

Data science for everyone

Prof. Jones-Rooy & Prof. Policastro

March 2, 2020

6.1:Visualizations

ANNOUNCEMENTS

Feedback on assignments

Lab 1 & 2 are returned

HW 1 to be reviewed in section this week

Lab 3 & 4 will be returned ASAP

HW2 code will be returned March 10

We will be discussing questions similar to HW2 conceptual questions

I. Assignments

- Lab 3 is out; due Wed., March 4, 8p
- Homework 2 is out; due Mon., March 9, 8p
- Lab 4 out Wed., March 4; due Mon. March 9, 8p
- Project description released tonight; due May 4



Notice earlier due date

2. Midterm

- Midterm exam in lecture, Wed. March 11
- Review in lecture, Mon., March 9
- TAs will hold review session office hours Tuesday, March 10, 2-5p in CDS
- No sections next week
- Sections this week: Go over HW1, Lab 4, review (bring your Qs!)

Why Are American Women Running Faster Than Ever? We Asked Them — Hundreds of Them

By Talya Minsberg and Kevin Quealy Feb. 28, 2020



Read the article [here!](#)

More than 450 women will race in the U.S. Olympic Trials marathon in Atlanta on Saturday. That's a huge number, reflecting changes in rules, innovative shoe technology and a sea change in women's running.

Amateur women are running faster than ever, and, through communal networks online, showing others how to do so.

To qualify to race on Saturday, a woman had to complete a marathon in 2 hours 45 minutes or faster sometime in the last three years, roughly a pace of 6 minutes 17 seconds per mile.

The women represent a range of backgrounds, and many have pushed themselves to results they once thought unachievable. We know this, because we talked to them — hundreds of them.

Selecting on the DV

Another possible systematic bias in their responses?

Ok, but you can't claim anything causal

Through social media networks, running clubs, online message boards and the runners' own personal networks, The New York Times reached about two-thirds of the qualifiers.

DATA SCIENCE IN THE NEWS!

U.S. women's marathons at 2:45 or faster

Title could be clearer

Is the unit of analysis women or marathons?

Is number of women running overall a confounder?

Why is this highlighted?

(As a personal preference, pink is a gross choice)

Athletes may have more than one OTQ result per year ▷ Sources: World Athletics, Racing Statisticians

Selection bias?!



IT'S NOT JUST
THE NYT

When Women Run

One hundred years after women were granted the right to vote, the U.S. has more women in political office than ever before. Yet gender has been a major theme of the 2020 campaign, as candidates, voters and the media debate whether a woman can win the presidency. To better understand what it's really like to try and win an election as a woman, we spoke to women from every state who have done it — 97 women in all.

These are their stories, in their words.

*Each chapter contains highlights from our interviews.
Interviews have been lightly edited for clarity.*

🎧 Put on your headphones

Did Voters Ask You Gendered Questions?

Did Your Appearance Become An Issue?

What's The Worst Sexism You Experienced In Politics?

Are American Voters Sexist?

How Did Media Talk About Your Gender?

What Advice Would You Give To Women Running Today?

 FiveThirtyEight

Find the project [here!](#)

Outline

1.What is a data visualization?

2.Types of visualizations

3.What makes a good visualization?

Lecture 6.2:

An example from start through visualization
Functions & more helpful code

DATA VISUALIZATION

- The graphical display of abstract information
- For two purposes:
 - Sense-making
 - Communication
- Using principles of data analysis, design, and storytelling (yes)

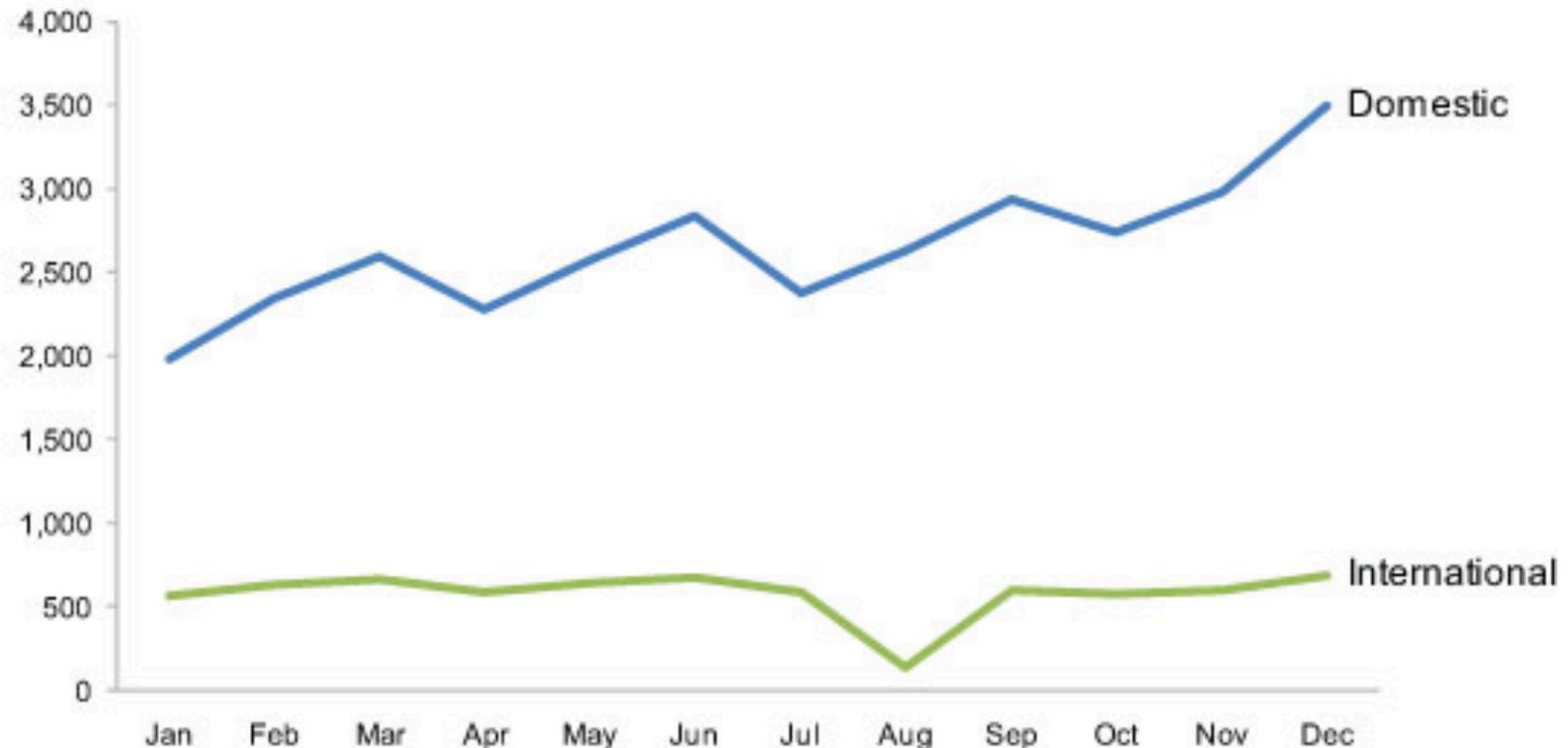


2009 Sales (thousands of U.S. \$)

Region	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Total
Domestic	1,983	2,343	2,593	2,283	2,574	2,838	2,382	2,634	2,938	2,739	2,983	3,493	31,783
International	574	636	673	593	644	679	593	139	599	583	602	690	7,005
Total	2,557	2,979	3,266	2,876	3,218	3,517	2,975	2,773	3,537	3,322	3,585	4,183	38,788

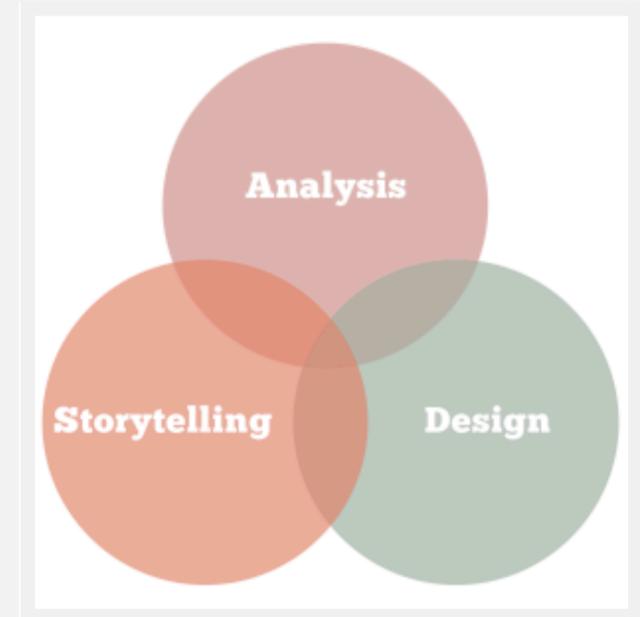
U.S. Dollars
(thousands)

2009 Sales



DATA VISUALIZATION

- Again: Science and art
 - **Science:** How to measure, how to model, what are the interesting patterns or findings in this data and research?
 - **Art:** How to measure, how to model, what is interesting, and **how to capture or share what's interesting?**
- Something humans have been doing for a (relatively) long time!

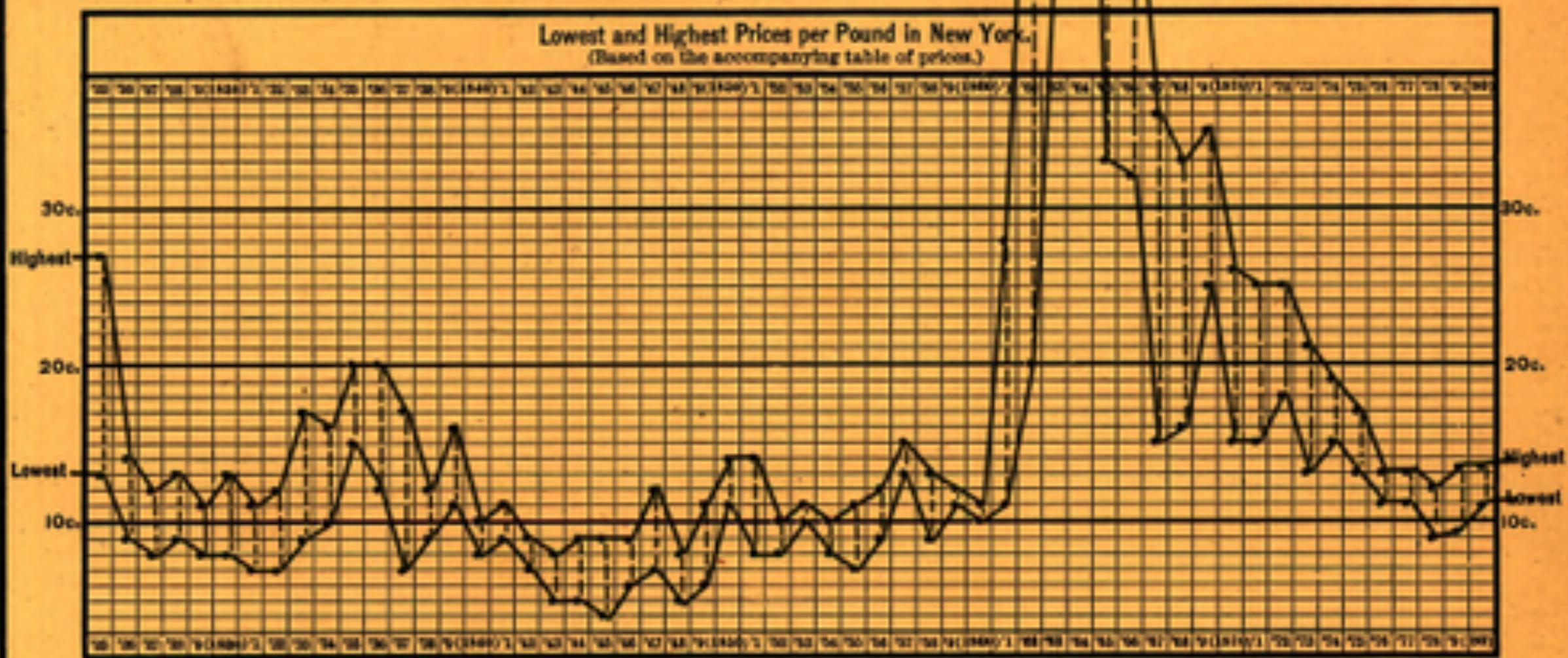


1825.

PRICE-CHART OF UPLAND COTTON FOR 56 YEARS.

1880.

Lowest and Highest Prices per Pound in New York.
(Based on the accompanying table of prices.)



SCALE.—Each space between the price lines represents one cent.

KEYS TO A GOOD VISUALIZATION

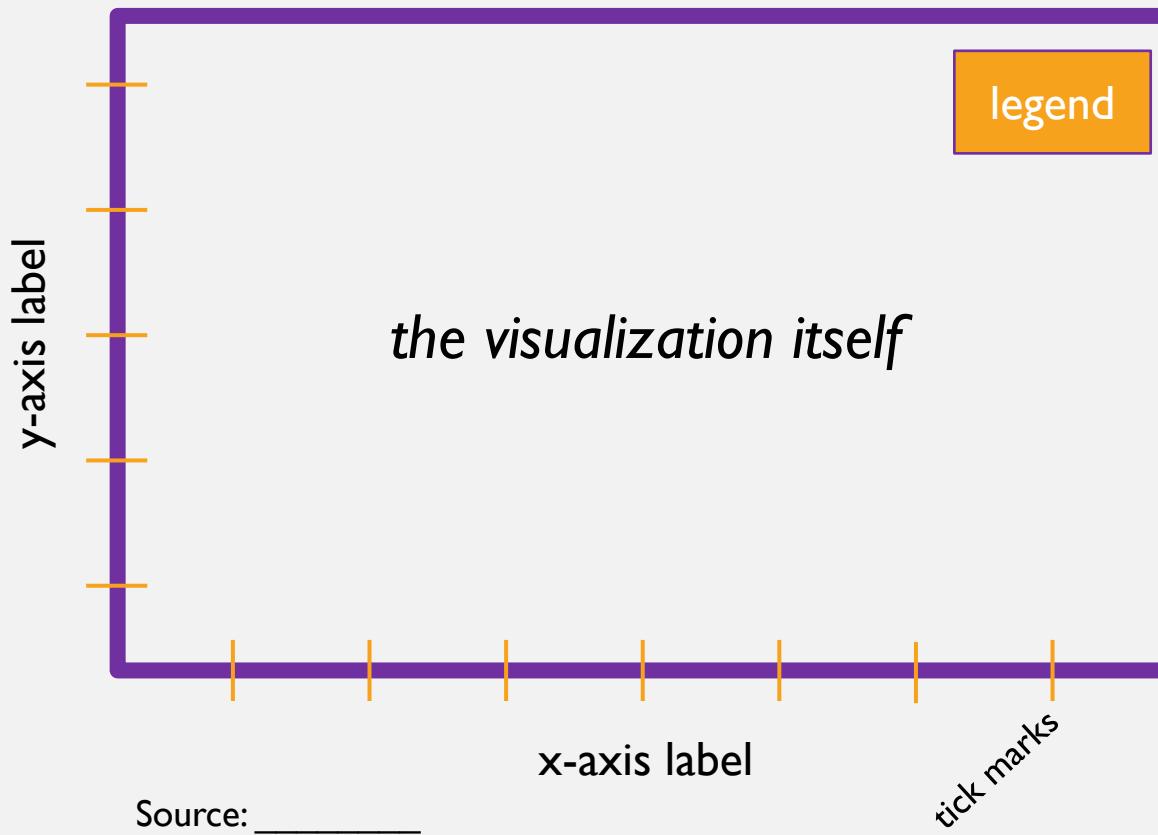
- Clear/understandable
- Ideally can stand on its own
- Well-labeled
- Draws attention to what's interesting/what the point is
- Honest: Doesn't mislead
- Transparent: about where the data comes from, how it's measured, what it really captures
- Replicable



