



INTRODUCTION TO DATA SCIENCE



Lecture 2

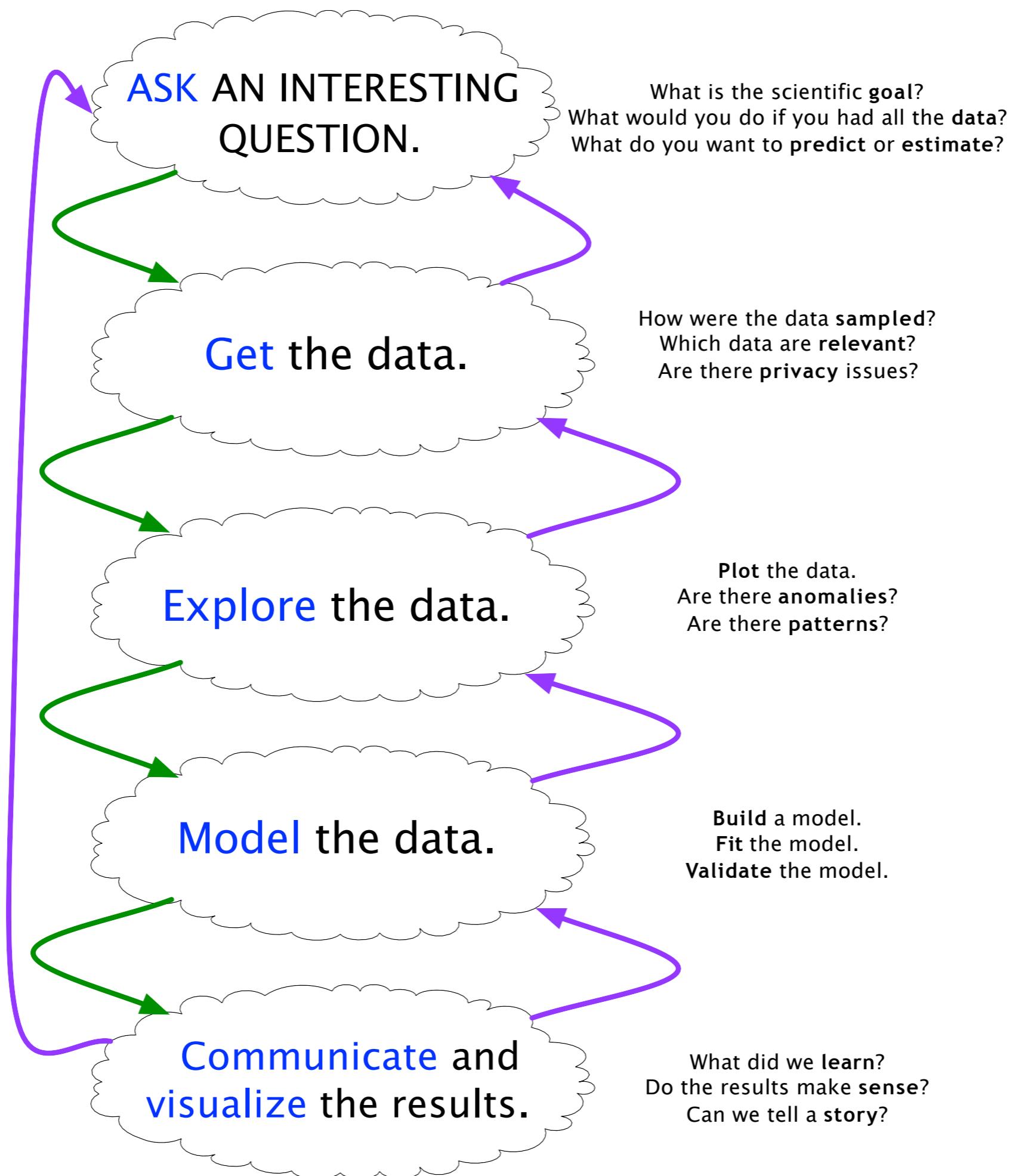
Effective Exploratory Data Analysis and Visualization

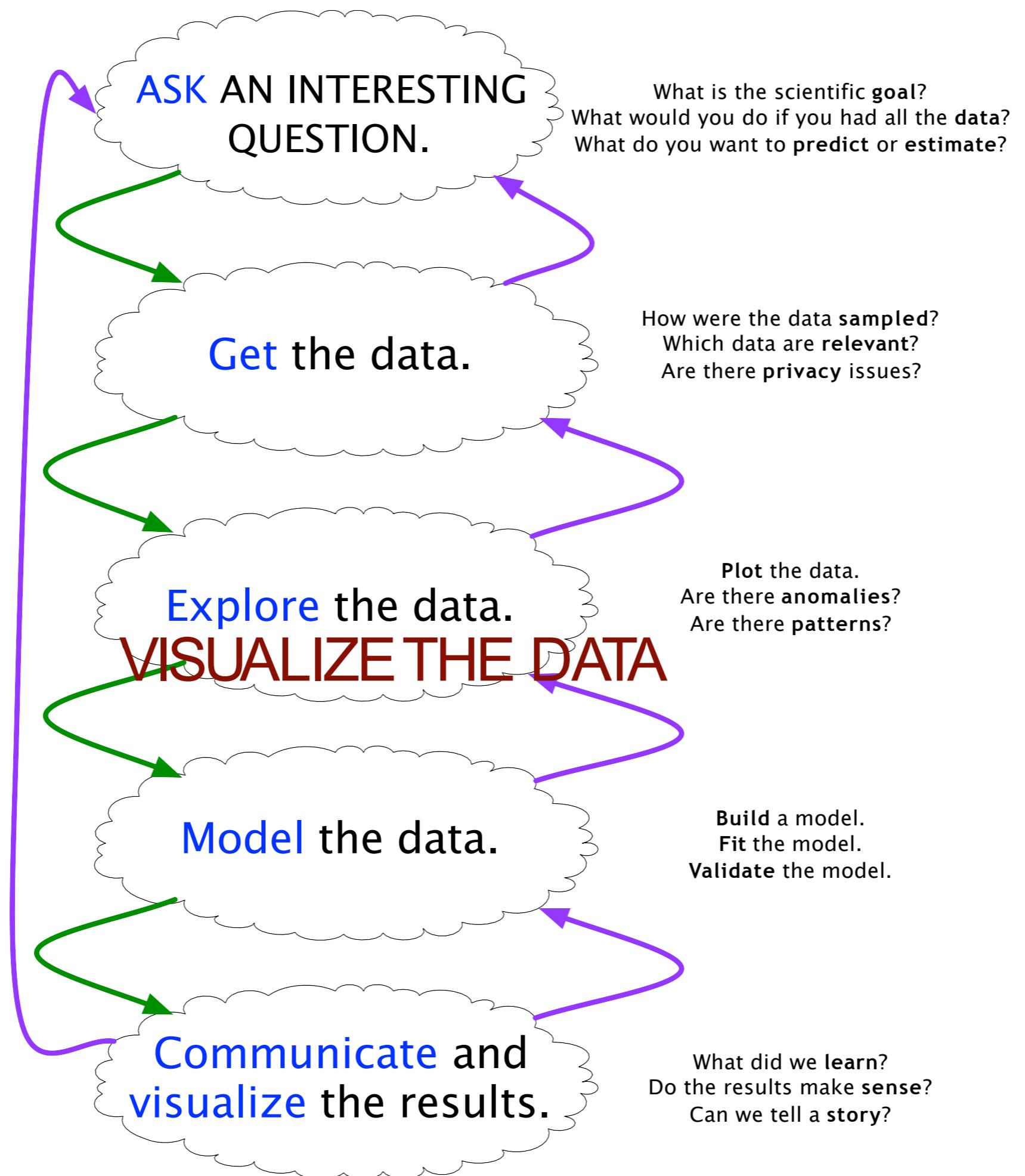
Data Science and Engineering Department
Faculty of Informatics
ELTE University

