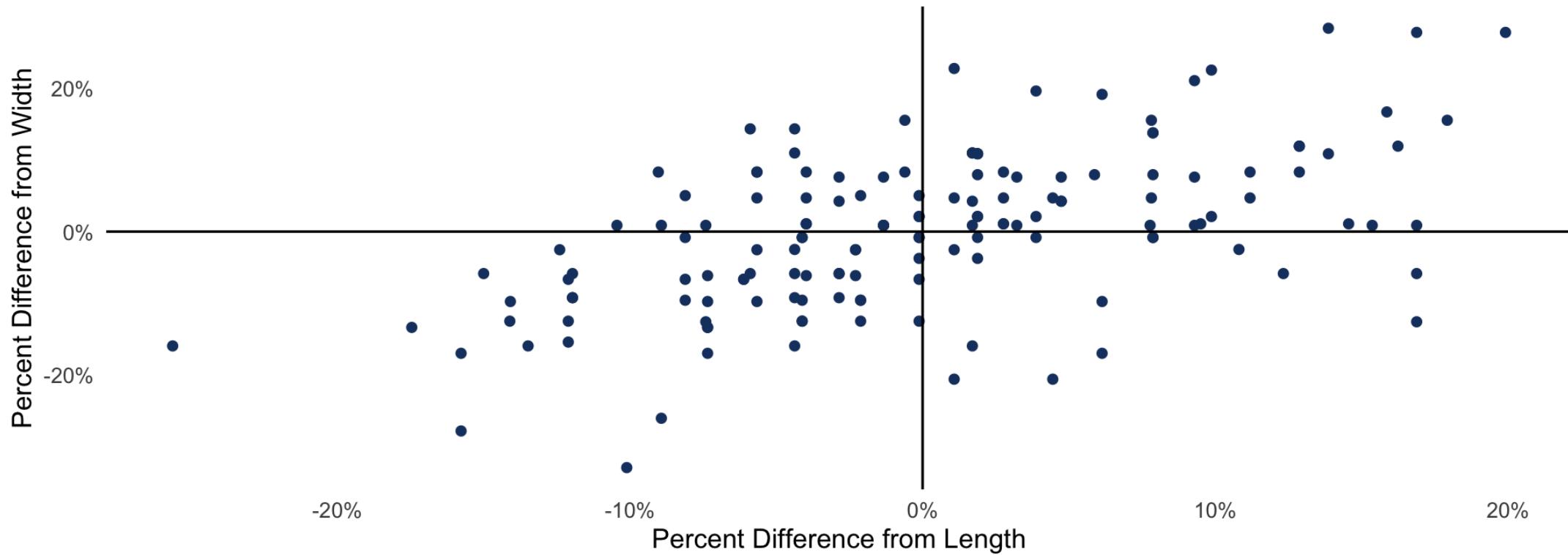
At this step, you are choosing the proper plot type, ensuring your plot axes are correct, title and labels are included, and color and typography are clear (more on that shortly).