KEY ELEMENTS

Title (ideally explaining key idea)

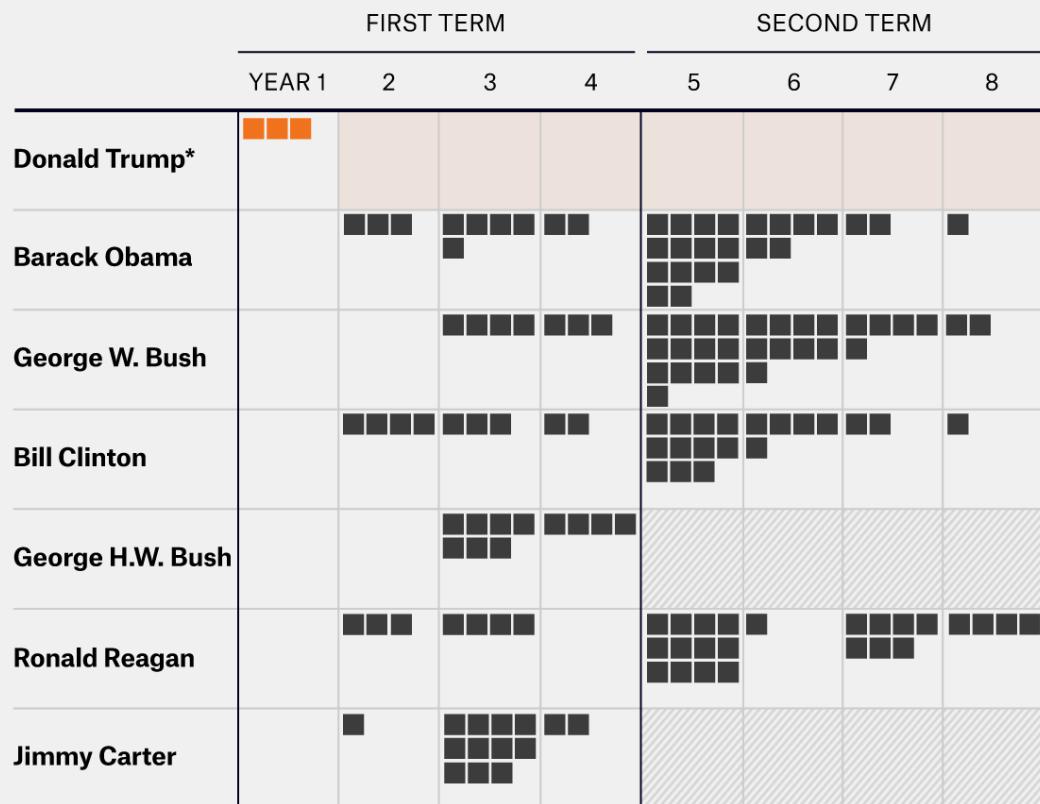
Subtitle (usually explaining specifics of the graphic)



Optional: Decoder
of any additional
information (e.g.,
colors)

Turnover among Trump's Cabinet has been high

Number of replacements in positions included in Trump's Cabinet, for recent presidents



*Trump's count is of departures because some replacements haven't happened yet.

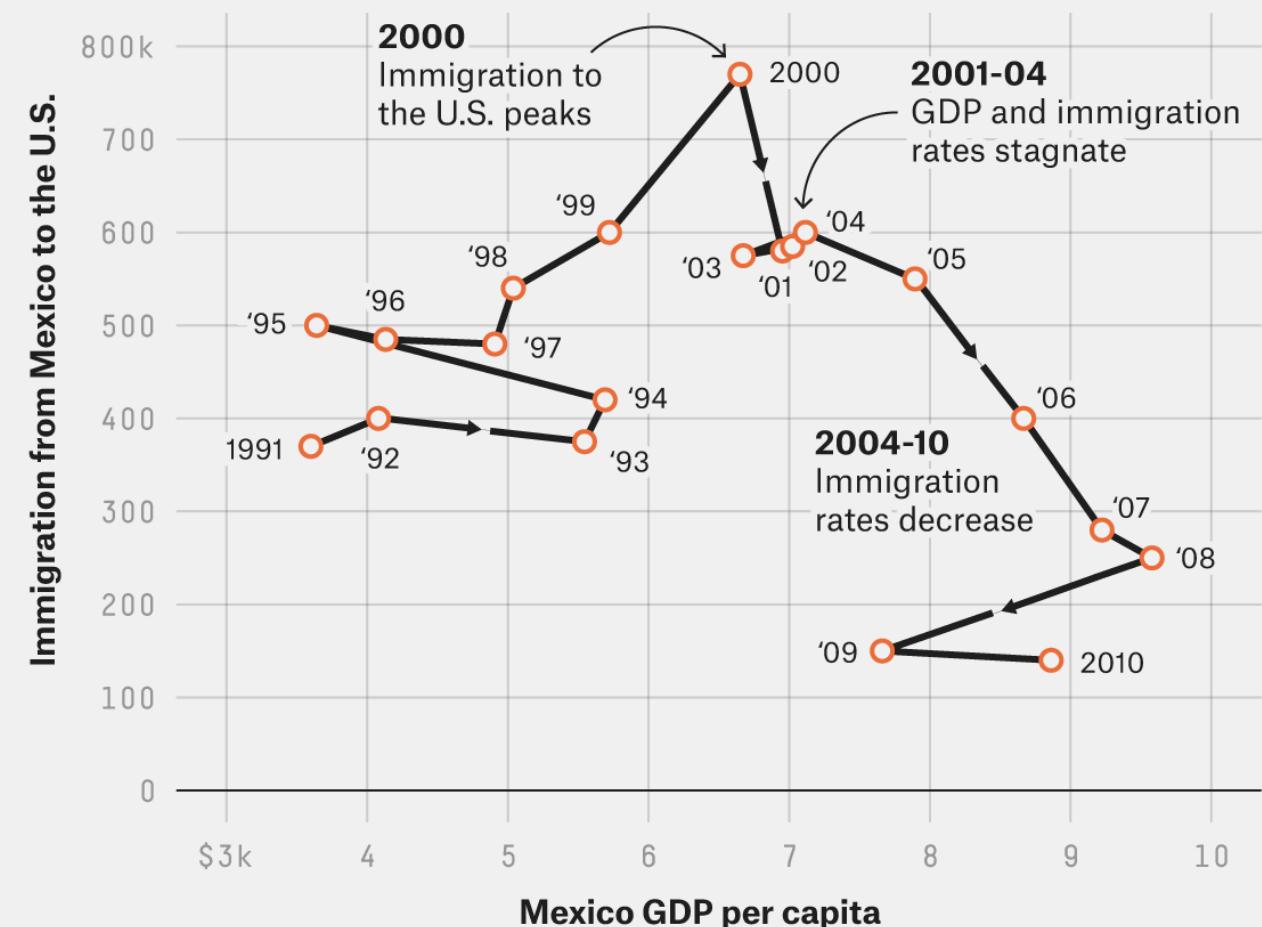
The number of total replaceable positions increases from 19 to 24 between Carter and Trump.

FiveThirtyEight

SOURCES: WHITE HOUSE, SENATE, DEPARTMENT WEBSITES

When the economy is up, immigration is down

Number of immigrants from Mexico to the United States versus Mexico's GDP per capita, 1991-2010



FiveThirtyEight

SOURCE: PEW RESEARCH CENTER, PEW HISPANIC CENTER, THE WORLD BANK

Outline

1.What is a data visualization?

2.Types of visualizations

3.What makes a good visualization?

FOUNDATIONAL VISUALIZATIONS

1. Bar charts
2. Line graphs
3. Scatter plots
4. Histograms
5. Others

Generating visualizations in Python
may be simple or technically
challenging

But a major challenge is often
**determining the right
visualization** in the first place

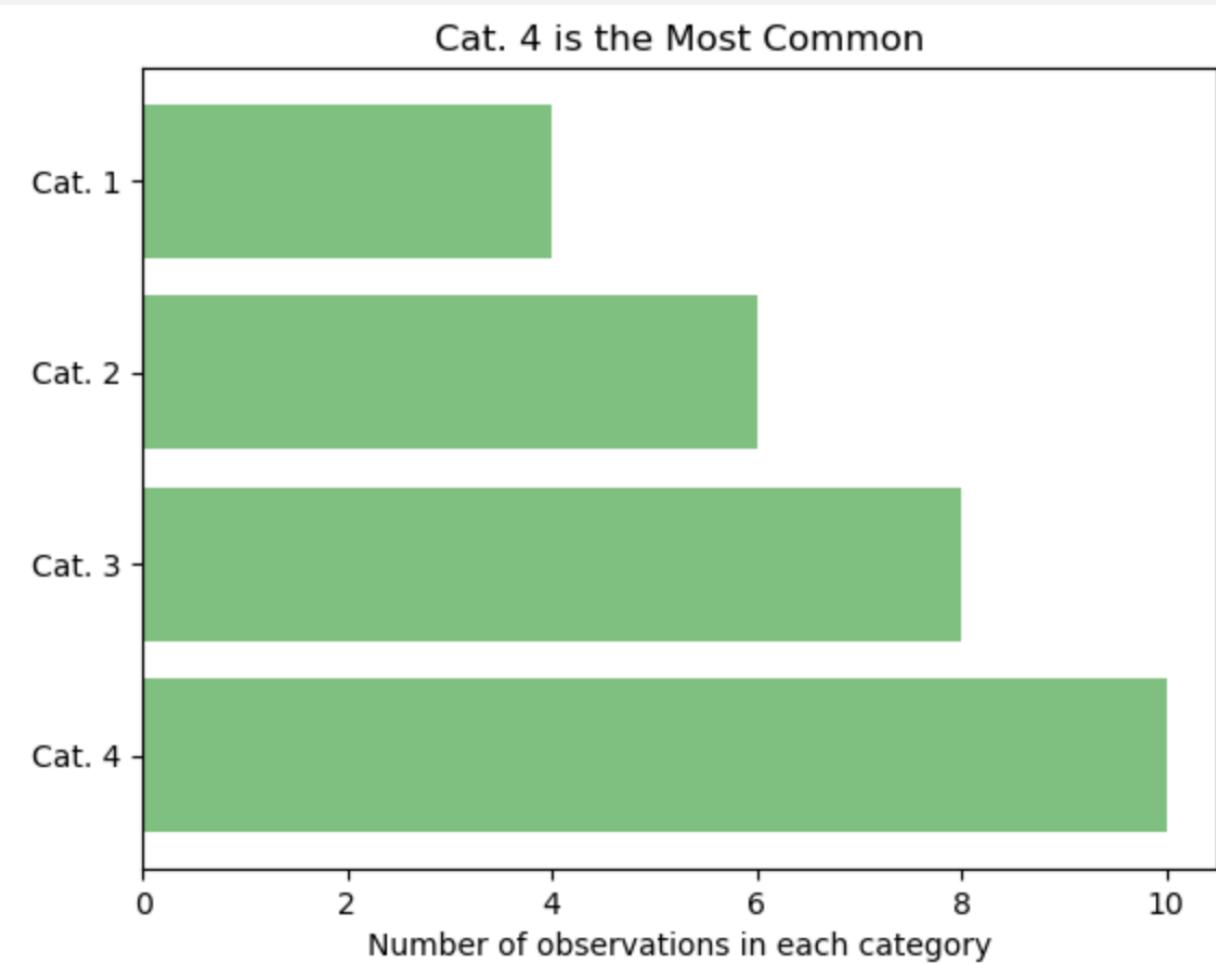
BAR CHARTS

Commonly used with categorical variables

Discrete, usually limited and finite, often qualitative property



Might also be useful for visualizing specific observations



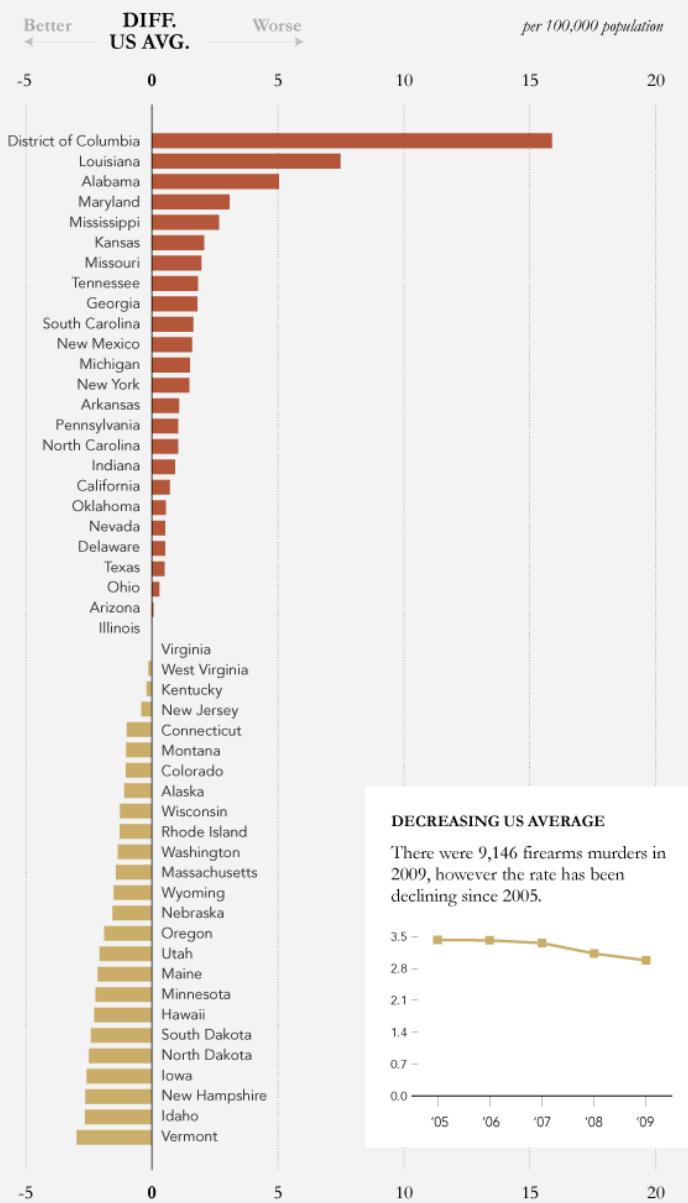
Examples

- Blood type: A, B, AB, O
- Hair color
- Type of book
- Political party
- Primary construction material
- Type of parole requirements
- Type of movie or award