Zakarya Farou

zakaryafarou@inf.elte.hu

THE DATA SCIENCE PROCESS





Anscombe's quartet

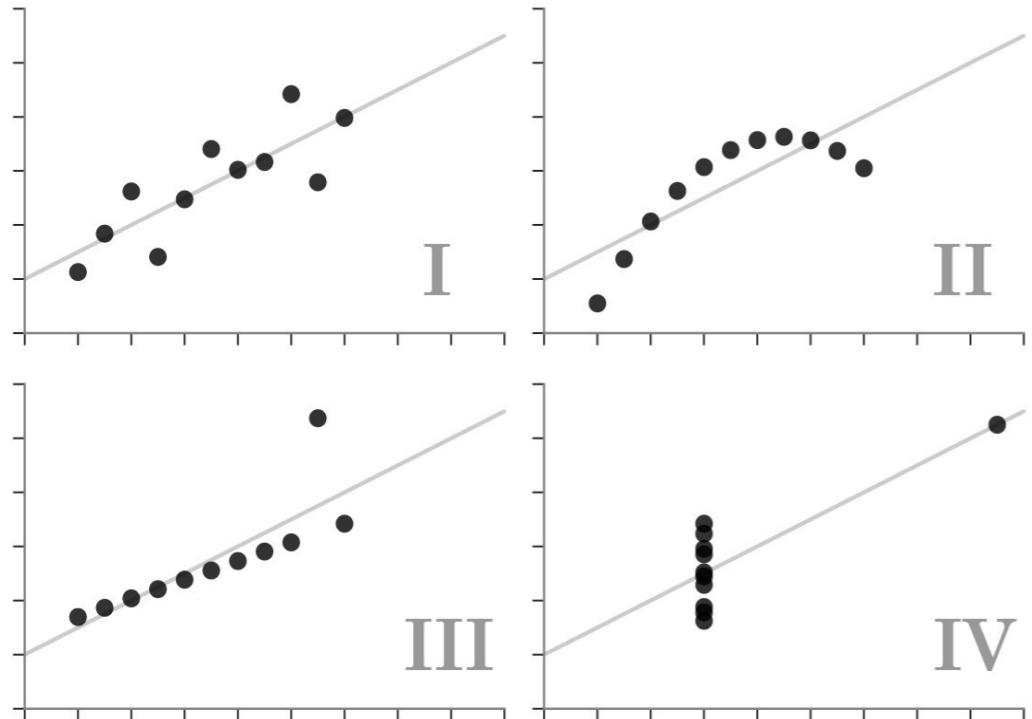
Anscombe's Data									
Observation	x1	y1	x2	y2	x3	y3	x4	y4	
1	10	8.04	10	9.14	10	7.46	8	6.58	
2	8	6.95	8	8.14	8	6.77	8	5.76	
3	13	7.58	13	8.74	13	12.74	8	7.71	
4	9	8.81	9	8.77	9	7.11	8	8.84	
5	11	8.33	11	9.26	11	7.81	8	8.47	
6	14	9.96	14	8.1	14	8.84	8	7.04	
7	6	7.24	6	6.13	6	6.08	8	5.25	
8	4	4.26	4	3.1	4	5.39	19	12.5	
9	12	10.84	12	9.13	12	8.15	8	5.56	
10	7	4.82	7	7.26	7	6.42	8	7.91	
11	5	5.68	5	4.74	5	5.73	8	6.89	
Summary Statistics									
N	11	11	11	11	11	11	11	11	11
mean	9.00	7.50	9.00	7.500909	9.00	7.50	9.00	7.50	
SD	3.16	1.94	3.16	1.94	3.16	1.94	3.16	1.94	
r	0.82		0.82		0.82		0.82		

The statistical information for these four data sets are approximately similar!

Same Stats, Different Graphs

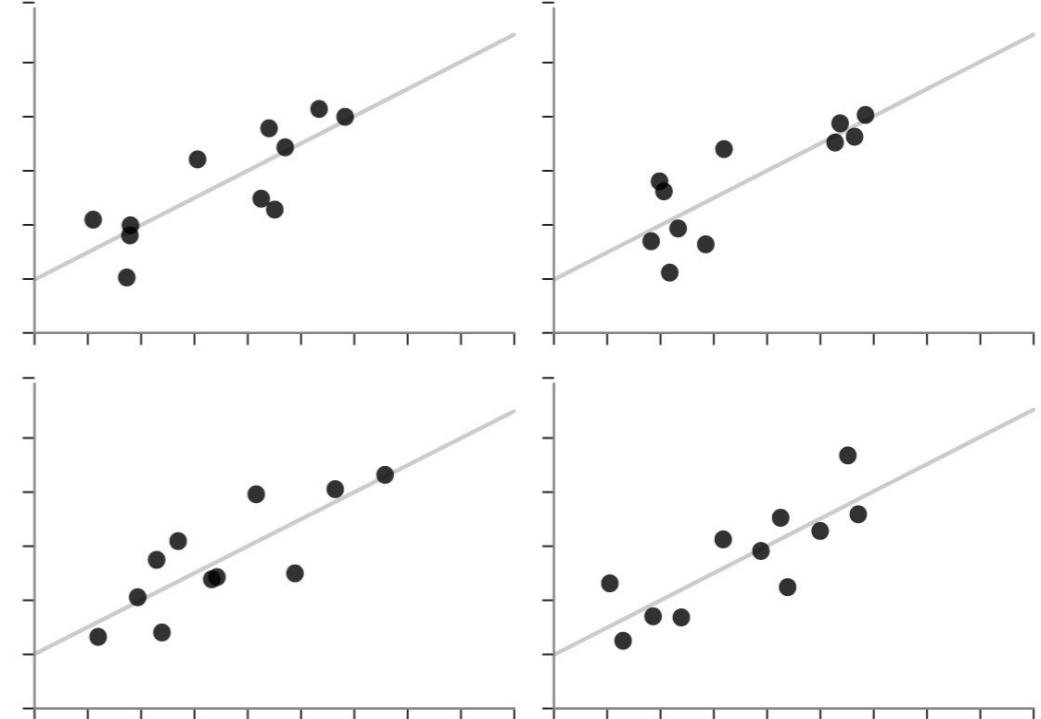
✓ Anscombe's Quartet

Each dataset has the same summary statistics (mean, standard deviation, correlation), and the datasets are *clearly different*, and *visually distinct*.



✗ Unstructured Quartet

Each dataset here also has the same summary statistics. However, they are not *clearly different* or *visually distinct*.