Mean Sepal Length for Iris Data Set

Stock Dataset for Data Science



Percent Difference from Averages

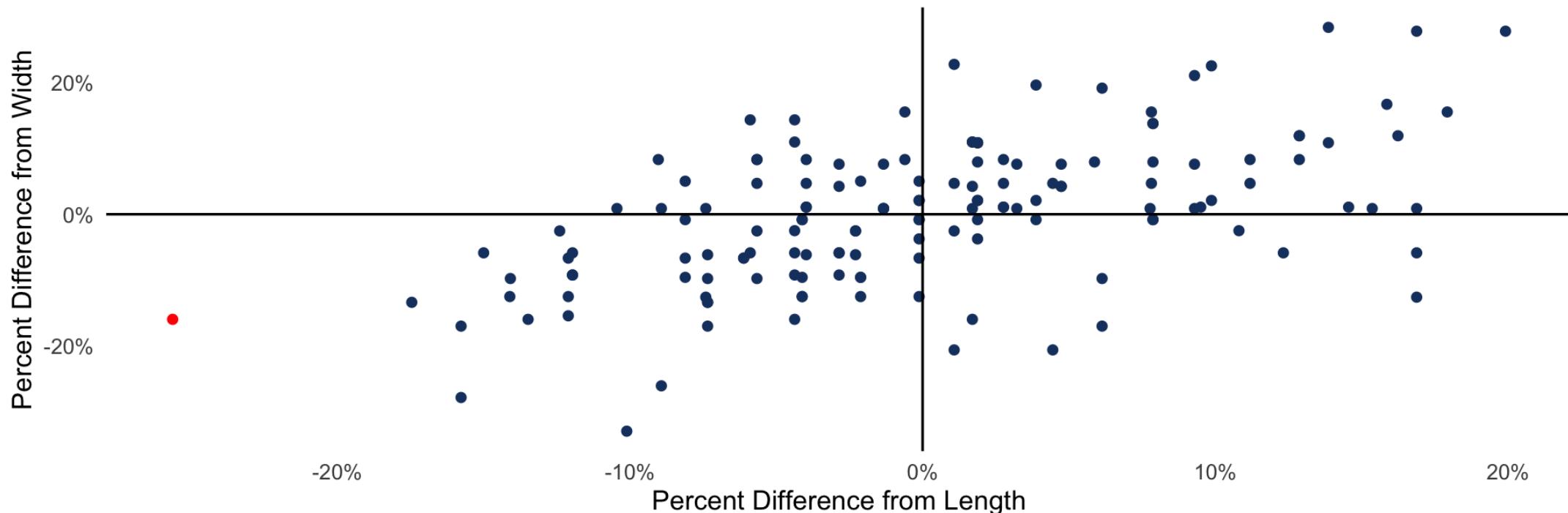


4. Review

Often overlooked due to expediency, ensuring the data is accurate will preempt questions and enhance trust in the data.

Percent Difference from Averages

Is the Red Point an Outlier or Bad Data?



**It is easy to slap something together. It is
HARD to put something clean, clear, and
meaningful together.**

A taxonomy for data graphics

Taxonomy, smaxonomy. That's just a fancy name for saying you got to know the finer parts of a data viz before you, too, can make fine data viz.

Data graphics can be understood in terms of four basic elements:

1. Visual cues
2. Coordinate systems
3. Scale
4. Context

And two bonus items:

1. Facets
2. Layers

Visual Cues

These are the building blocks of data viz. Visual cues are graphical elements that draw the eye to what we want our audience to focus upon. Human beings' ability to perceive difference in magnitude accurately decends in this order.

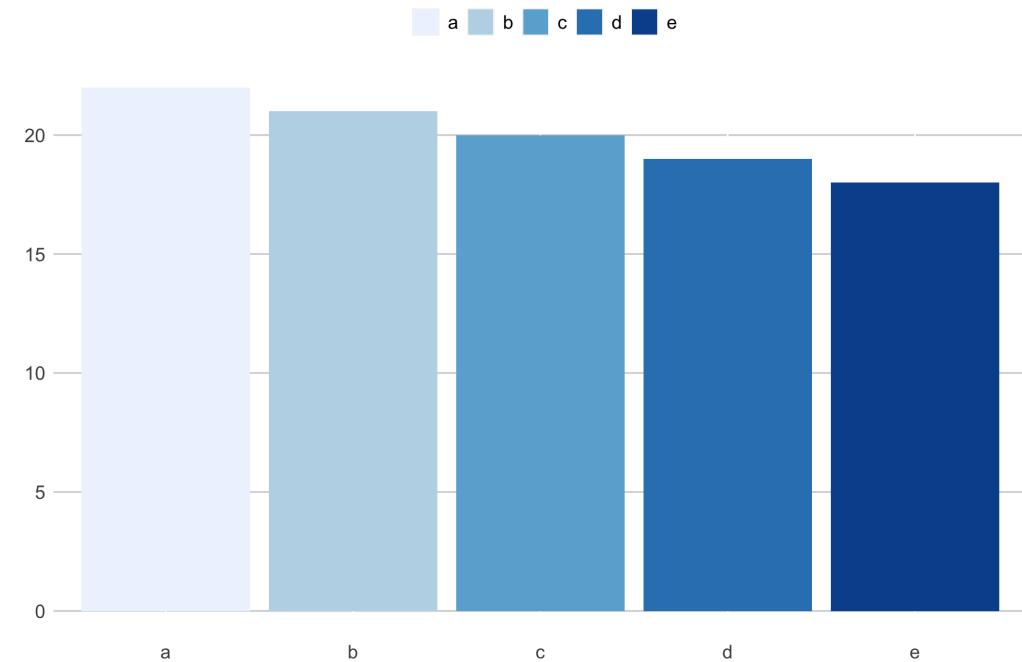
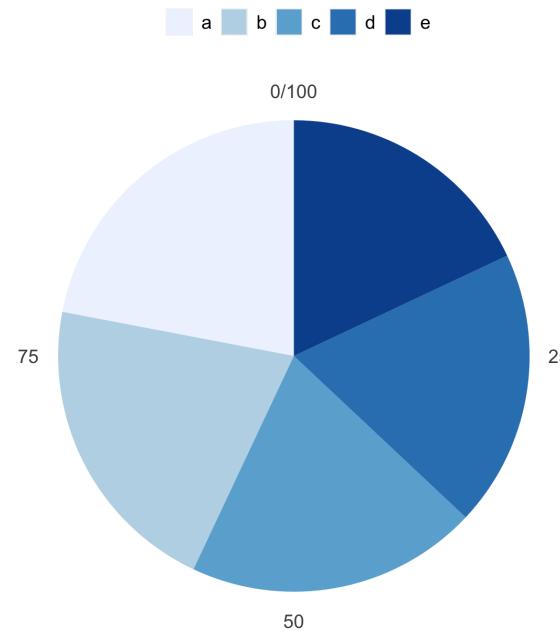
Visual cues and what they signify.