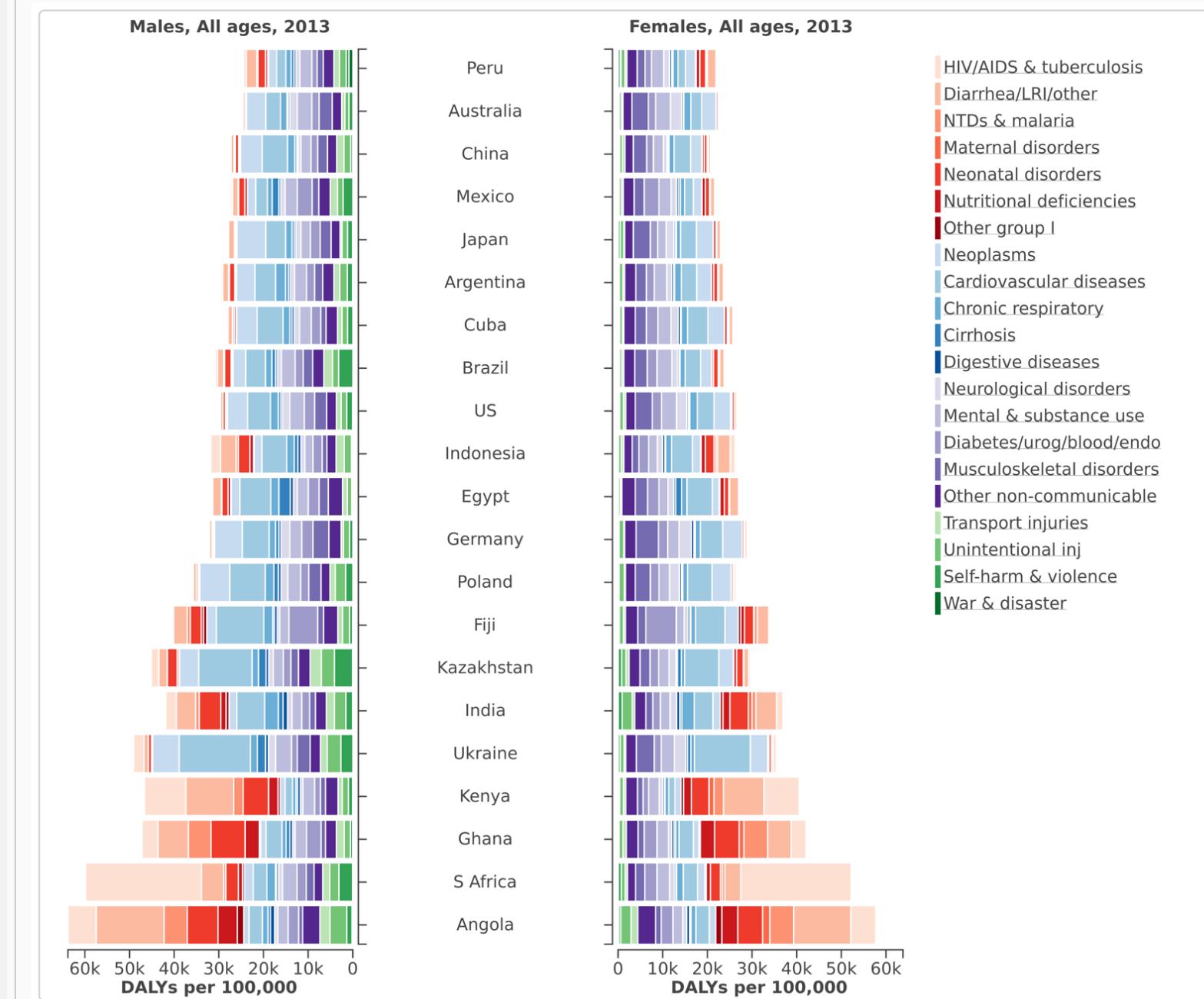
```
import matplotlib.pyplot as plt  
plt.bar(...) or plt.bart(...)  
plt.show()
```

MURDERS BY FIREARM

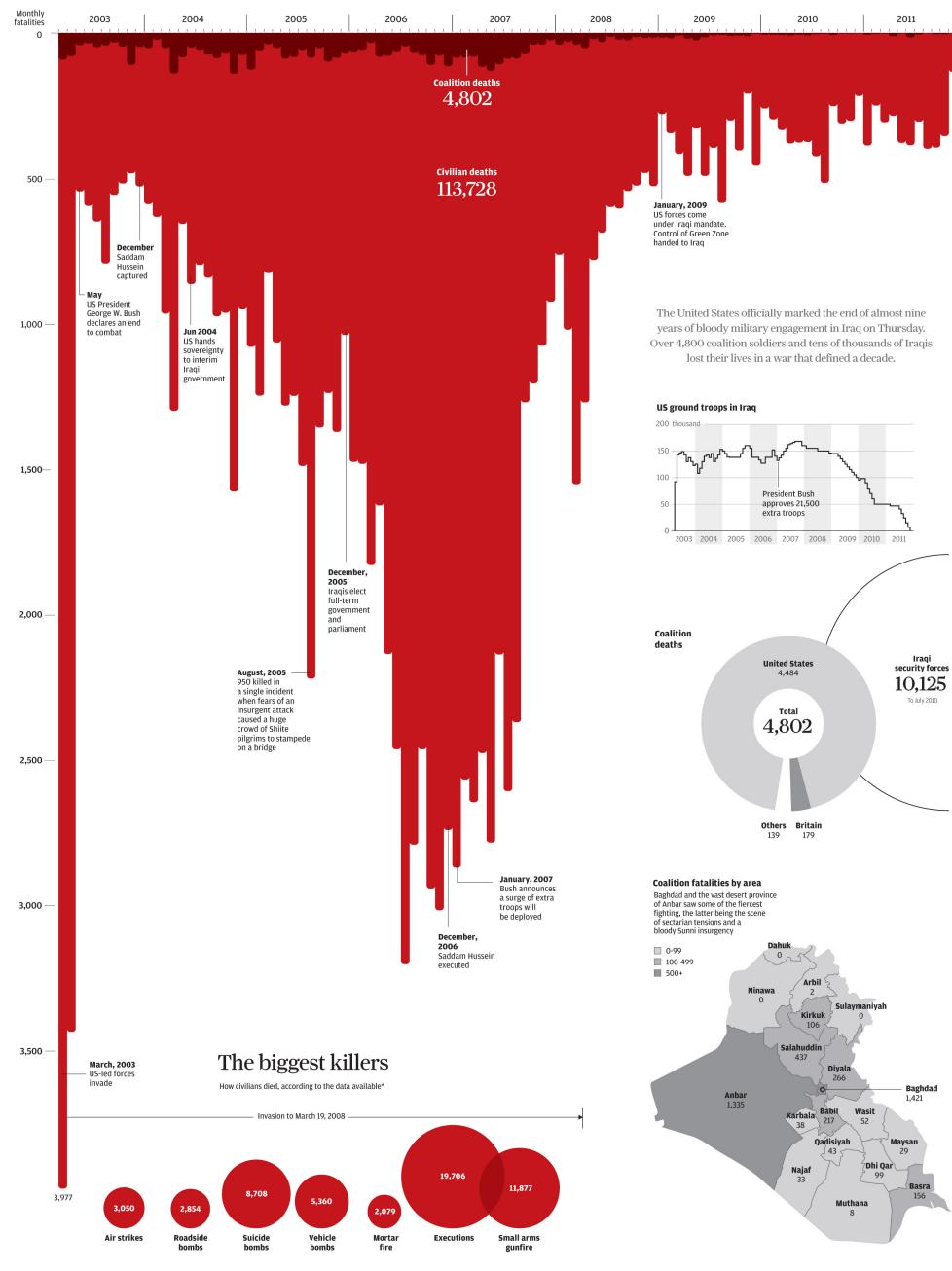
The average rate for the United States was 2.98 murders by firearm per 100,000 population. The below shows how states compared to the national average.



Burden of disease by cause, country, and gender (2013 estimates) – produced by IHME Viz Hub



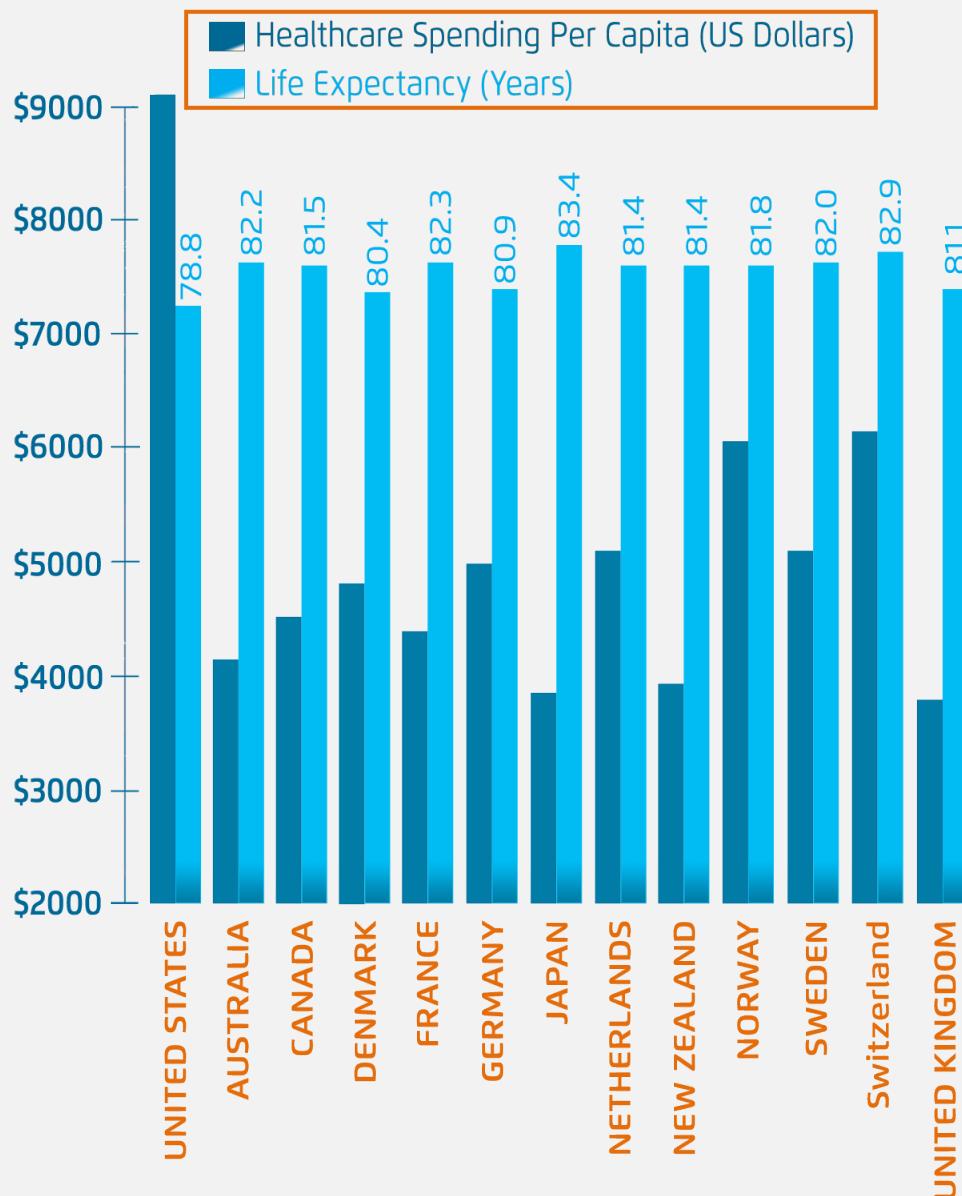
Iraq's bloody toll



Source: Washington Post

US Population Health

SPENDING ≠ POPULATION OUTCOMES



PERCENT GDP SPENT ON HEALTHCARE



PUBLIC \$PENDING



US public spending
\$4197 per Capita



UK public spending
\$2802 per Capita

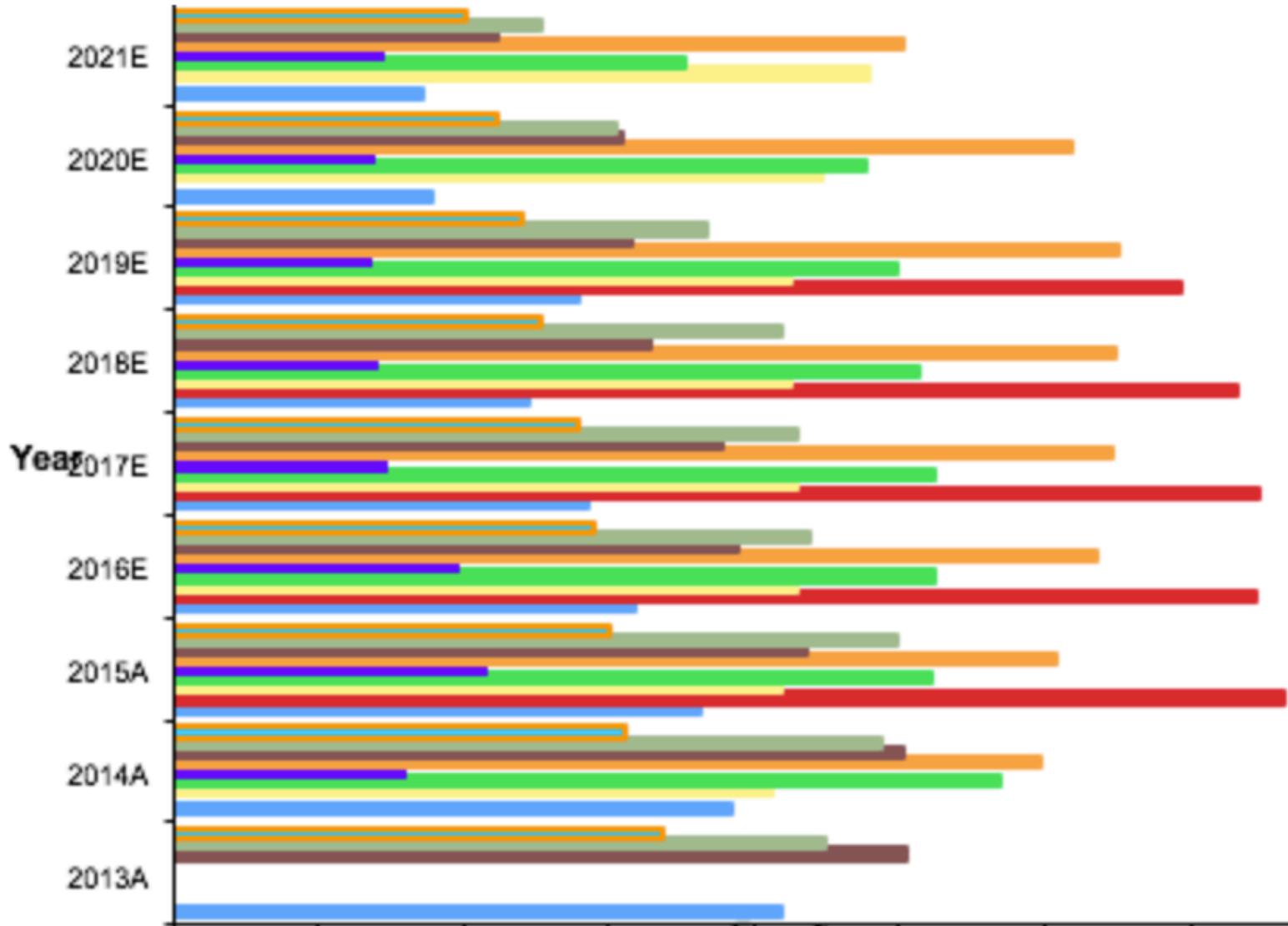


RESIDENTS COVERED BY PUBLIC SPENDING



Source: Health Catalyst

Please don't ever do this



Please explore these interactives!

