



Welcome to DSCI 101

Introduction to Data Science

Week 4-5 Topics Recap

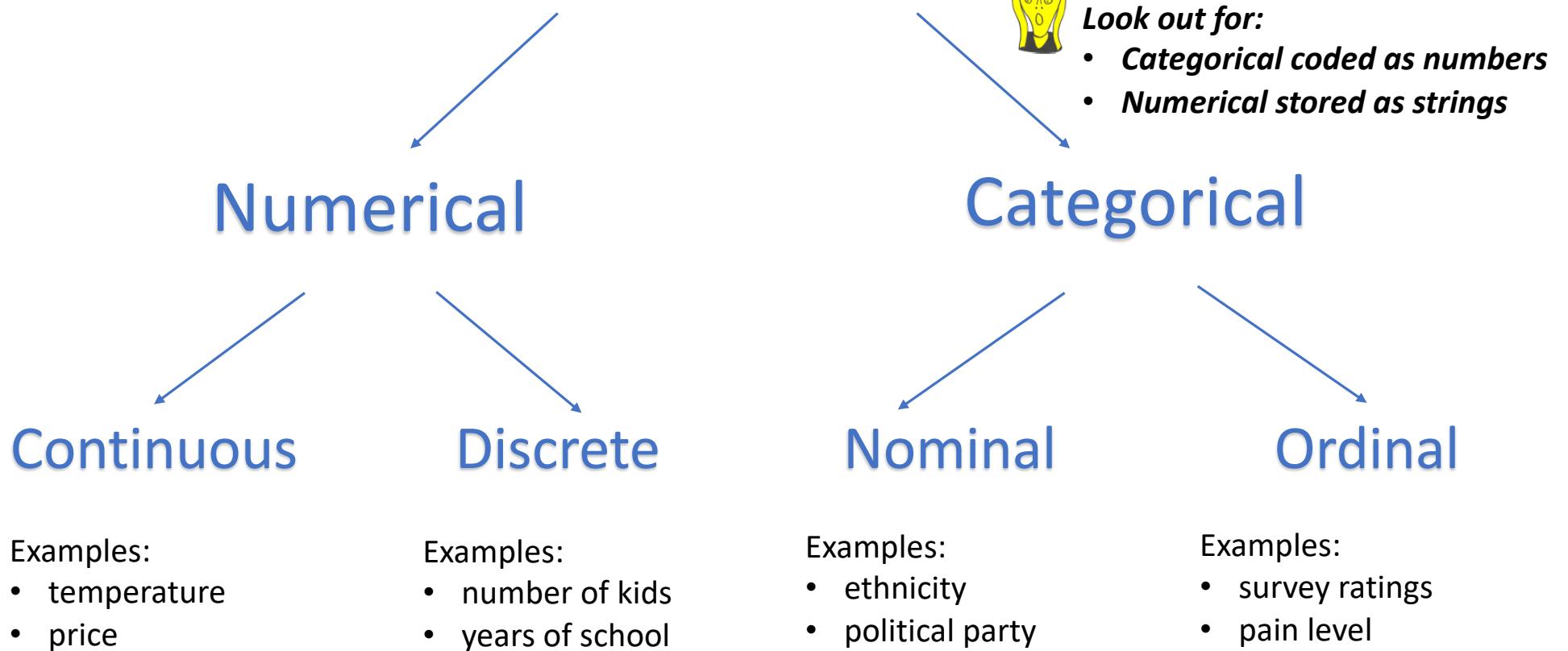
- Basic dataframe methods:
 - import with `pd.read_csv`
 - indexing with `df[], df.loc` and `df.iloc`
 - sort, rename, mutate, drop a column content
 - `df.sample`: get a random sample of rows
- Data summaries and manipulation:
 - `series.value_counts`, `df.groupby`, `df.pivot_table`, `df.describe`
 - `df.merge`: join two df by common column(s)
 - missing values

Week 6 Topics Preview

- Exploring data through visualization:
 - Variable types and basic plots: bar plot, histogram, boxplot, line plot and scatter plot
- Storytelling with data
 - Modern data and advanced plots
- Data visualization learning outcome
 - choose the right plots and implement in Python
 - read and interpret plots

A review

Variable



Categorical – bar plot

- X-axis:
 - categories
- Y-axis can be:
 - numerical
 - count or proportion
- Length of the bar encode values
- Width of the bar encode nothing!

app_type	homeownership			Total	
	rent	mortgage	own		
	individual	3496	3839	1170	8505
joint	362	950	183	1495	1495
Total	3858	4789	1353	10000	10000

Figure 2.17: A contingency table for app_type and homeownership.

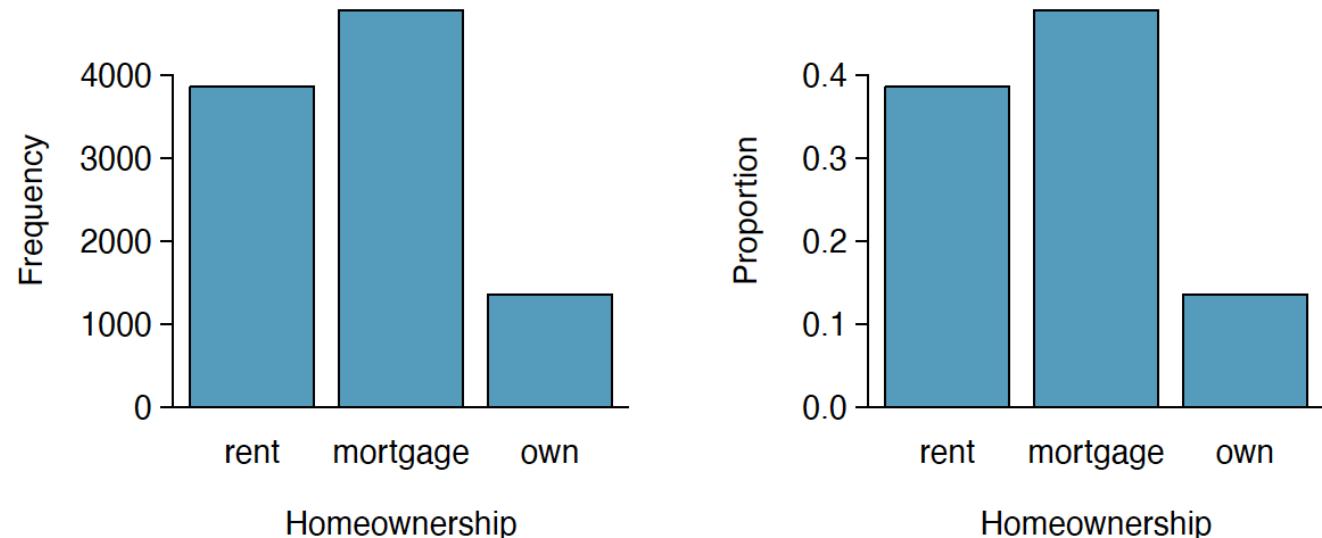


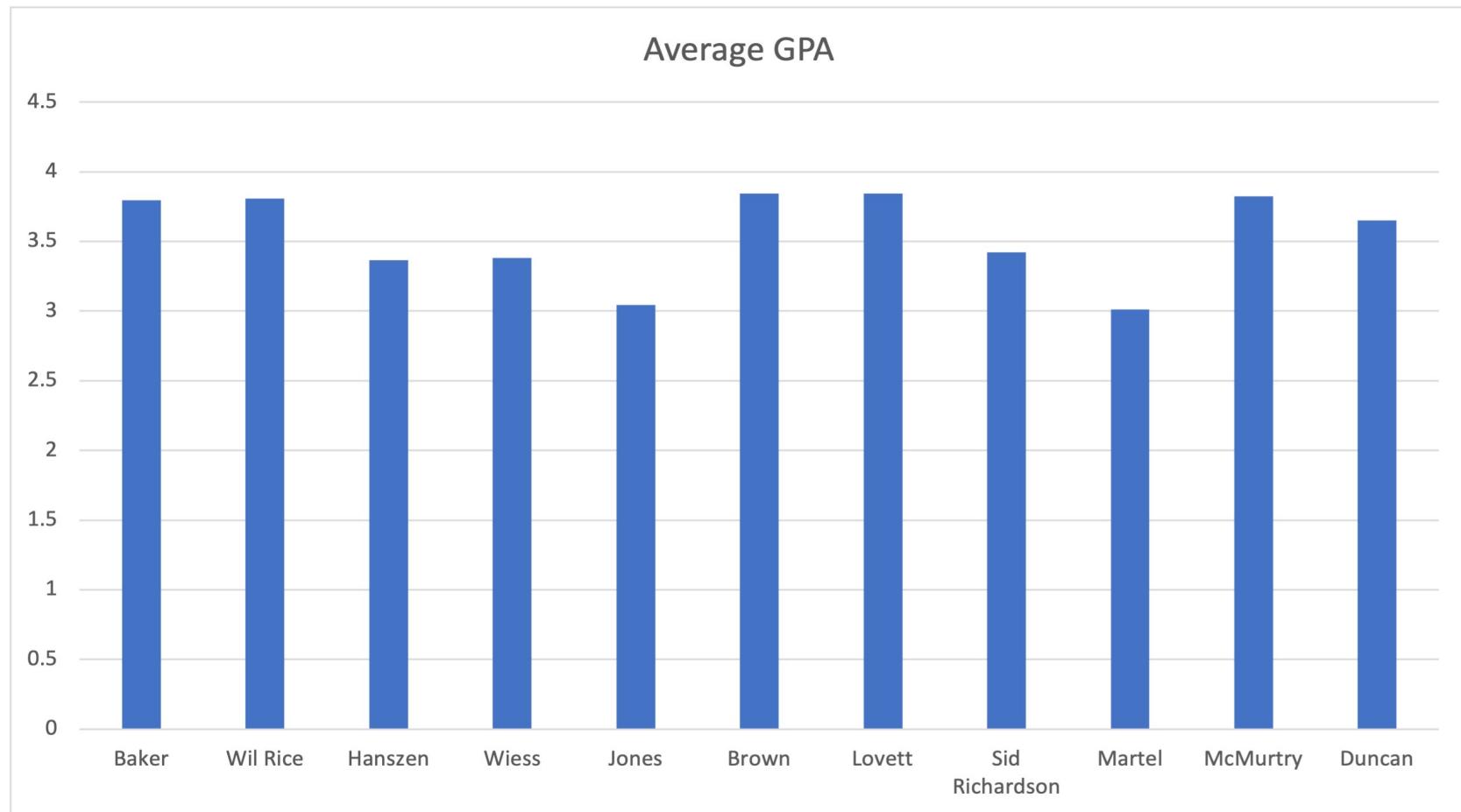
Figure 2.19: Two bar plots of number. The left panel shows the counts, and the right panel shows the proportions in each group.