| Visual Cue | Variable Type | Question |
|------------|---------------|--|
| Position | numerical | where in relation to other things? |
| Length | numerical | how big (in one dimension)? |
| Angle | numerical | how wide? parallel to something else? |
| Direction | numerical | at what slope? in a time series, going up or down? |
| Shape | categorical | belonging to which group? |
| Area | numerical | how big (in two dimensions)? |
| Volume | numerical | how big (in three dimensions)? |
| Shade | either | to what extent? how severely? |
| Color | either | to what extent? how severely? |

Position, Length, & Area

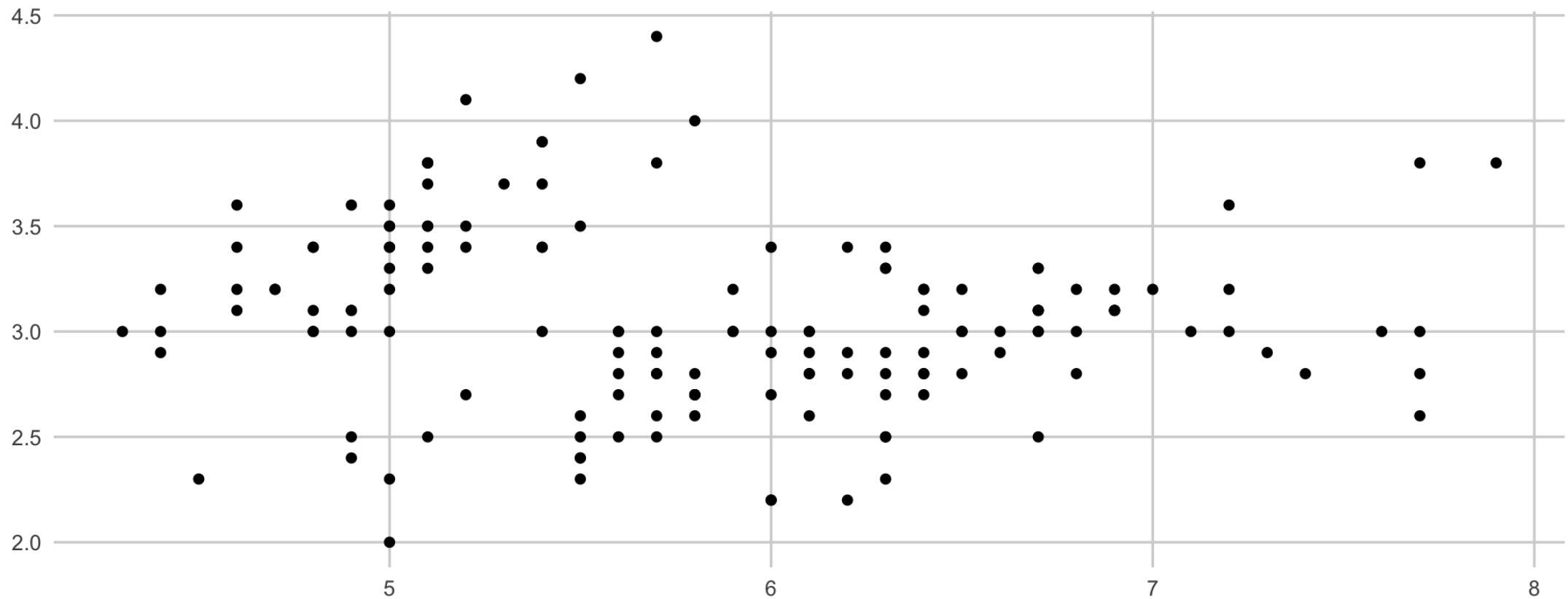
Because I feel like this is getting a little too much me talking at you, let us do an example together.

Someone tell me what they see in this plot, especially as it relates to **position**, **length**, and **area**? Which letter has the most?



