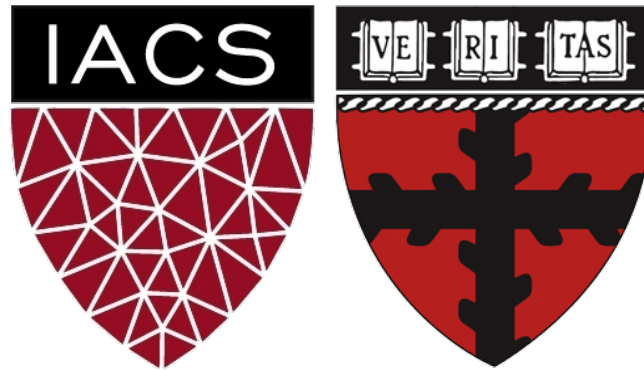# Visualization

## CS109A Introduction to Data Science

Pavlos Protopapas, Natesh Pillai

As the matplotlib thickens …

# ANNOUNCEMENTS

- Homework 4 will be released Oct 20 and is due Oct 27 (Wed) @ 11:59pm

# Learning Objectives

- Understand why visualization/plotting is important

- Learn aspects that tend to make visualizes effective and ineffective

- Feel comfortable designing plots that best convey your message

- Gain experience in producing plots with Python

**Extra Goal** be more cognizant of broader design choices (e.g., typography, s p a c i n g, colors)

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

# Agenda

- EDA Refresher
- Effective Visualization
    - Graphical Integrity
    - Scope
    - Displays
    - Sensible Design

- Communication
    - Motivation
    - Key Considerations

Assume you know a given dataset is credible, complete with the info you want, and has no missing values.
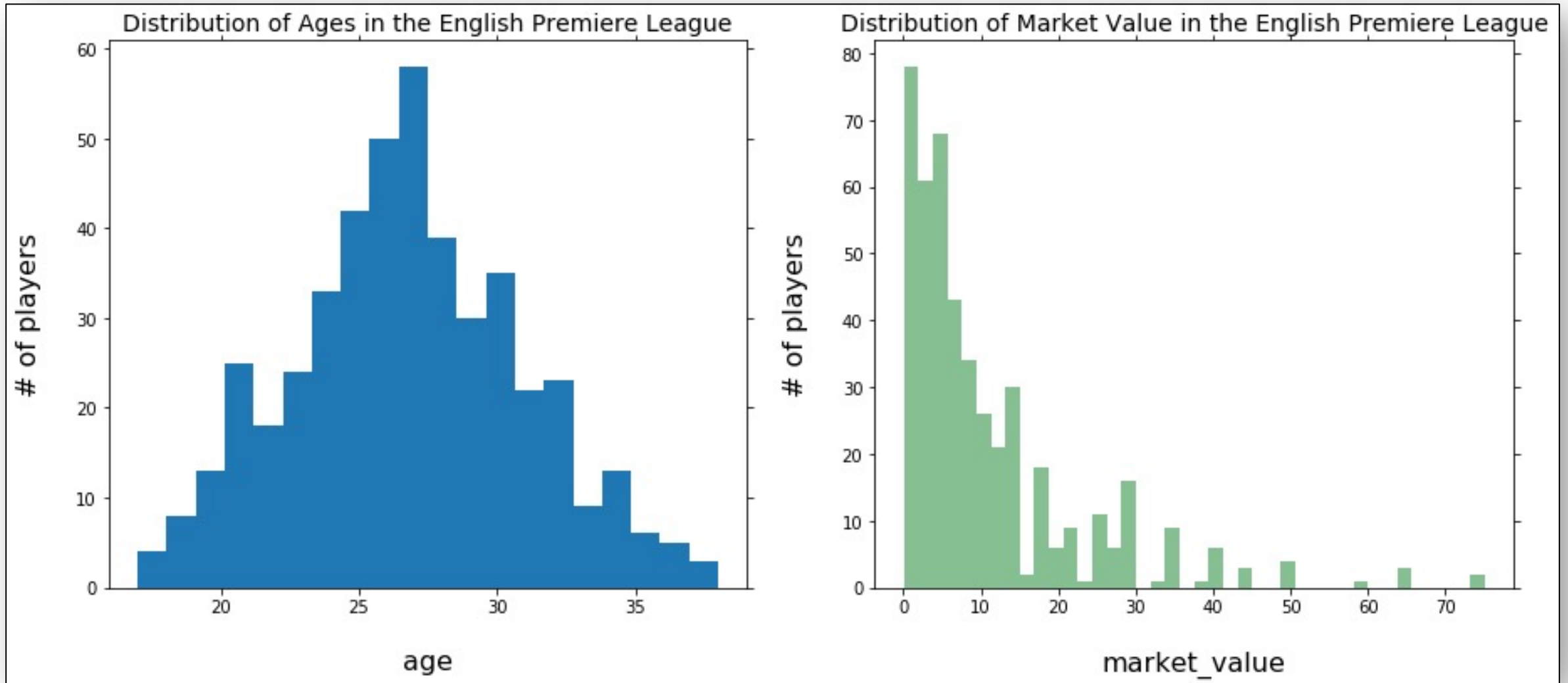
Why do further EDA?

# Purposes of EDA:

- Maximize insight into a dataset

- Uncover underlying structure

- Detect outliers

- Test underlying assumptions

- Develop parsimonious models
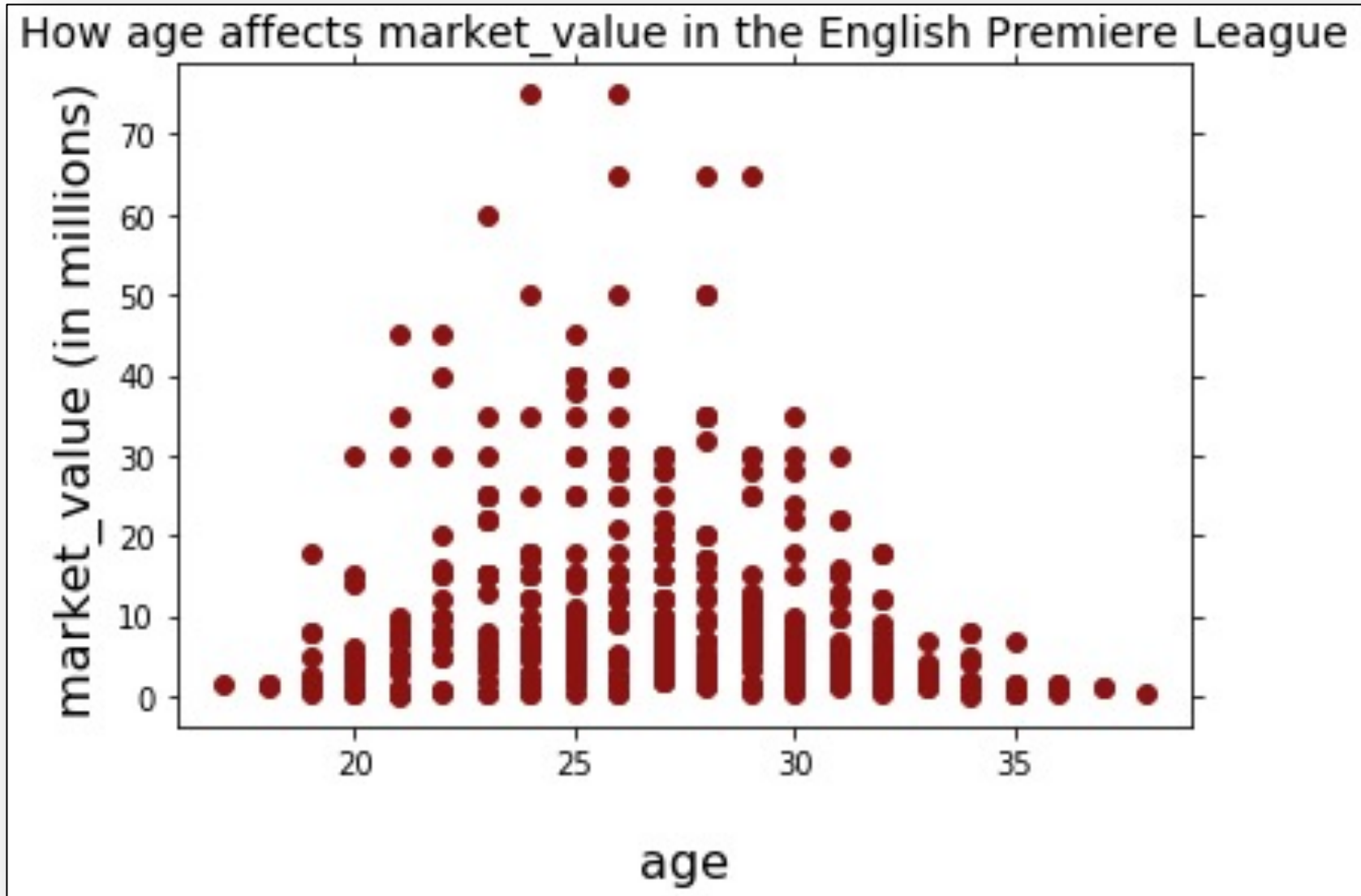
# EDA Refresher: English Premier League

| name | club | age | position | market value |
| --- | --- | --- | --- | --- |
| Alexis Sanchez | Arsenal | 28 | LW | 65 |
| Mesut Ozil | Arsenal | 28 | AM | 50 |
| Petr Cech | Arsenal | 35 | GK | 7 |
| Theo Walcott | Arsenal | 28 | RW | 20 |
| Laurent Koscielny | Arsenal | 31 | CB | 22 |

from www.transfermarkt.us

# EDA Refresher: English Premier League

# EDA Refresher: English Premier League



How age affects market_value in the English Premiere League

# EDA Refresher: English Premier League



How age affects market_value in the English Premiere League

Are the outliers legit?