More bar plots

*Be mindful of
color choices!*

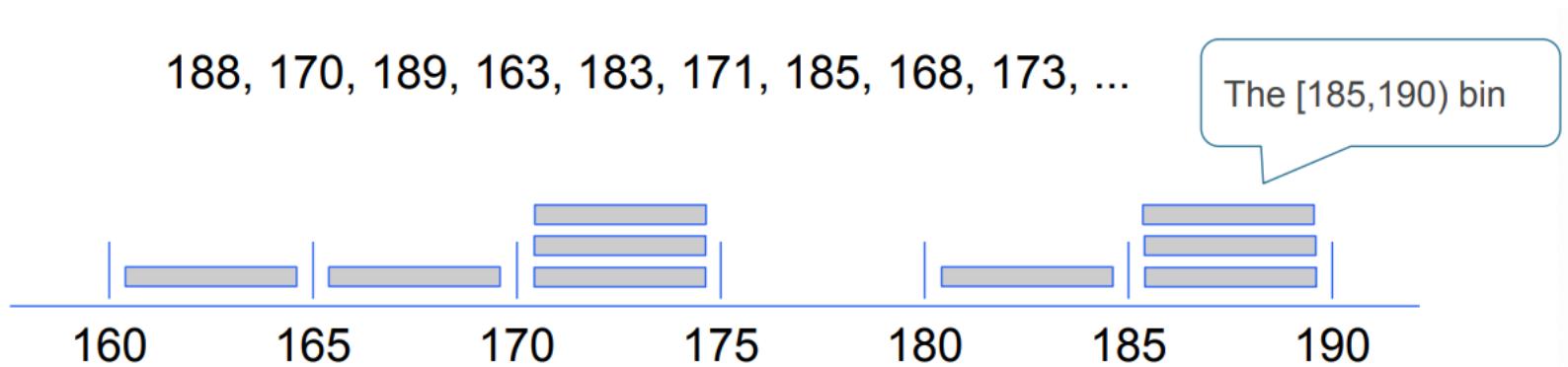


More bar plots – categorical and numerical



Numerical - histogram

- Binning a numerical variable:
 - count the number of numerical values that lie within ranges, called bins
 - Bins are defined by their lower bounds (inclusive), and upper bounds (exclusive)
 - The upper bound is the lower bound of the next bin



Visualizing distributions

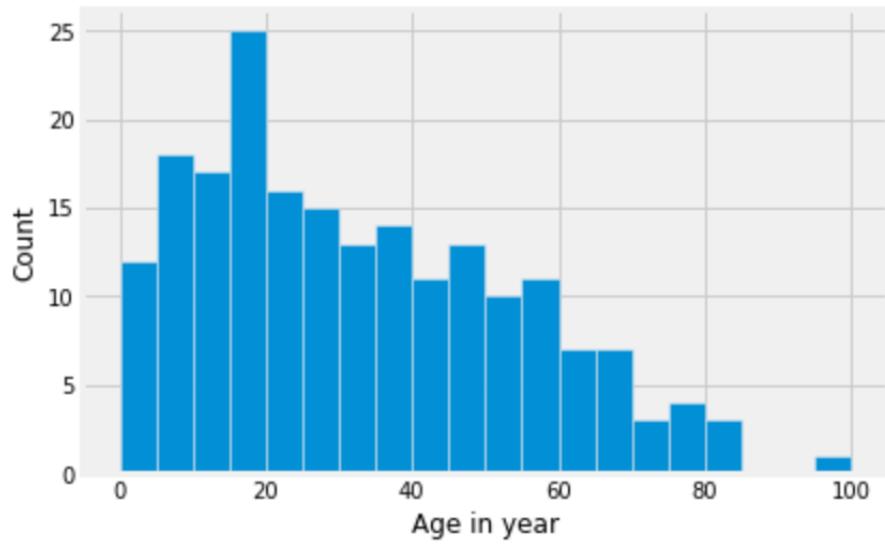
Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.



Histogram – two ways

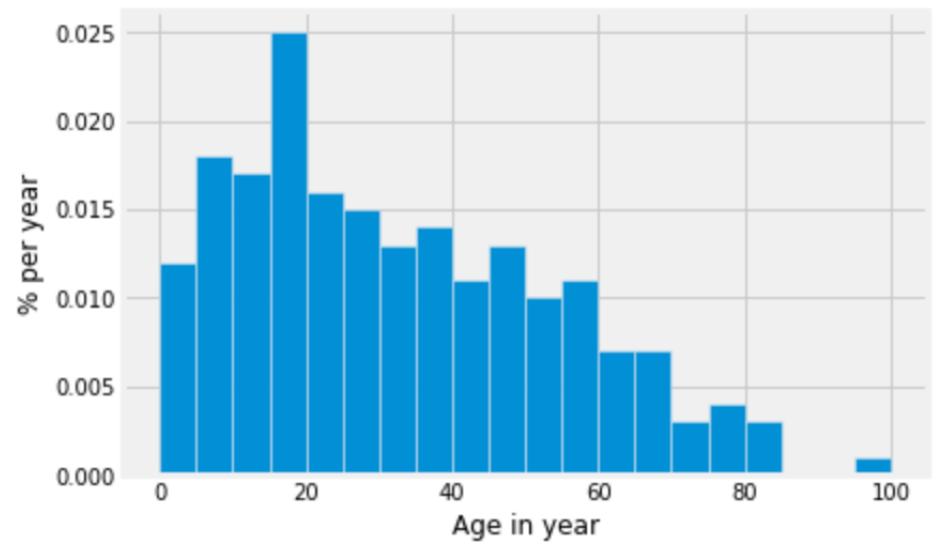
- Y-axis show counts

- intuitive
- most common (default)
- even sized bins only!!!



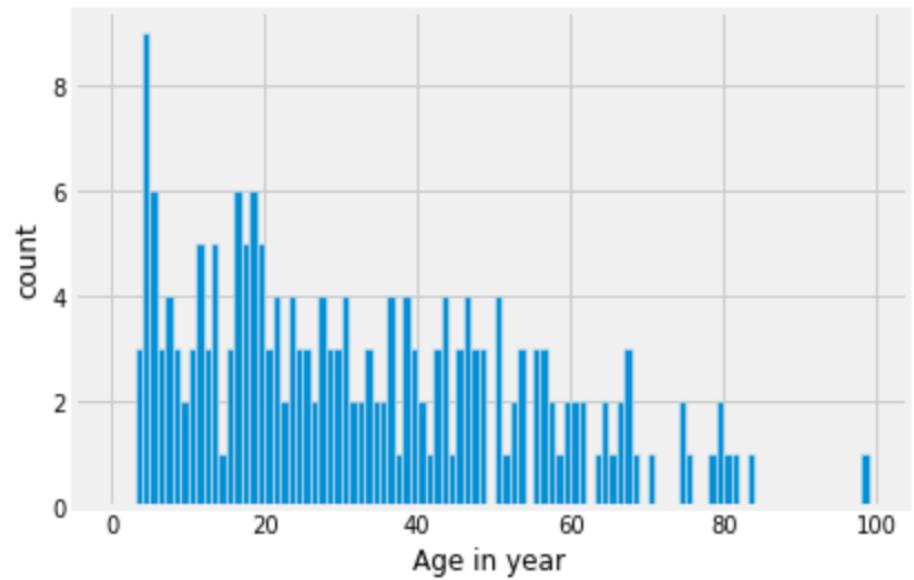
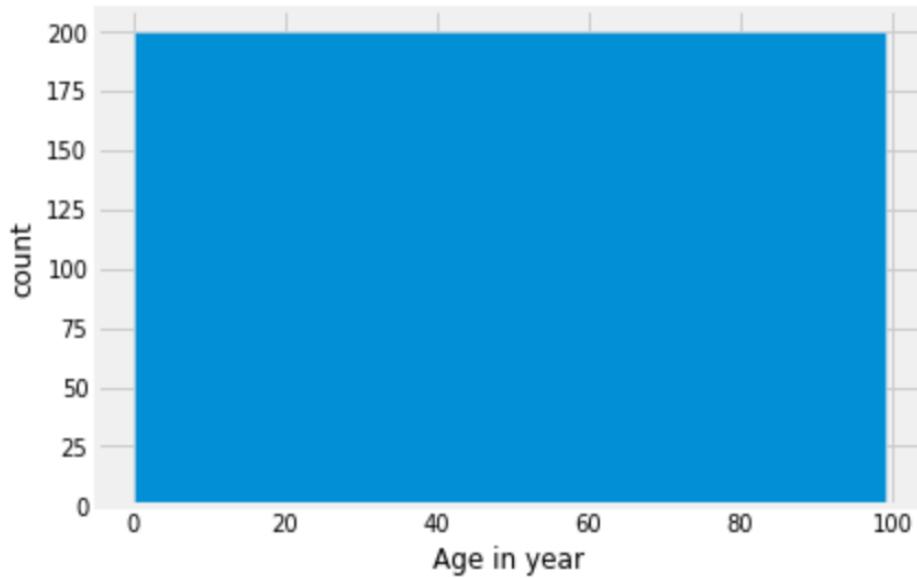
- Y-axis show density

- area principle
- height = % / bin width
- can use uneven bins

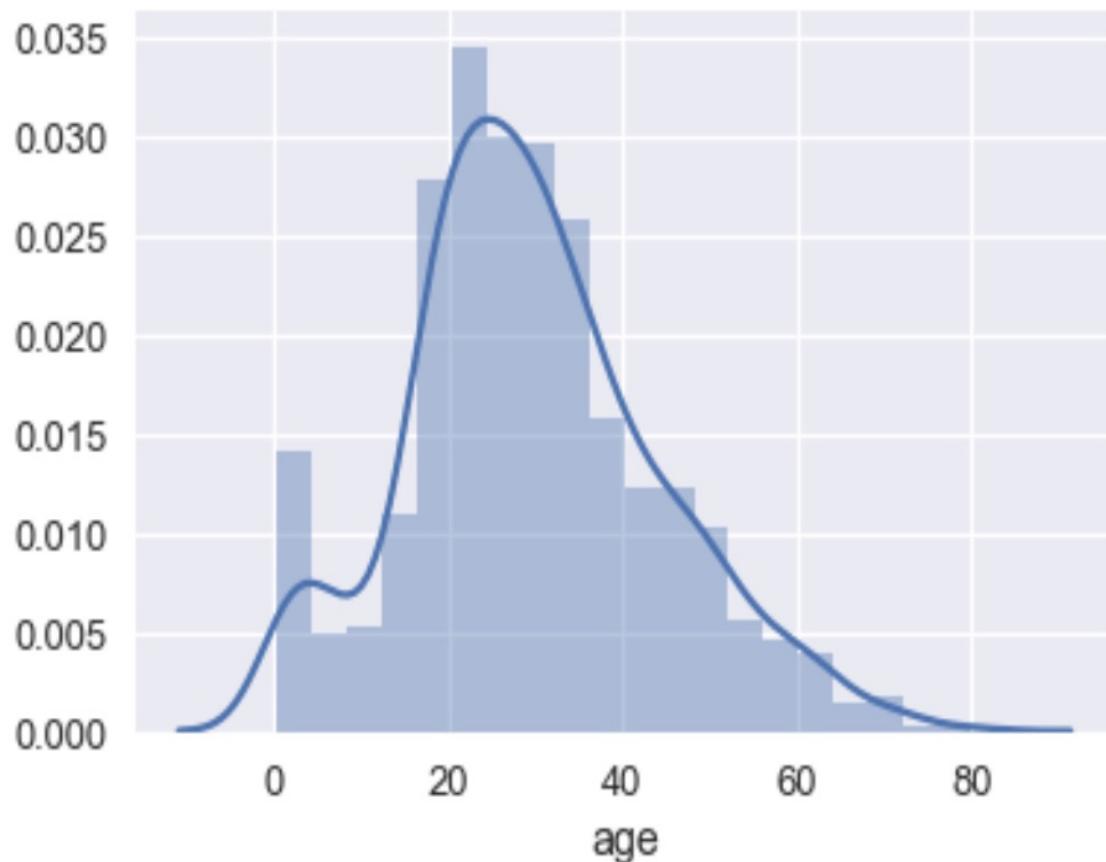


How to choose bin size?