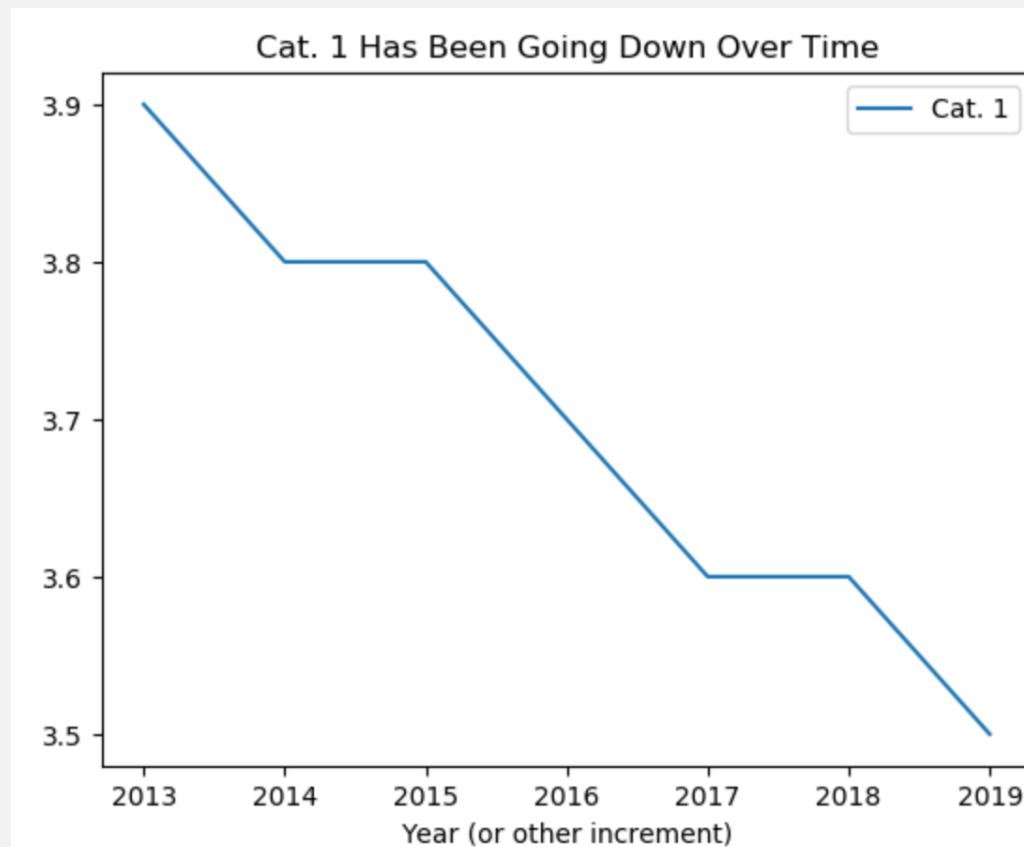
- [Daily routines of famous creative people](#)
- [Distributions of annual incomes](#)

LINE GRAPHS

Often used for chronological trends and patterns

Discrete or continuous, doesn't have to be over time, but is usually implied

Other continuous trends – temperature, speed, pressure readings



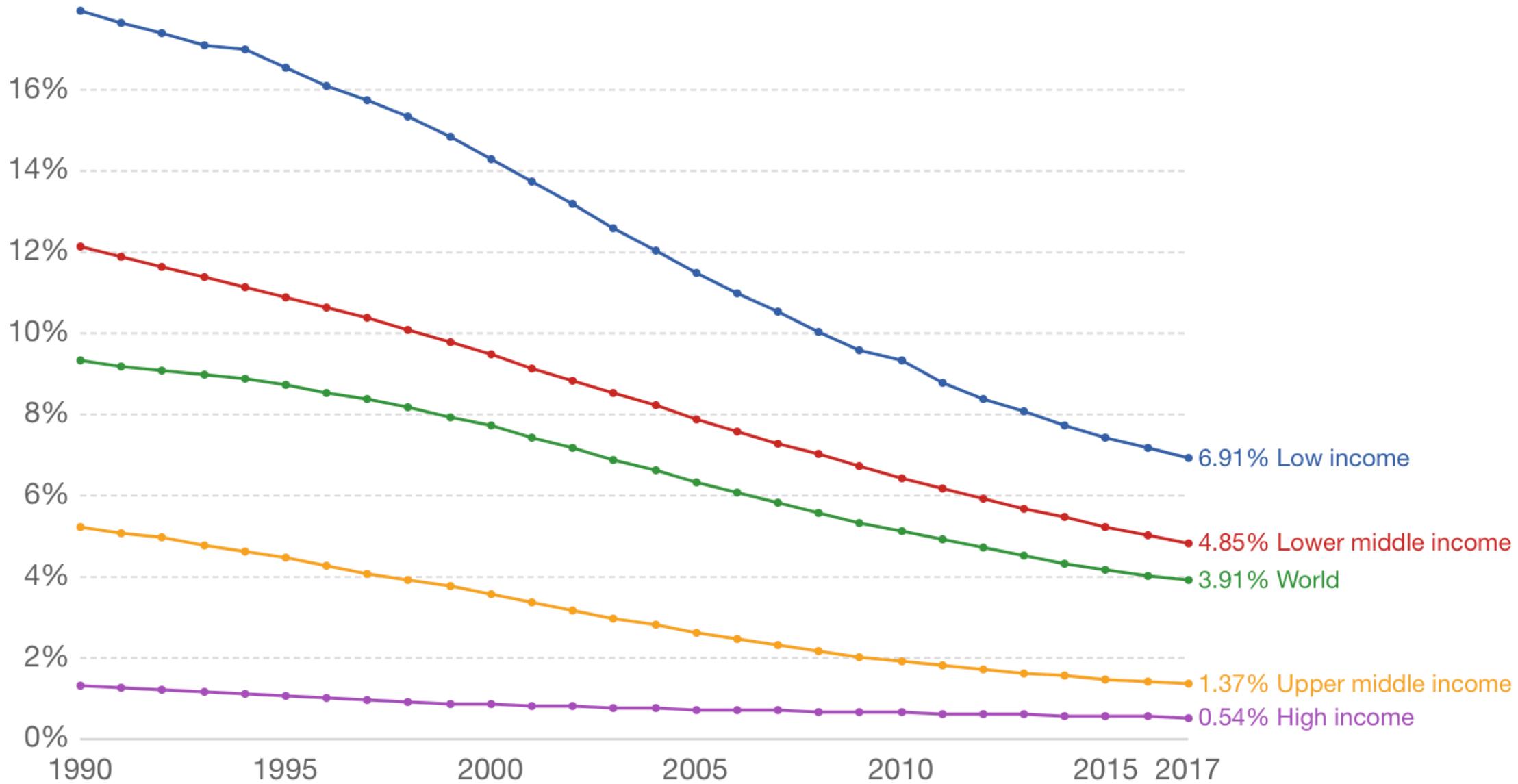
Examples

- GPA over time
- Voter turnout in elections historically
- Average rainfall per year in a location
- Effects of a change in temperature on a substance
- Pollution levels of a body of water over time

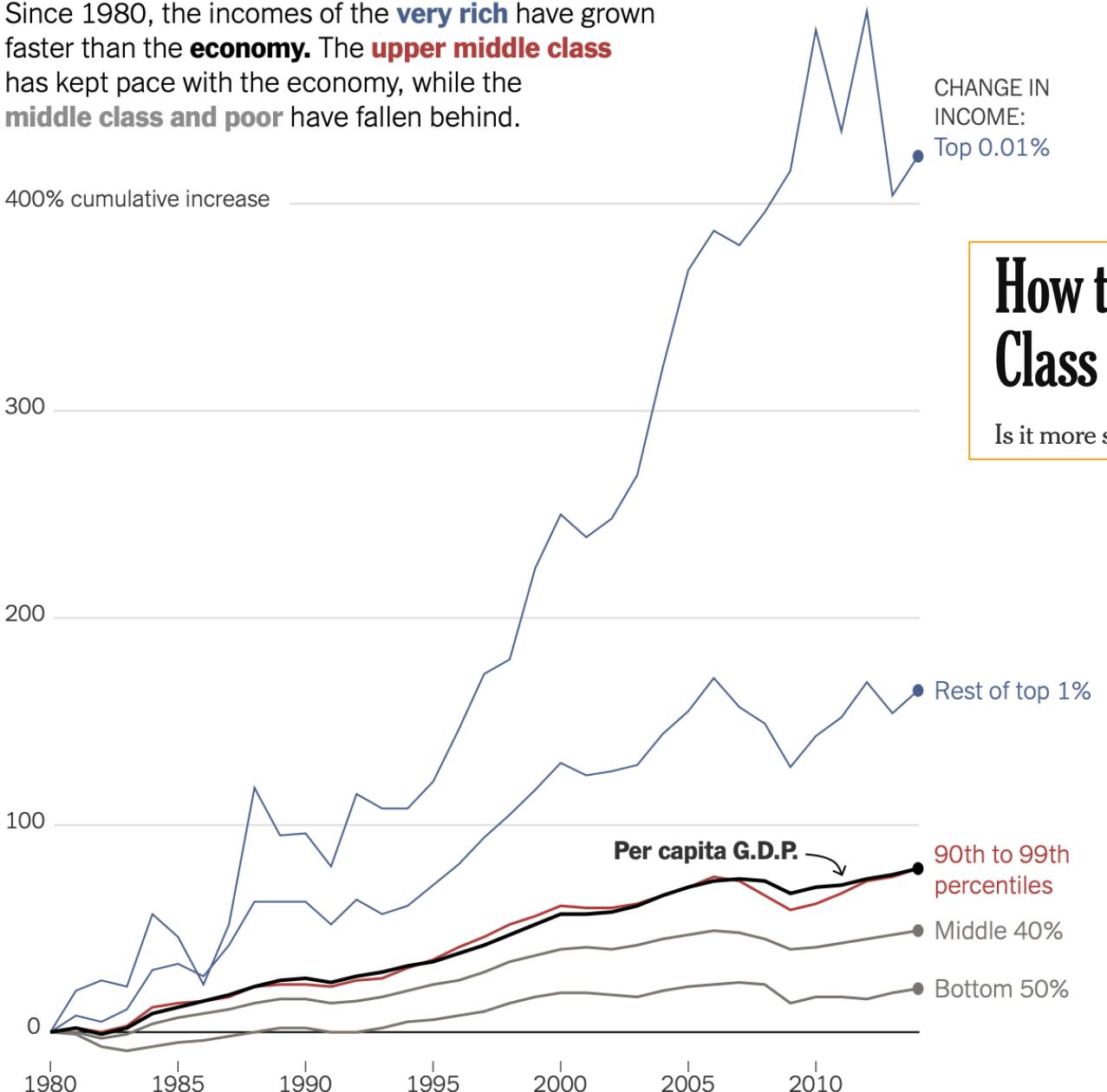
```
import matplotlib.pyplot as plt  
plt.plot(...)  
plt.show()
```

Child mortality by income level of country

The child mortality rate measures the share of children that die before reaching the age of 5.



Since 1980, the incomes of the **very rich** have grown faster than the **economy**. The **upper middle class** has kept pace with the economy, while the **middle class and poor** have fallen behind.



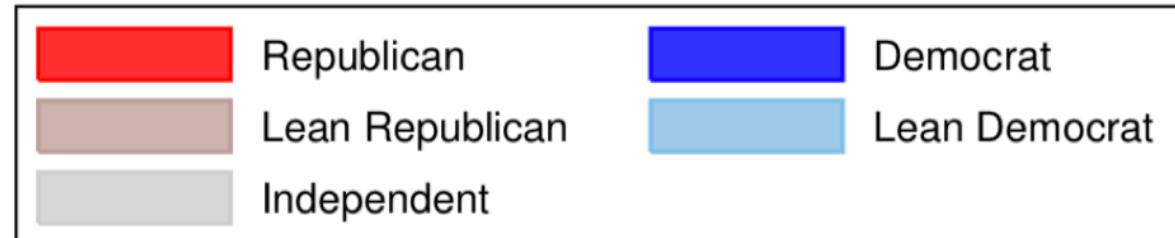
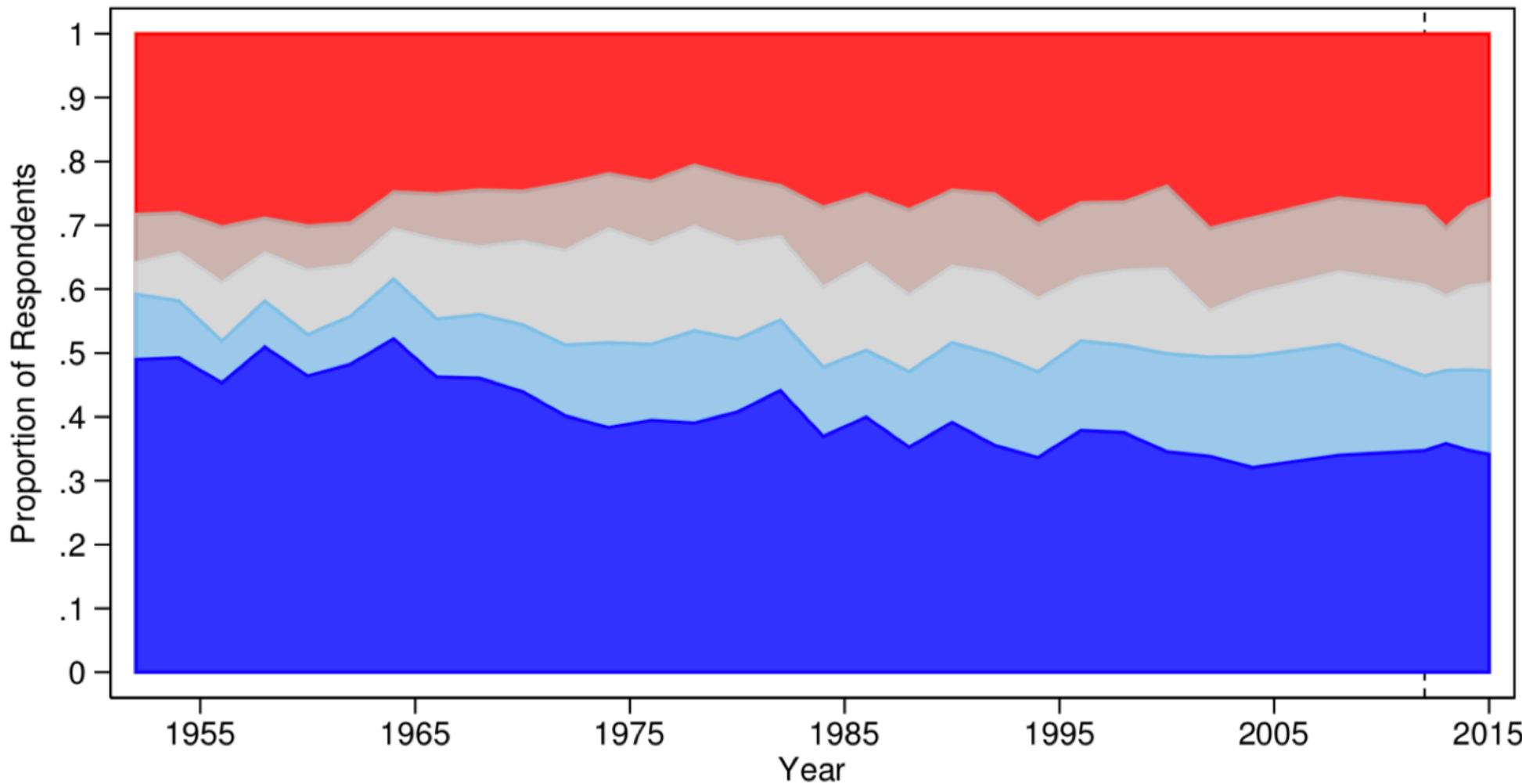
How the Upper Middle Class Is Really Doing

Is it more similar to the top 1 percent or the working class?