# EDA Refresher: English Premier League

```python
league_df.loc[league_df['age']<20][['name', 'club', 'age', 'position', \
    'market_value']].sort_values(by="age")
```

|     | name | club | age | position | market_value |
|-----|------|------|-----|----------|--------------|
| 233 | Ben Woodburn | Liverpool | 17 | LW | 1.50 |
| 231 | Trent Alexander-Arnold | Liverpool | 18 | RB | 1.50 |
| 350 | Josh Tymon | Stoke+City | 18 | LB | 1.00 |
| 435 | Jonathan Leko | West+Brom | 18 | RW | 1.50 |
| 147 | Tom Davies | Everton | 19 | CM | 8.00 |
| 155 | Ademola Lookman | Everton | 19 | LW | 5.00 |
| 239 | Dominic Solanke | Liverpool | 19 | CF | 2.00 |
| 270 | Marcus Rashford | Manchester+United | 19 | CF | 18.00 |
| 281 | Axel Tuanzebe | Manchester+United | 19 | CB | 1.00 |
| 282 | Timothy Fosu-Mensah | Manchester+United | 19 | DM | 2.50 |
| 375 | Tammy Abraham | Swansea | 19 | CF | 8.00 |
| 436 | Sam Field | West+Brom | 19 | CM | 0.25 |

# EDA Refresher: English Premier League

```python
league_df.loc[league_df['market_value']>=60][['name', 'club', \
    'age', 'position', 'market_value']].sort_values(by="age")
```

|     | name | club | age | position | market_value |
|-----|------|------|-----|----------|--------------|
| 377 | Harry Kane | Tottenham | 23 | CF | 60.0 |
| 263 | Paul Pogba | Manchester+United | 24 | CM | 75.0 |
| 92  | Eden Hazard | Chelsea | 26 | LW | 75.0 |
| 240 | Kevin De Bruyne | Manchester+City | 26 | AM | 65.0 |
| 0   | Alexis Sanchez | Arsenal | 28 | LW | 65.0 |
| 241 | Sergio Aguero | Manchester+City | 29 | CF | 65.0 |

# EDA Refresher: English Premier League

```python
league_df.loc[league_df['market_value']>=60][['name', 'club', \
    'age', 'position', 'market_value']].sort_values(by="age")
```

|     | name | club | age | position | market_value |
| --- | --- | --- | --- | --- | --- |
| 377 | Harry Kane | Tottenham | 23 | CF | 60.0 |
| 263 | Paul Pogba | Manchester+United | 24 | CM | 75.0 |
| 92 | Eden Hazard | Chelsea | 26 | LW | 75.0 |
| 240 | Kevin De Bruyne | Manchester+City | 26 | AM | 65.0 |
| 0 | Alexis Sanchez | Arsenal | 28 | LW | 65.0 |
| 241 | Sergio Aguero | Manchester+City | 29 | CF | 65.0 |

# Agenda

- **EDA Refresher**
- Effective Visualization
  - Graphical Integrity
  - Scope
  - Displays
  - Sensible Design
- Communication
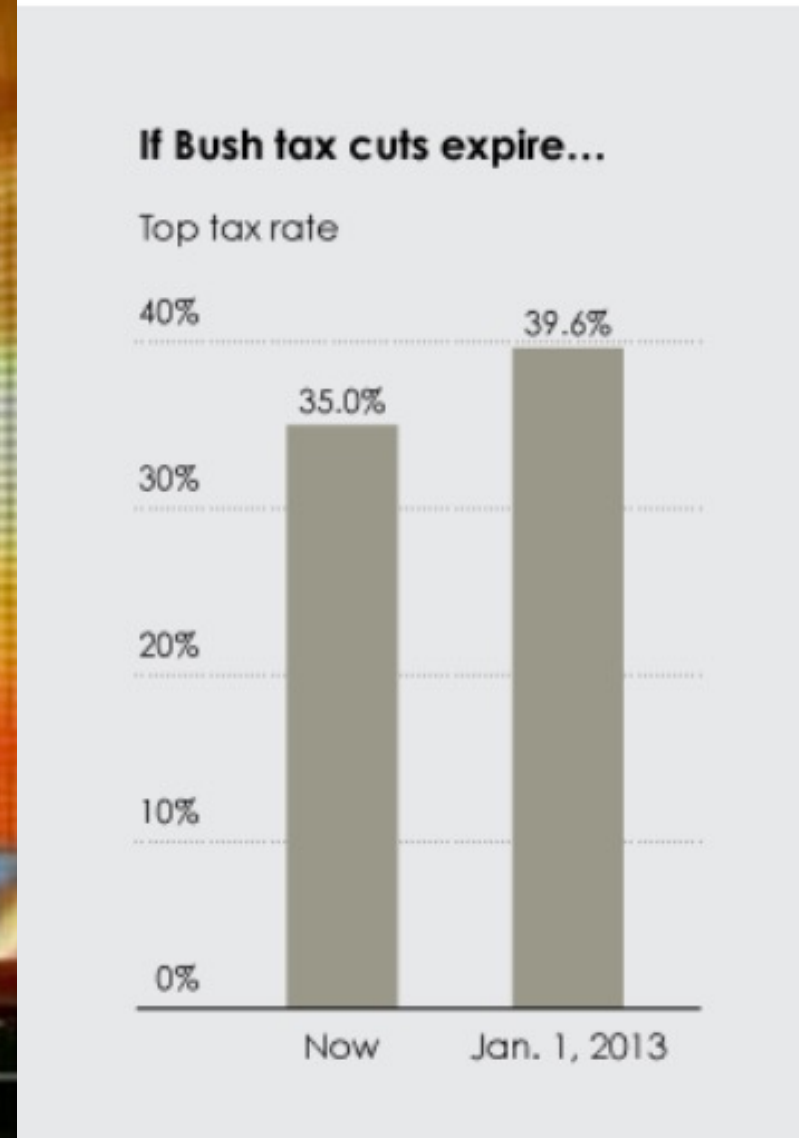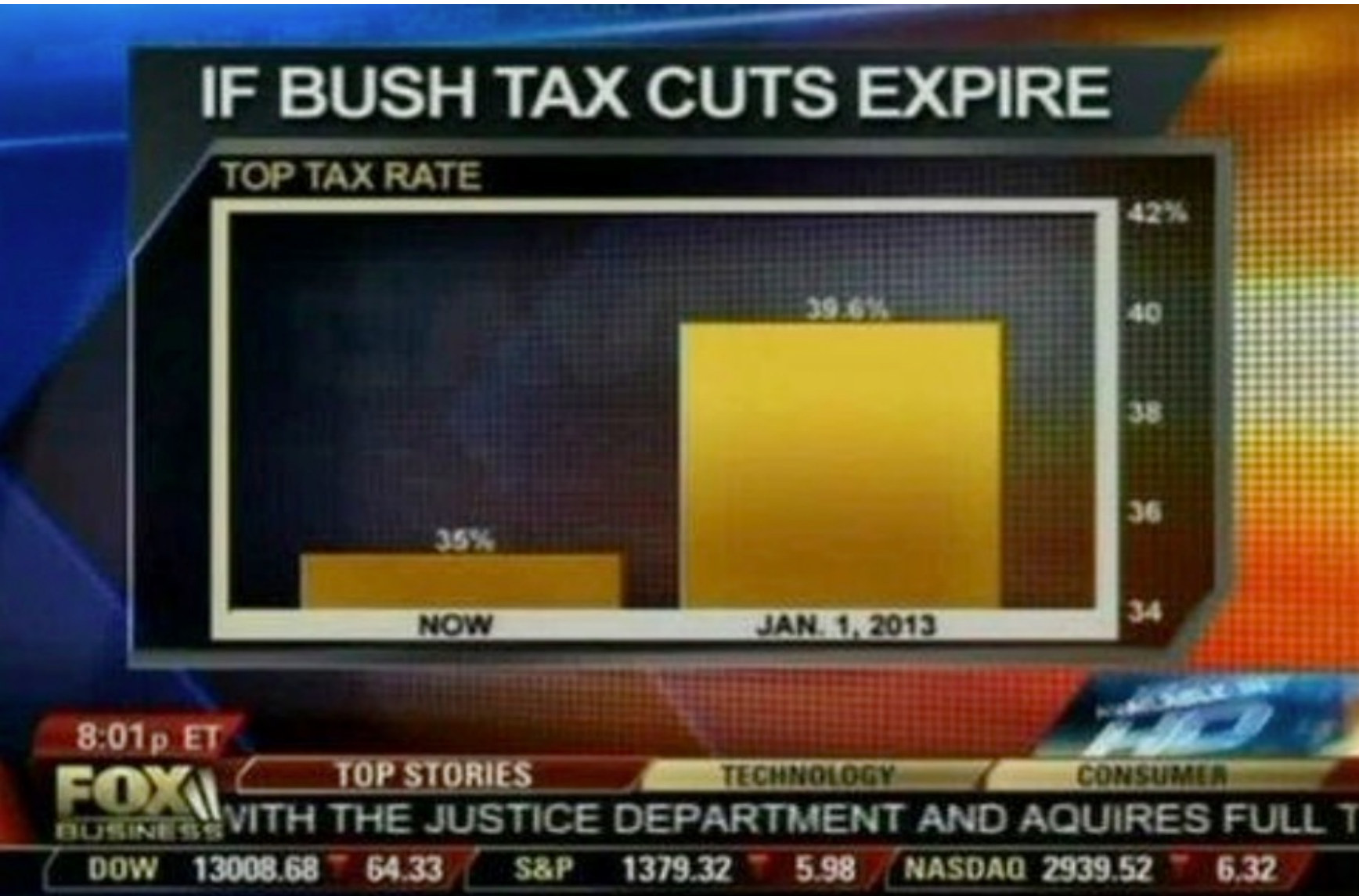  - Motivation
  - Key Considerations