- Think about what bin size encodes?
 - really large bin size v.s. really small bin size
 - trade off: granularity v.s. interpretability



Histogram and kernel density estimate



- Remember height of the bars in histogram measure density
- “Smooth out” the height
- KDE
 - estimate a **probability density function** (pdf) from data
 - DSCI 301 / 303

Bar plot

vs.

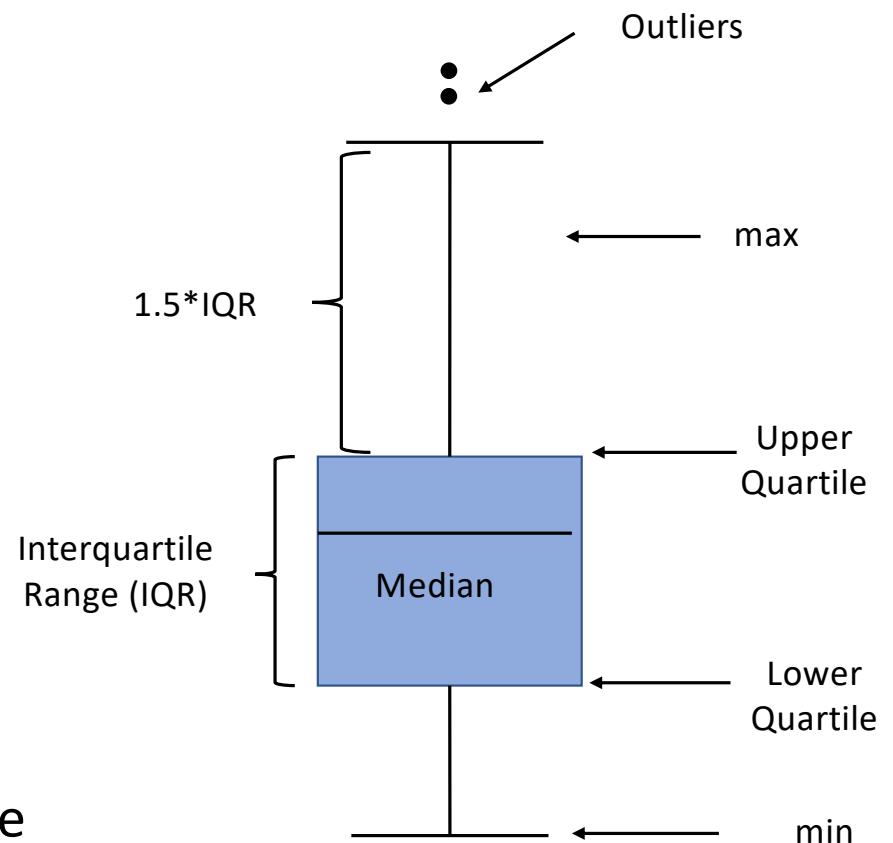
Histogram

- Distribution of one categorical variable
- Bars have arbitrary width
- Bars usually do not touch

- Distribution of one numerical variable
- Bar width defined by bins
- No gap between bars

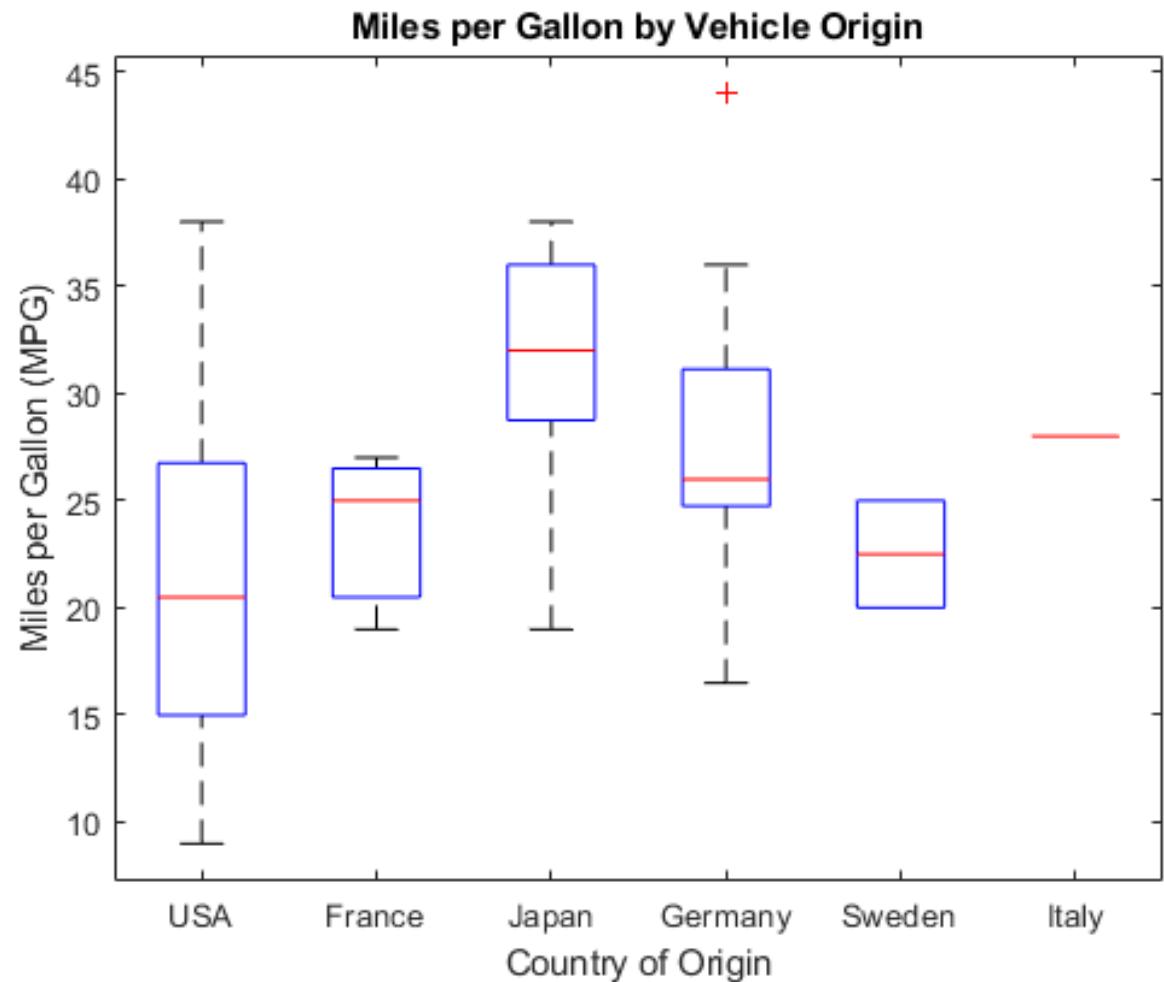
Box plot

- The pth percentile:
 - Smallest value that has at least p% of the data at or below it
- Quartile: 25% & 75% percentile
- Median: 50% percentile
- IQR: 75% percentile – 25% percentile
(the middle 50% of the data)
- Outlier: $1.5 * \text{IQR}$ away from lower and upper quartiles



More box plot

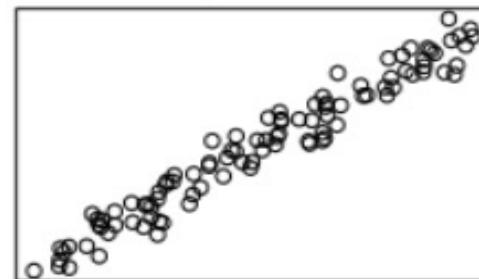
- side-by-side box plots
- interaction between a numerical and categorical variables
- comparison of the same measure among different groups



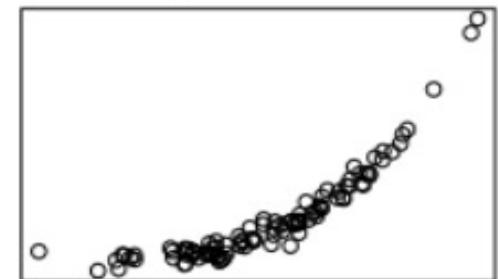
Scatter plot

- Two numerical variables
- Relationship between the pair
- Identify data pattern to help inform model choices