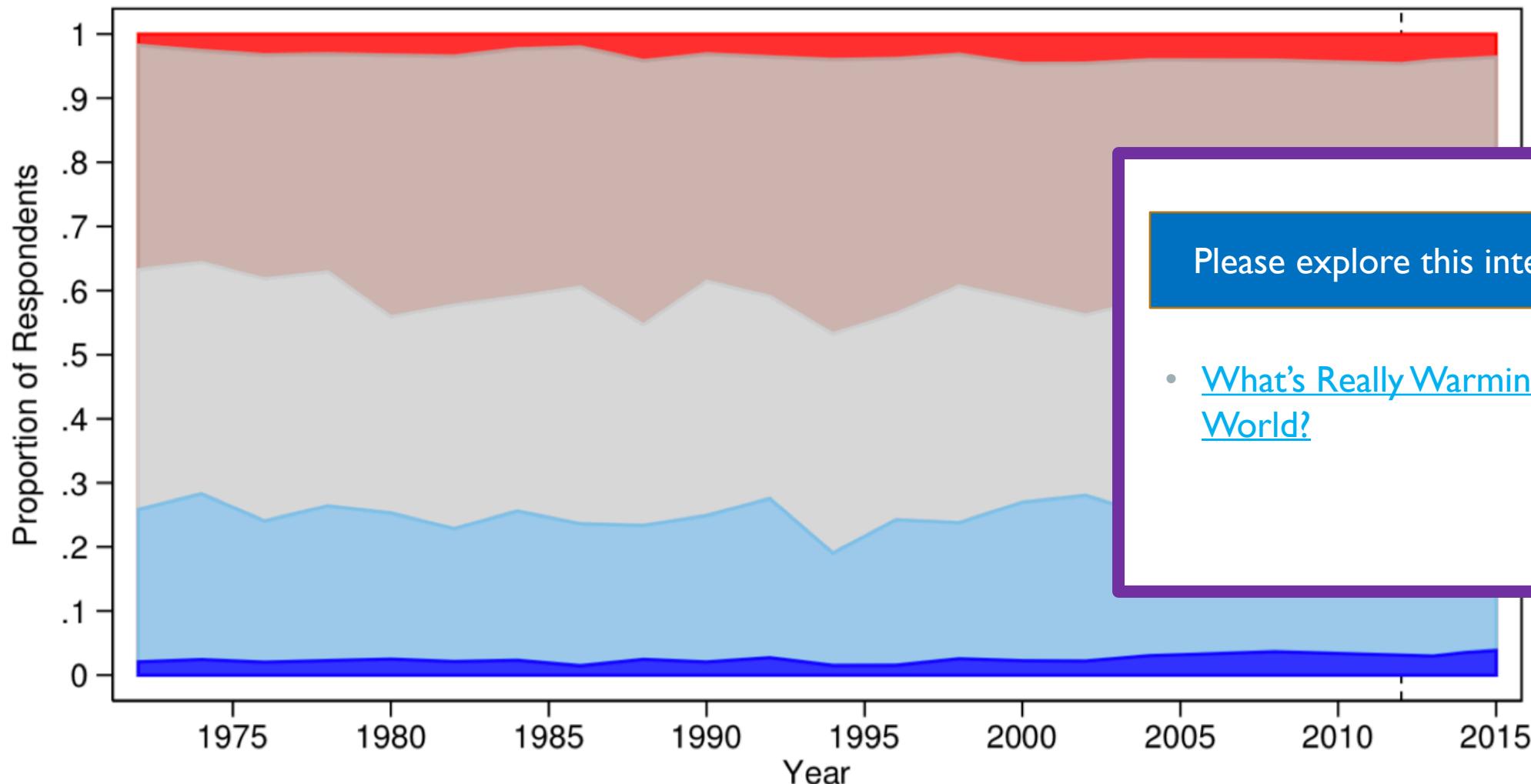
Note: Incomes are after taxes and include government transfers. • Sources: Thomas Piketty, Emmanuel Saez and Gabriel Zucman (incomes); Bureau of Economic Analysis (G.D.P.) • By The New York Times

Source: New York Times

Political Party

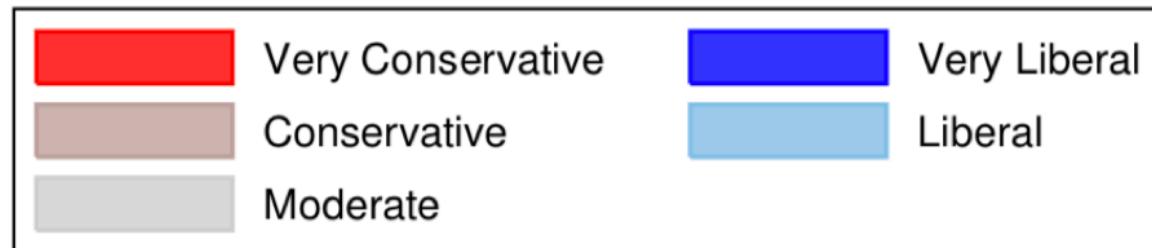


Political Ideology



Please explore this interactive!

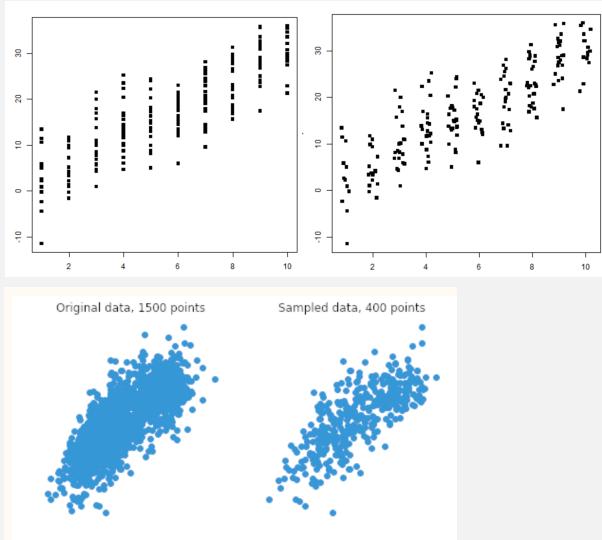
- [What's Really Warming the World?](#)



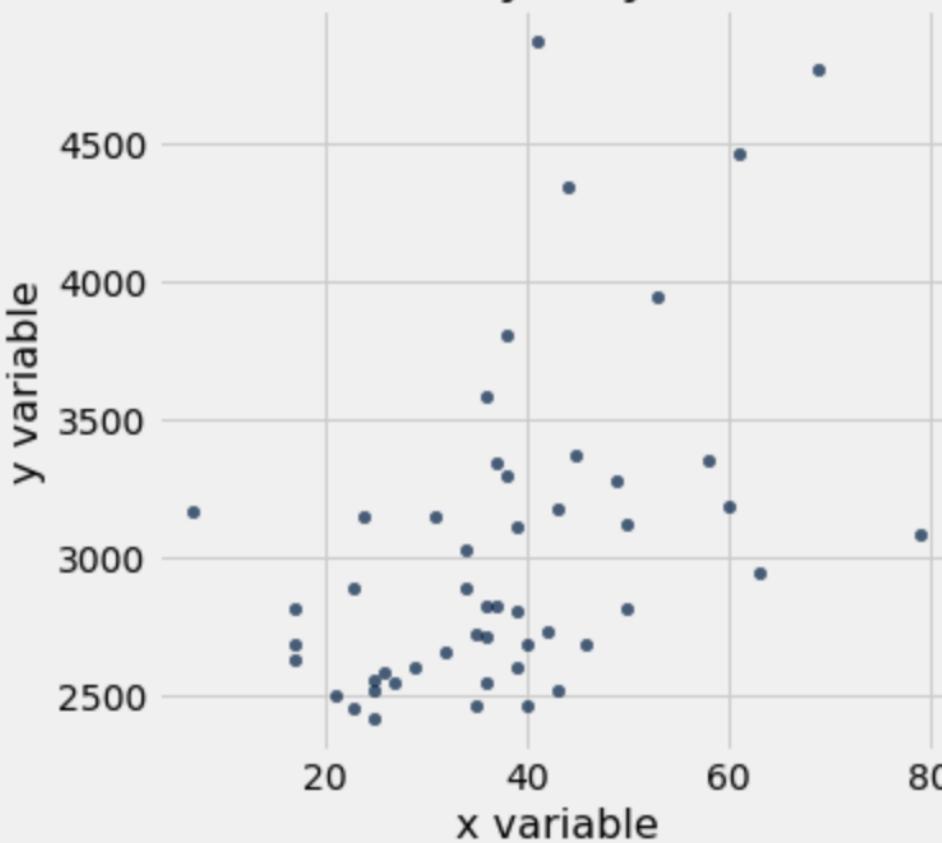
SCATTERPLOTS

Best for visualizing **association** between two variables

Continuous numeric variables tend to work best, though there are ways to work with others (e.g., jitter, select sub-sample)



Variables x and y May Be Associated

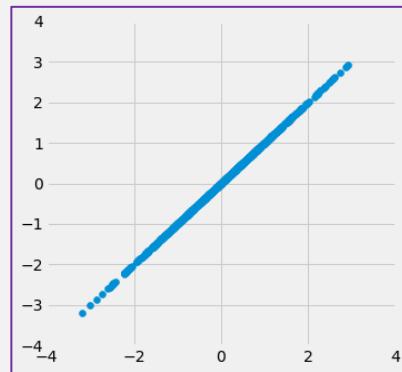
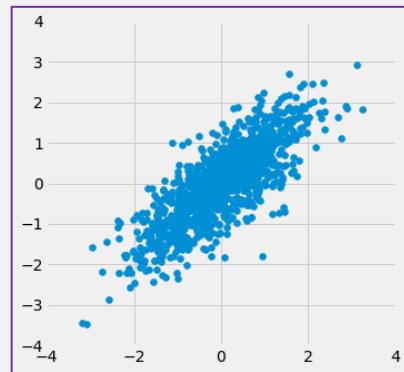
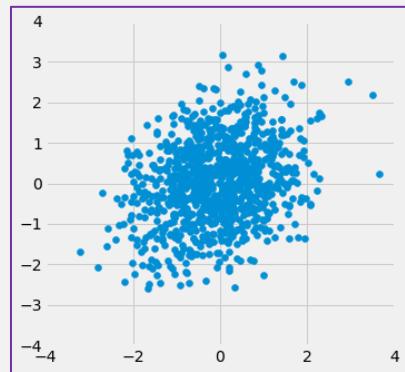
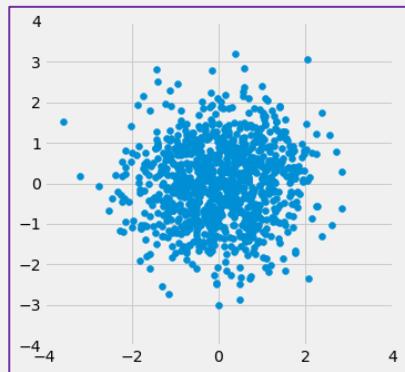
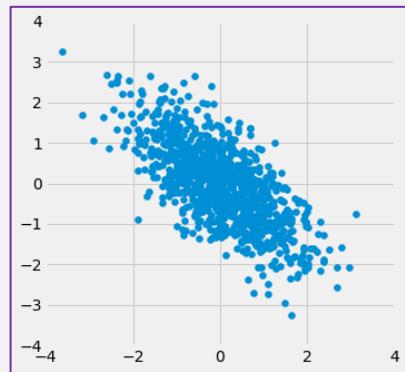
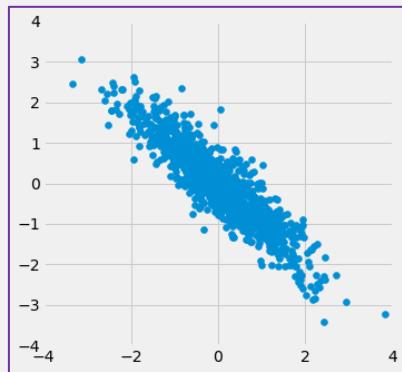


Examples

- Education and income
- Political stability and national GDP
- Coffee consumption and exam performance
- Social media usage and happiness
- Air pollution and hospital trips in a city

```
import matplotlib.pyplot as plt  
plt.scatter(...)  
plt.show()
```

VISUAL INSPECTION OF ASSOCIATION BETWEEN TWO VARIABLES



After the midterm, we will assign numbers to these relationships to quantify their correlation (r):