ANScombe'S QUARTET FOUR DATASETS

- **DS1:** fits the linear regression model pretty well.
- **DS2:** cannot fit the linear regression model because the data is non-linear.
- **DS3 – DS4:** shows the outliers involved in the data set, which the linear regression model cannot handle.
- Anscombe's quartet helps us to understand the importance of data visualization and how easy it is to fool a regression algorithm.
- Before attempting to interpret and model the data or implement any machine learning algorithm, we first need to visualize the data set to help build a well-fit model.

**EXAMPLE: ANTIBIOTICS
WILL BURTON, 1951**

DATA

Table 1: Burtin's data.

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus faecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

DATA

Genus, Species

Table 1: Burtin's data.

Bacteria		Antibiotic			
		Penicillin	Streptomycin	Neomycin	Gram Staining
<i>Aerobacter aerogenes</i>		870	1	1.6	negative
<i>Brucella abortus</i>		1	2	0.02	negative
<i>Brucella anthracis</i>		0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>		0.005	11	10	positive
<i>Escherichia coli</i>		100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>		850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>		800	5	2	negative
<i>Proteus vulgaris</i>		3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>		850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>		1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>		10	0.8	0.09	negative
<i>Staphylococcus albus</i>		0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>		0.03	0.03	0.001	positive
<i>Streptococcus faecalis</i>		1	1	0.1	positive
<i>Streptococcus hemolyticus</i>		0.001	14	10	positive
<i>Streptococcus viridans</i>		0.005	10	40	positive

DATA

Genus, Species

Table 1: Burtin's data.

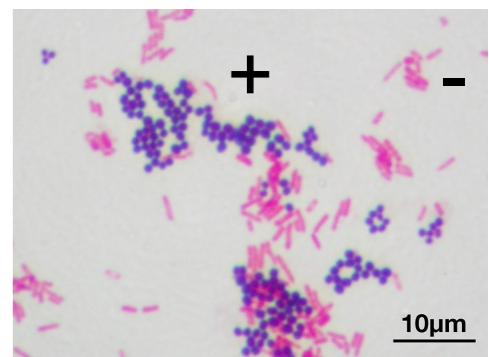
Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus fecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

DATA

Genus, Species

Table 1: Burtin's data.

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus fecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

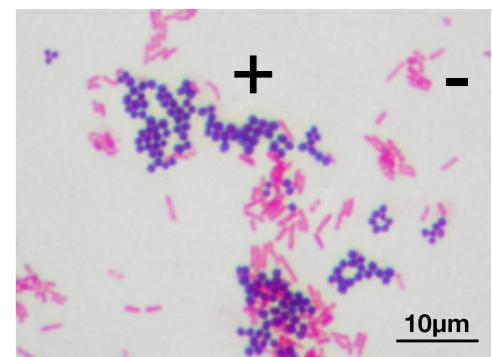


DATA

Genus, Species

Table 1: Burtin's data.

Bacteria	Min. Inhibitory Concentration [ml/g]	Antibiotic			Gram Staining
		Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6		negative
<i>Brucella abortus</i>	1	2	0.02		negative
<i>Brucella anthracis</i>	0.001	0.01	0.007		positive
<i>Diplococcus pneumoniae</i>	0.005	11	10		positive
<i>Escherichia coli</i>	100	0.4	0.1		negative
<i>Klebsiella pneumoniae</i>	850	1.2	1		negative
<i>Mycobacterium tuberculosis</i>	800	5	2		negative
<i>Proteus vulgaris</i>	3	0.1	0.1		negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4		negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008		negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09		negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001		positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001		positive
<i>Streptococcus faecalis</i>	1	1	0.1		positive
<i>Streptococcus hemolyticus</i>	0.001	14	10		positive
<i>Streptococcus viridans</i>	0.005	10	40		positive



WHAT QUESTIONS ?

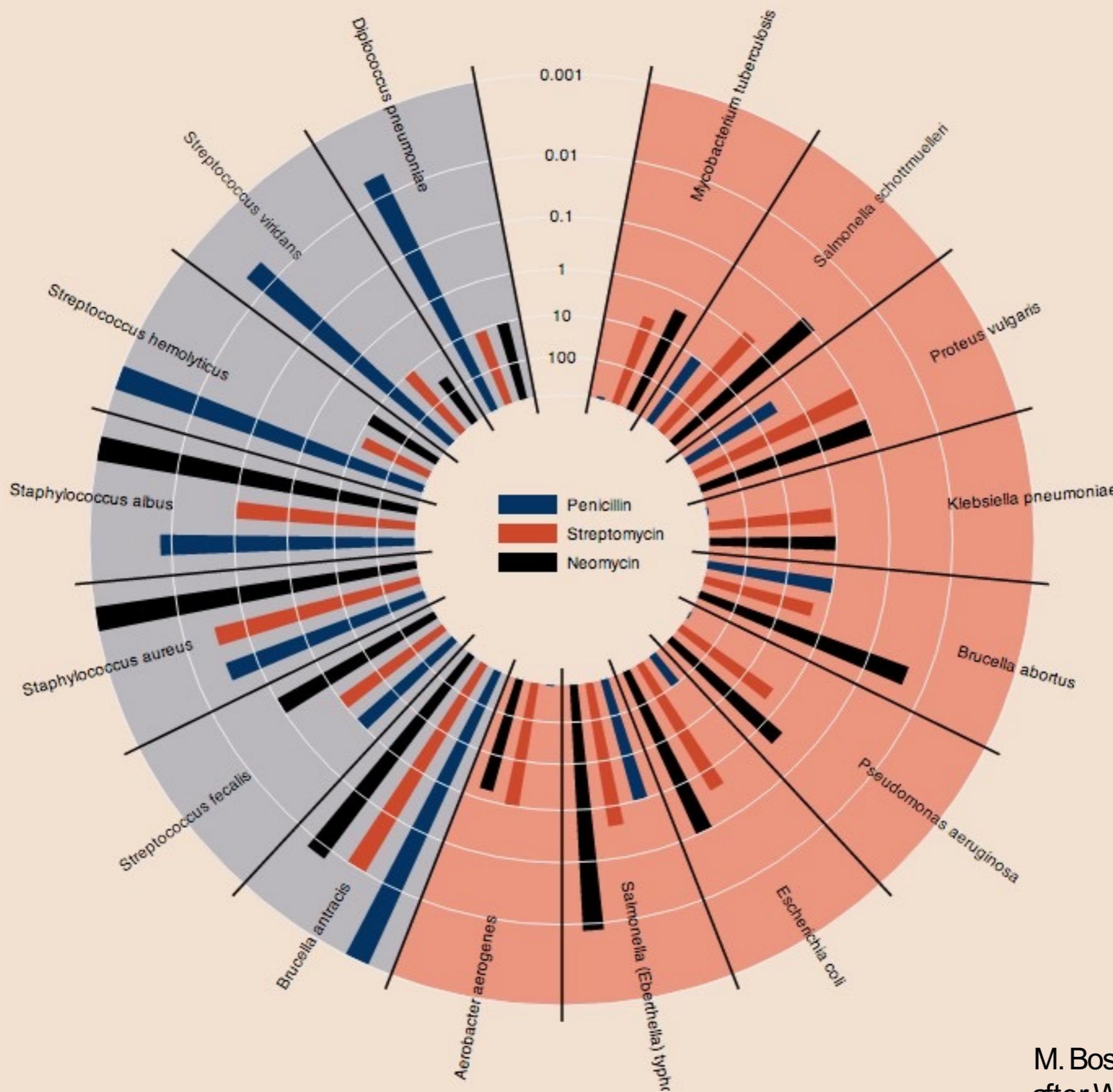
Table 1: Burtin's data.