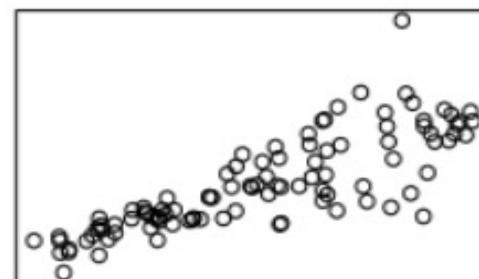
simple linear



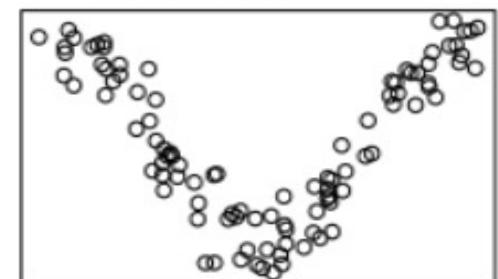
simple nonlinear



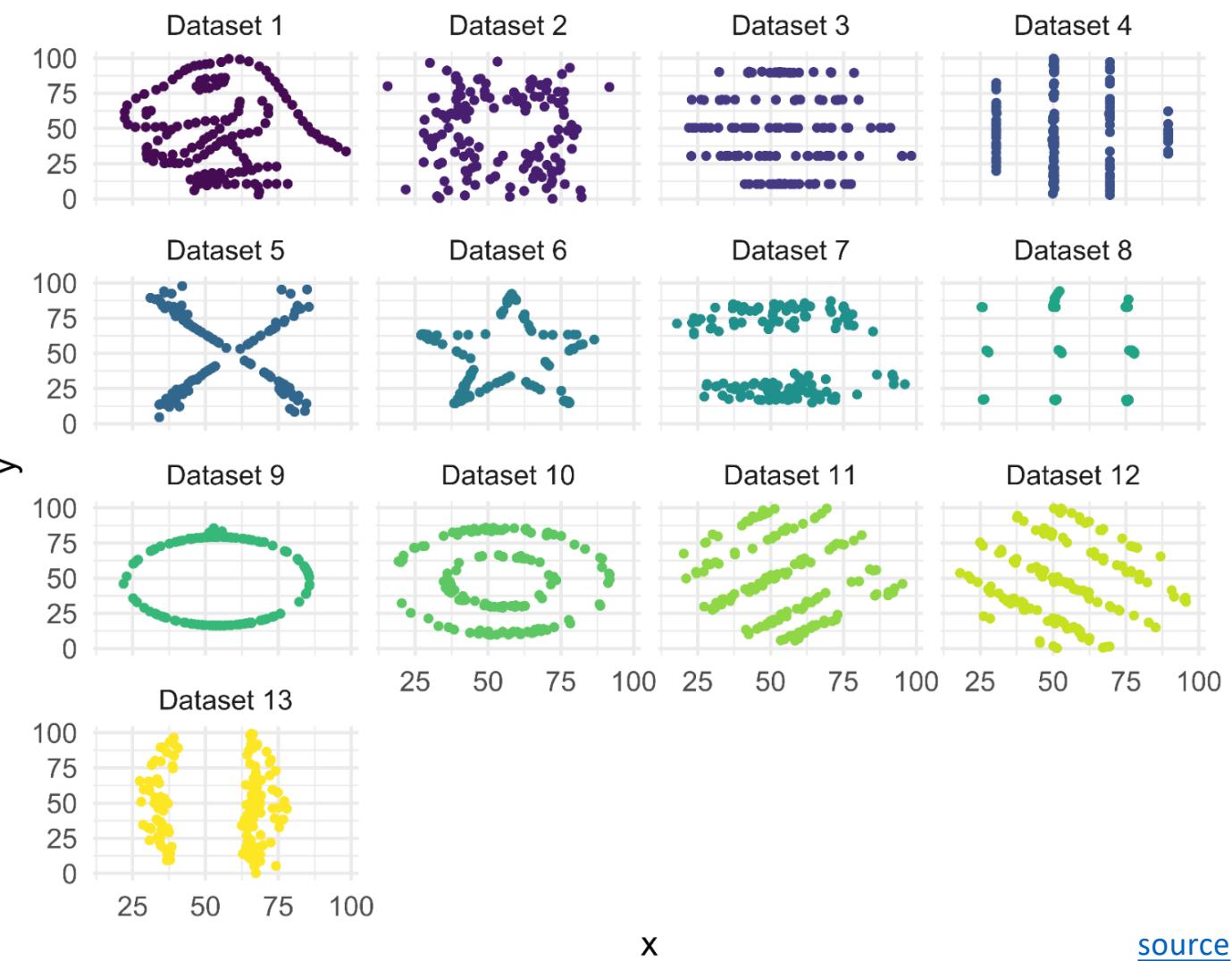
unequal spread



complex nonlinear



dataset	n	Average x	Average y	St Dev x	St Dev y	Correlation
Dataset 1	142	54.3	47.8	16.8	26.9	-0.1
Dataset 2	142	54.3	47.8	16.8	26.9	-0.1
Dataset 3	142	54.3	47.8	16.8	26.9	-0.1
Dataset 4	142	54.3	47.8	16.8	26.9	-0.1
Dataset 5	142	54.3	47.8	16.8	26.9	-0.1
Dataset 6	142	54.3	47.8	16.8	26.9	-0.1
Dataset 7	142	54.3	47.8	16.8	26.9	-0.1
Dataset 8	142	54.3	47.8	16.8	26.9	-0.1
Dataset 9	142	54.3	47.8	16.8	26.9	-0.1
Dataset 10	142	54.3	47.8	16.8	26.9	-0.1
Dataset 11	142	54.3	47.8	16.8	26.9	-0.1
Dataset 12	142	54.3	47.8	16.8	26.9	-0.1
Dataset 13	142	54.3	47.8	16.8	26.9	-0.1



"The simple graph has brought more information to the data analyst's mind than any other device."

John Tukey

[source](#)

Line plot

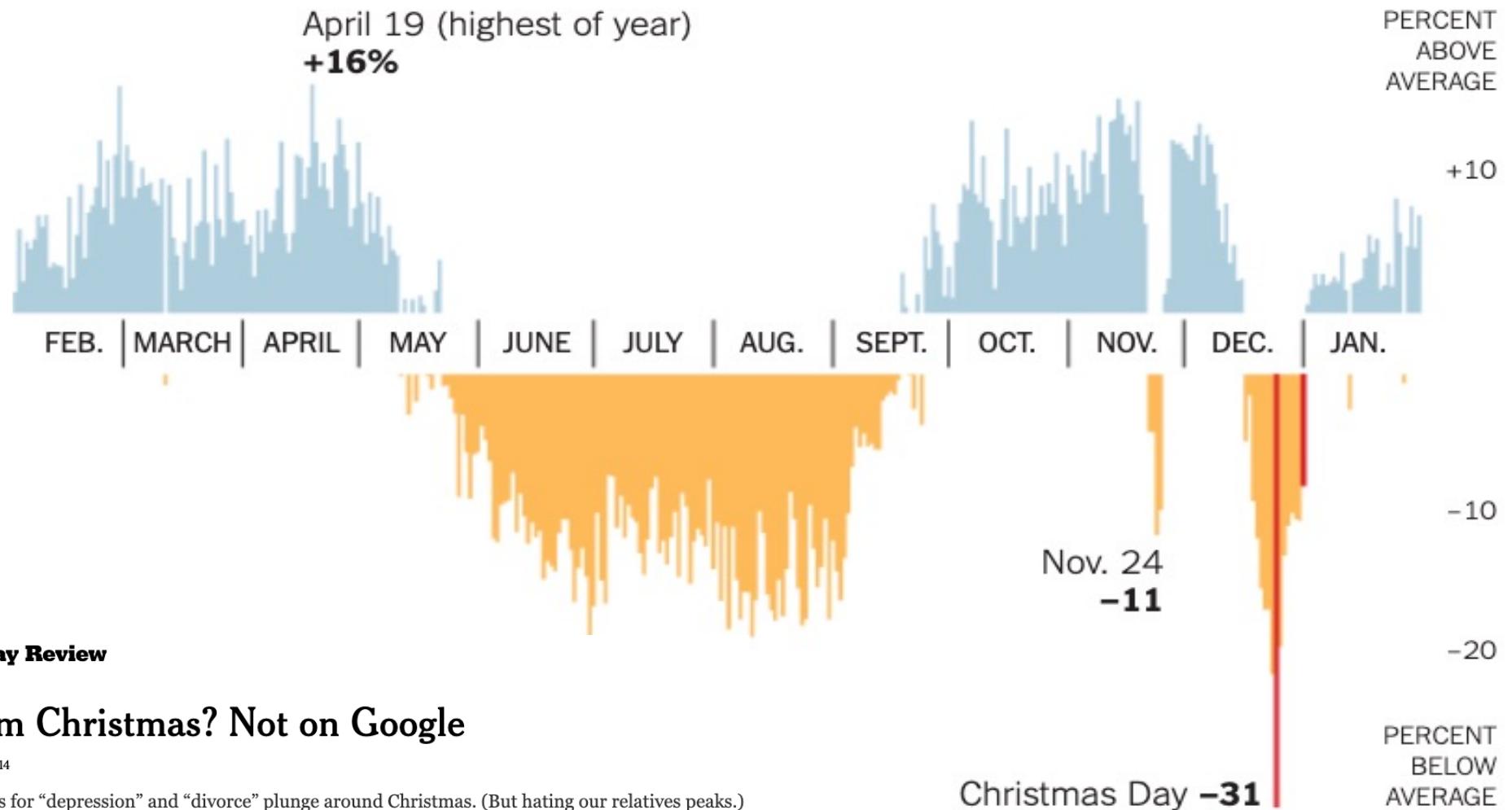
- Two numerical variables
- Often the variable on x-axis is time
- Best to show changes over time and predict future trend
- Let's look at some plots of life expectancy and fertility rate per women over time for some countries.
- Can you guess which country?

More plots of real data to look at

- Can you tell what kind of plot it is?
- Can you read the plot and explain what's going on?
- What do you think makes a good plot?

depression

<https://www.nytimes.com/interactive/2014/12/21/sunday-review/glum-christmas-not-on-google.html?mtrref=www.google.com&assetType=REGIWALL>



The emotional arcs of stories are dominated by six basic shapes

Andrew J Reagan^{1*}, Lewis Mitchell², Dilan Kiley¹, Christopher M Danforth¹ and Peter Sheridan Dodds¹

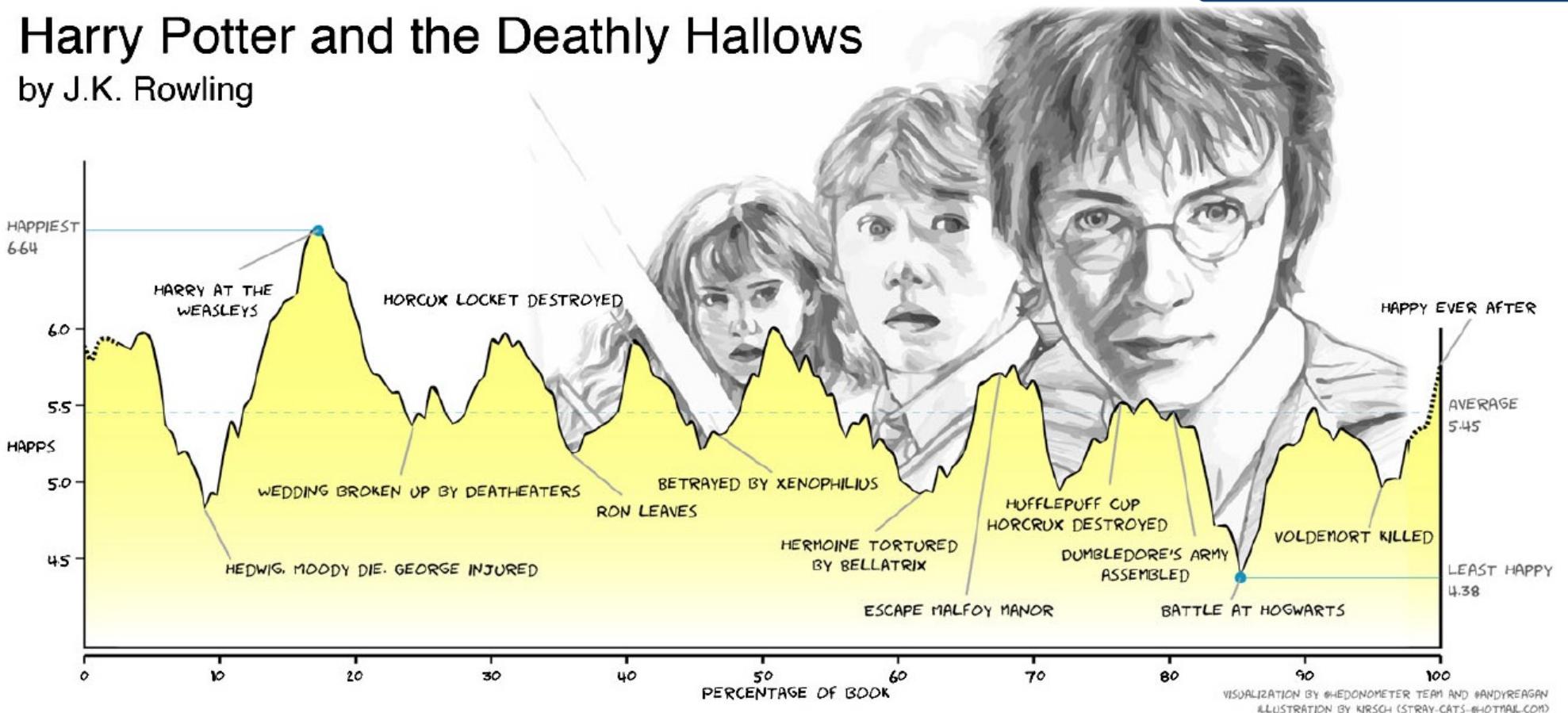
Reagan et al. *EPJ Data Science* (2016) 5:31
DOI 10.1140/epjds/s13688-016-0093-1