$$r = -0.9$$

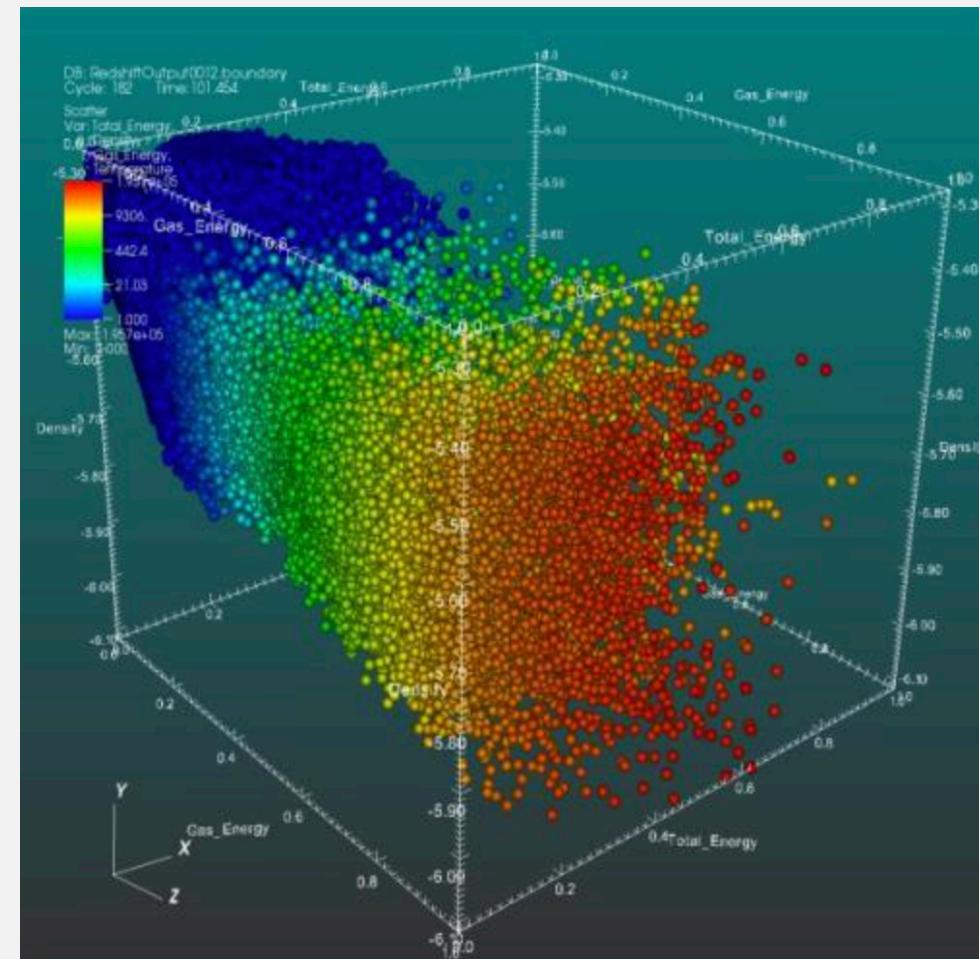
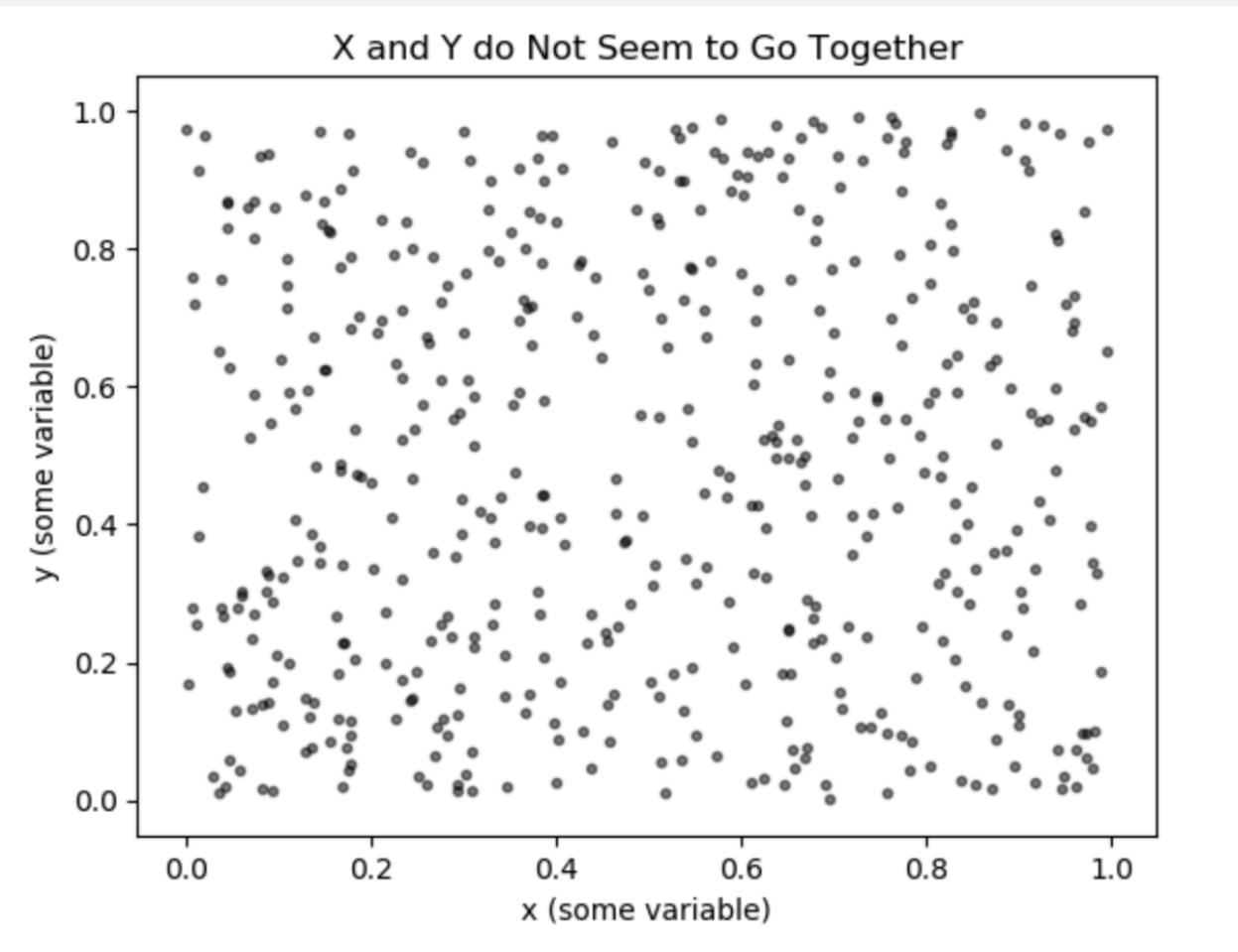
$$r = -0.7$$

$$r = 0$$

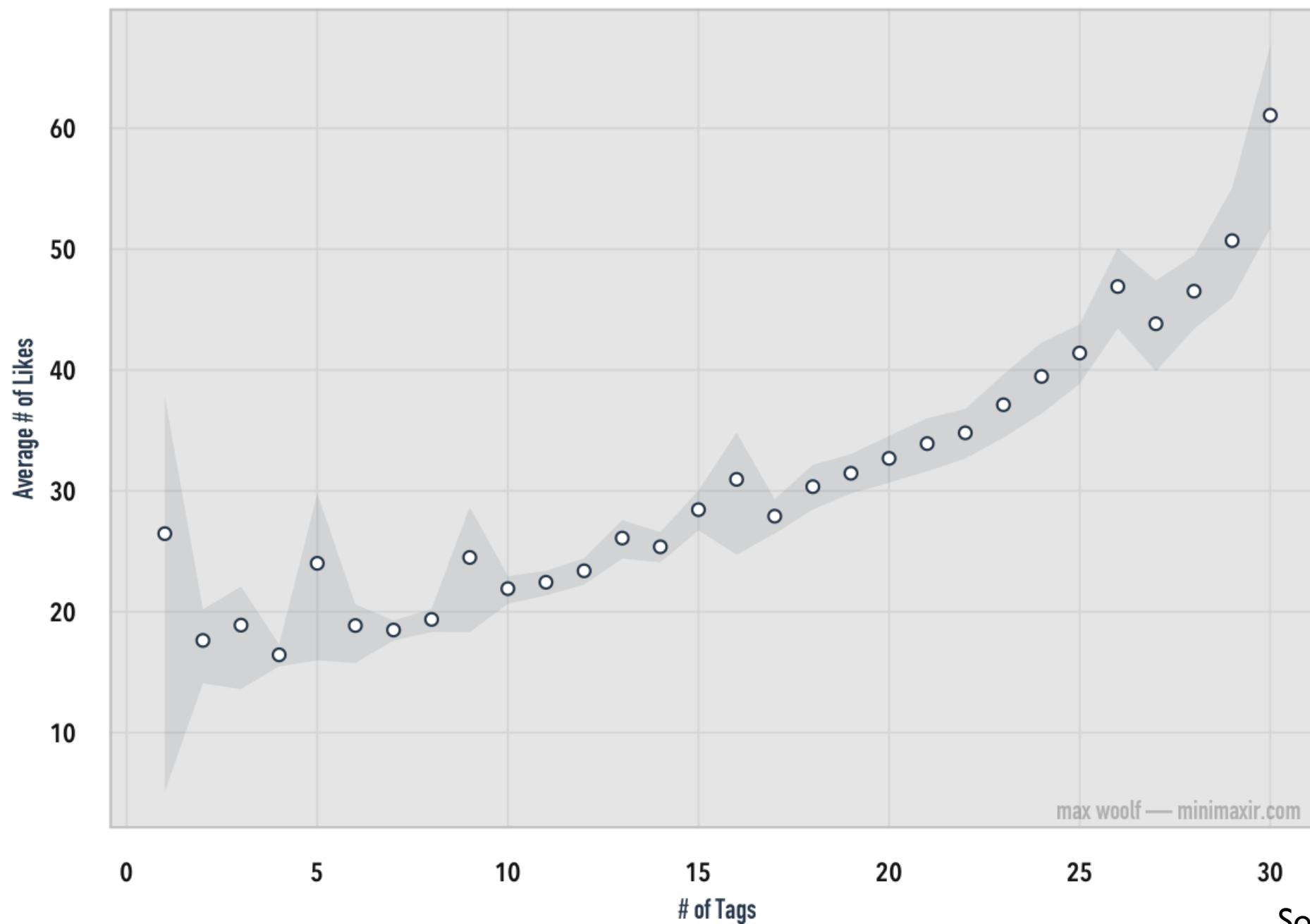
$$r = 0.25$$

$$r = 0.8$$

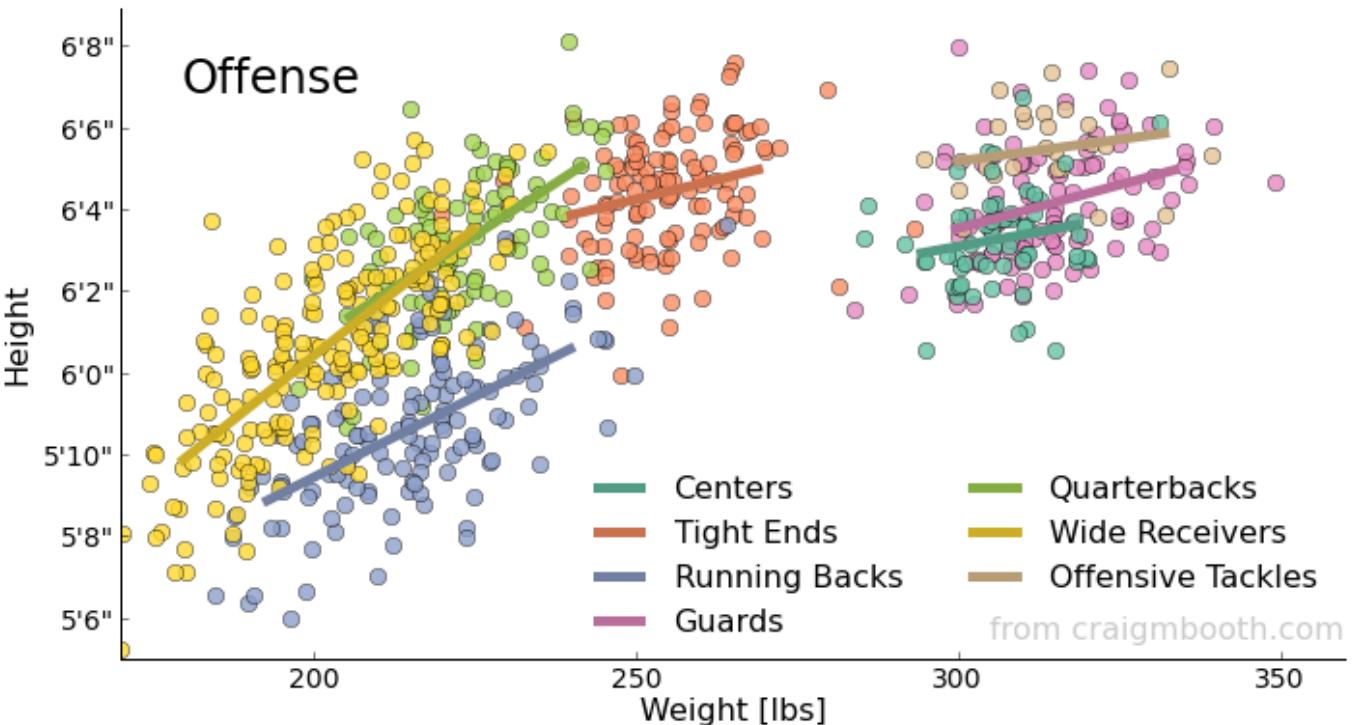
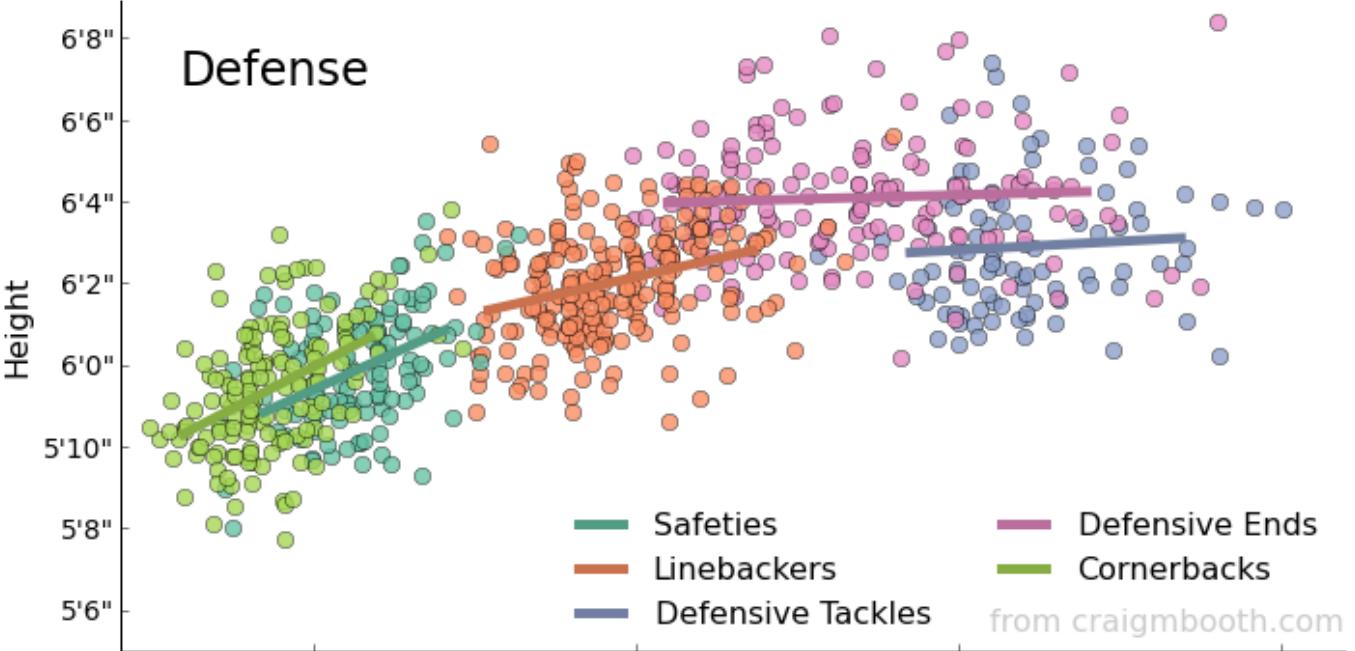
$$r = 1$$



OF TAGS VS. AVERAGE # OF LIKES ON 120,346 INSTAGRAM PHOTOS BY # OF TAGS



Height and weight of the players on all 2013 NFL rosters, broken down by position. Data scraped from nfl.com



Source: Craig Booth

Healthcare Expenditure vs. GDP, 2014

Total healthcare expenditure per capita and GDP per capita in PPP-adjusted constant 2011 international dollar.

LINEAR



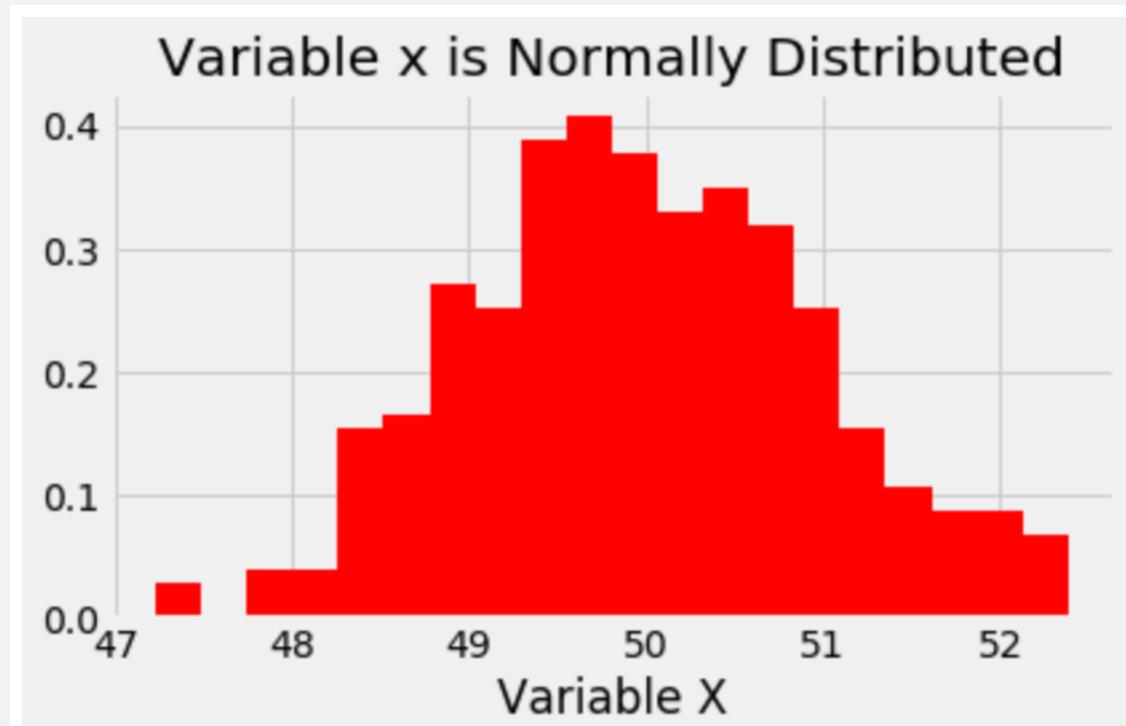
HISTOGRAMS

Used to understand the **distribution** of numerical variables (we already know this, too!)