**DISCLAIMER.** Some of these examples involve political data. In no way should this be taken as a signal of my support or endorsement in any beliefs; the point is merely to convey good and bad choices when it comes to effective visualization.
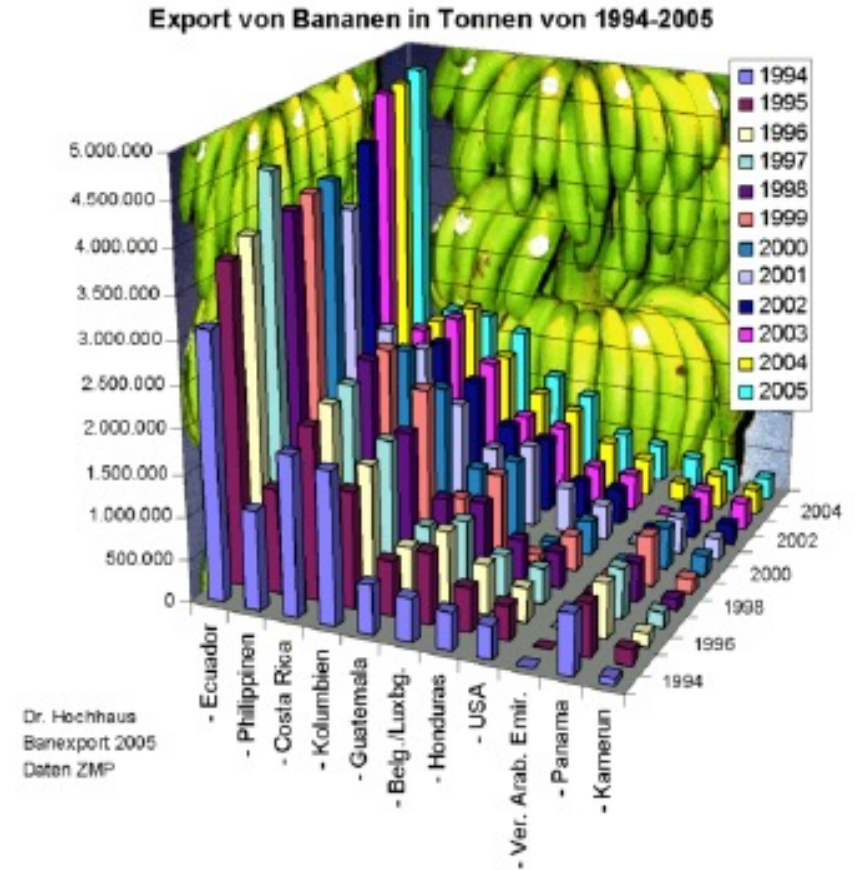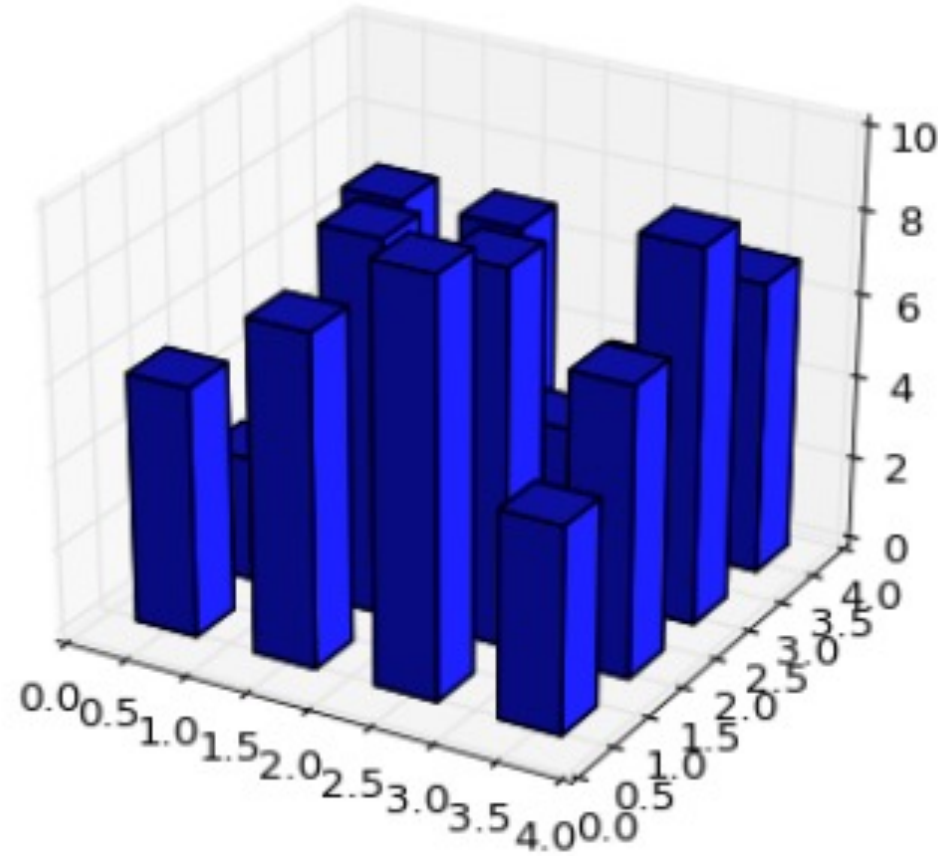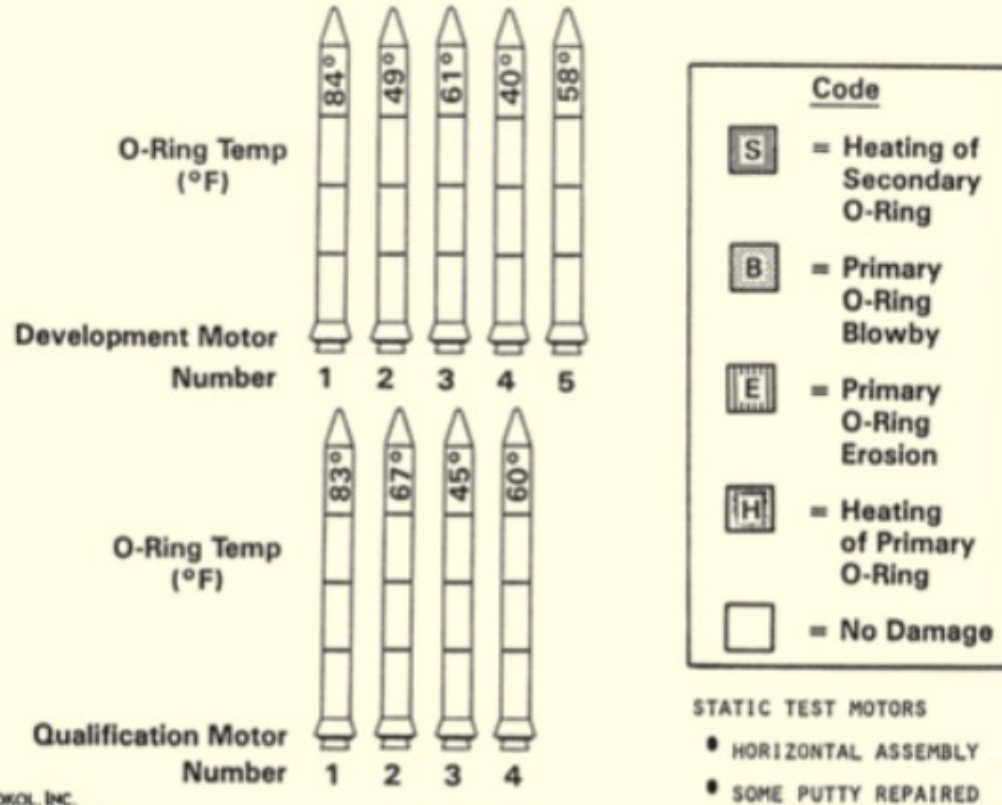
# Agenda

- EDA Refresher
- Effective Visualization
  - Graphical Integrity
  - Scope
  - Displays
  - Sensible Design

- Communication
  - Motivation
  - Key Considerations

**DISCLAIMER.** Some of these examples involve political data. In no way should this be taken as a signal of my support or endorsement in any beliefs; the point is merely to convey good and bad choices when it comes to effective visualization.