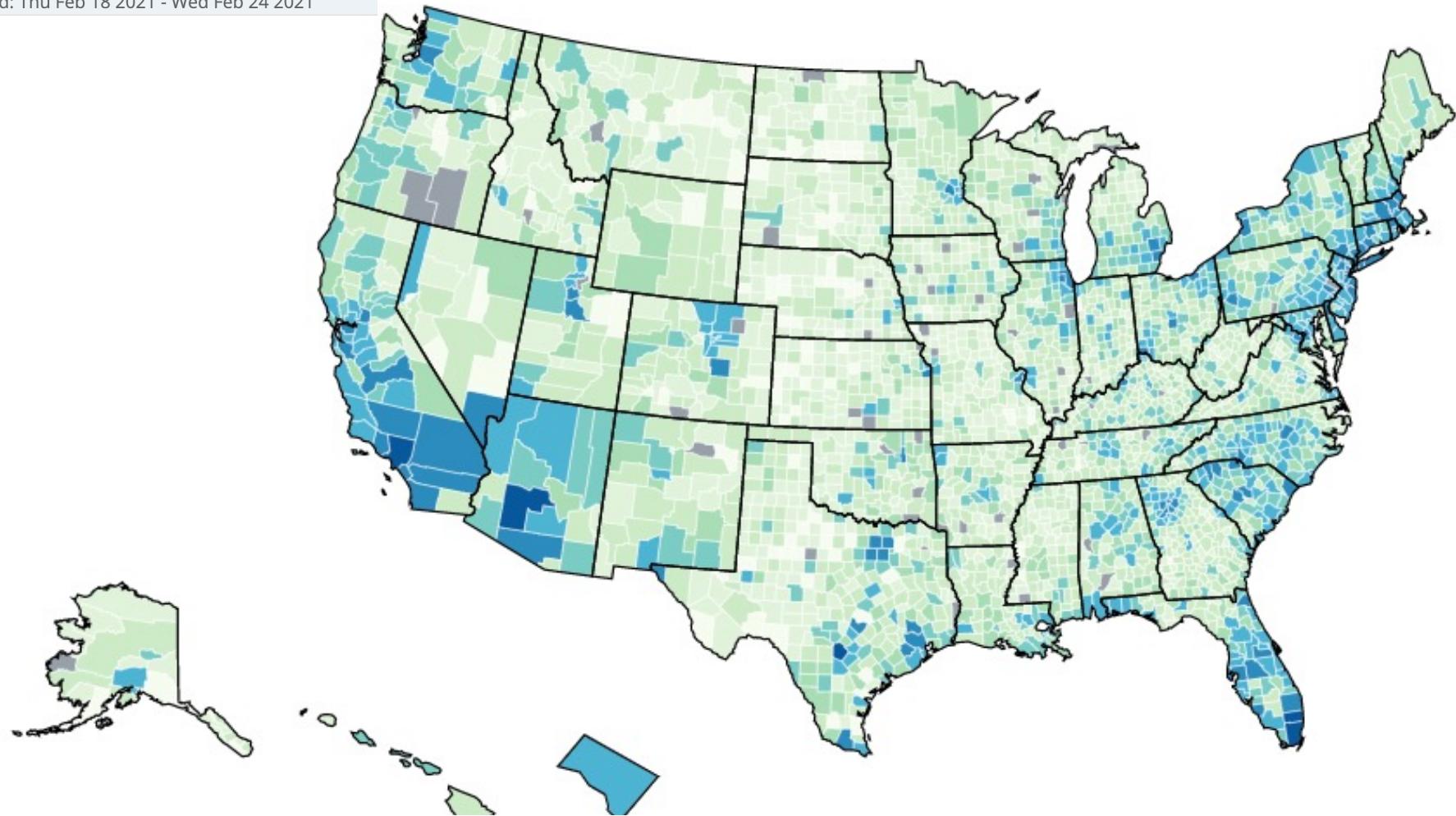
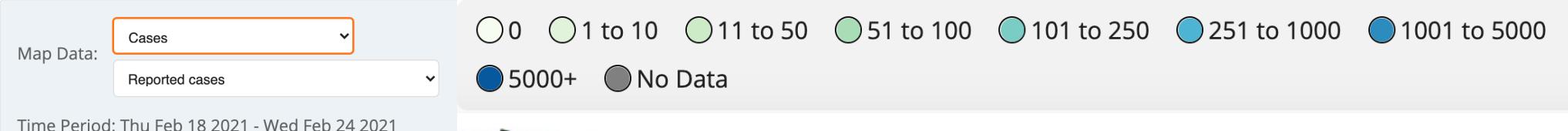


REGULAR ARTICLE

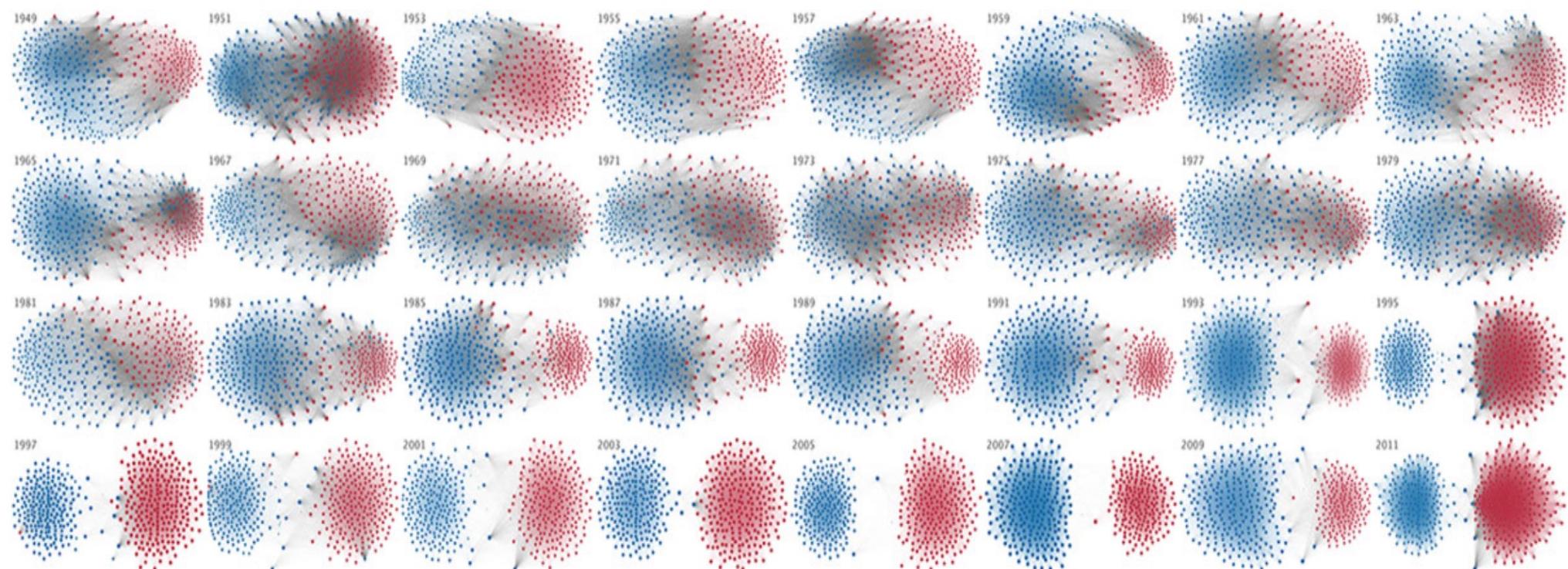
Harry Potter and the Deathly Hallows

by J.K. Rowling





The rise of partisanship in the US House of Representatives



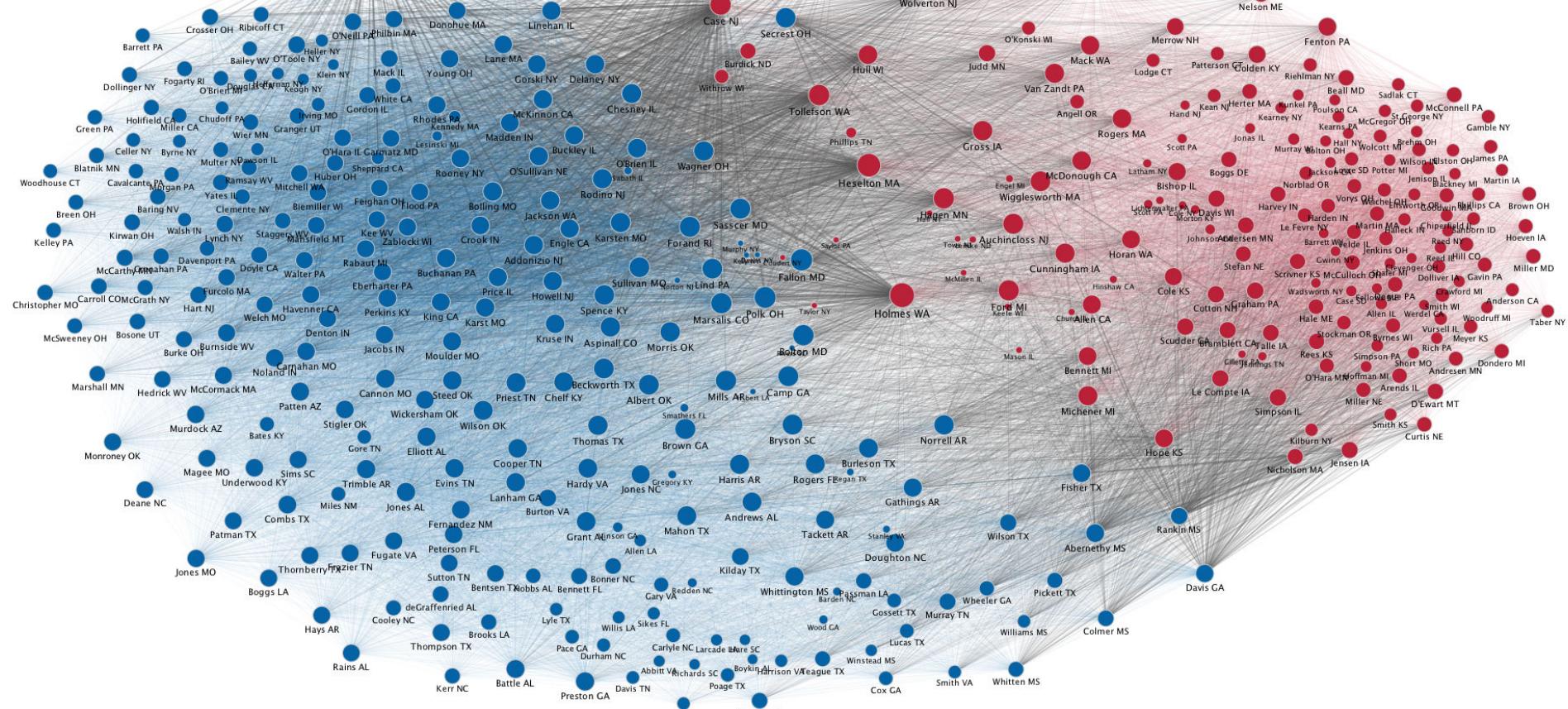
https://www.mamartino.com/projects/rise_of_partisanship/index.html