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus faecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

HOW EFFECTIVE ARE THE DRUGS?

Gram
Positive

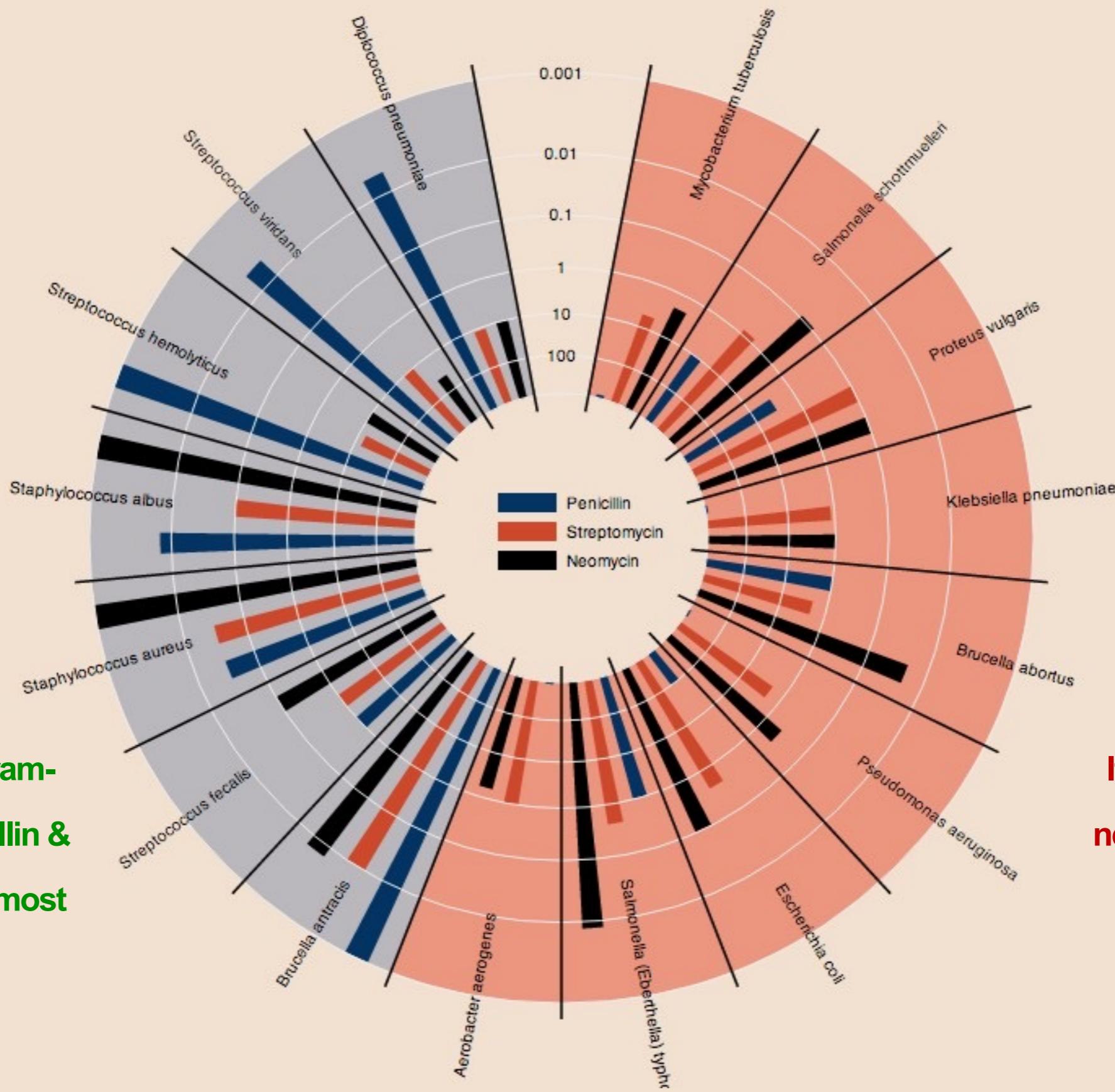
Gram
Negative



HOW EFFECTIVE ARE THE DRUGS?

Gram
Positive

If bacteria is gram-
positive, Penicillin &
Neomycin are most
effective.



Gram
Negative

If bacteria is gram-
negative, Neomycin
is most effective.

EXPLORATORY DATA ANALYSIS

“The greatest value of a picture is when it forces us to notice what we never expected to see.”



John Tukey

VISUALIZATION GOALS

Communicate (Explanatory)

- Present data and ideas
- Explain and inform
- Provide evidence and support
- Influence and persuade

Analyze (Exploratory)

- Explore the data
- Assess a situation
- Determine how to proceed
- Decide what to do

COMMUNICATE

755

Steroids or Not, the Pursuit Is On

Babe Ruth is taking aim at the career home run record. He needs only six more to tie Babe Ruth and 47 to equal Hank Aaron.

Lines are cumulative home runs



Hank Aaron
755 homers
23 seasons

Babe Ruth
714 homers
22 seasons

Barry Bonds
708 homers
20 seasons



Bonds takes lead
Home runs after 16 seasons
Bonds: 567
Aaron: 554
Ruth: 516

000

400

200

14th season

According to allegations in a book about Bonds, he began taking steroids before the 1999 season, his 14th in the league. Two seasons later, he hit 73 home runs, surpassing Aaron's career pace.

755
714
23 seasons
22 seasons
20 seasons
Bonds was injured last season. He played 14 games and hit 5 homers

Homer Pace After Age 34

If the accusations are correct, Bonds was 34 in his first season on steroids. Here are projected home run paces for each player after age 34.

PROJECTED PACE BASED ON AVERAGE OF PREVIOUS FIVE SEASONS

Aaron
Actual homers slightly outpace projected homers for five seasons.

Ruth
Averaged 46.4 homers a season from age 30 to 34. Averaged 42.5 for next four seasons.

Bonds
From age 35 to 39, he averaged 14 more homers a season than projected.

Note: Ages as of July 1 of each season.

Others Taking Aim

Alex Rodriguez

Is ahead of the pace set by all three home run leaders.
429 HR
Aaron, Ruth and Bonds
12 SEASONS

Albert Pujols

Averaging 40 homers a season, he has started stronger than the three leaders did.
261 HR
16 SEASONS

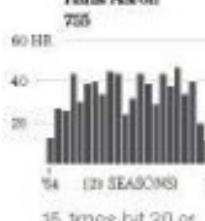
Ken Griffey Jr.

Many thought he would be the first to catch Ruth and Aaron until injuries limited his output.
536 HR
17 SEASONS

Differing Paths to the Top of the Charts

The top seven players on the career home run list, along with a look at Griffey (12th), Rodriguez (37th) and Pujols (tied 257th)

Hank Aaron
755



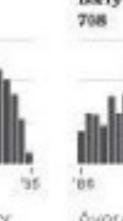
15 times hit 30 or more (M.L. most).

