

UC San Diego

Discussion Section A01 & A02 **Nov 2, 2022**

Ashwin Mishra, TA

Jiesen Zhang, IA

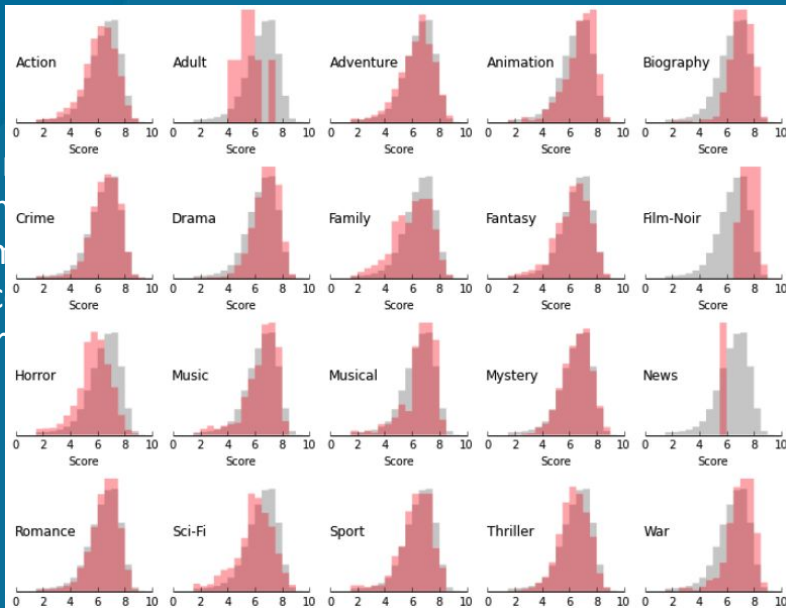
Lindsey Gu, IA

Visualization Goals

- Explanatory (Communicate)
 - Present data
 - Explain and inform
 - Provide evidence
 - Influence and persuade
- Exploratory (Analyze)
 - Explore the data
 - Assess a situation
 - Determine how to proceed
 - Decide what to do

Exploratory Data Analysis (EDA)

- Exa
- Har
- Ren
- Enc
- Nor

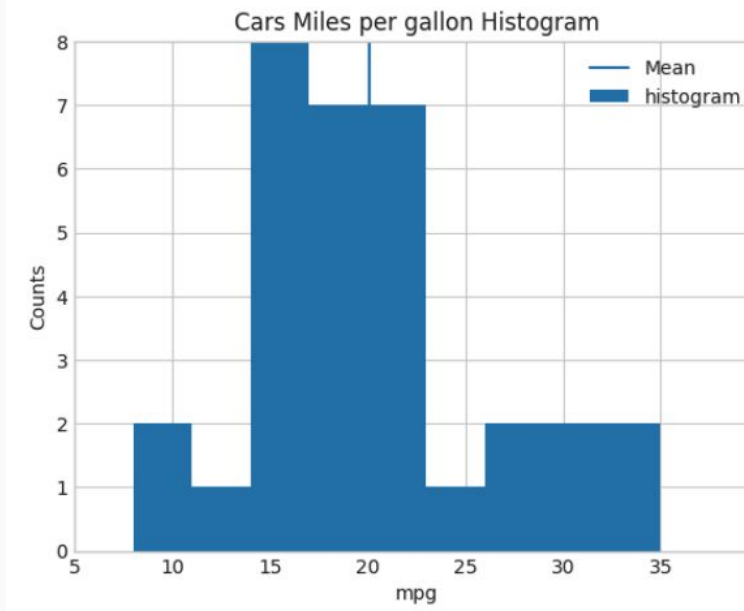


But how?

- **Build** a DataFrame (pandas) from the data
- **Clean** the DataFrame, i.e.,
 - Each row describes a single object
 - Each column describes a property of that object
- Explore **summary** of the data (histograms, scatterplots, aggregation functions)
- Explore **subsets** of data (groupBy)

But how?

	name	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb	maker
0	Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4	Mazda
1	Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4	Mazda
2	Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1	Datsun
3	Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1	Hornet
4	Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2	Hornet



Not “Effective”