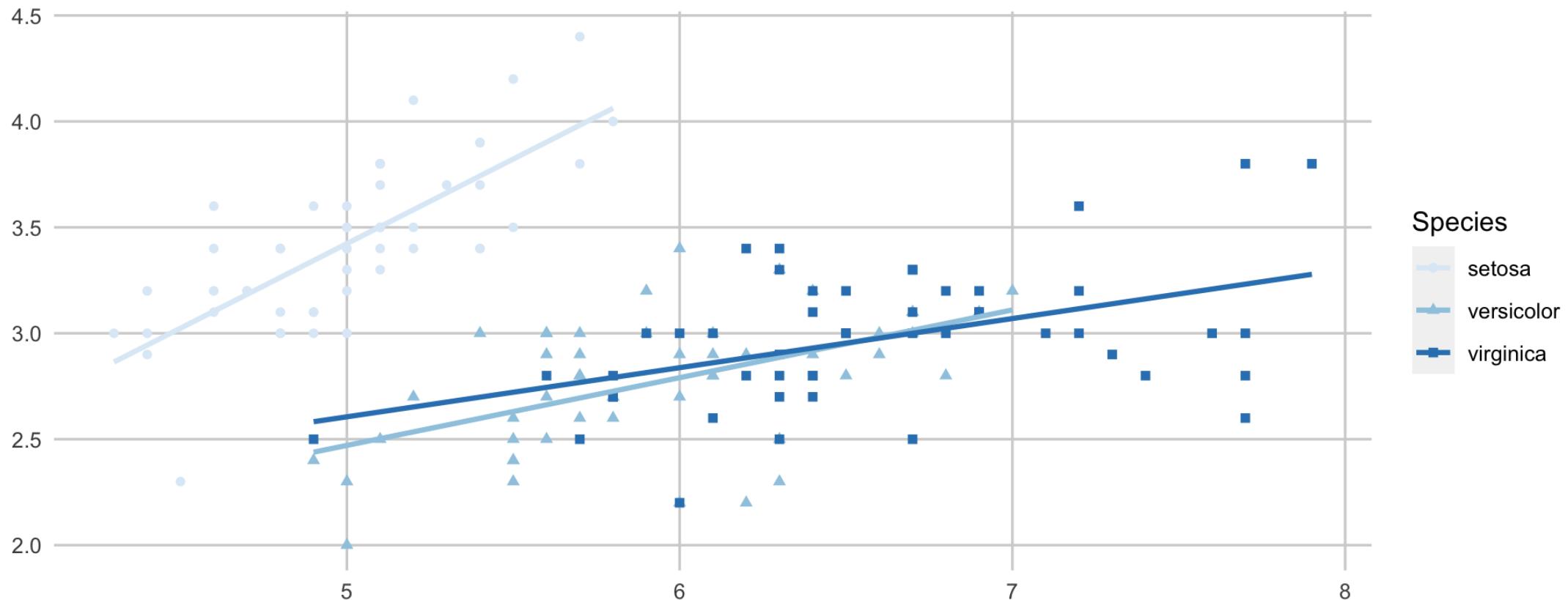
Direction

What **direction** do we see in this plot? Anything else?



Angle, Color, Shape

Talk to me about the **angle, color, and shape**.



Plots Speak Louder than Words

Do NOT!

Company A had revenues of \$1MM and
Company B had revenues of \$750k in
FY2021.

Do!

Annual Company Earnings (\$)
Visualization Gives a Fuller Understanding of the Data.



Let the Data Speak for Itself

Do not use any bs to make the graph "pretty." A clean crisp chart allows the reader to focus on the data and the message of the story.

This means:

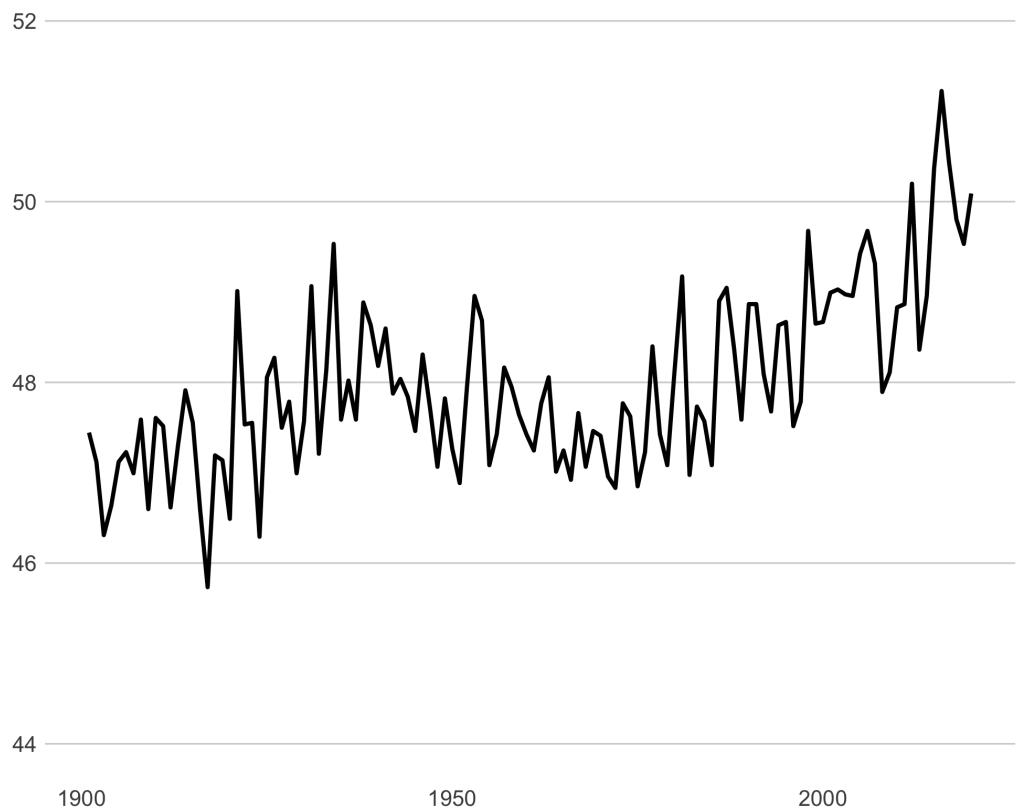
1. No 3D rendering EVER
2. No "pretty" colors (more on this soon)
3. No heavy grid lines
4. NEVER NEVER NEVER a pie chart (though the WSJ doesn't agree with me)

Misrepresenting the Trend

Misrepresenting the Trend

- The data should fill about 2/3 of the y-axis height
- The line should not be too thick or thin
- The grid shouldn't be thicker than the line
- The y-axis increments shouldn't be awkward
 - Natural increments: (1s, 2s, 5s, 10s, etc.)
 - Awkward increments: (3s, 6s, 8s, 12s, etc.)