# Agenda

- **EDA Refresher**
- **Effective Visualization**
  - Graphical Integrity
  - Scope
  - Displays
  - Sensible Design
- **Communication**
  - Motivation
  - Key Considerations

**DISCLAIMER.** Some of these examples involve political data. In no way should this be taken as a signal of my support or endorsement in any beliefs; the point is merely to convey good and bad choices when it comes to effective visualization.

**Treasury Quarterly Net Marketable Borrowing**
"Net Cash"
Fiscal Quarter

**Figure 10**

Figure 1 **Mortality and disease burden (DALYs) in women aged 15–44 years by region and broad causes, 2004**

* High-income countries are excluded from the regional groups.
Source: World Health Organization.[1]

Coupons Maturing*
February 15, 2010-November 15, 2039

*Based on coupon securities outstanding as of January 21, 2010

**Figure 5.2** Mean prevalence rates of *Cryptosporidium* oocysts by animal species.
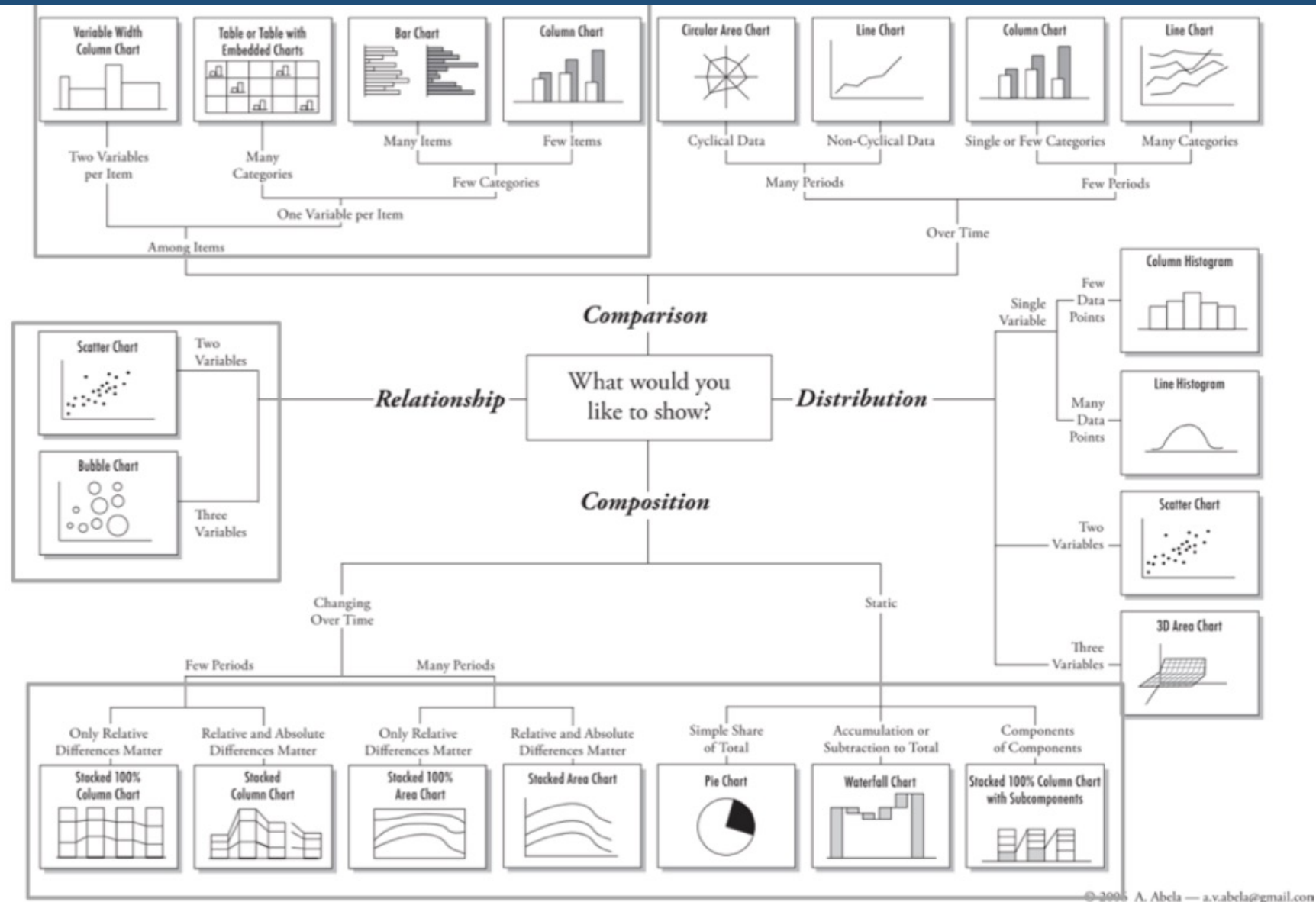
# Agenda

- 🟥 EDA Refresher
- 🟦 Effective Visualization
  - ■ Graphical Integrity
  - ■ Scope
  - ■ Displays
  - ■ Sensible Design

- 🟪 Communication
  - ■ Motivation
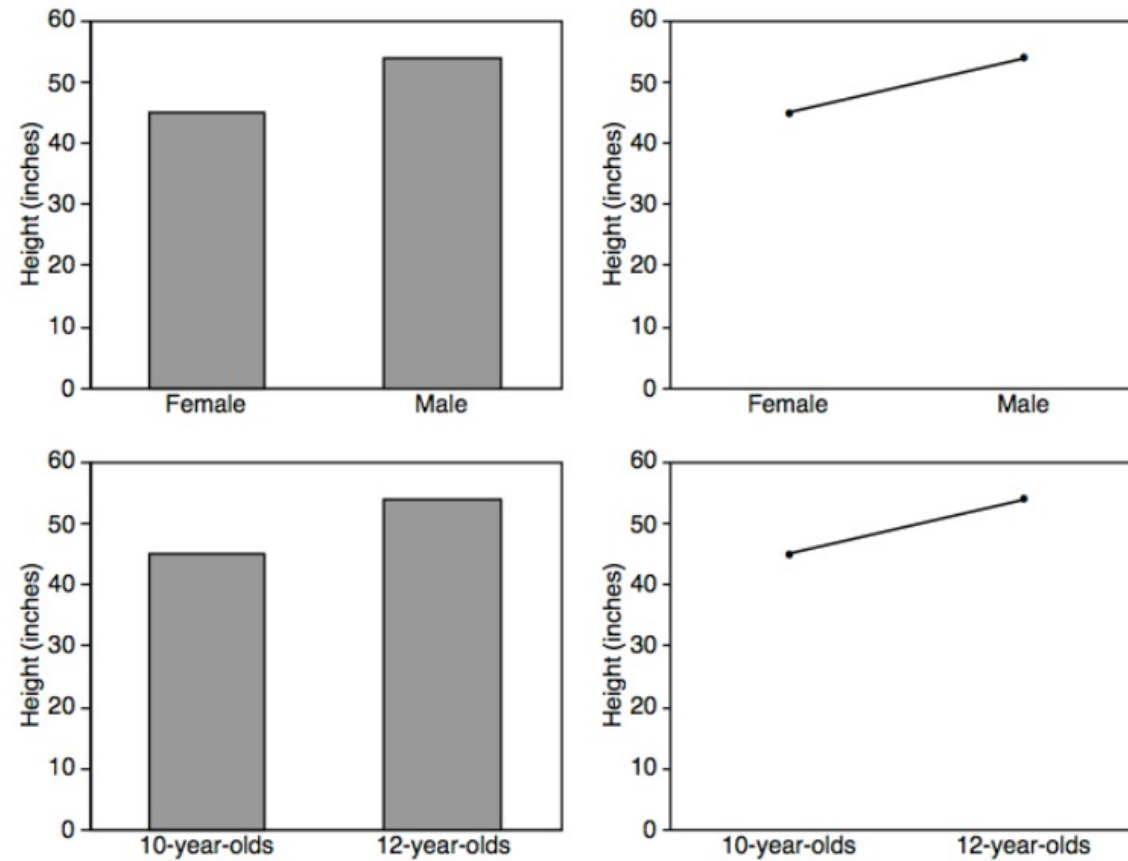  - ■ Key Considerations