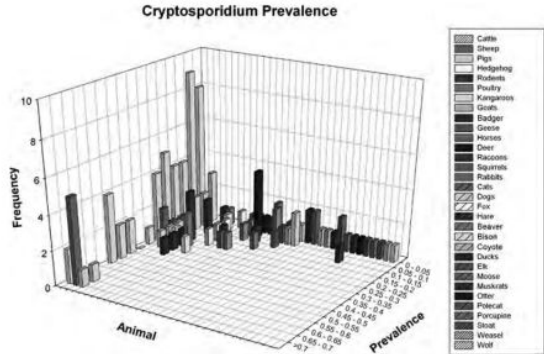
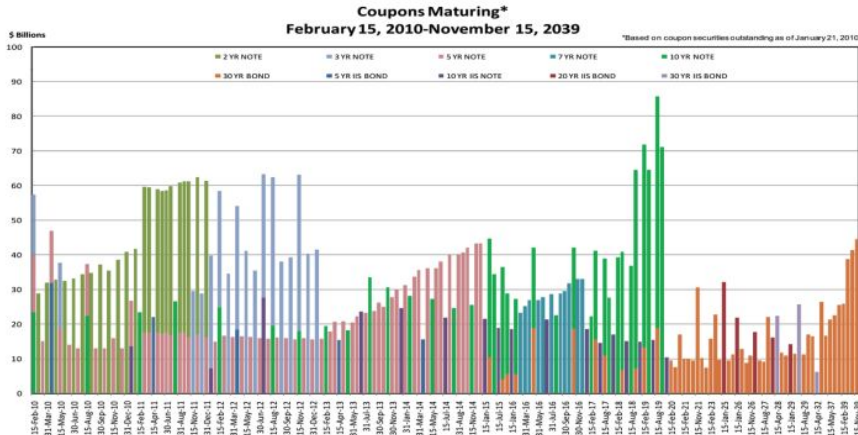
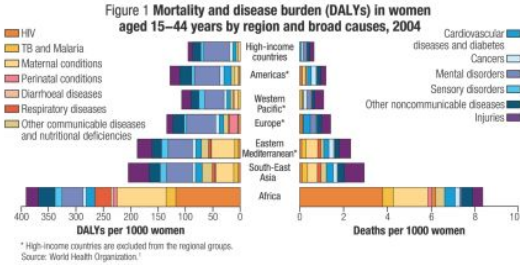
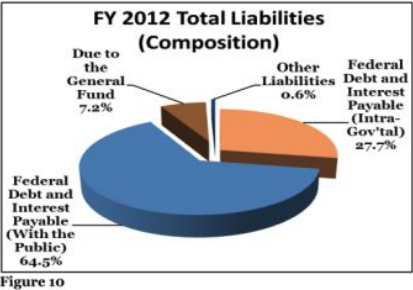
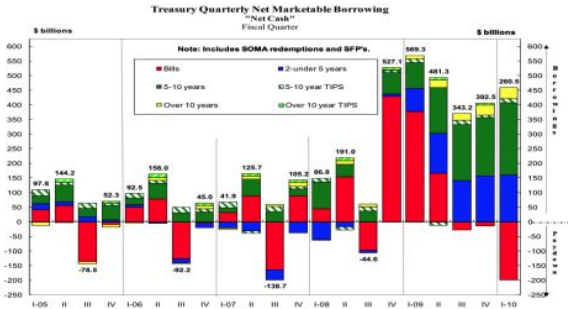
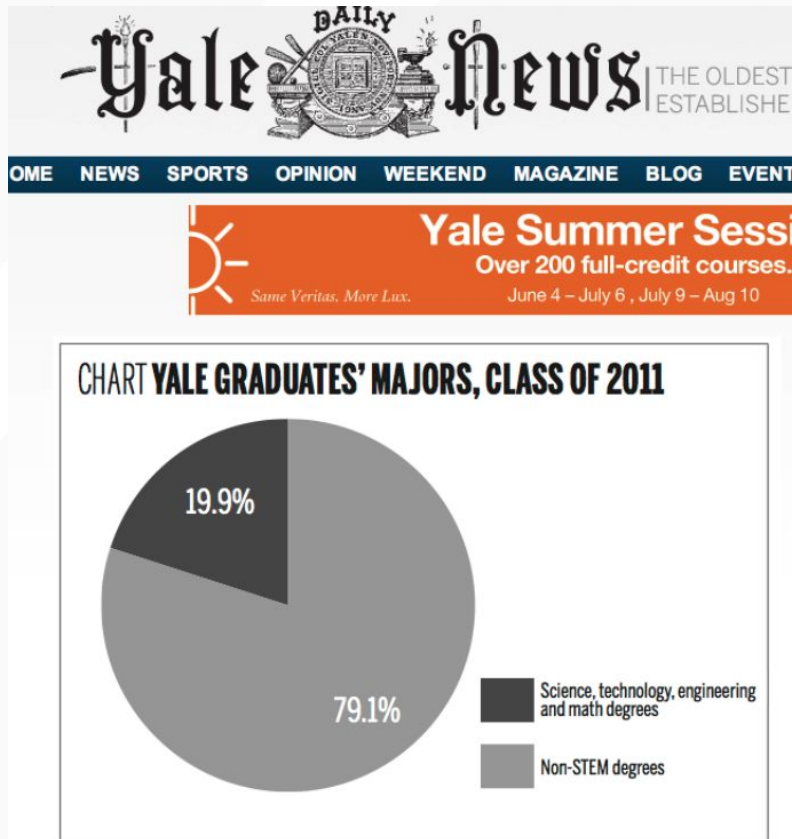


Figure 5.2 Mean prevalence rates of *Cryptosporidium* oocysts by animal species.