Year: 1949

D-D ●● D-R ●● R-R ●●

Degree 1 o 400

Few ||||| Many

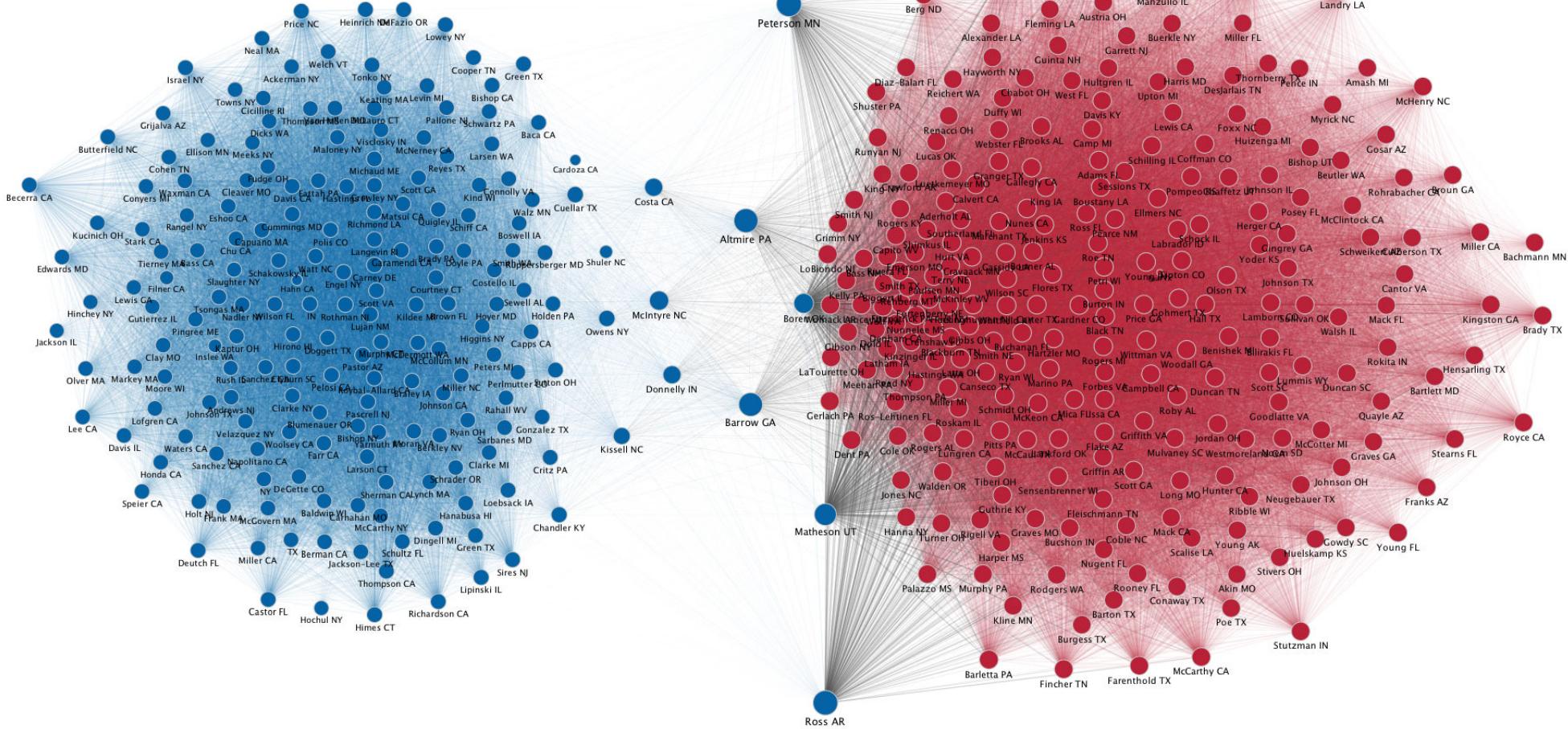


Year: 2011

D-D ●● D-R ●● R-R ●●

Degree 1 ○ 400

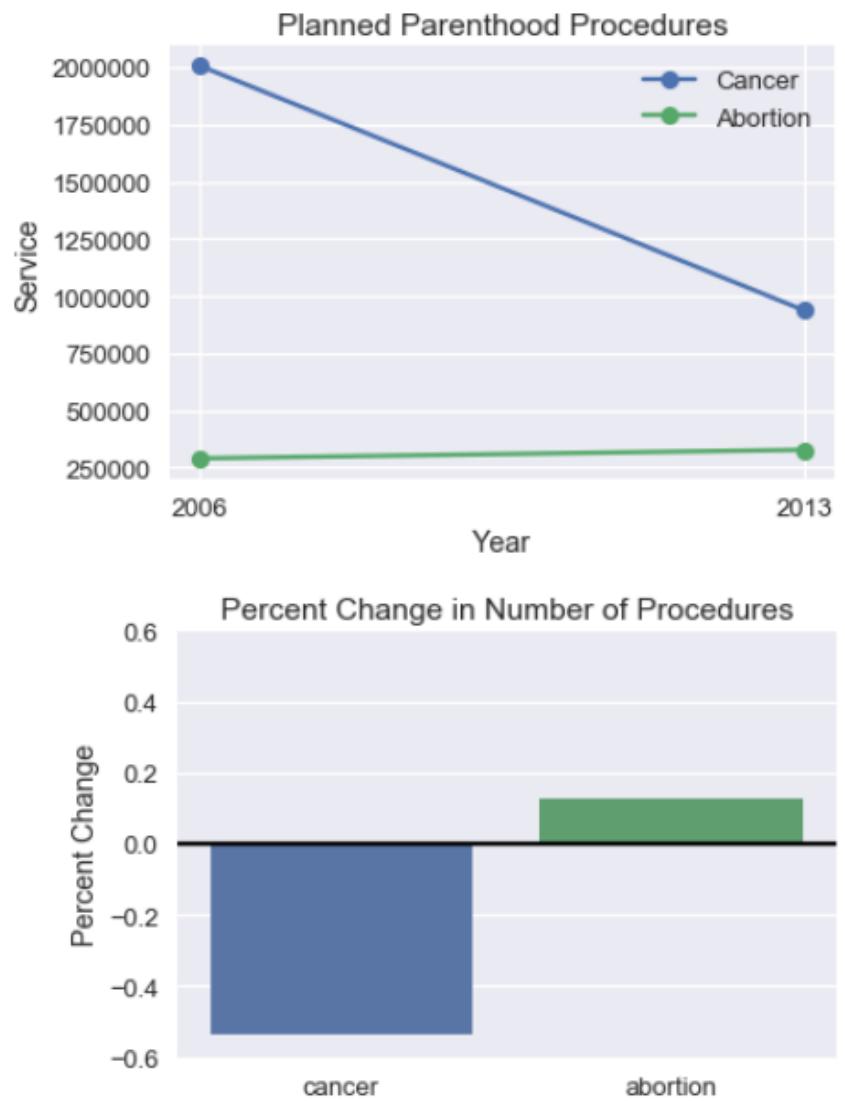
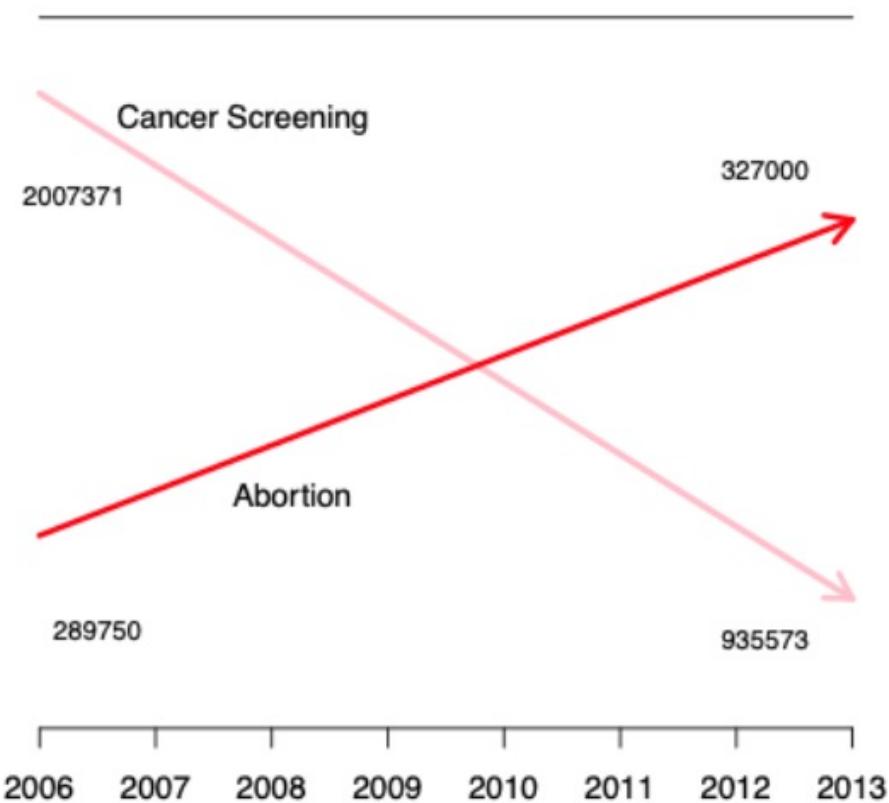
Few ||||| Many



Recipe for good data visualization

- Good understanding of your data and your purpose!
 - What does your data tell?
 - What do you want to tell?
- Some principles
- A dash of creativity