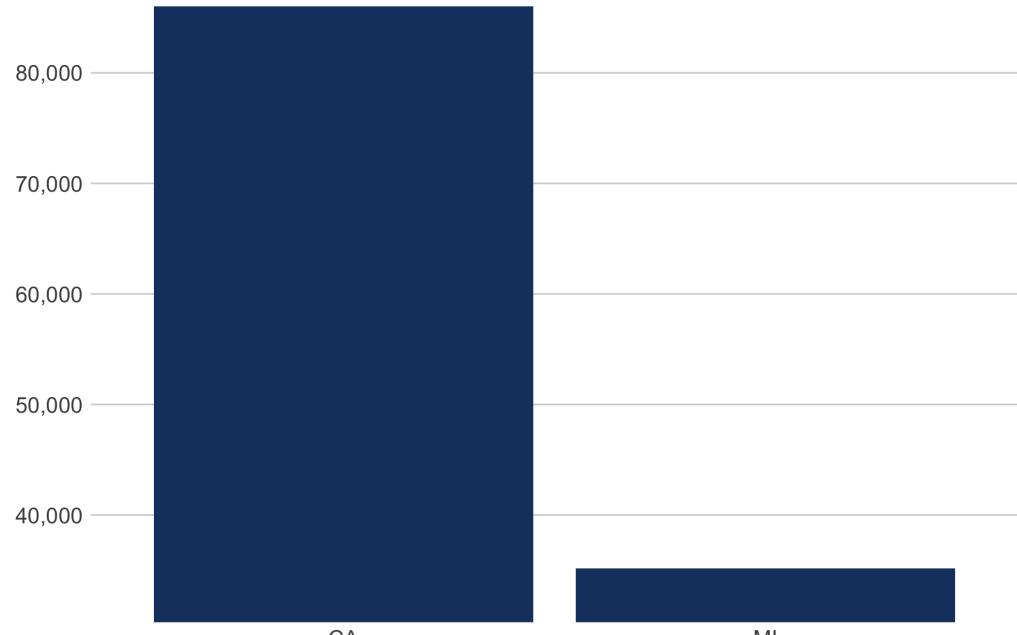
Average US Temperature in Fahrenheit
1901 - 2020



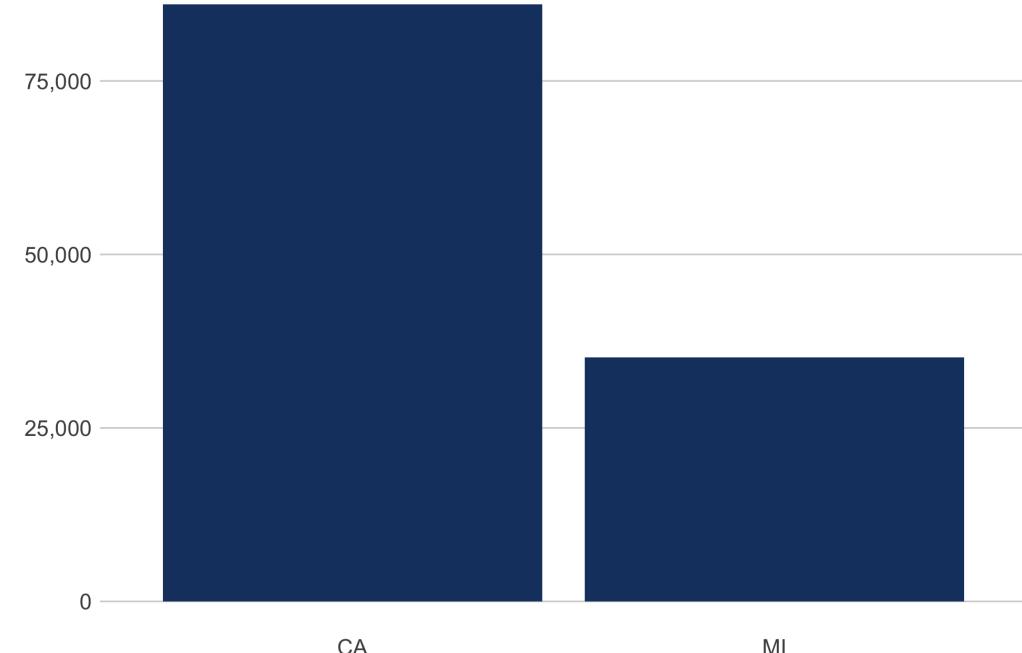
Zero Baseline

- With line charts, as we just saw, we want the y-axis to be the right size
- With bar charts, we are interpreting area; because of this, we always start at 0