Babe Ruth
714



Hit only 20 over first five seasons.

Barry Bonds
708



Averaged 52 from 2000 to 2004.

Willie Mays
660



No one hit more from 1950-69.

Sammy Sosa
588



Three 60-homer seasons is record.

Frank Robinson
586



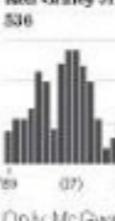
Triple Crown in '66 (49, 122, 316)

Mark McGwire
583



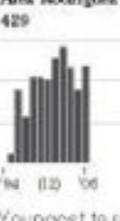
First to hit 70 in a season.

Ken Griffey Jr.
536



Only McGwire had more in the 90's.

Alex Rodriguez
429



Youngest to reach 400 homers.

Albert Pujols
301

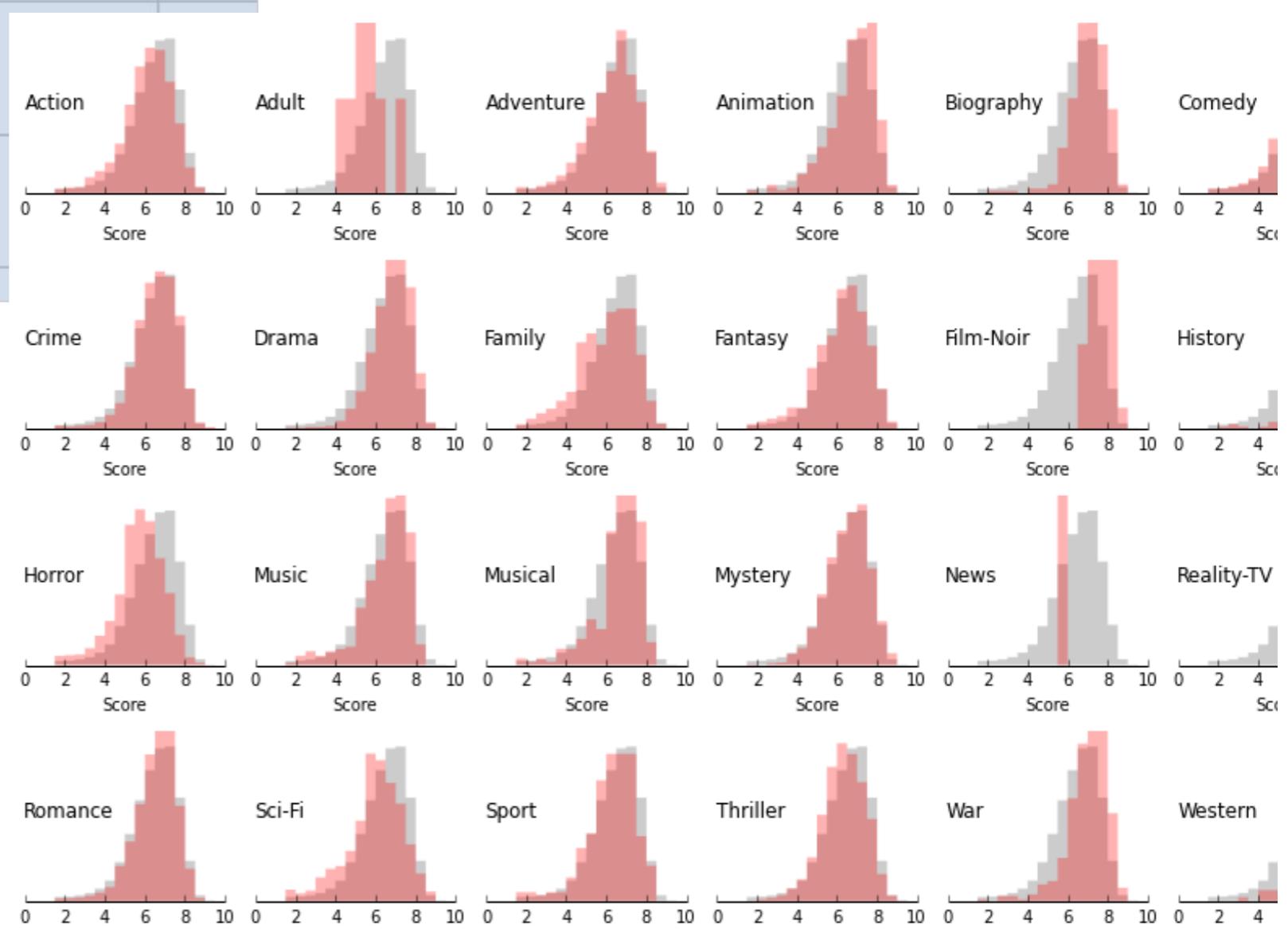
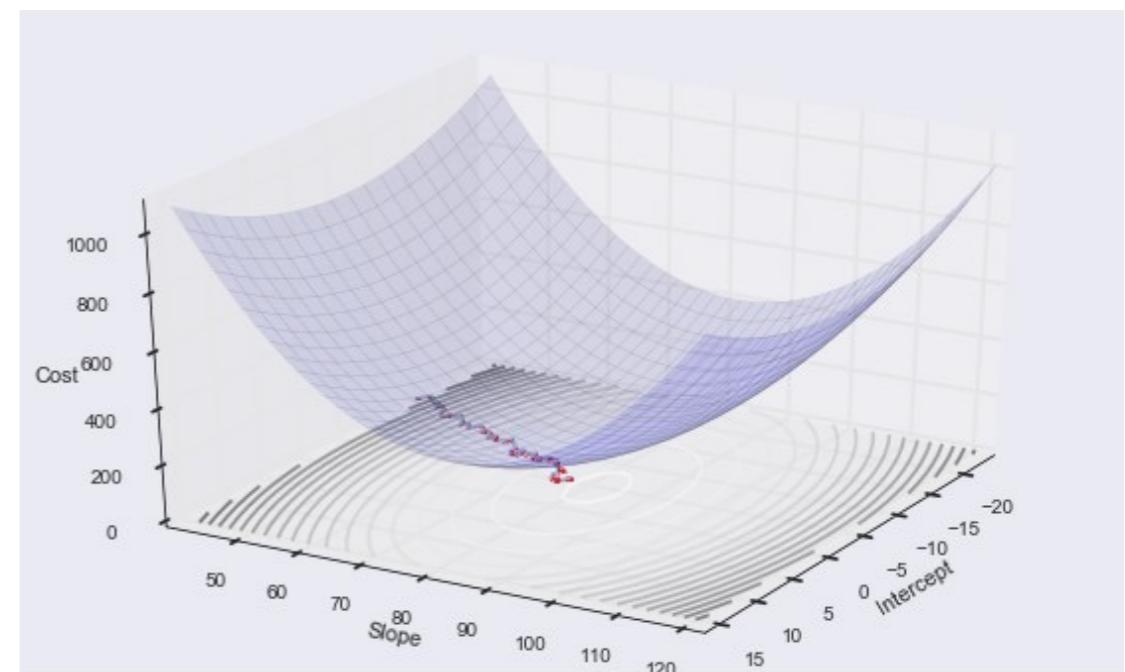
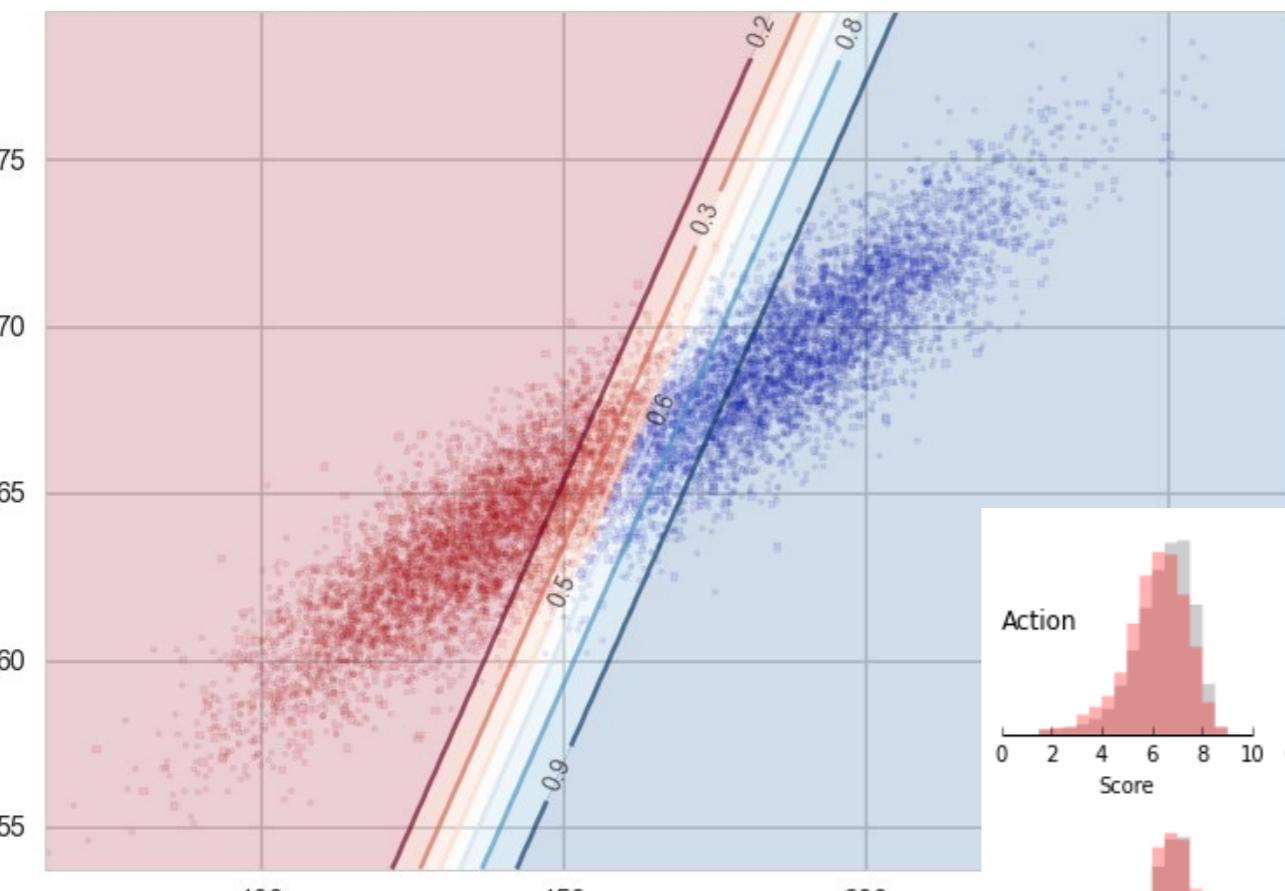