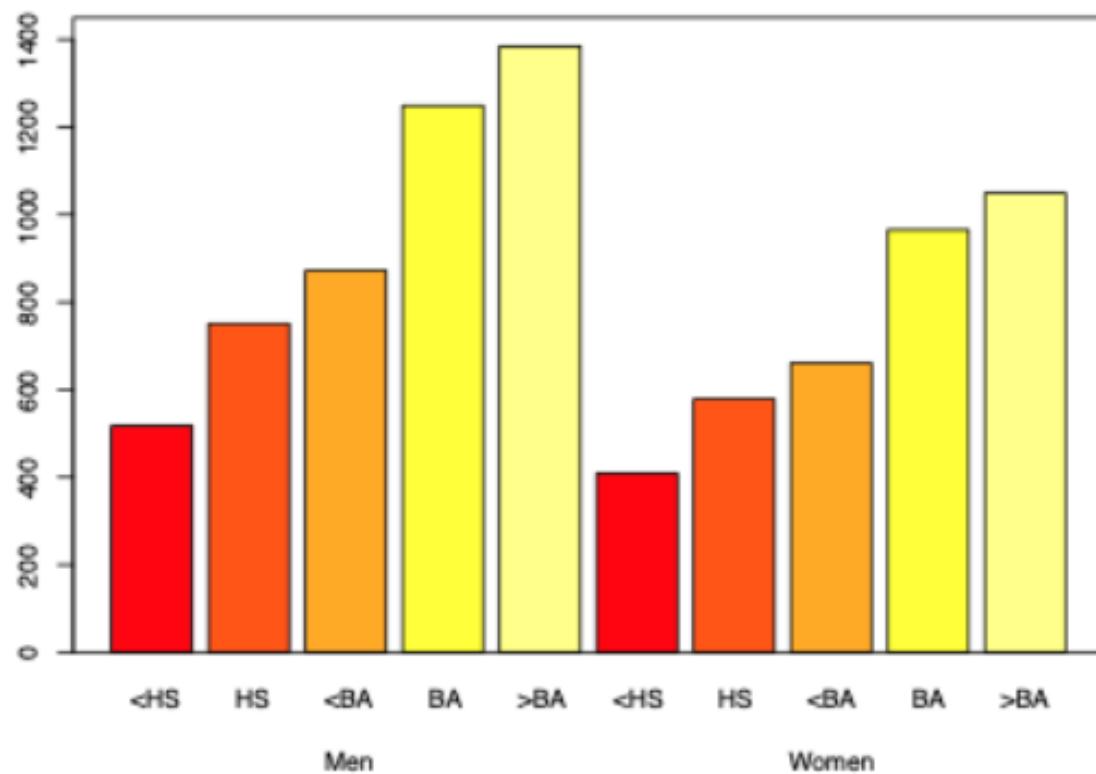
Scale



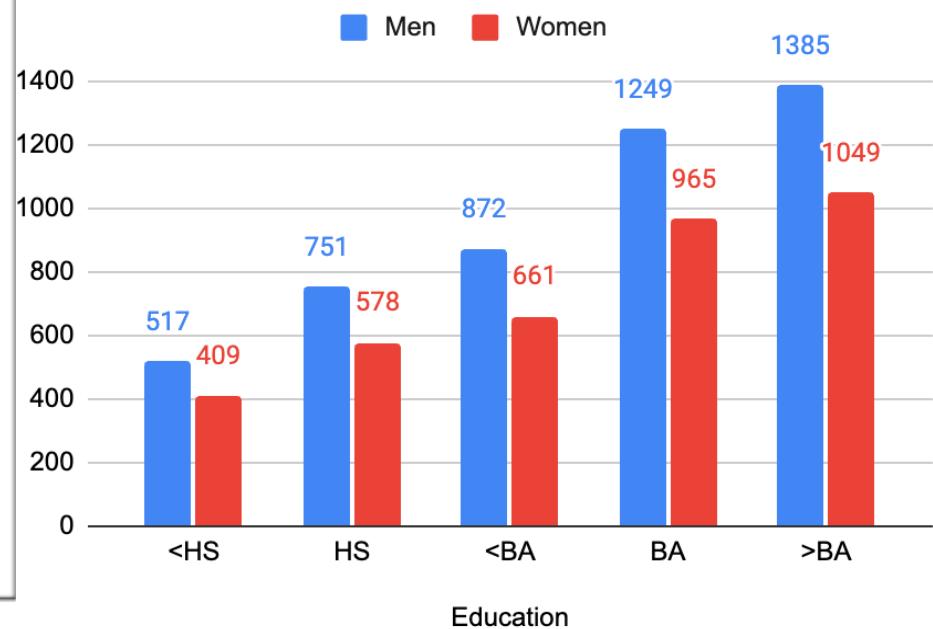
<https://oversight.house.gov/interactivepage/plannedparenthood/>

Comparison

2014 Median Weekly Earnings
Full-Time Workers over 25 years old



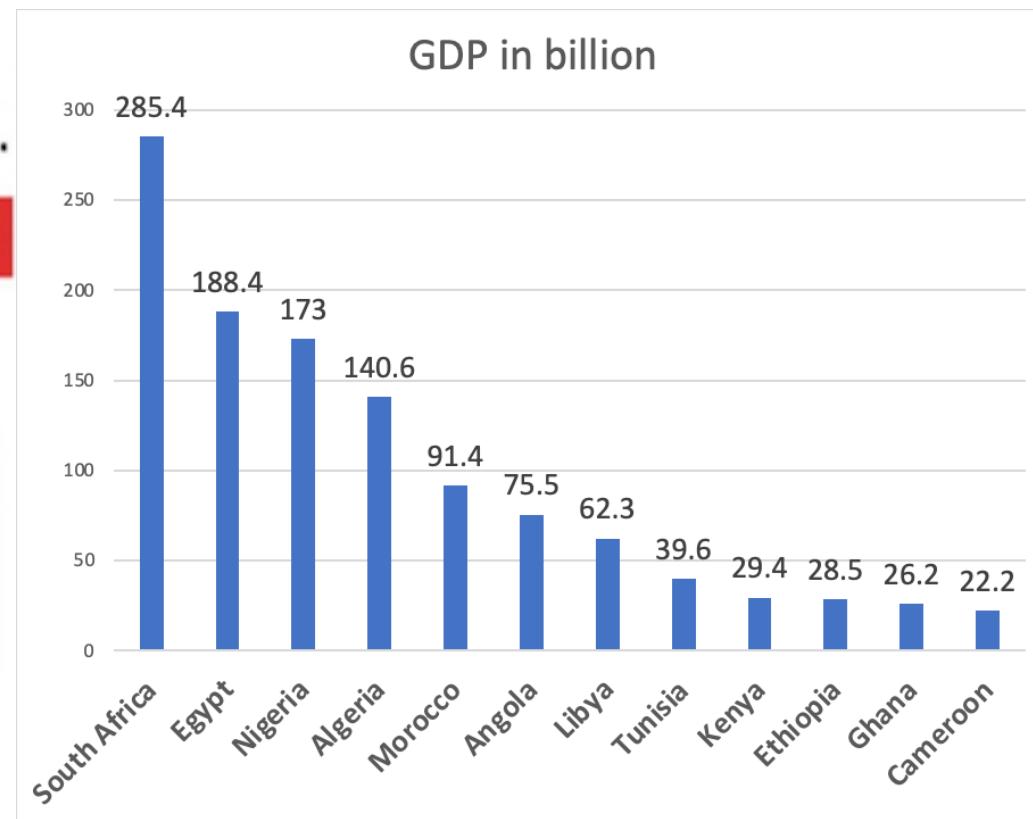
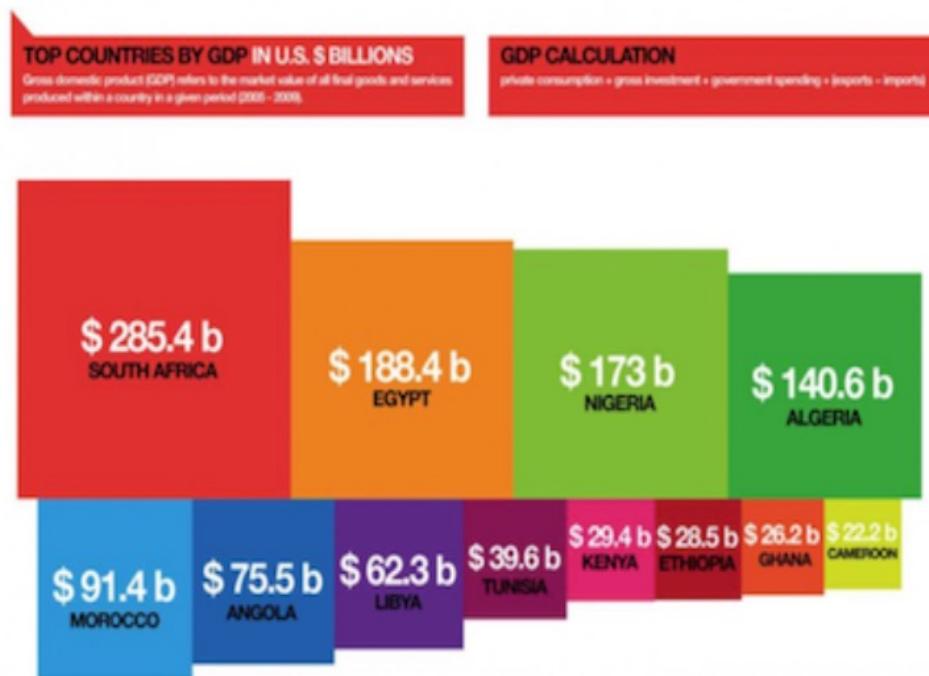
2014 median weekly earnings for full-time workers over 25 years old



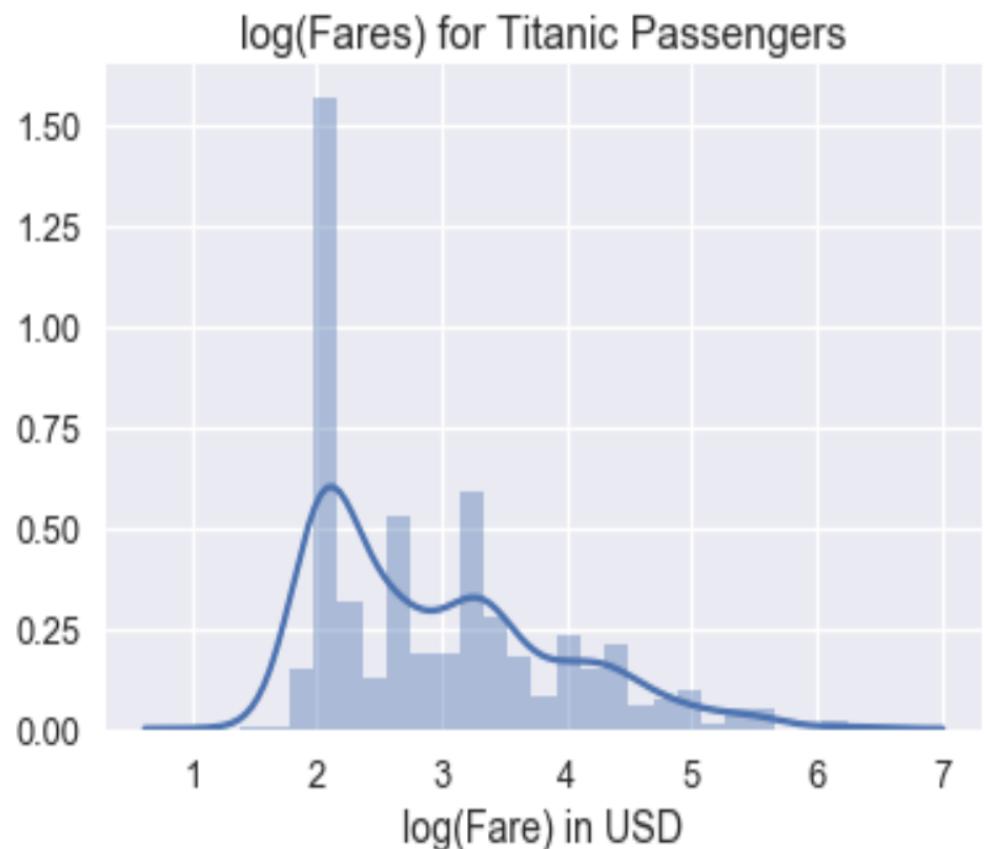
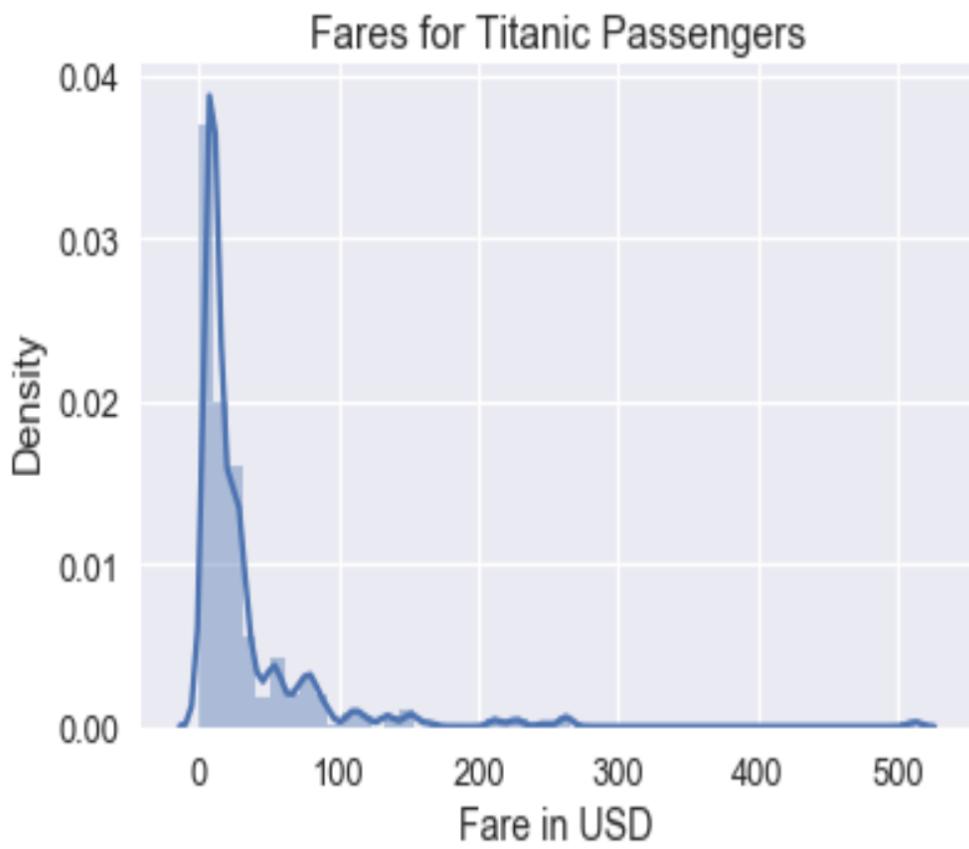
[Bureau of Labor Statistics](#)