Total Deaths due to COVID-19
CA vs MI
Data as of March 2022



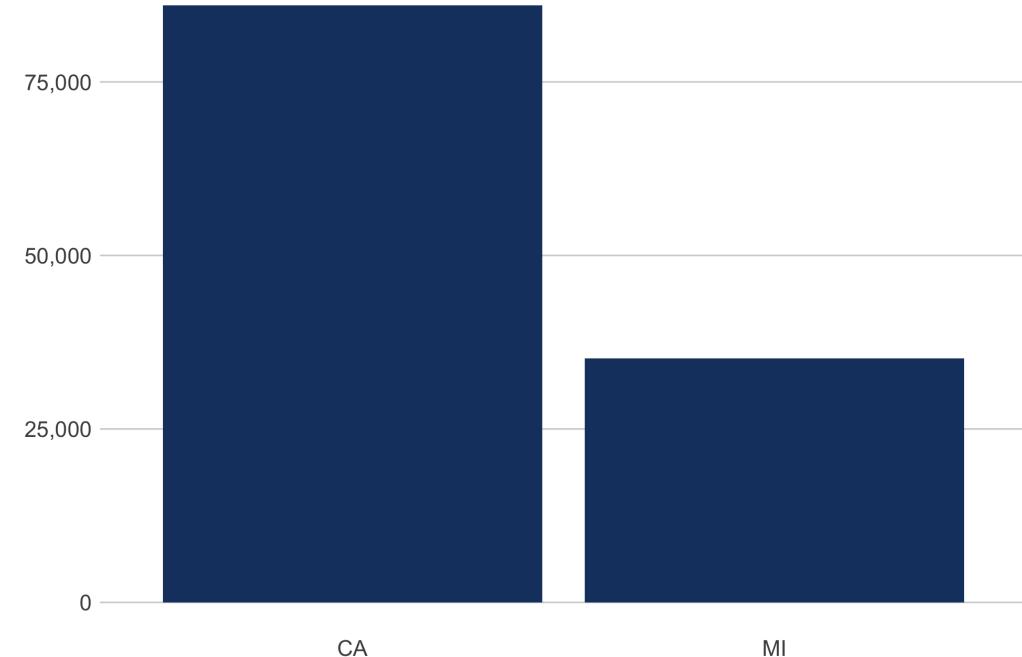
Total Deaths due to COVID-19
CA vs MI
Data as of March 2022



Zero Baseline

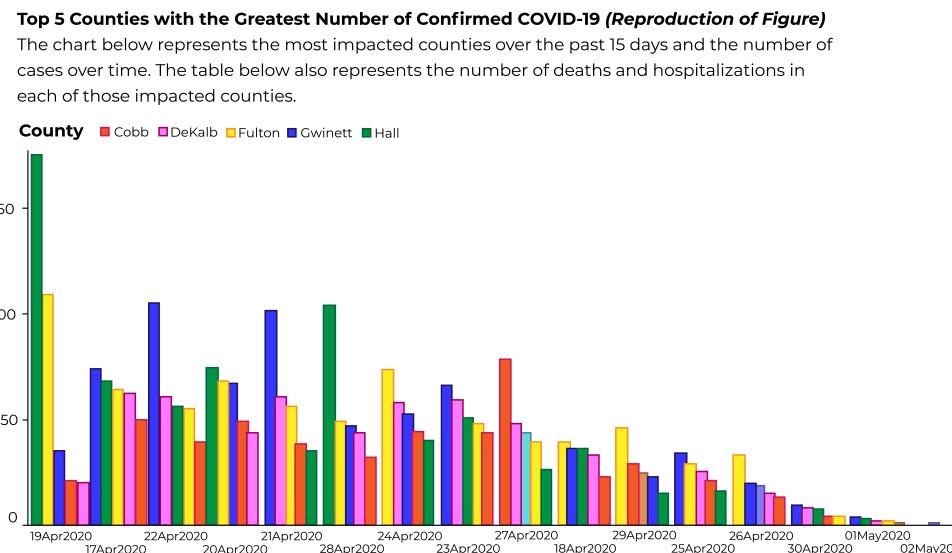
- We can clearly see that MI has about half the total COVID Deaths as CA
- Even this plot may still be misleading or hard to interpret: CA and MI have vastly different populations (should probably be a percentage of population)

Total Deaths due to COVID-19
CA vs MI
Data as of March 2022



Order Bars in the Correct Order

This is not good (look at the dates).