Second-most ever in franchise history.

Annika Cox and Joe Ward/The New York Times

New York Times

EXPLORE



EDA WORKFLOW

1. **Build** a DataFrame from the data (ideally, put all data in this object)
2. **Clean** the DataFrame. It should have the following properties
 - Each row describes a single object
 - Each column describes a property of that object
 - Columns are numeric whenever appropriate
 - Columns contain atomic properties that cannot be further decomposed
3. Explore **global properties**. Use histograms, scatter plots, and aggregation functions to summarize the data.
4. Explore **group properties**. Use groupby and small multiples to compare subsets of the data.

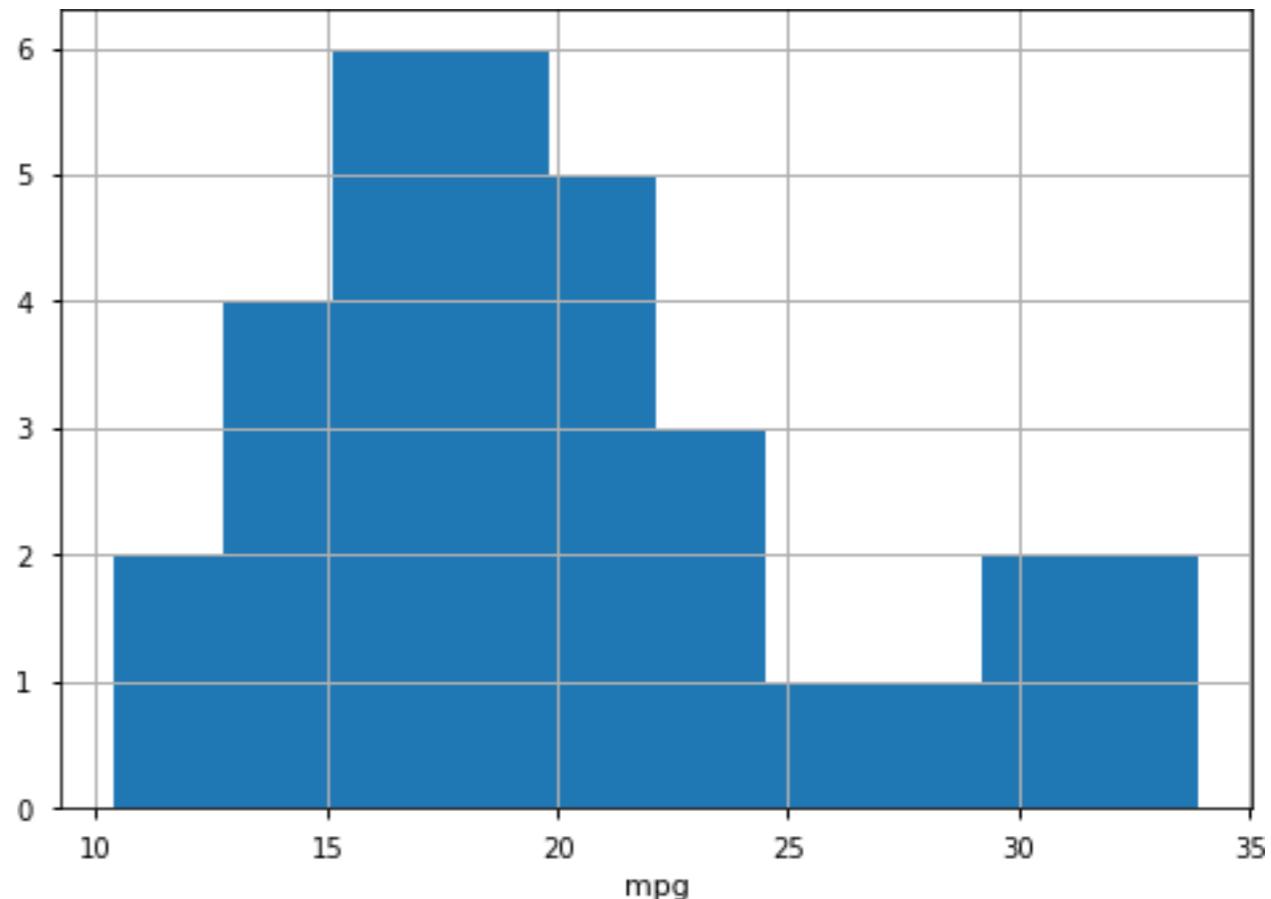
VISUALIZATION OPTIONS

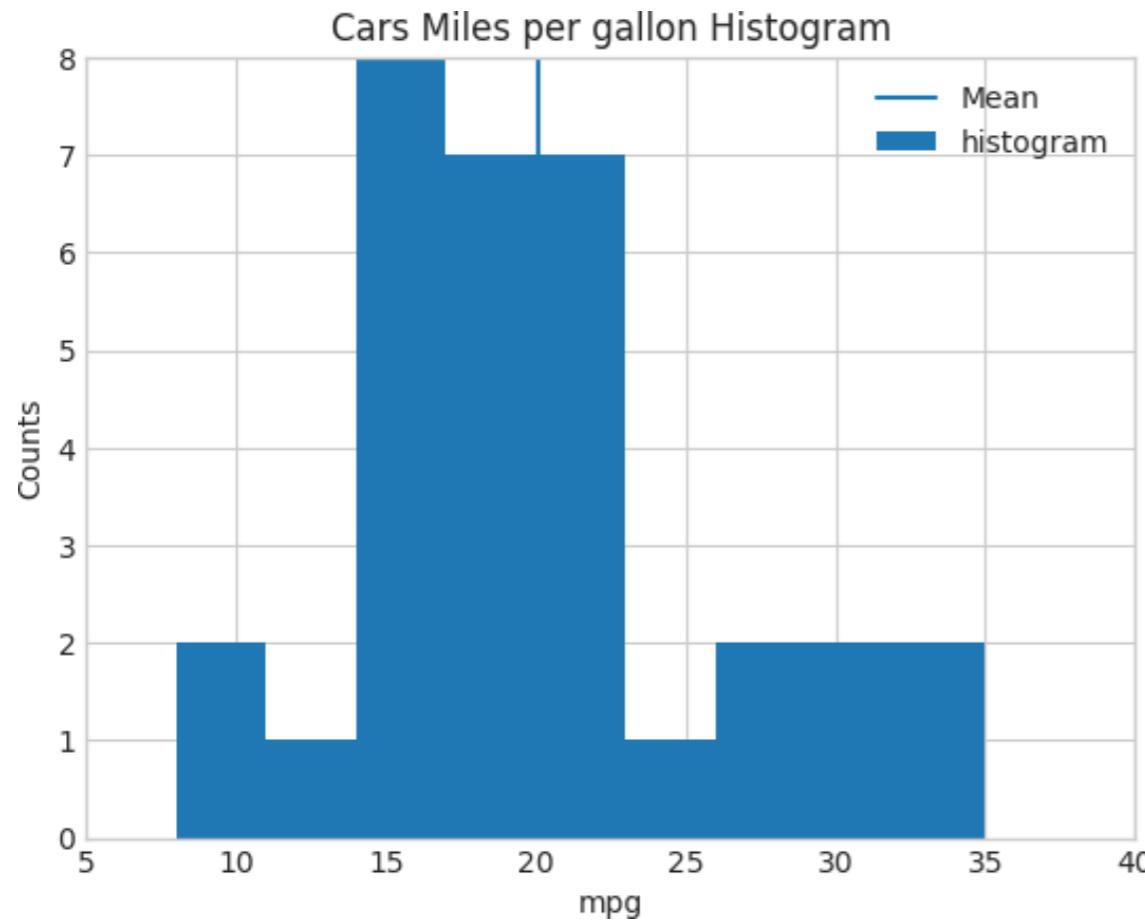
- Pandas Visualization module
- Matplotlib
- Seaborn
- The above three are inter-mixable
- Other options: Bokeh, Vega, Vincent, Altair

CARS DATASET

	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

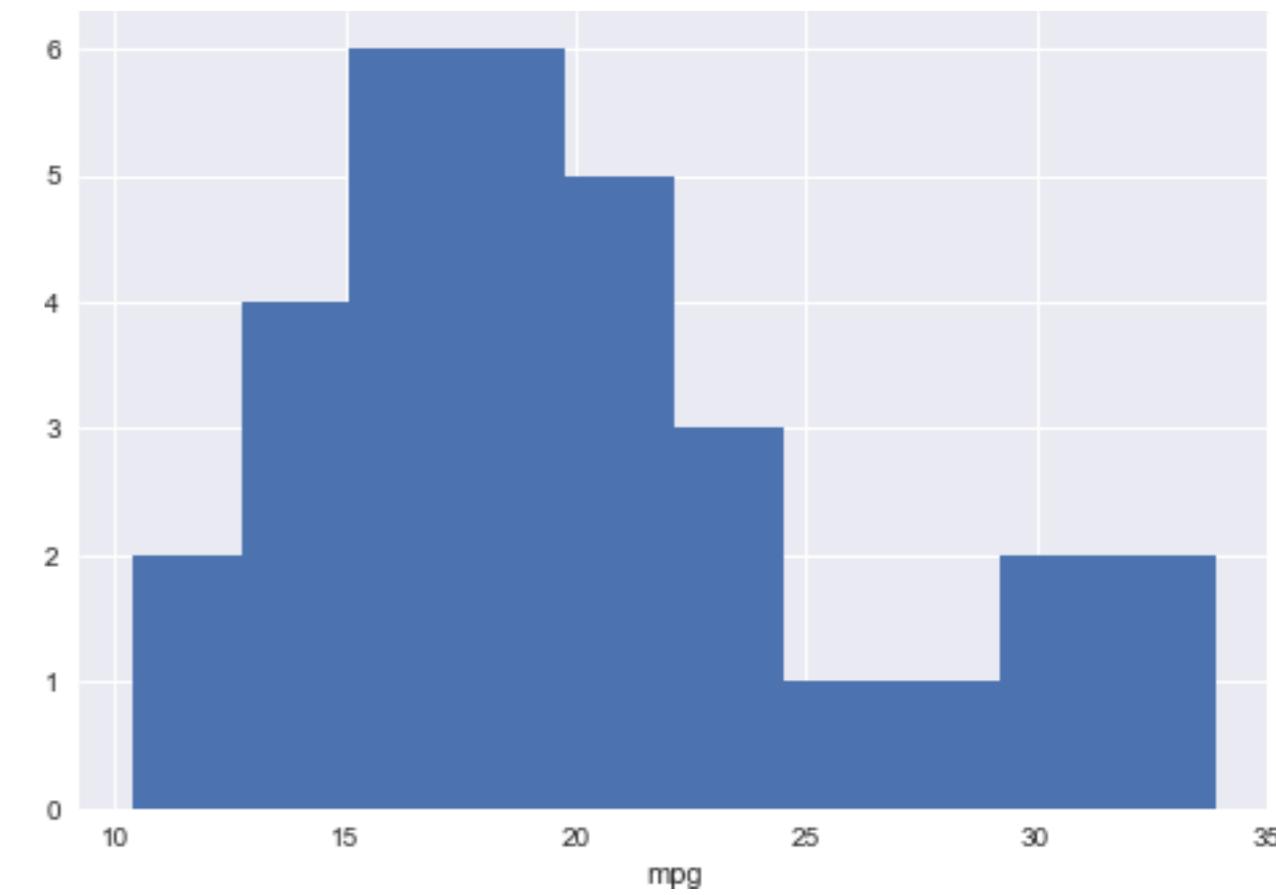
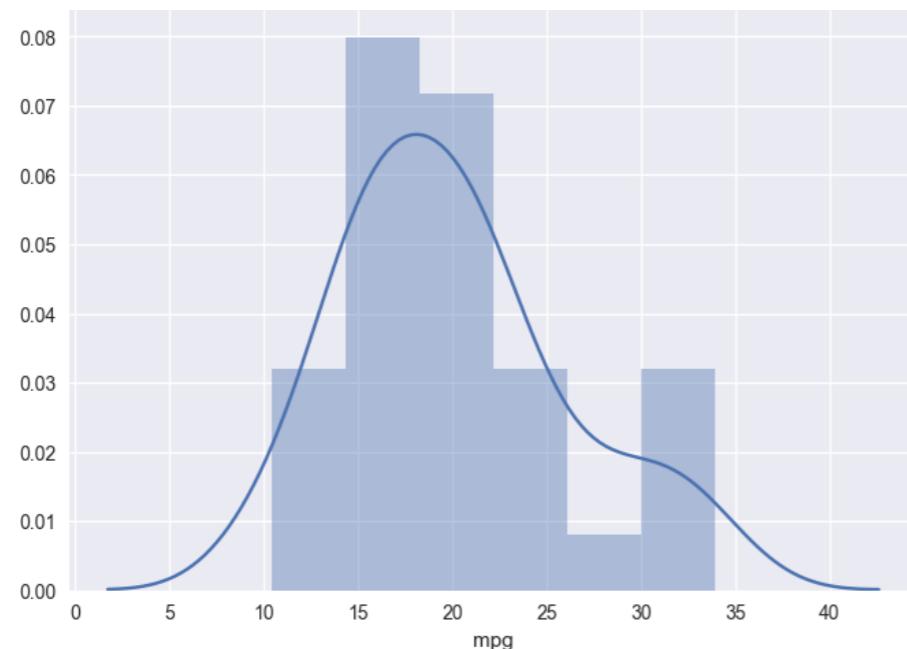
Basic Pandas/matplotlib