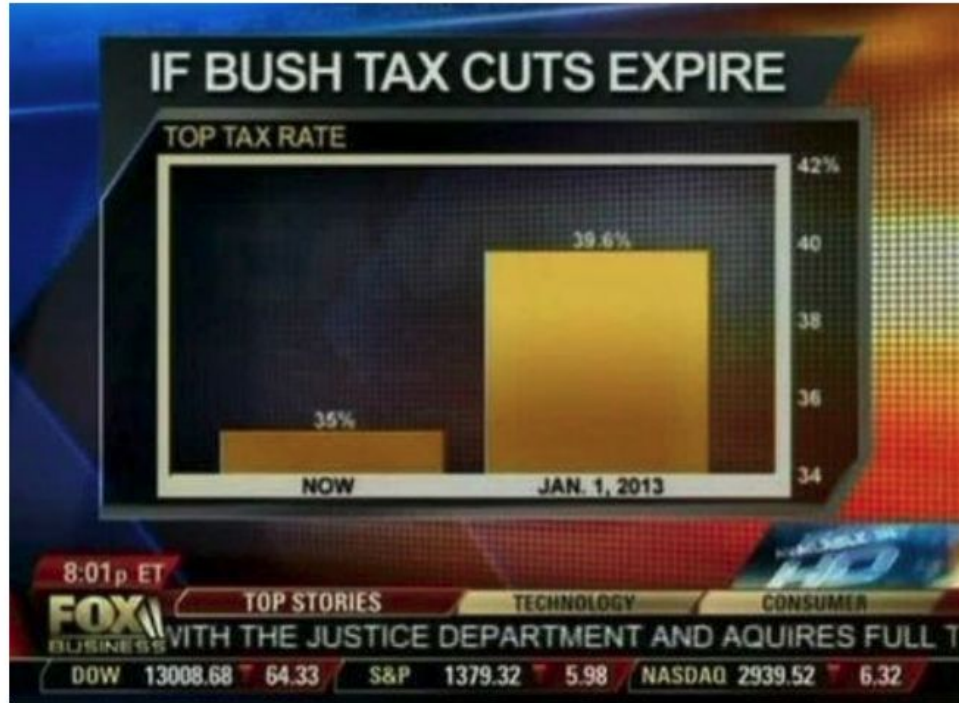
Effective EDA Viz

- Have graphical integrity
- Keep it simple
- Use the right chart
- Use the right colors

Graphical Integrity

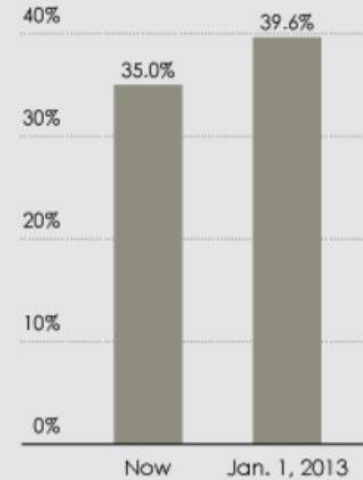


Graphical Integrity



If Bush tax cuts expire...