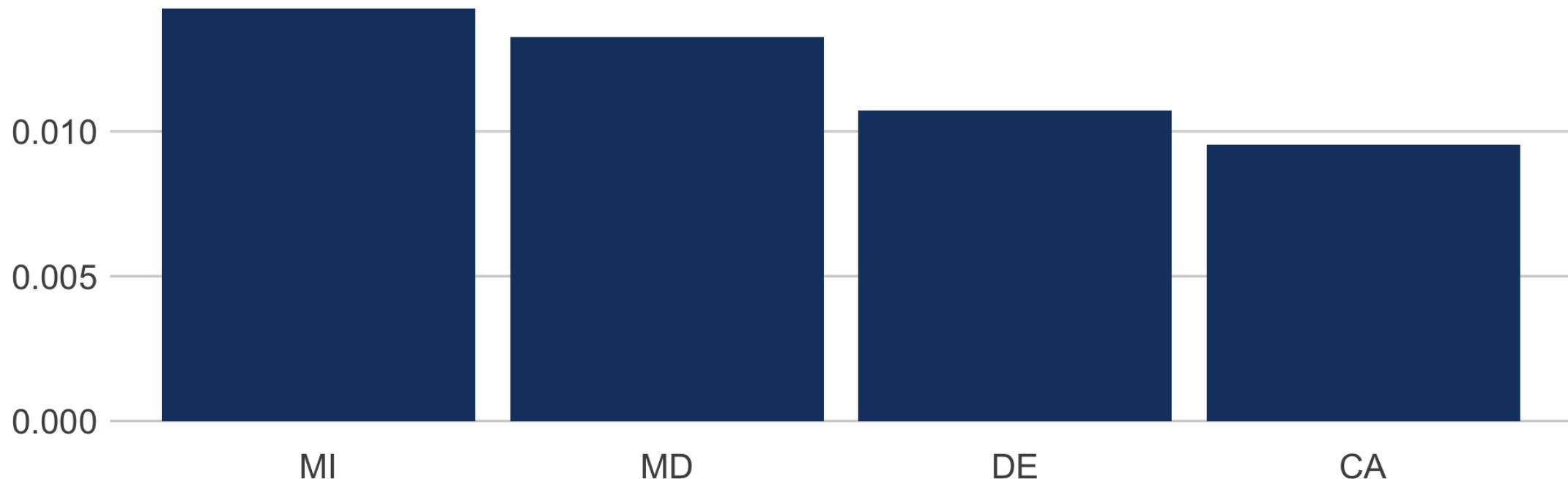
- Ordinal data should stay in order
 - Grades
 - Dates
 - Age groups (Child, Teen, Young Adult, Adult, etc)

Order Bars in the Correct Order

Other categorical data should be reordered so the "highest" is on the left and "lowest" is on the right.

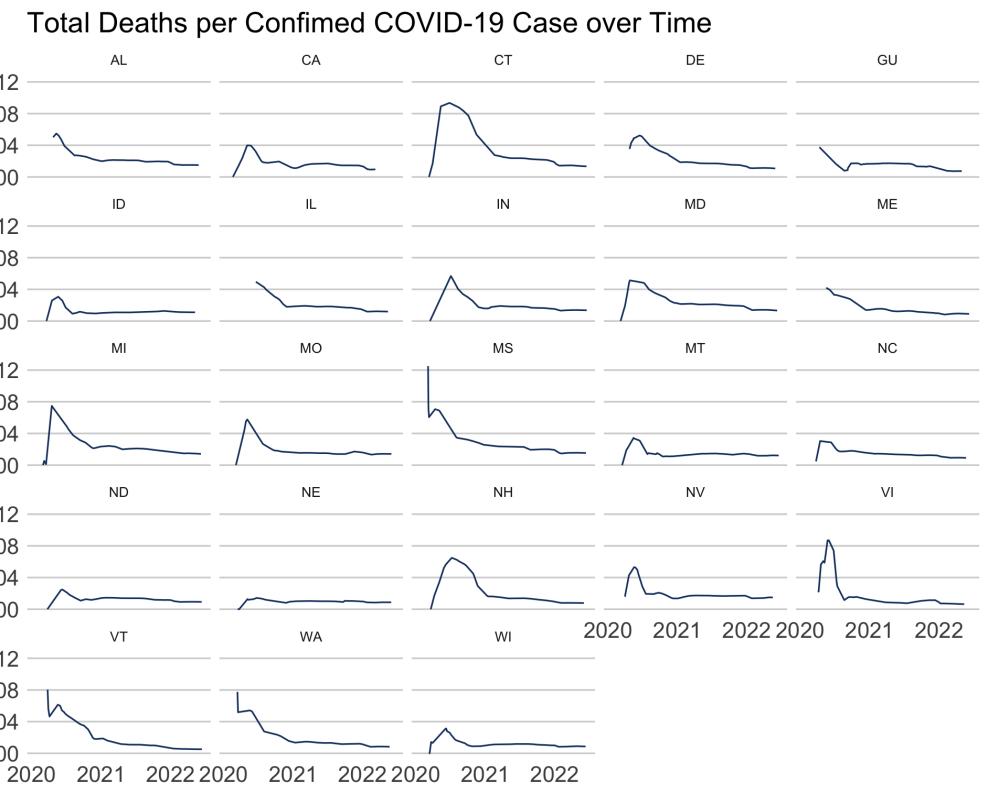
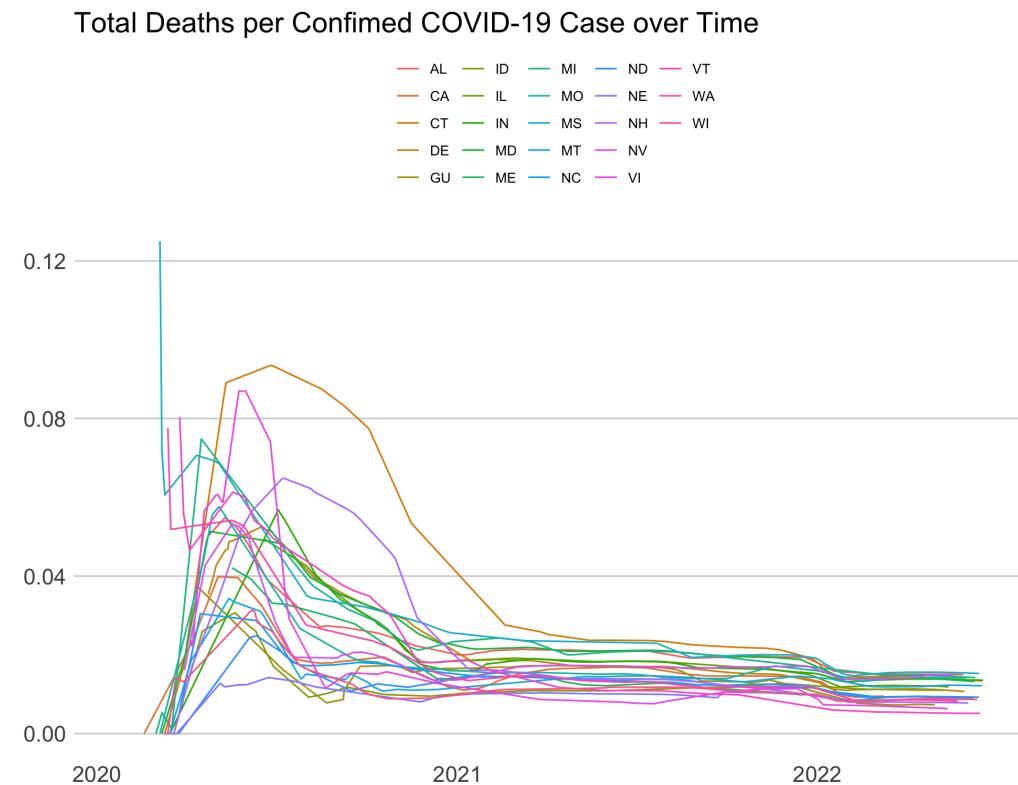
Total Deaths per Confimed COVID-19 Case