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

UNEMPLOYMENT LEVEL BY RANDOM QUARTER

**October 1, 2019**

Alberto Cairo • University of Miami • www.thefunctionalart.com • Twitter: @albertocairo

**Lesson:** be proportional

What non-scientists are not aware of (cone is just 66% probability)

Alberto Cairo • University of Miami • www.thefunctionalart.com • Twitter: @albertocairo

Hurricane
**CAIRO**
(category 5)

What we could be showing instead

Alberto Cairo • University of Miami • www.thefunctionalart.com • Twitter: @albertocairo

**Lesson II:** include uncertainty when possible

Counties with the LOWEST
kidney cancer death rates
(1980-1989)

Counties with the HIGHEST
kidney cancer death rates
(1980-1989)

**Lesson III:** plot all the data

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

# Agenda

- EDA Refresher
- Effective Visualization
    - Graphical Integrity
    - Scope
    - Displays
    - Sensible Design

- Communication
    - Motivation
    - Key Considerations

Export von Bananen in Tonnen von 1994-2005

History of O-Ring Damage in Field Joints

History of O-Ring Damage in Field Joints (Cont)

**Lesson IV: keep it simple… enough**

*"You should have stayed with the soup question. The object of a question is to obtain information that matters only to us"*

-- Sean Connery in Finding Forrester (movie)

- Making plots is effectively providing an answer to an implicit question

- You get to pick the way to express the answer

- Ensure the answer doesn't leave the viewer with uncertainty as to **what it's answering** or the **completeness of the answer**

- A good plot should invoke and inspire new questions

# Agenda

- EDA Refresher
- Effective Visualization
    - Graphical Integrity
    - Scope
    - Displays
    - Sensible Design

- Communication
    - Motivation
    - Key Considerations

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

http://extremepresentation.typepad.com/blog/files/choosing_a_good_chart.pdf

Bars vs. Lines

http://xkcd.com/388/

**London Cholera Epidemic**

-- Edward Tufte,
Visual and Statistical Thinking

friday @ 10:30am

sunday @ 5pm

saturday @ 2pm

sunday @ 2pm

friday @ 11:30am

1st and 3rd wed

2nd and 4th wed

2nd and 4th wed

saturday @ 3pm

sunday @ 4:30pm

friday @ 7:30pm

saturday @ 8pm

friday @ 8pm

saturday @ 6pm

saturday @ 7

saturday @ 6

Bin Width

binwidth = 0.1

binwidth = 0.01

https://www.autodeskresearch.com/publications/samestats

GROUP

THE TRILOGY METER

#1 In A Series of Pop Cultural Charts      DANMETH.COM

# Displays: scatter plot matrix



Scatterplot Matrix (SPLOM) for Diabetes Dataset
Data source: [1]

# How do you feel about doing science?

Table

| Interest | Before | After |
|---|---|---|
| Excited | 19 | 38 |
| Kind of interested | 25 | 30 |
| OK | 40 | 14 |
| Not great | 5 | 6 |
| Bored | 11 | 12 |

Data courtesy of Cole Nussbaumer

Stacked bar, not very useful

Data Transposed Bar Chart

Difference Bar Chart

How do you feel about doing science?

After the pilot program,

# 68%

of kids expressed interest towards science,
compared to 44% going into the program.

# Time for a game

🎉

# How much longer?

# How much longer?



4x

# How much steeper slope?

A          B

How much steeper slope?

A          B

4x

# How much larger area?

A

B

How much larger area?

A

B

4.5x

# How much darker?

A          B

# How much darker?

**2x**

A          B

# How much bigger value?
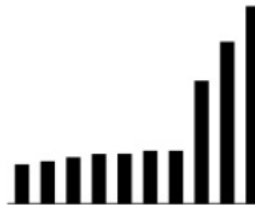


A          B

2                              16

# How much bigger value?



A          B

4x

2          16

C. Mulbrandon
VisualizingEconomics.com

Most Effective

VisualizingEconomics.com

# Agenda

- EDA Refresher
- Effective Visualization
  - Graphical Integrity
  - Scope
  - Displays
  - Sensible Design