Top tax rate

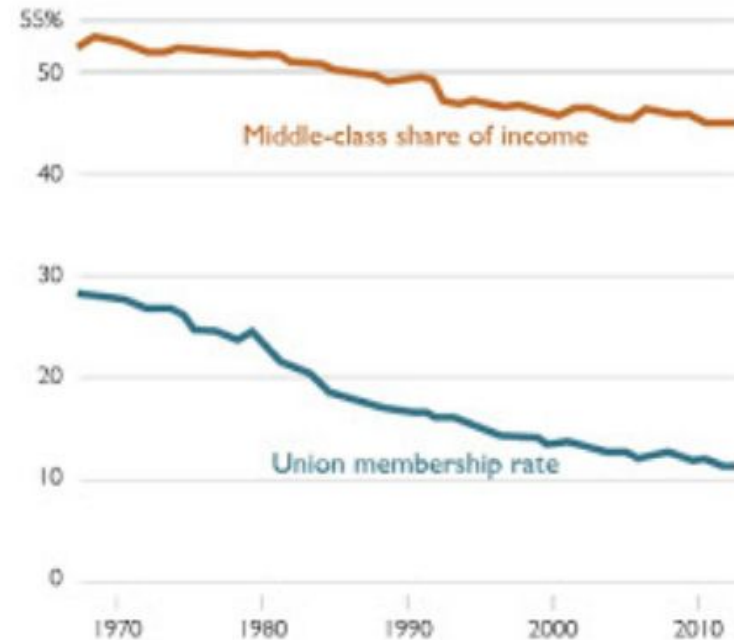


Graphical Integrity

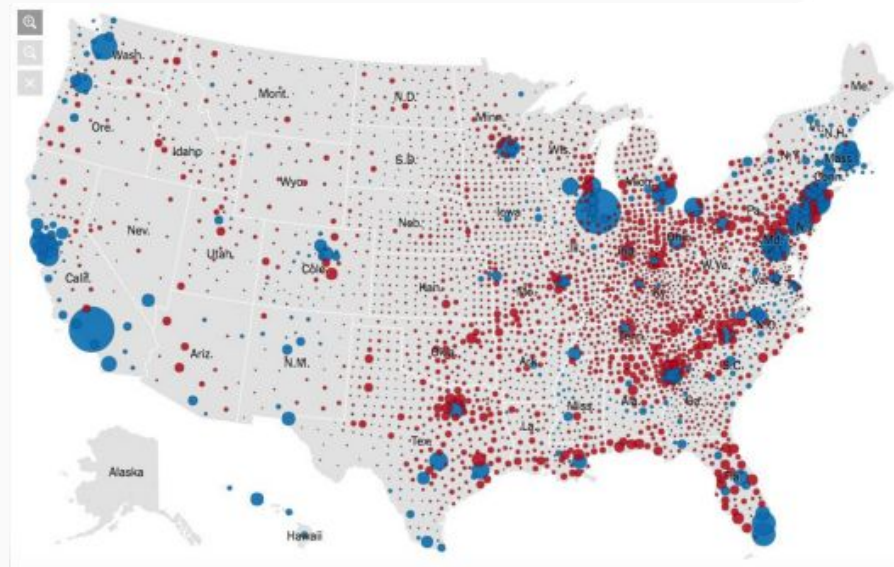
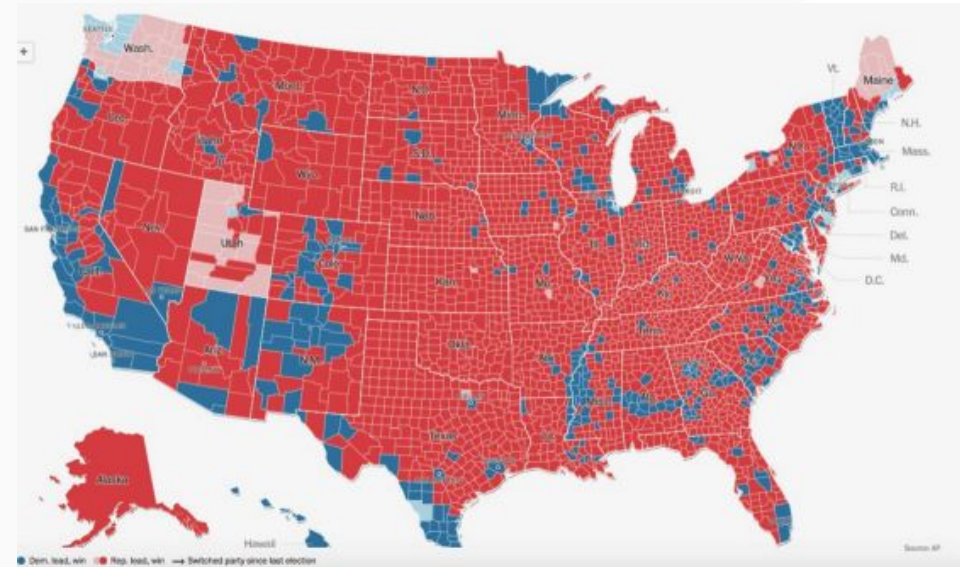
FIGURE 7. AS UNION MEMBERSHIP DECLINES, THE SHARE OF INCOME GOING TO THE MIDDLE CLASS SHRINKS