Often for **continuous** data but discrete works, too

Data is broken up into groups called **bins**

The distribution of a variable is useful for understanding and predicting

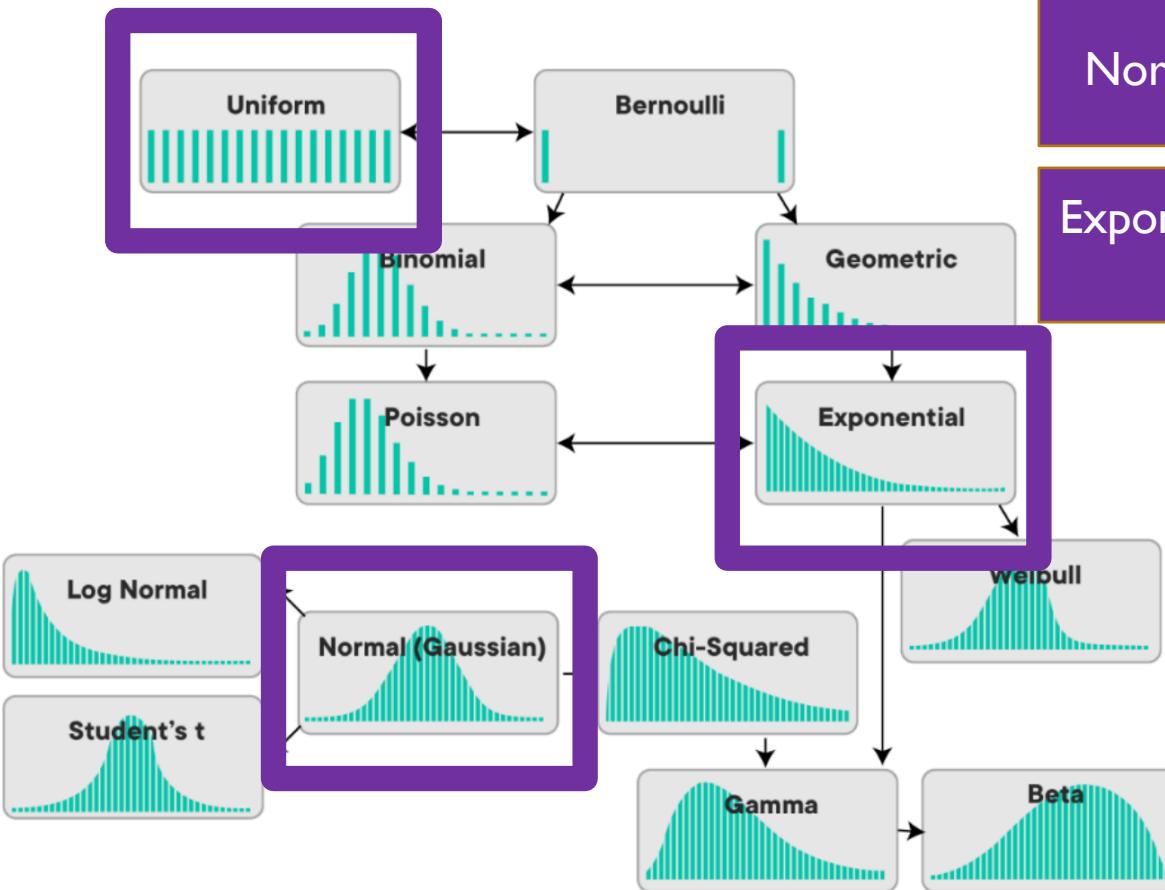


Examples

- Height
- Income
- Age
- How many copies of books sold
- Size of earthquakes

```
import matplotlib.pyplot as plt  
plt.hist(...)  
plt.show()
```

Understanding Different Types of Distributions You Will Encounter As A Data Scientist



Uniform: All outcomes are equally likely

Normal: Middle values are most likely

Exponential: Low values most likely, high values very unlikely



Power law: Even more extreme

Read more [here!](#)

Distribution of CEO Compensation, 2013



Source: Business Insider

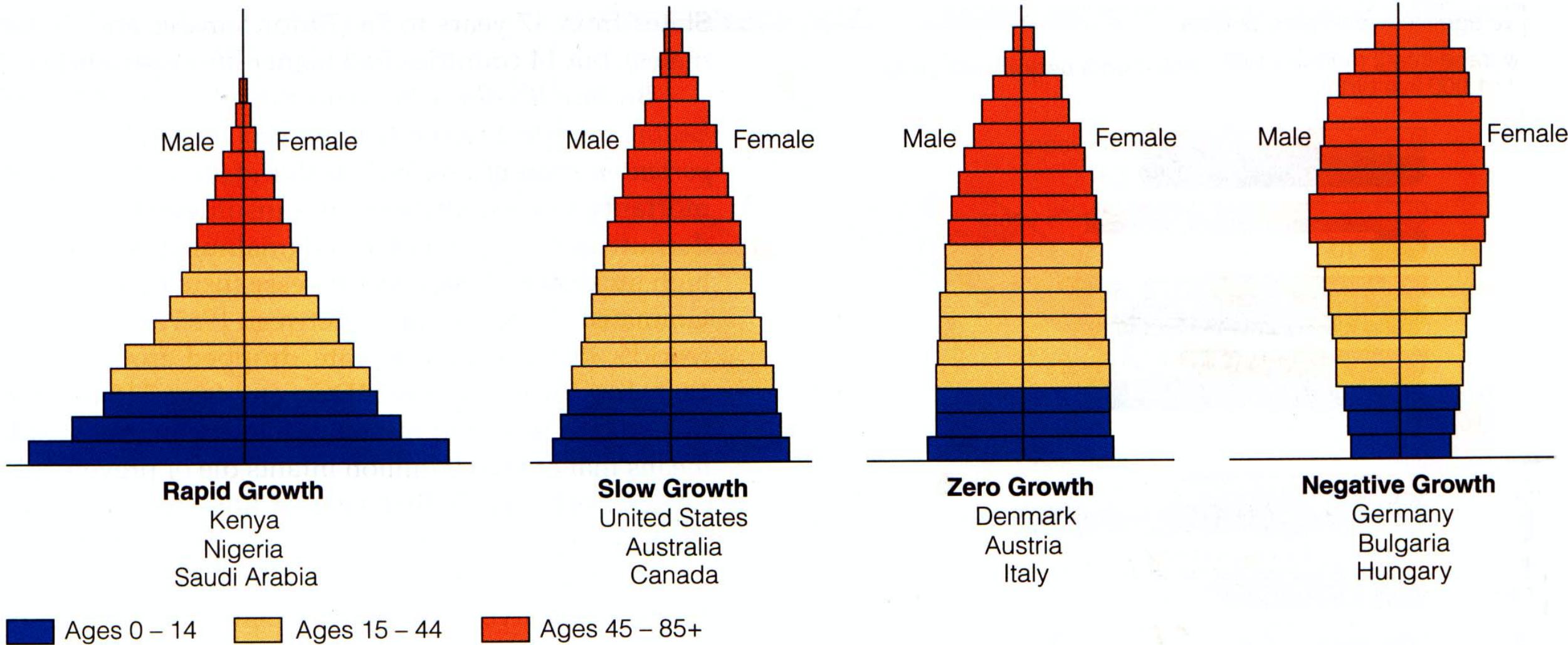
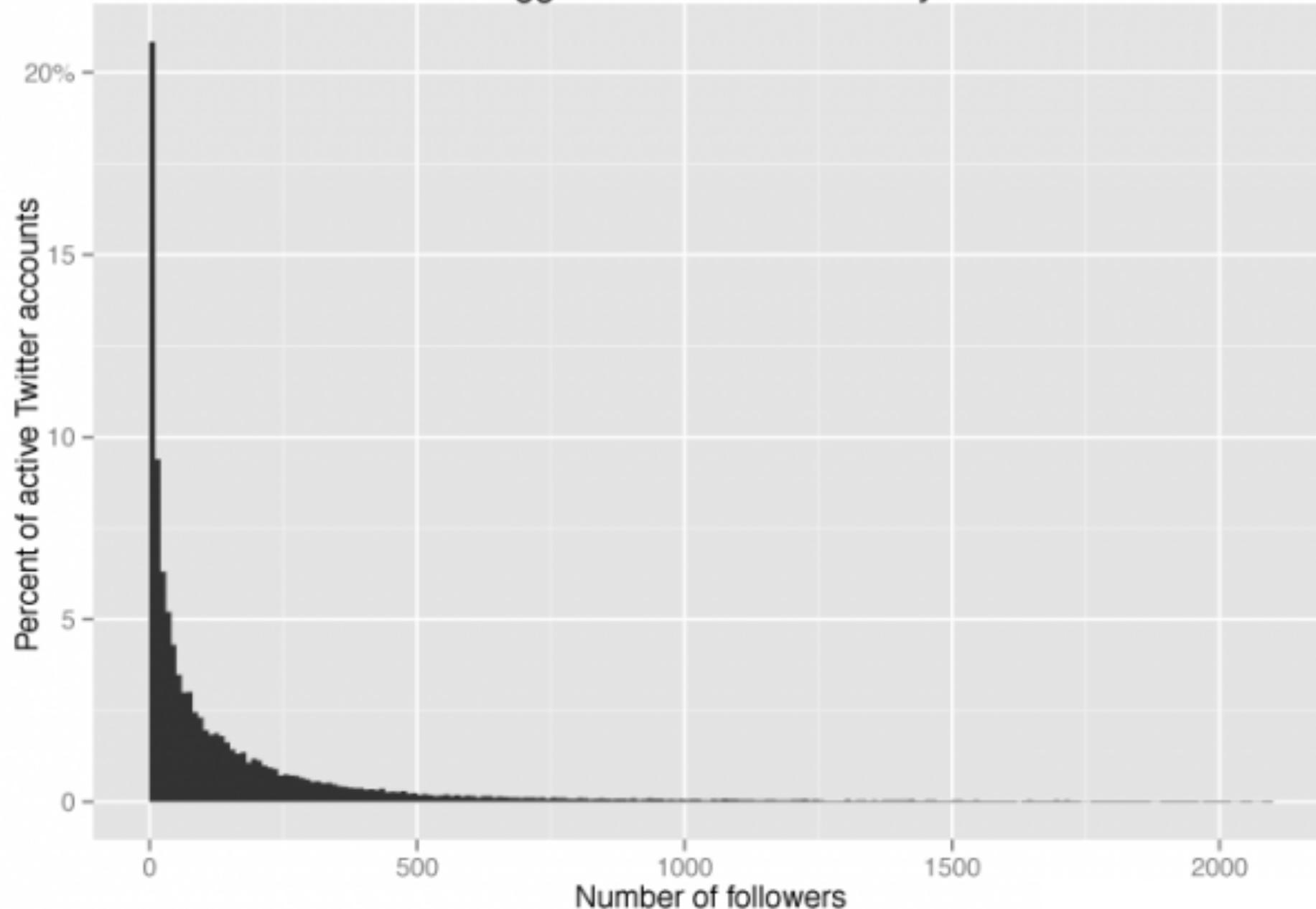


Figure 0 1 Population age structure diagrams for countries with rapid, slow, zero, and negative population growth rates. (Data from Population Reference Bureau)

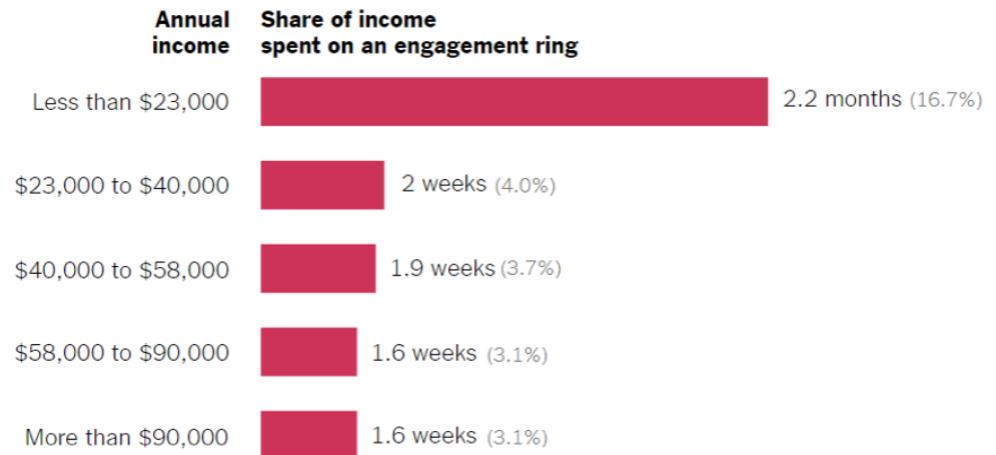
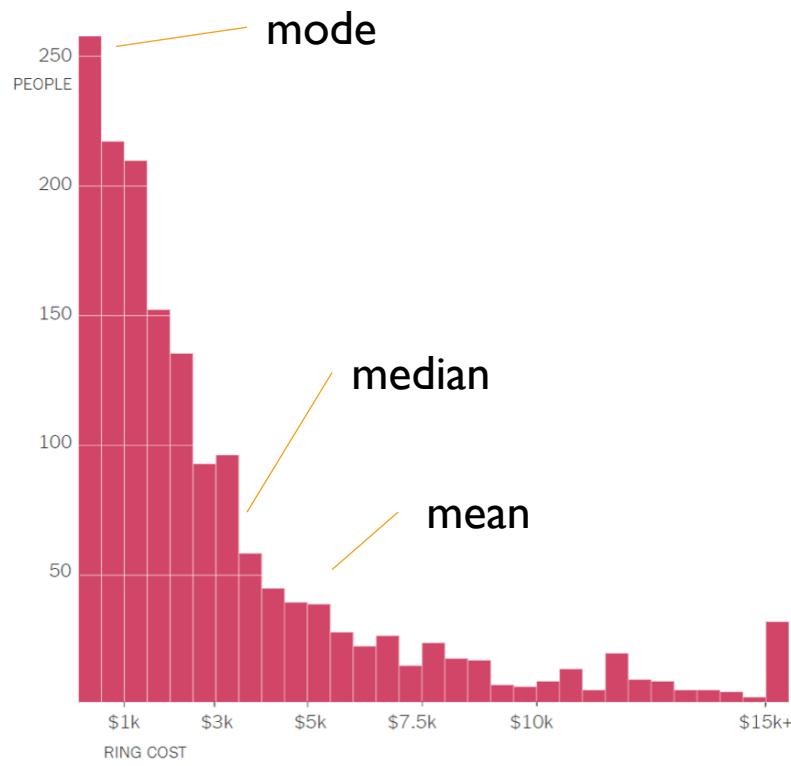
You're a bigger deal on Twitter than you think



Source: BHIVE Social Media Labs

Engagement Ring Costs

How much do people spend on engagement rings?