- Communication
  - Motivation
  - Key Considerations

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

# Colors for Categories
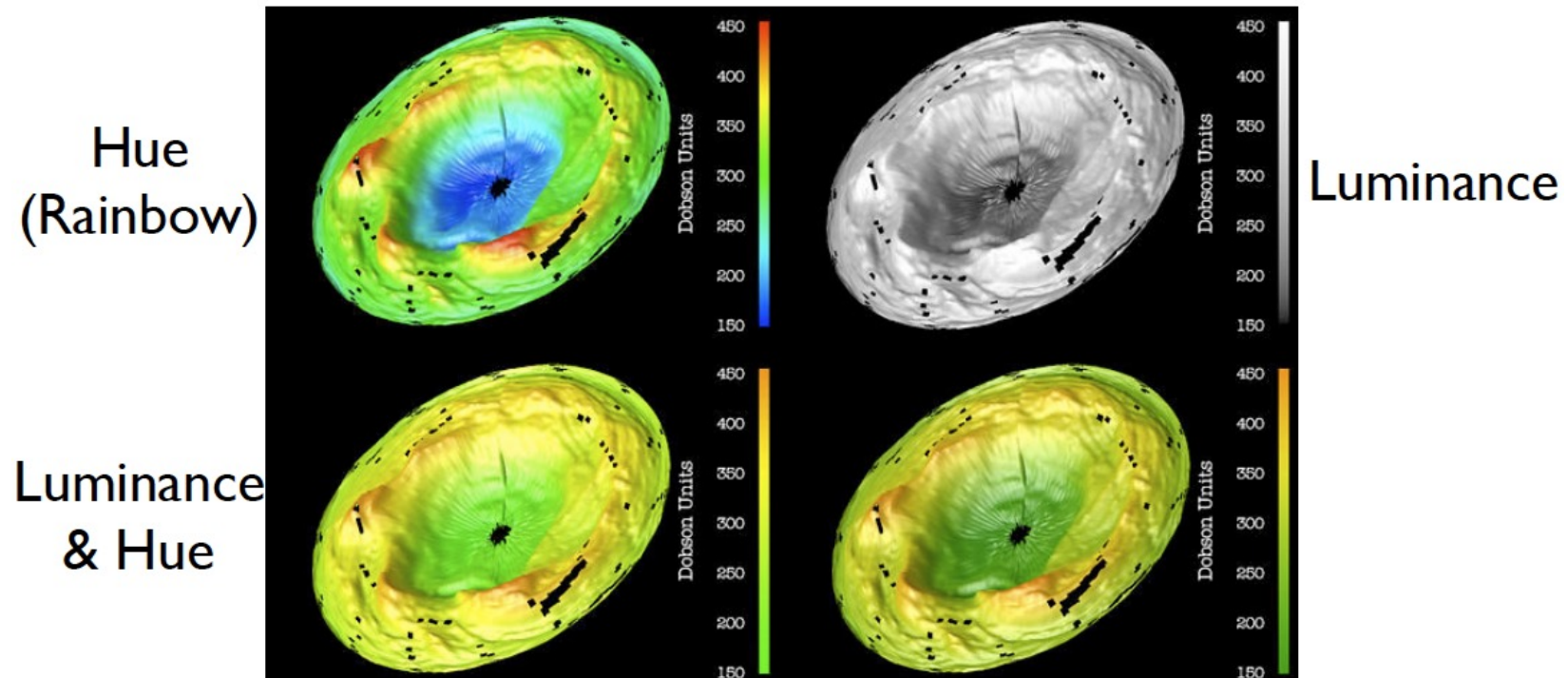
## Do not use more than 5-8 colors at once

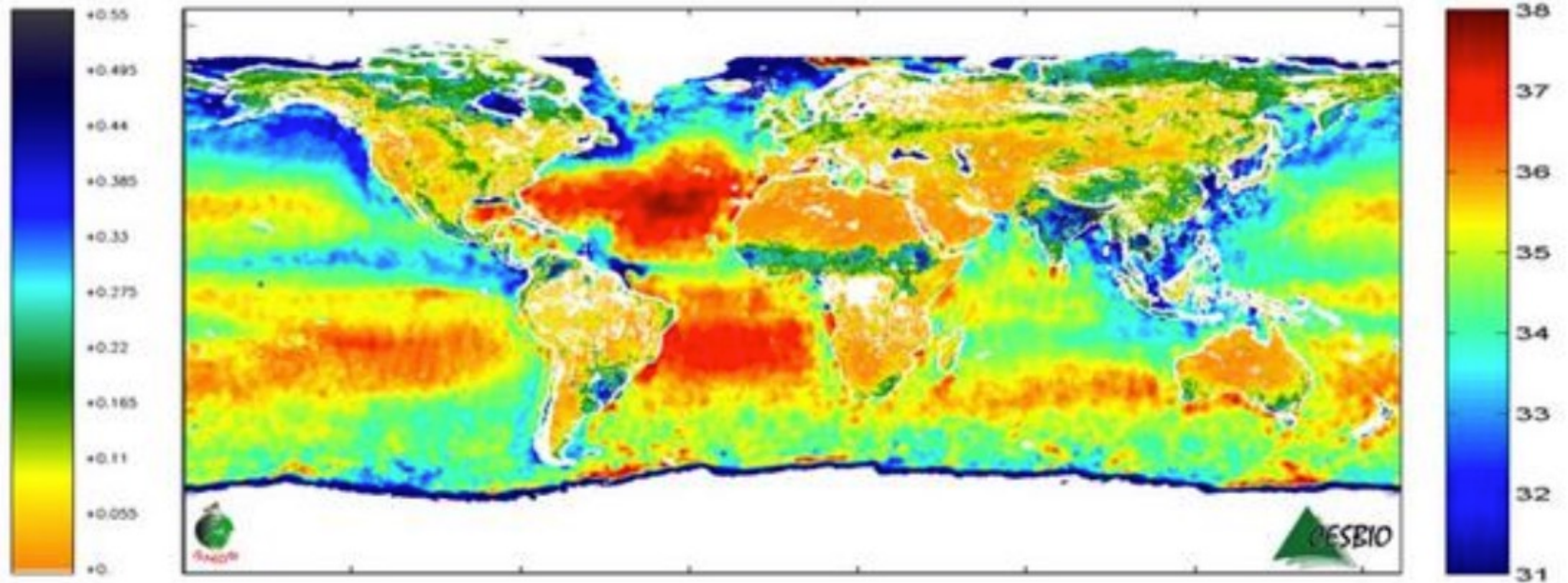Colors for Ordinal Data

Vary luminance and saturation

Zeilis et al, 2009, "Escaping RGBland: Selecting Colors for Statistical Graphics"
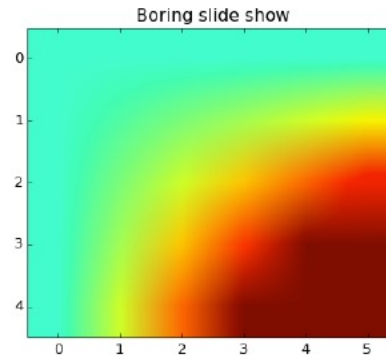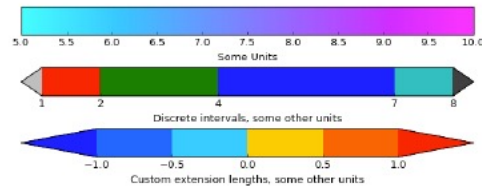
Colors for Quantitative Data

Rogowitz and Treinish, Why should engineers and scientists be worried about color?
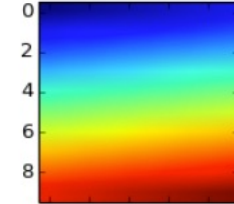
Avoid Rainbow Colors!

Diverging Palette for Quantitative or Ordinal

Sequential Palette for Densities

Color Blindness

Protanope     Deuteranope     Tritanope

Red / green deficiencies     Blue / Yellow deficiency

Color Blindness

Normal — Protanope — Deuteranope — Lightness

Great sites for selecting color schemes:

- http://colorbrewer2.org
- https://coolors.co

## Color Brewer



## Coolors.co

How much do you trust this text?

How much do you trust this text?

How much do you trust this text?

How much do you trust this text?

The everyday items that are designed the best are the ones that we never have to think about how to use/interact with it.

Can you think of examples?



Highly recommended

# Lesson V: design matters

# Agenda

- EDA Refresher
- Effective Visualization
  - Graphical Integrity
  - Scope
  - Displays
  - Sensible Design

- Communication
  - Motivation
  - Key Considerations

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

## Analyze (Exploratory)

- Explore the data

- Assess a situation

- Determine how to proceed

- Decide what to do

## Communicate (Explanatory)

- Present data and ideas

- Explain and inform

- Provide evidence and support

- Influence and persuade

# The Persuasive Power of Data Visualization

Anshul Vikram Pandey
*New York University*

Anjali Manivannan
*New York University*

Oded Nov
*New York University*

Margaret L. Satterthwaite
*NYU School of Law*, satterth@exchange.law.nyu.edu

Enrico Bertini
*New York University*

After looking into common effects in attitude formation and change we searched for specific mentions to the graphical appearance of charts as a driver for persuasion. Some of the comments we collected seem to back up the findings we found in our results. Some participants explicitly mention the charts as being the main reason for their change: *"I already knew that increased incarceration didn't lower crime, but I wasn't sure of the statistics. To see it on the graphs is really eye opening."; "I was influenced by the bar graph showing the reasons why the survey respondents played video games."*; "I would not know exact numbers on this issue - the graphs gave a visual and helped identify the numbers"; "Seeing the graphs conflicted with my previous opinion, so I feel like I need to reevaluate my stance in a way."

It is also important to mention that the graphical appearance of charts is not the only factor that has a strong impact on people's attitude. In our collected feedback, we found numerous references to statistics and numbers, suggesting that mere exposure to data does have a persuasive effect – maybe at least partially due to the increased sense of objectivity evidence supported by numbers carries. We found comments like: *"It was concrete data that seemed compelling.; "Seeing numbers is a good indicator of change rather than just reading what someone has to say"; "It showed a large amount of different sources, which made it more credible"*. More research is needed to disentangle what kind of specific effects each of these components have on persuasion.

http://lsr.nellco.org/cgi/viewcontent.cgi?article=1476&context=nyu_plltwp

# Minard's Graphic on Napoleon's Russia Campaign



Figurative Map of the successive losses in men of the French Army in the Russian campaign 1812 ~ 1813

Drawn by M. Minard, Inspector General of Bridges and Roads (retired).

Paris, November 20, 1869.

The numbers of men present are represented by the widths of the colored zones at a rate of one millimeter for every ten thousand men; they are further written across the zones. 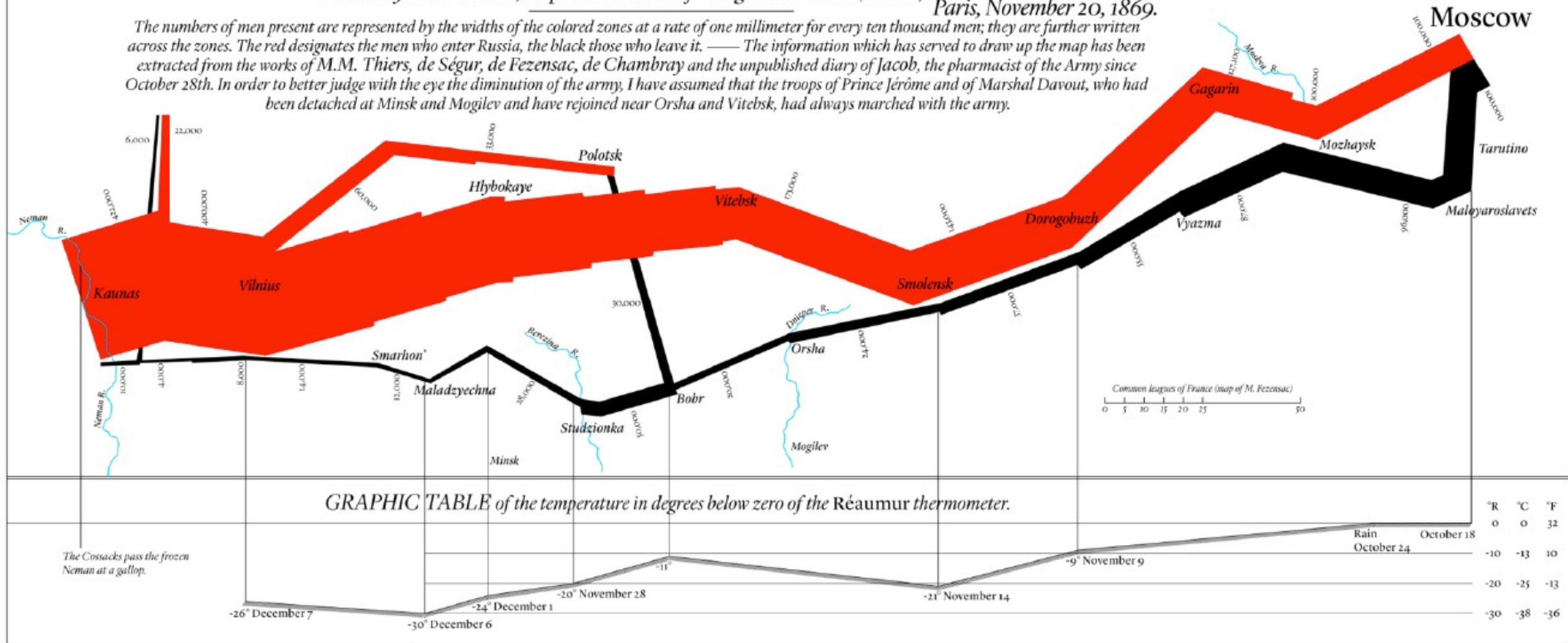The red designates the men who enter Russia, the black those who leave it. —— The information which has served to draw up the map has been extracted from the works of M.M. Thiers, de Ségur, de Fezensac, de Chambray and the unpublished diary of Jacob, the pharmacist of the Army since October 28th. In order to better judge with the eye the diminution of the army, I have assumed that the troops of Prince Jérôme and of Marshal Davout, who had been detached at Minsk and Mogilev and have rejoined near Orsha and Vitebsk, had always marched with the army.