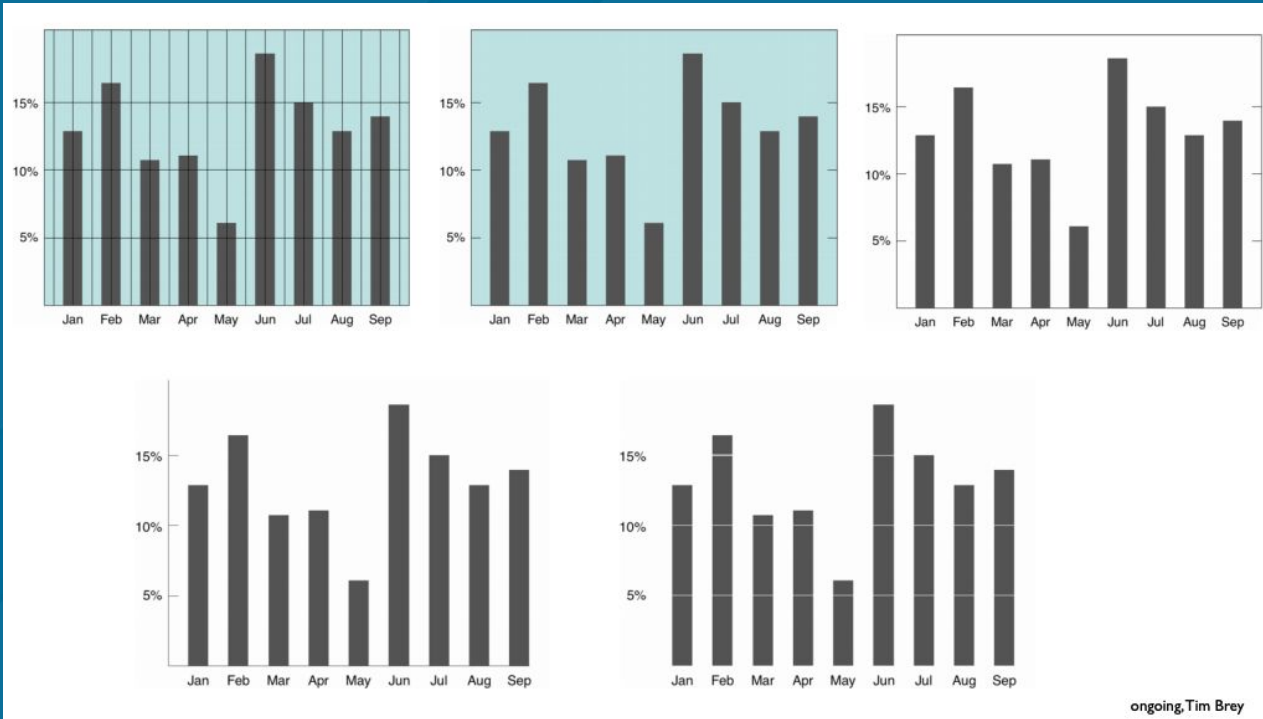
NEW VERSION



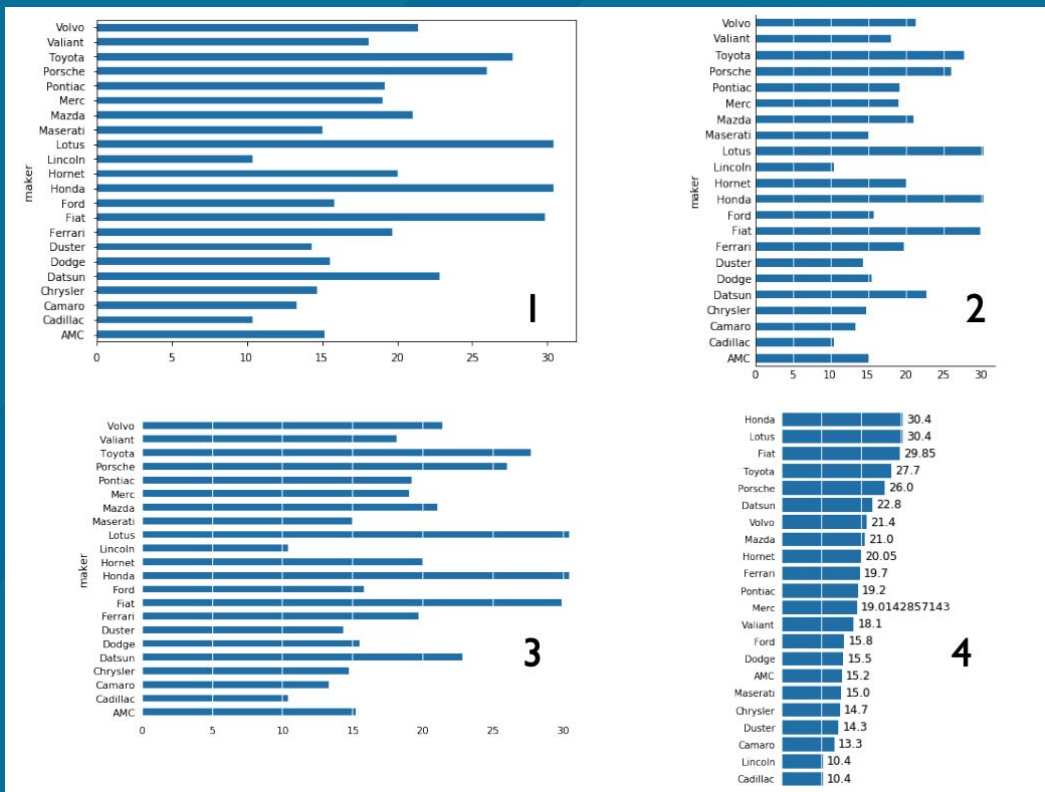
Graphical Integrity



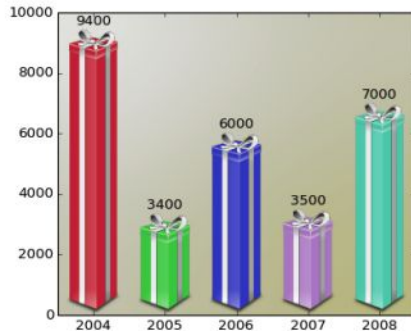
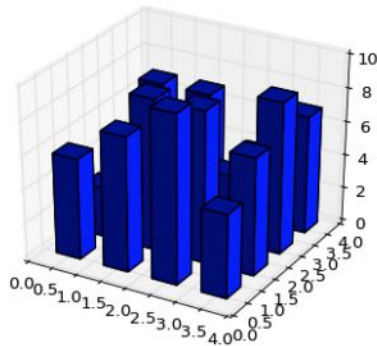
Keep it simple



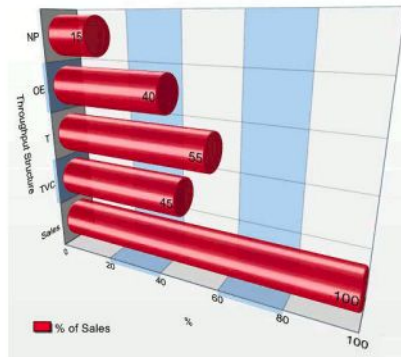
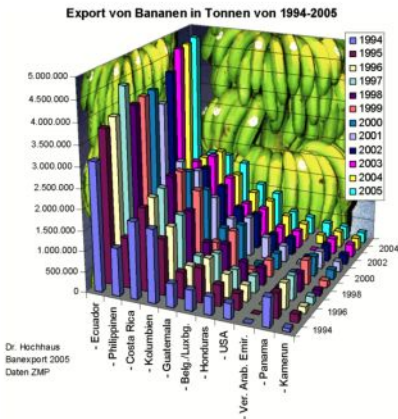
Keep it simple



Nope...