Data as of March 2022

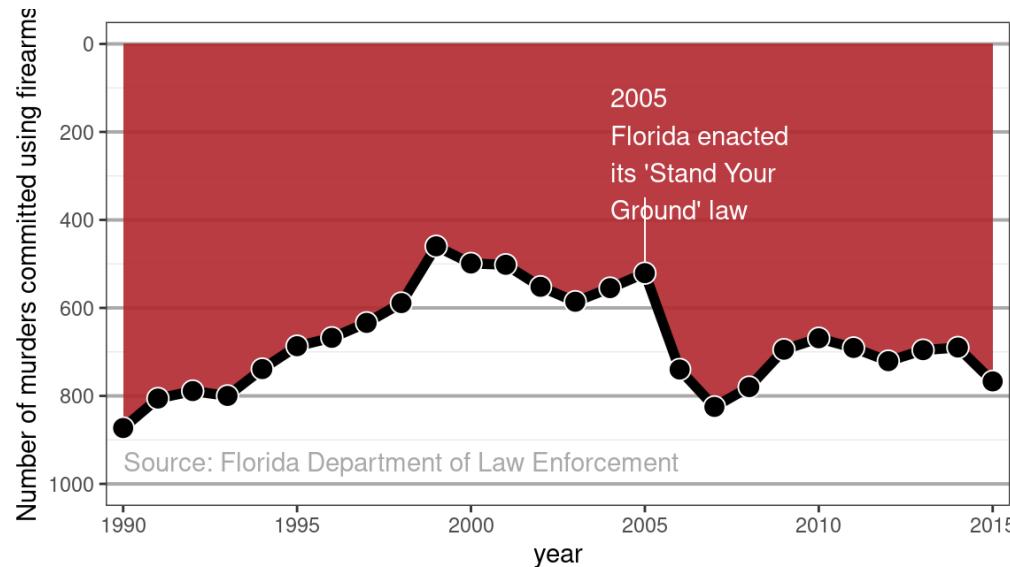


Too many lines on a plot

When you have too many categories, you should split your plot into separate "facets."



Negative Values



This is another example of what not to do.

Any values under the x-axis or left of the y-axis are treated naturally as negative values. It is misleading to flip the axes.

Another Exercise Together!

<https://tinyurl.com/GFA-DV-Final>

Learning More

1. SQL
2. Tableau or other Data Visualization Tool
3. Data analysis in R