Perception

African Countries by GDP



Transformation

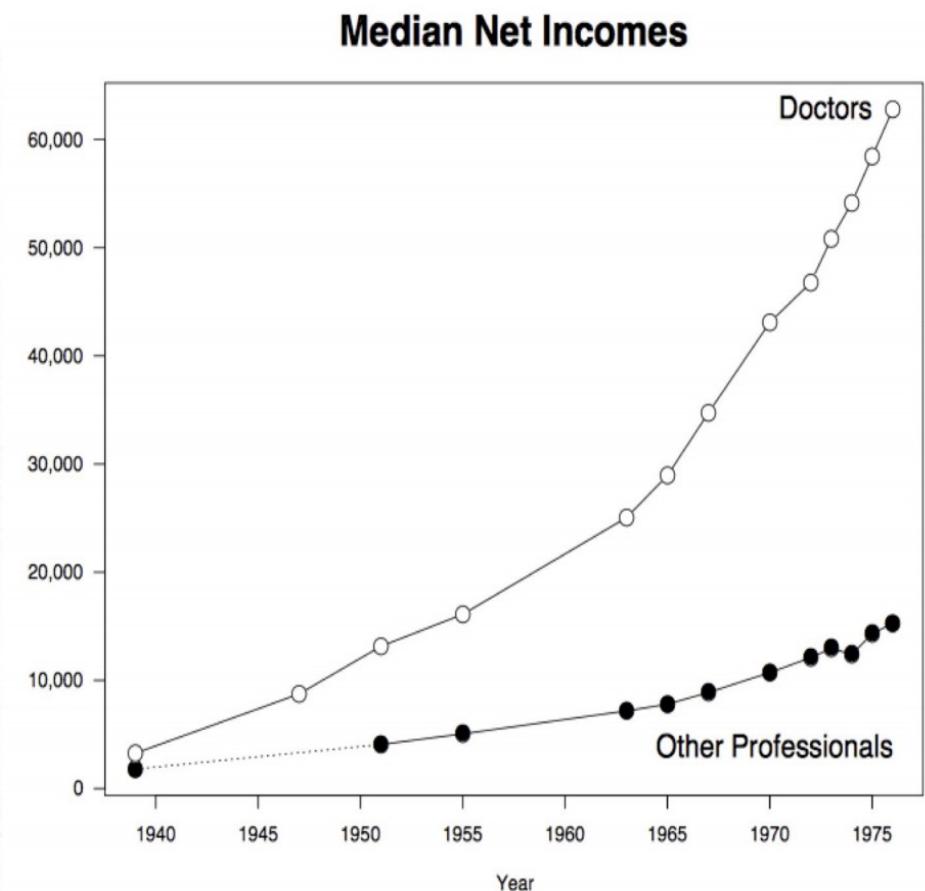
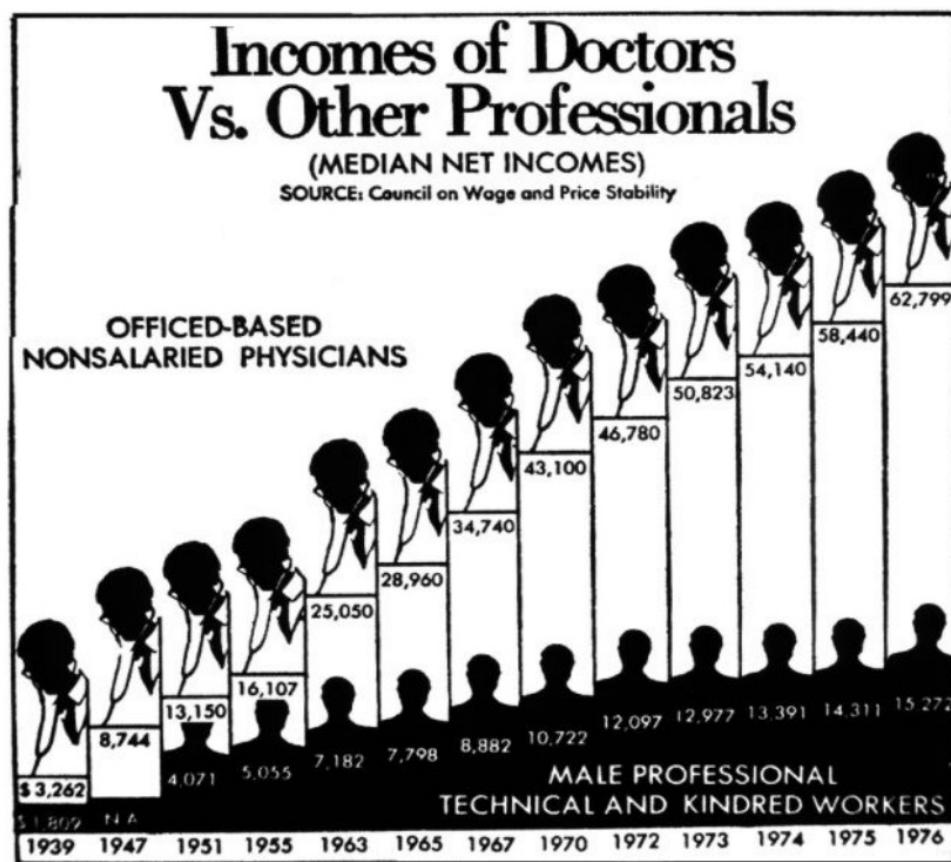


Context

- The most common mistake: **No Label**
- In general, we provide context for a plot through:
 - Plot title
 - Axes labels
 - Reference lines and markers for important values
 - Labels for interesting points and unusual observations
 - Captions that describe the data and its important features



Don't be over creative...



Storytelling with data

- “the greatest graph of all time”
- Visualized Napoleon's 1812 invasion of Russia, including
 - the number of soldiers
 - the direction of the march
 - the latitude and longitude of each city
 - the temperature on the return journey
 - dates in November and December



Charles Joseph Minard
1781 - 1870
French civil engineer

Carte Figurative des pertes successives en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite

Paris, le 20 Novembre 1869.

Les nombres d'hommes perdus sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en lettres des zones. Le rouge désigne les hommes qui ont péri en Russie; le noir ceux qui en sortent. Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. Chier, de Segur, de Fézencac, de Chambray et le journal intérieur de Jacoby, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout, qui avaient été détachés sur Minsk et Mogilow et qui rejoignirent Orléans en Wilno, avaient toujours marché avec l'armée.

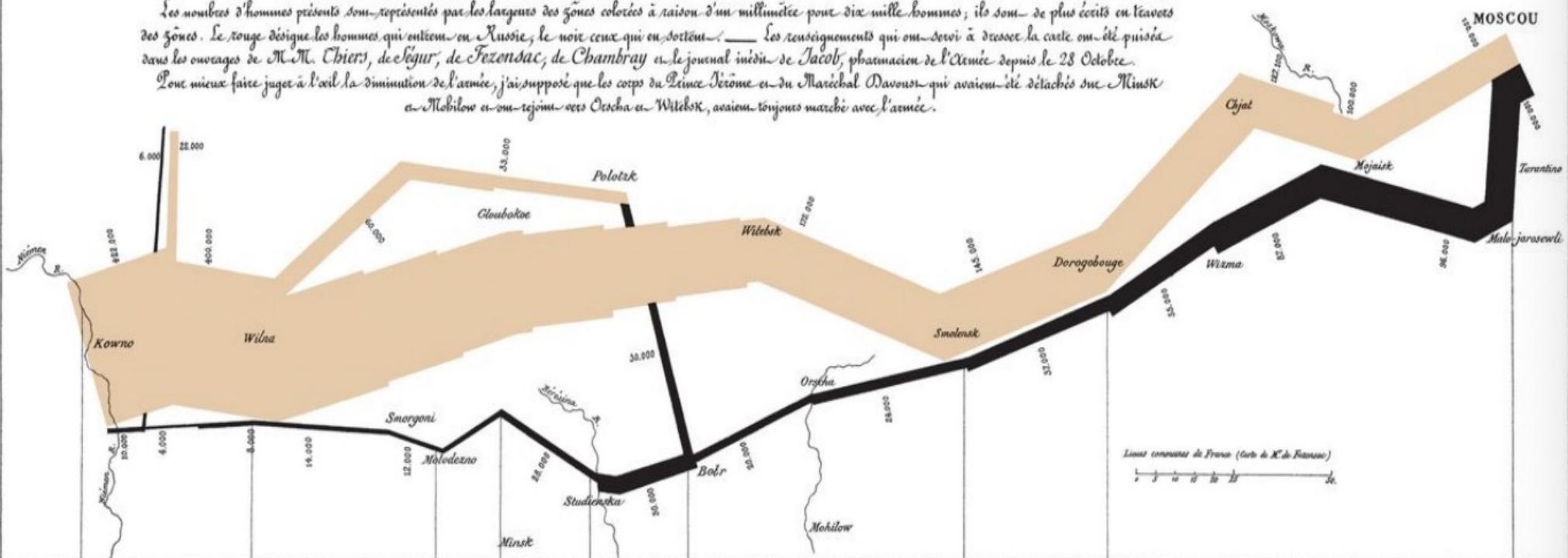
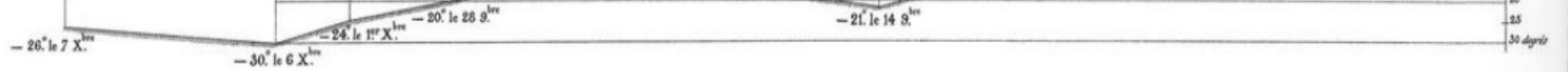


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Cosaques passent au galop
le Niemen gelé.



3 ways to plot in Python

- Matplotlib
 - the underlying plotting library powering all 3 of these
- Pandas.plot()
 - knows how to make some default plots for you
- Seaborn
 - allows to create sophisticated visualizations quickly
 - good default setting
- Documentation is your friend!
 - [Matplotlib Gallery](#)
 - [Seaborn Example Gallery](#)
- DSCI 304: data visualization in R

Demo