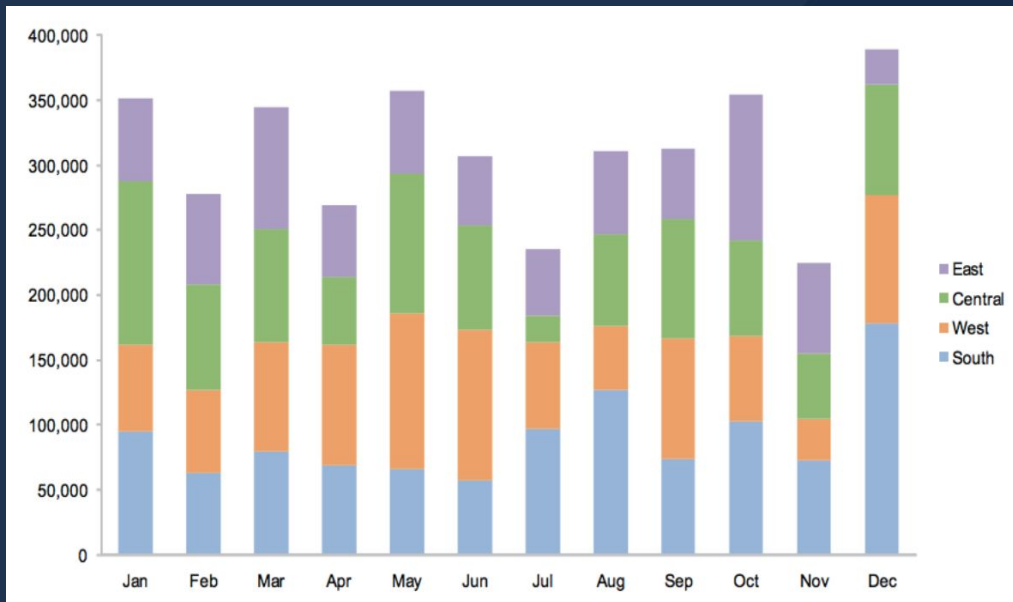
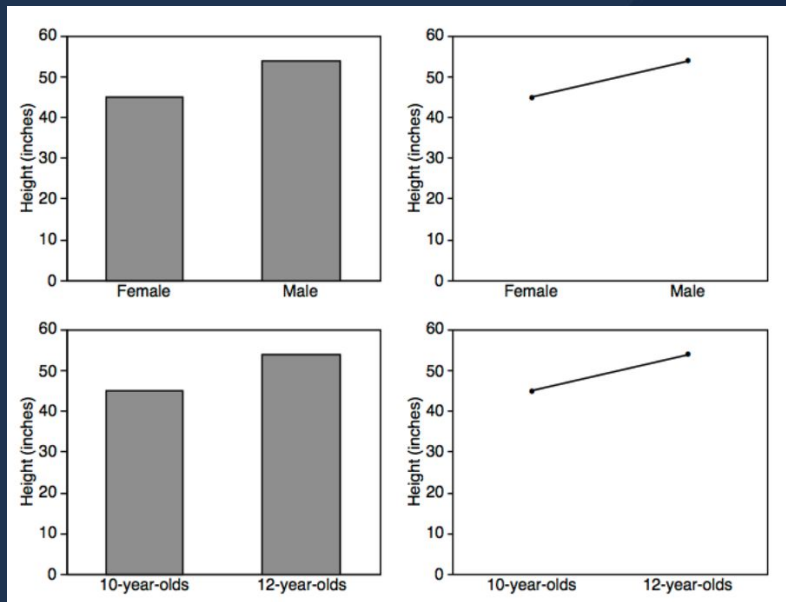
matplotlib gallery



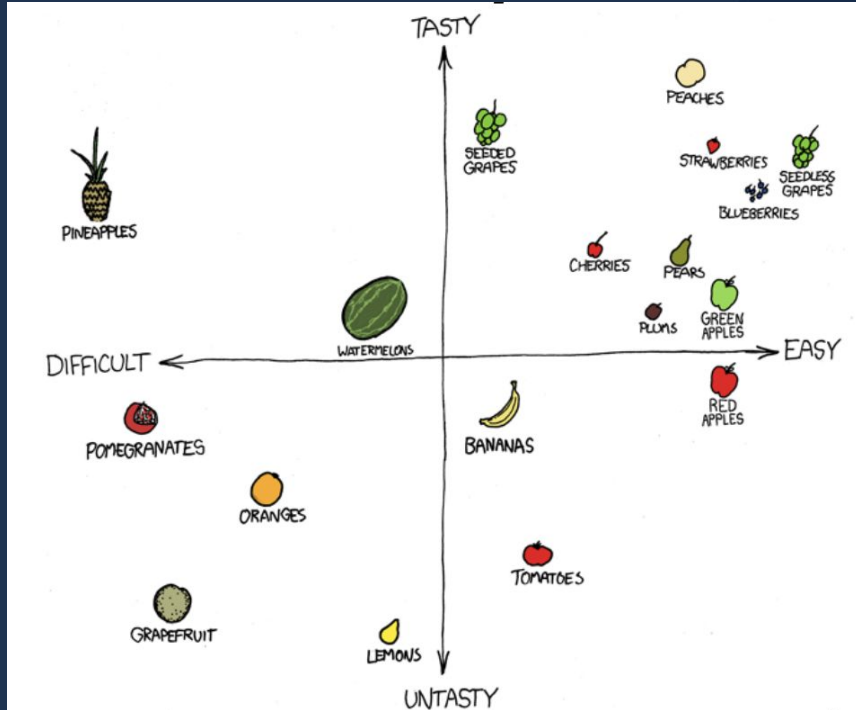
Excel Charts Blog

Use the right chart

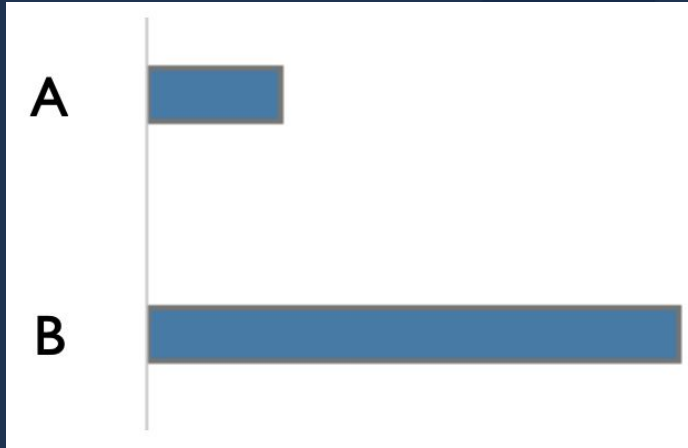
https://extremepresentation.typepad.com/blog/files/choosing_a_good_chart.pdf