Source: National poll of 1,640 adults between Jan. 24 – 30, 2019 by Morning Consult

Interactives

- [Population mountains](#)
- [Human terrain](#)
- [Summers are getting hotter](#)

Learn about the article & graph [here!](#)

Outline

1.What is a data visualization?

2.Types of visualizations

3.What makes a good visualization?

MATPLOTLIB FOR VISUALIZATIONS

CHART TYPES ▾

- Boxplot
- Scatterplot
- Histogram
- Network
- Barplot
- Area chart
- Wordcloud
- Density
- Violin
- Heatmap
- Other ..

TOOLS



THE PYTHON
GRAPH GALLERY

Matplotlib



Matplotlib is a plotting library for the Python programming language. It allows to make quality charts in few lines of code. Most of the other python plotting library are build on top of Matplotlib. It makes that a basic understanding

of matplotlib is probably needed to make any chart with python. I highly advise you to have a look to the [matplotlib homepage](#) and have a look to this [general concept page](#). This page aims to give a few tip concerning the general

usage of Matplotlib. It gives examples showing how to custom your title, the colors of your chart, how to annotate it etc. If you need to make a chart in particular, visit the [welcome page](#) of the gallery to find the one you need!

Note that datacamp offers a good and free [online course](#) on Matplotlib. Worth it.

<https://python-graph-gallery.com/matplotlib/>

This is a great resource for information on aesthetics as well as code for each element (e.g., ticks, titles, colors)

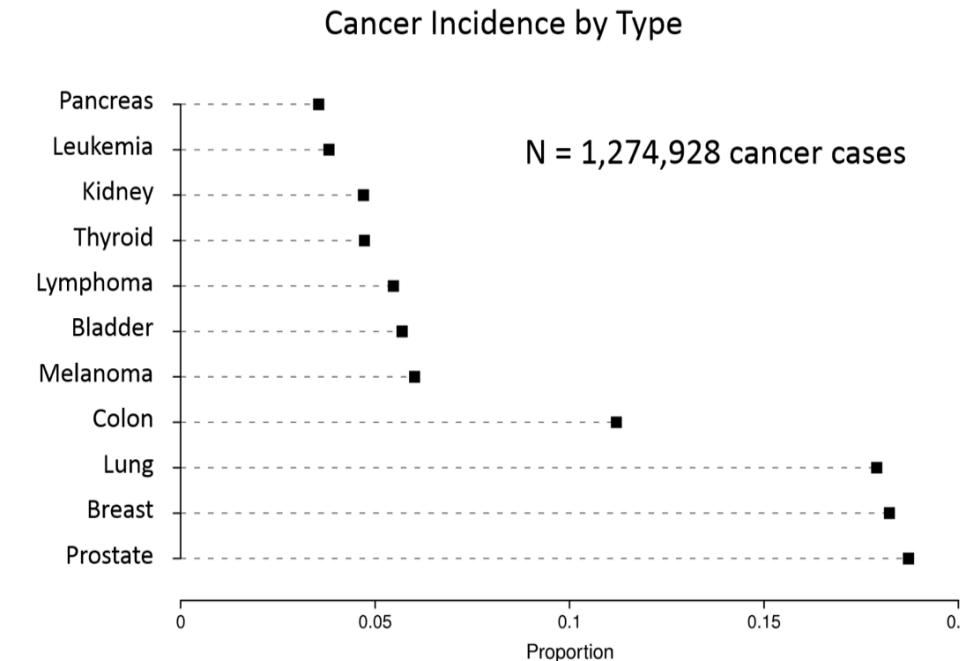
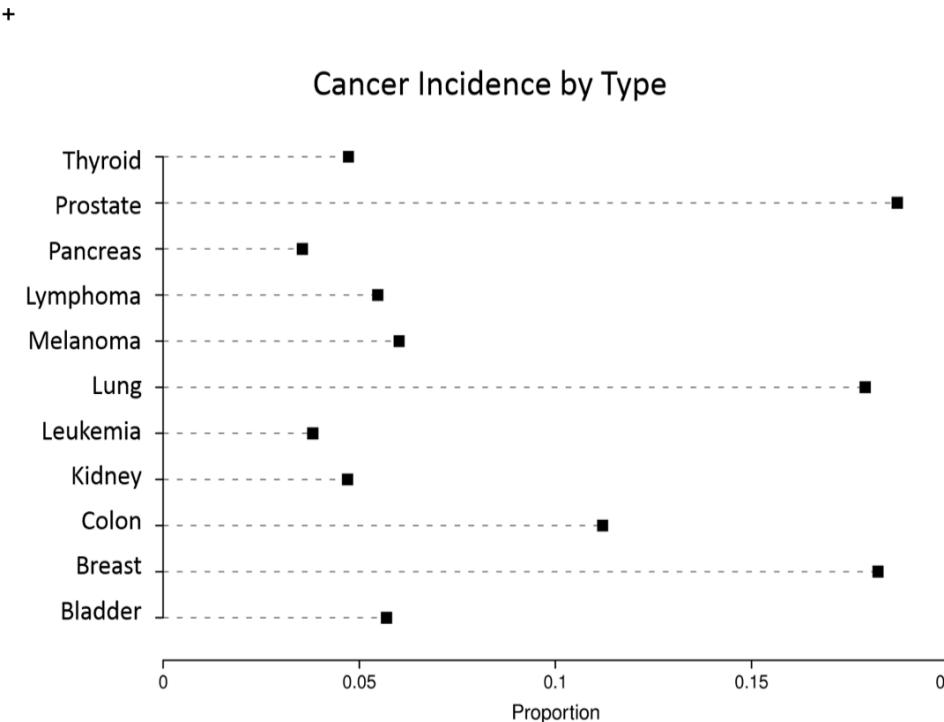
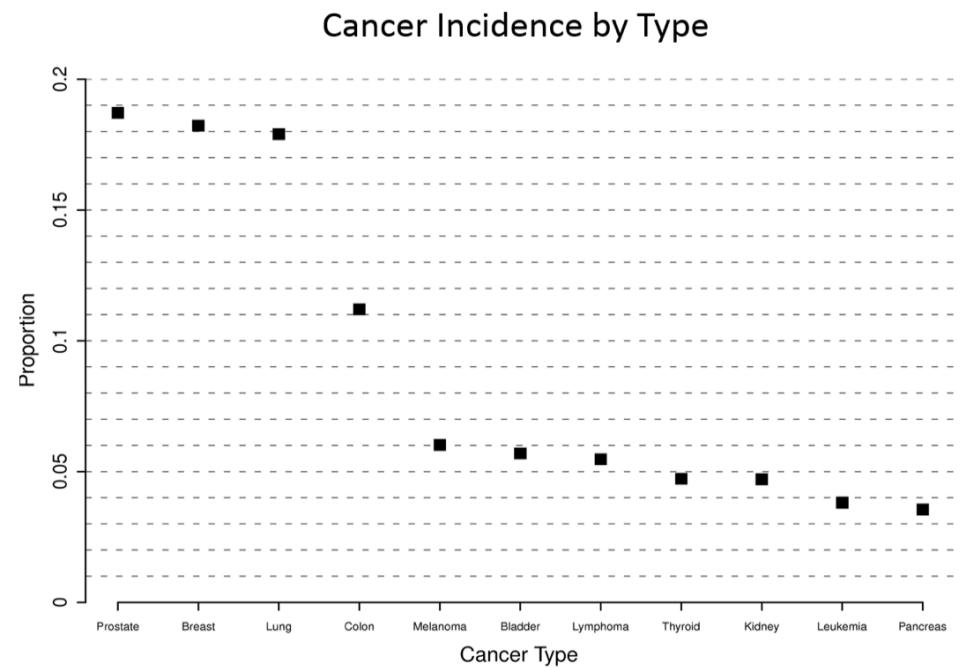
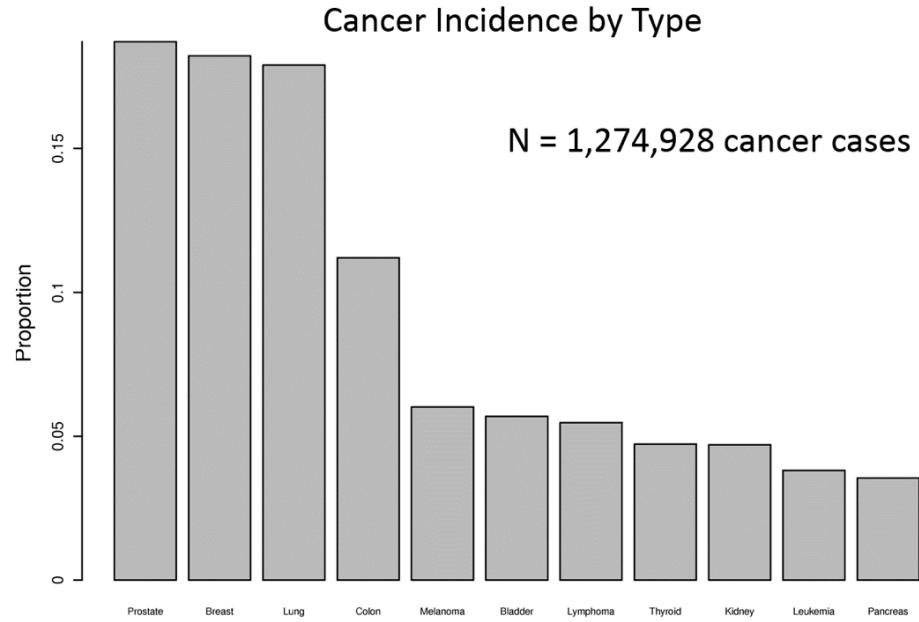
CHOICE OF TYPE OF VISUALIZATION IS IMPORTANT

Given your data type,
what is the type of
graph that is the best
fit?

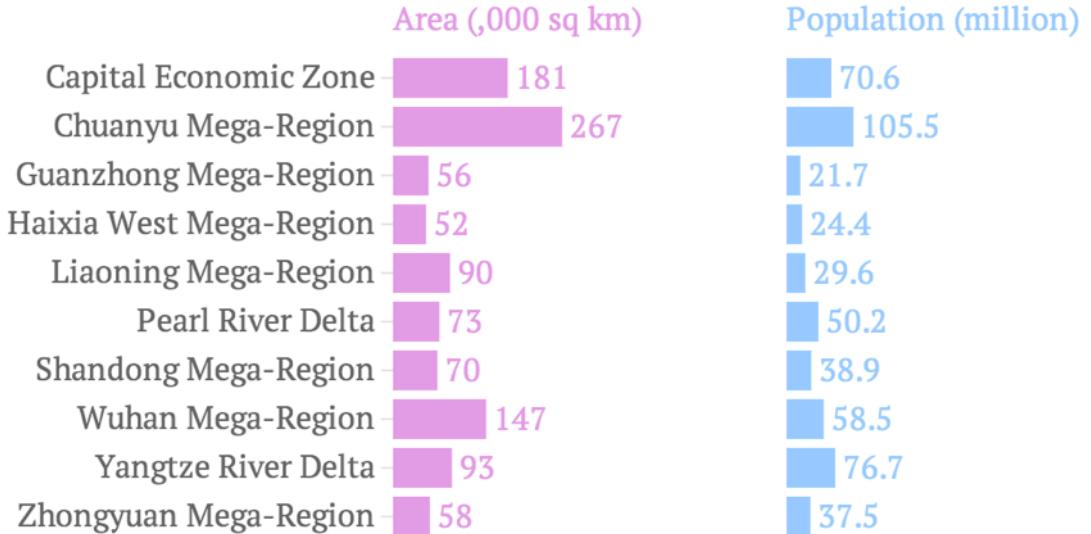
Sense-making

What story are you
trying to tell?
(What point are you
trying to make?)

Communication

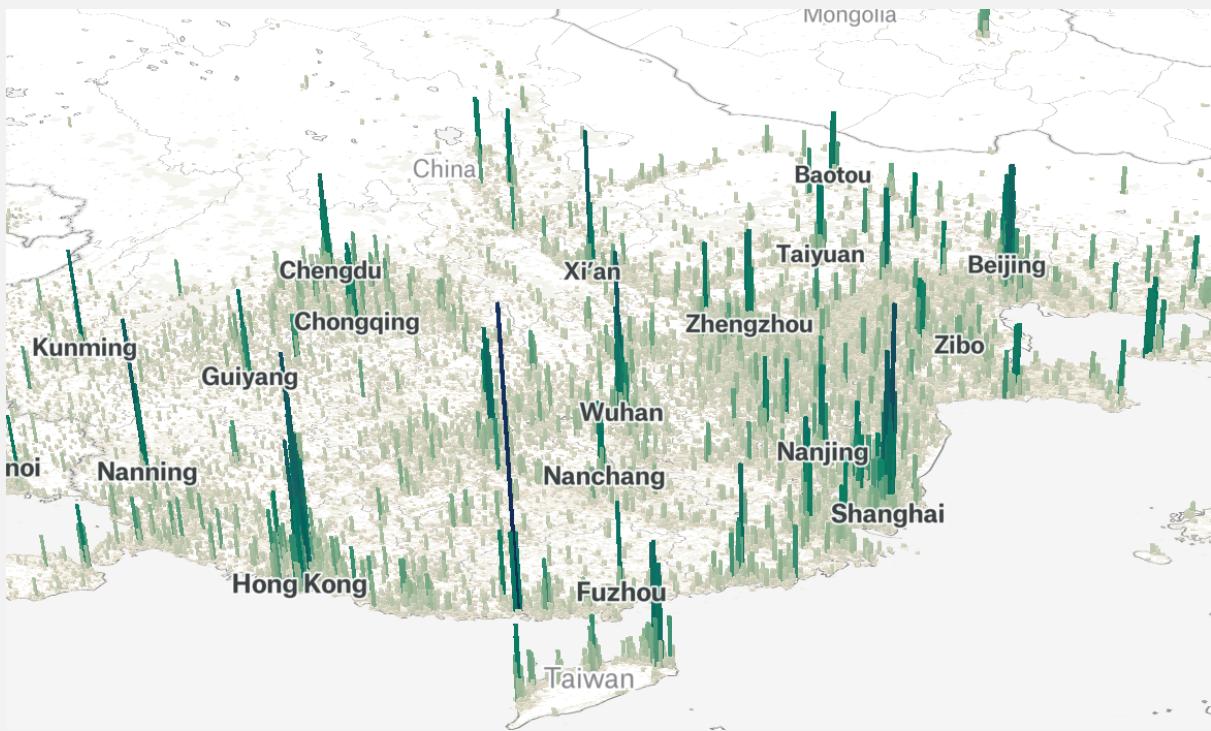


China's mega-regions



Quartz | qz.com

Data: Lincoln Institute of Land Policy



THE BASICS OF GOOD VISUALIZATIONS

1. Title that states the discovery and/or explains what's happening
 - Subtitle: If not in the title, exactly how to interpret what this is
2. Label everything, especially both axes, multiple variables, units, and sources
 - Ticks – Not too many, not too few!
3. All graphs should be able to stand on their own
 - I.e., they should make sense without surrounding text
4. Colors signify information.
 1. Should be used only to clarify or distinguish. Where possible use shapes and/or textures, too.
5. Consistency: Scale, color scheme, size
6. Simplicity!
 - Only include what's necessary. One major point per graphic. Be careful with multiple dimensions

Outline

- 1.What is a data visualization?
- 2.Types of visualizations
- 3.What makes a good visualization?