GRAPHIC TABLE of the temperature in degrees below zero of the Réaumur thermometer.

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

  - Motivation

  - Key Considerations

# Agenda

- EDA Refresher

- Effective Visualization

  - Graphical Integrity

  - Scope

  - Displays

  - Sensible Design

- Communication

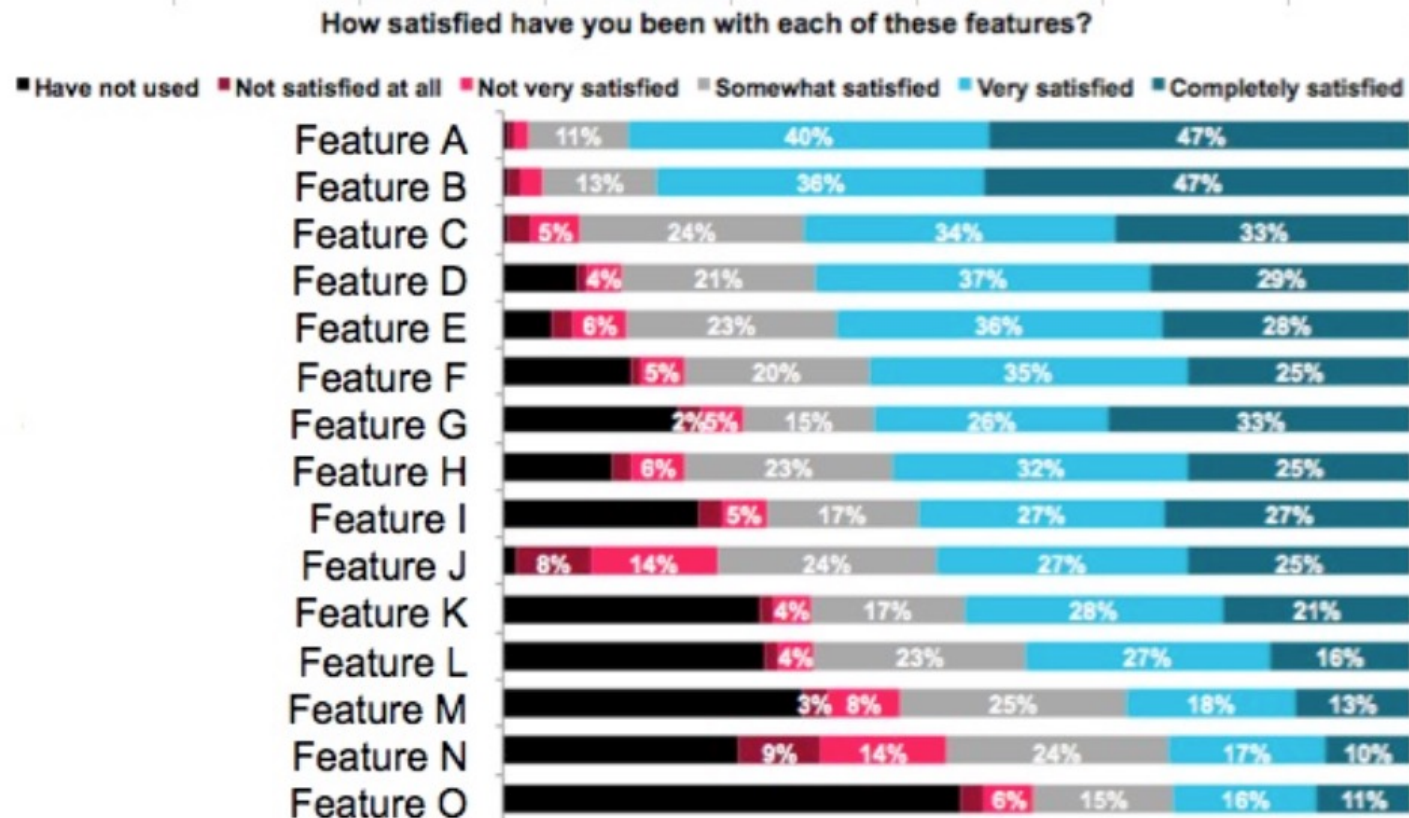  - Motivation

  - Key Considerations

## Key Considerations

- Who is your **audience**

- What questions are you answering?

- Why should the audience care?

- What are your major insights and surprises?

- What change do you want to affect?

# Don't make them think!

- The audience does not want to spend cognitive energy on dissecting and decoding your intended message.

- Lead them through the **major** steps of your story.

- Point out interesting key facts and insights using **captions** and **annotations**

# Don't Bury the Lead



How satisfied have you been with each of these features?

■ Have not used  ■ Not satisfied at all  ■ Not very satisfied  ■ Somewhat satisfied  ■ Very satisfied  ■ Completely satisfied

# Communication

Andy Cotgreave, Tableau

# Final Takeaways

- How you choose to display your data greatly influences how people interpret the data

- Humans are visual, *emotional* creations; make graphs that don't make others **feel** confused, insulted, etc.

- Your graphs should illicit good feelings and effectively convey your narrative

# Suggested Python Packages

- Matplotlib

- Seaborn

- plotly

- ggplot

# Further Good Examples

- [https://www.nytimes.com/](https://www.nytimes.com/) tends to have incredibly high-quality visualizations that convey information seamlessly

- [https://www.reddit.com/r/dataisbeautiful/](https://www.reddit.com/r/dataisbeautiful/)

- [fivethirtyeight.com](fivethirtyeight.com)

## Cases and deaths by state and county

This table is sorted by places with the most cases per 100,000 residents in the last seven days. Charts are colored to reveal when outbreaks emerged.

**Cases** | Deaths | Search counties

| | TOTAL CASES | PER 100,000 | CASES IN LAST 7 DAYS | ▼ PER 100,000 | WEEKLY CASES PER CAPITA |
|---|---|---|---|---|---|
| + North Dakota MAP » | 23,553 | 3,091 | 2,754 | 361 | |
| + South Dakota MAP » | 24,418 | 2,760 | 2,853 | 322 | |
| + Wisconsin MAP » | 139,941 | 2,403 | 17,769 | 305 | |
| + Montana MAP » | 14,738 | 1,379 | 2,547 | 238 | |
| + Utah MAP » | 77,618 | 2,421 | 6,675 | 208 | |
| + Iowa MAP » | 92,584 | 2,934 | 6,394 | 203 | |
| + Nebraska MAP » | 47,807 | 2,471 | 3,807 | 197 | |
| Guam | 3,586 | 2,128 | 331 | 196 | |
| + Arkansas MAP » | 87,013 | 2,883 | 5,770 | 191 | |
| + Idaho MAP » | 43,964 | 2,460 | 3,291 | 184 | |