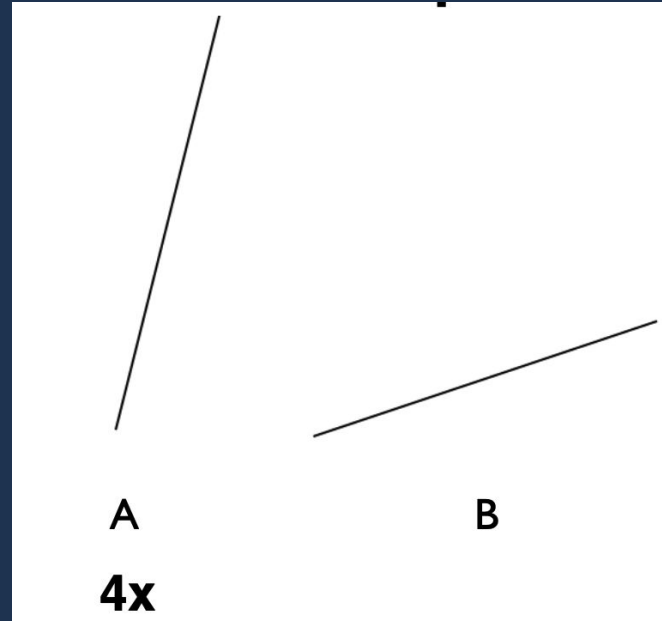
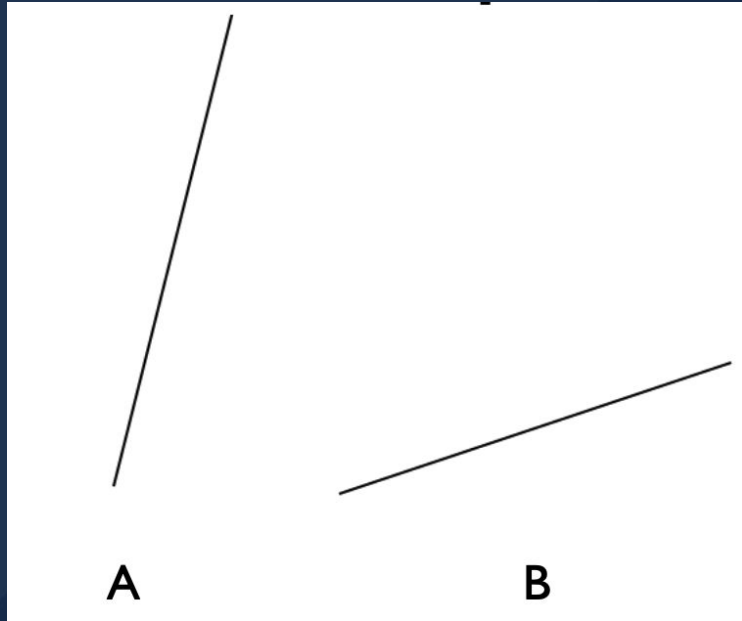
Use the right chart



How much longer?



How much steeper?



Most
Efficient



Least
Efficient

Position



Length



Slope



Angle



Area



Intensity



Color



Shape

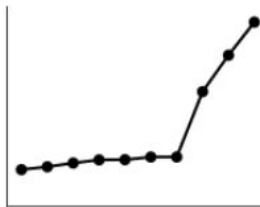


Quantitative

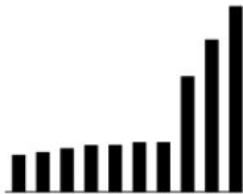
Ordered

Categories

Most Effective



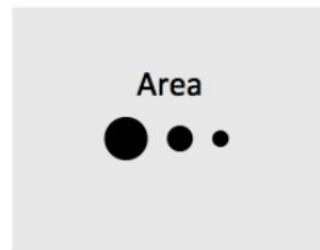
Position



Length



Less Effective

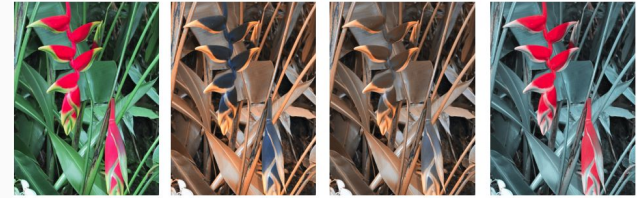


Use the right colors

Do not use more than 5-8 colors at once



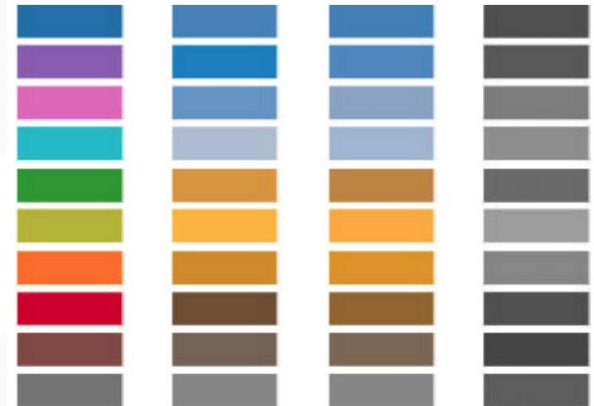
Color Blindness



Protanope Deuteranope Tritanope

Red / green
deficiencies

Blue / Yellow
deficiency



Normal

Protanope

Deuteranope

Lightness

Communicating

- Who:
 - Who is your audience?
 - What is your relationship to your audience?
- What:
 - What do you need your audience to know?
 - What do you want your audience to do?
 - What will your tone be?
- How:
 - How will you communicate this to your audience? Live or written?

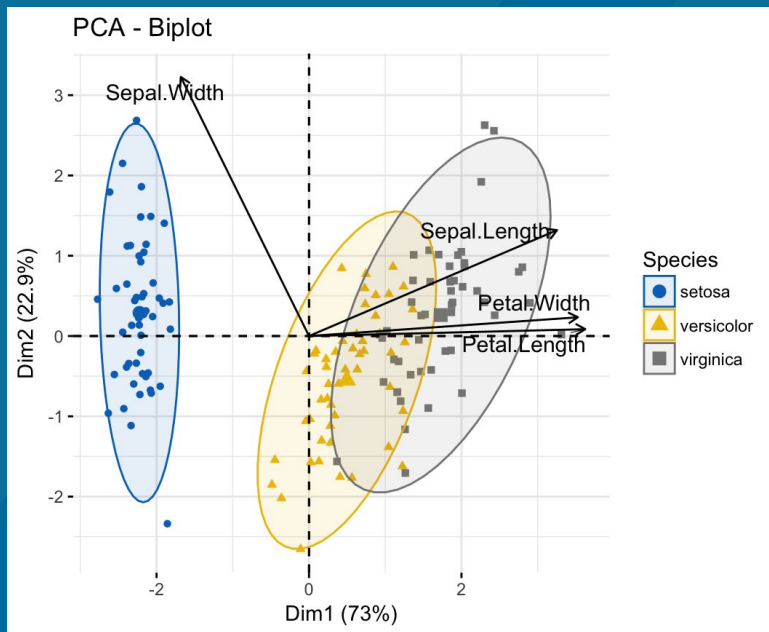
When NOT to do EDA?

- Identify/remove samples after analyzing data
- After running statistical test and obtaining p-value
- After getting an answer you don't like
- To improve the correlation between variables

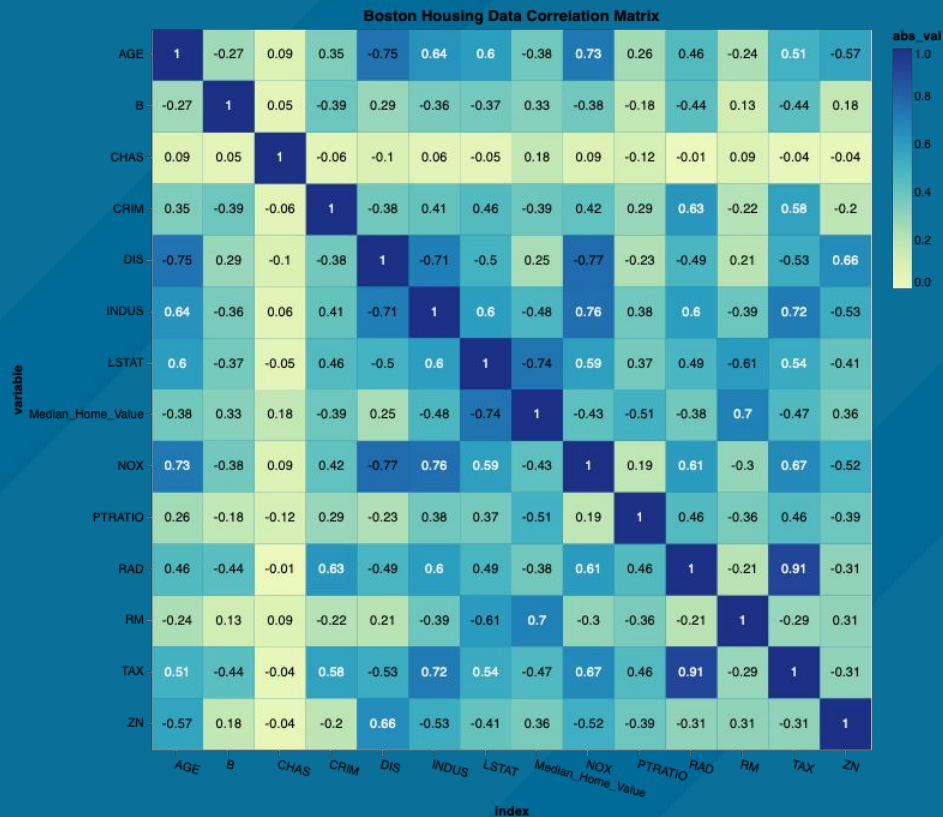
EDA is NOT a tool to get your data analysis to give you the results you want.

PCA Analysis

- Summarize and visualize the most important part of the data
- Pick the variables/features most closely related to the data



Correlation Matrix



Source: <https://harvard-iacs.github.io/2018-CS109A/lectures/lecture-3/presentation/lecture3.pdf>

UC San Diego