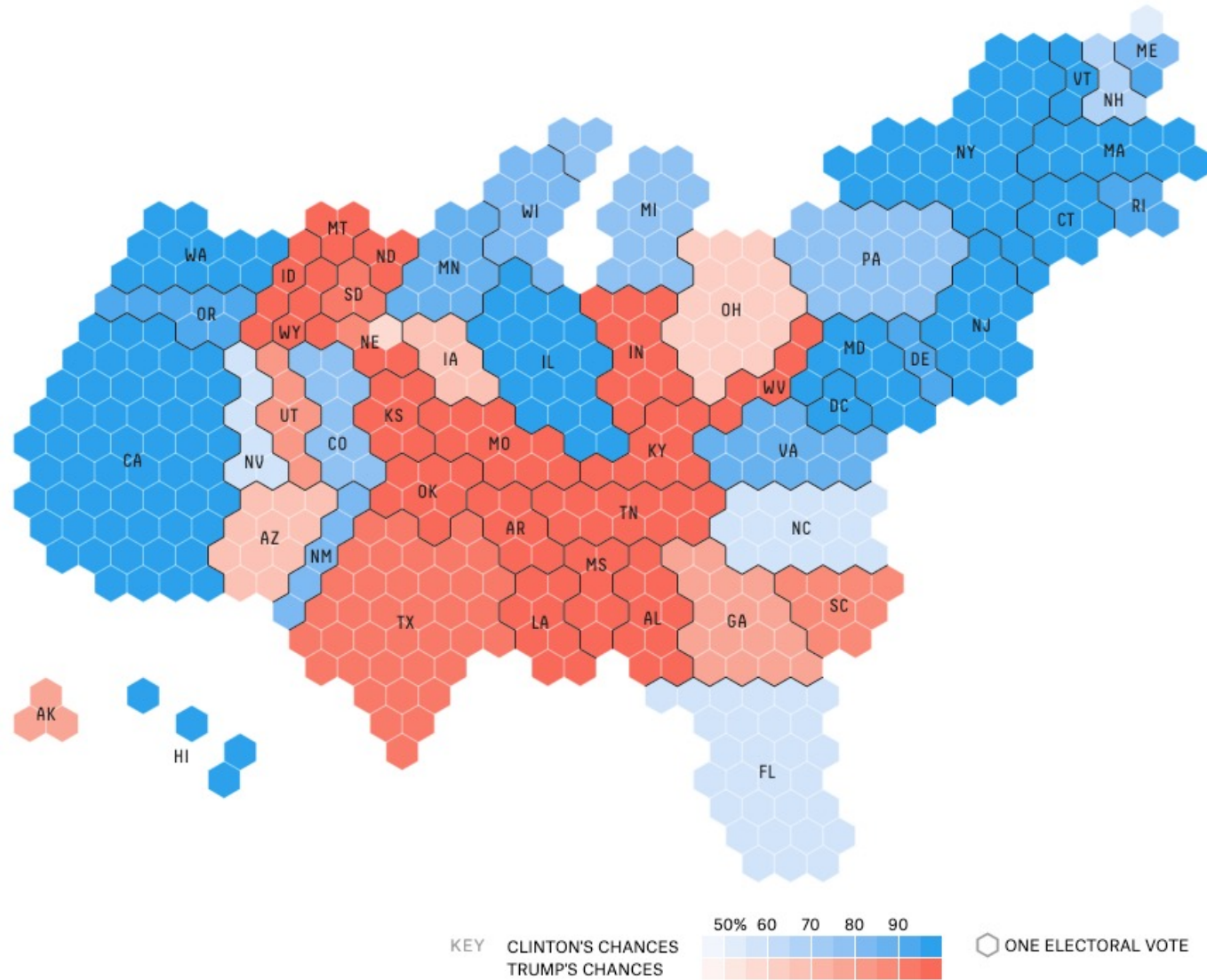WEEKLY CASES PER CAPITA: FEWER — MORE (March 1 — Oct. 3)

Show all

## Where new cases are higher and staying high

States where new cases are higher had a daily average of at least 15 new cases per 100,000 people over the past week. Charts show daily cases per capita and are on the same scale. Tap a state to see detailed map page.

North Dakota | South Dakota | Wisconsin | Montana | Utah
(7-day average, March 1 — Oct. 3, Last 14 days)

Iowa | Nebraska | Guam | Arkansas | Idaho

Oklahoma | Missouri | Kansas | Wyoming | Puerto Rico

Alabama | Minnesota | Kentucky | Texas | Alaska

Mississippi | Indiana | Illinois | Nevada
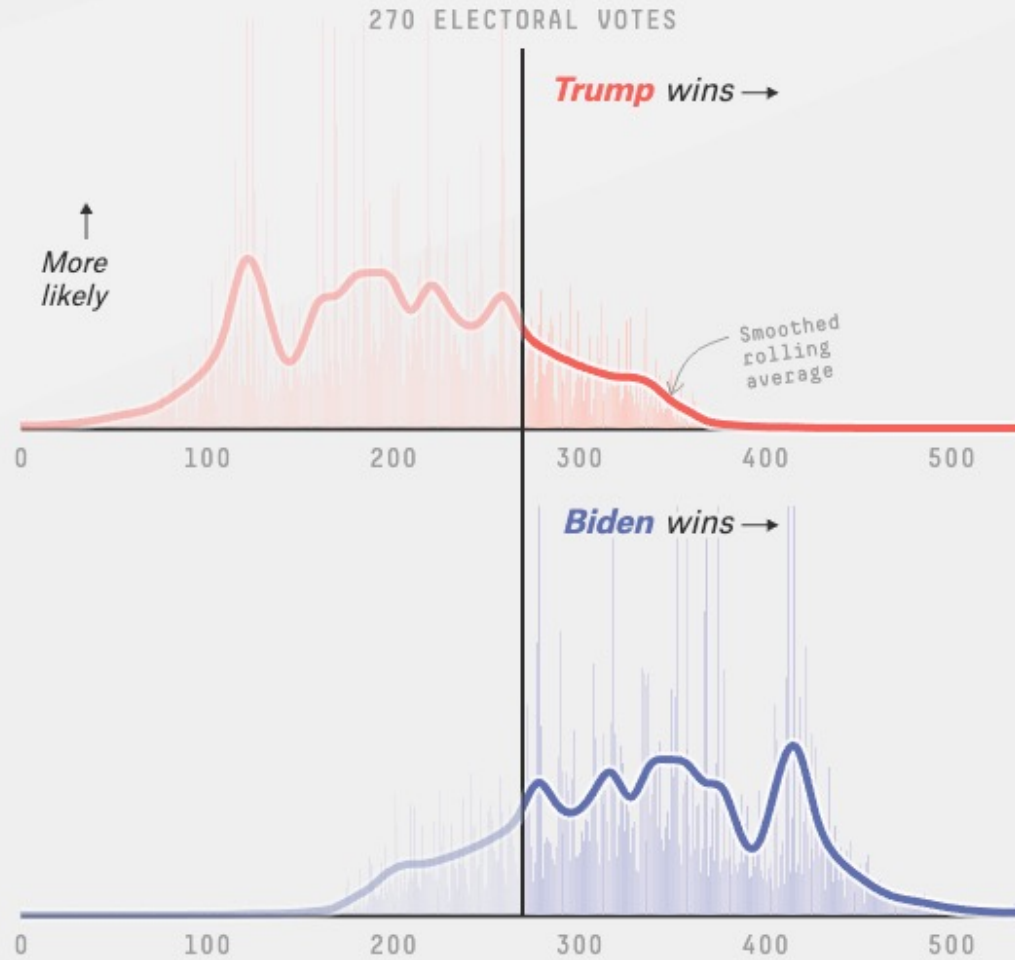
**It's all about the 538 Electoral College votes**

Here's a map of the country, with each state sized by its number of electoral votes and shaded by the leading candidate's chance of winning it.

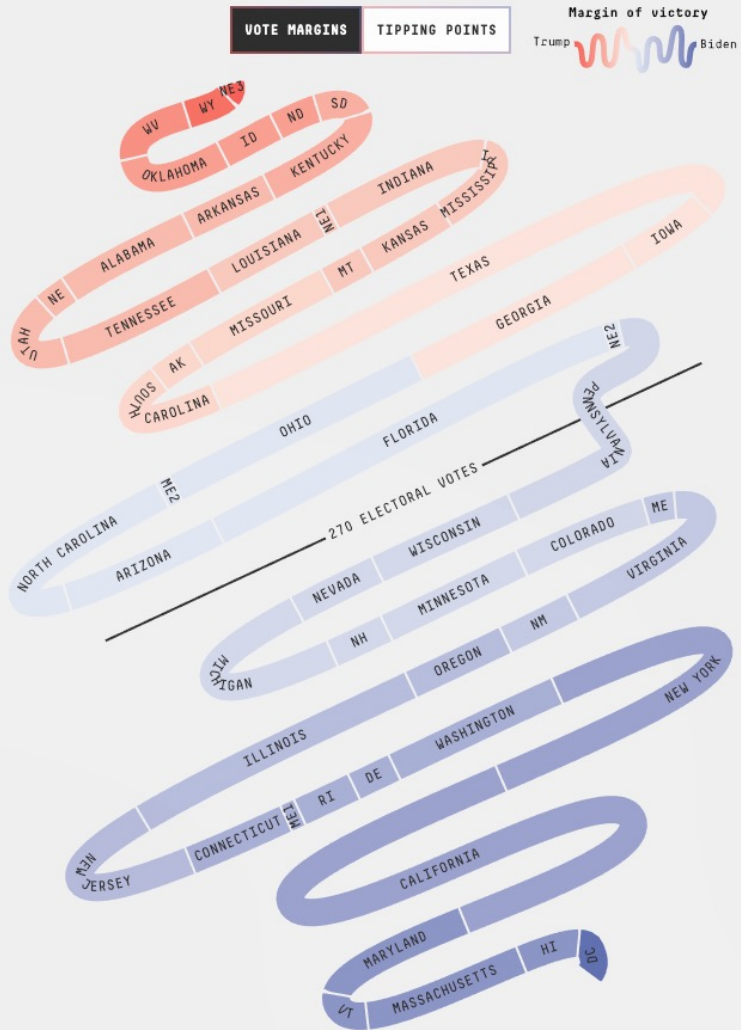KEY — CLINTON'S CHANCES / TRUMP'S CHANCES — 50% 60 70 80 90 — ONE ELECTORAL VOTE

https://projects.fivethirtyeight.com/2016-election-forecast/

https://projects.fivethirtyeight.com/2020-election-forecast/

# Exercise time!