Can set limits, tick styles, scales, add lines, annotations, titles, legends

Seaborn provides a different visual style and lots of canned plots.



EFFECTIVE VISUALIZATIONS

NOT EFFECTIVE...

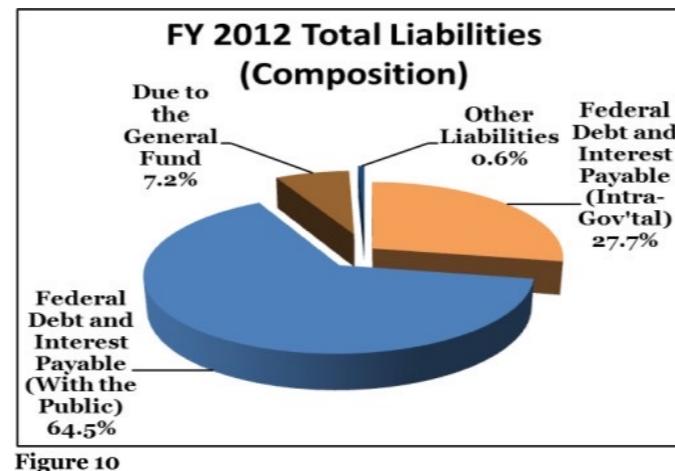
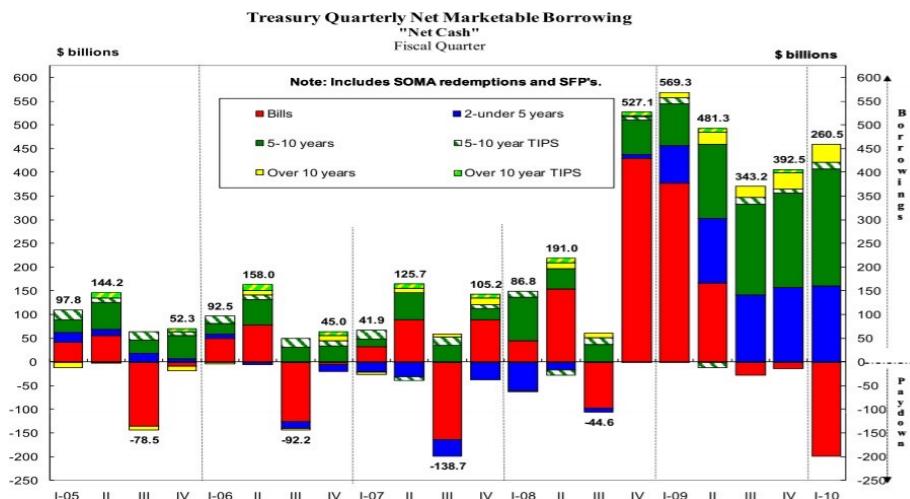


Figure 10

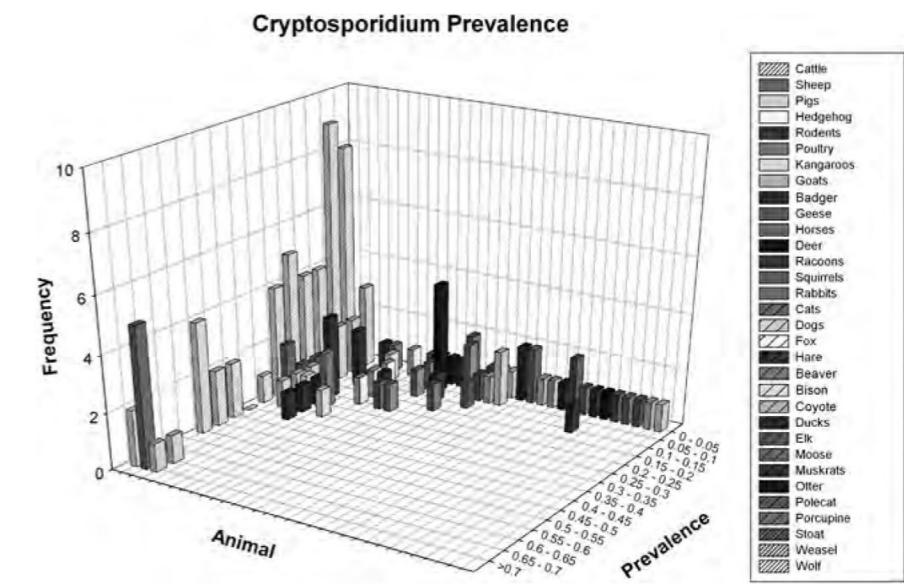
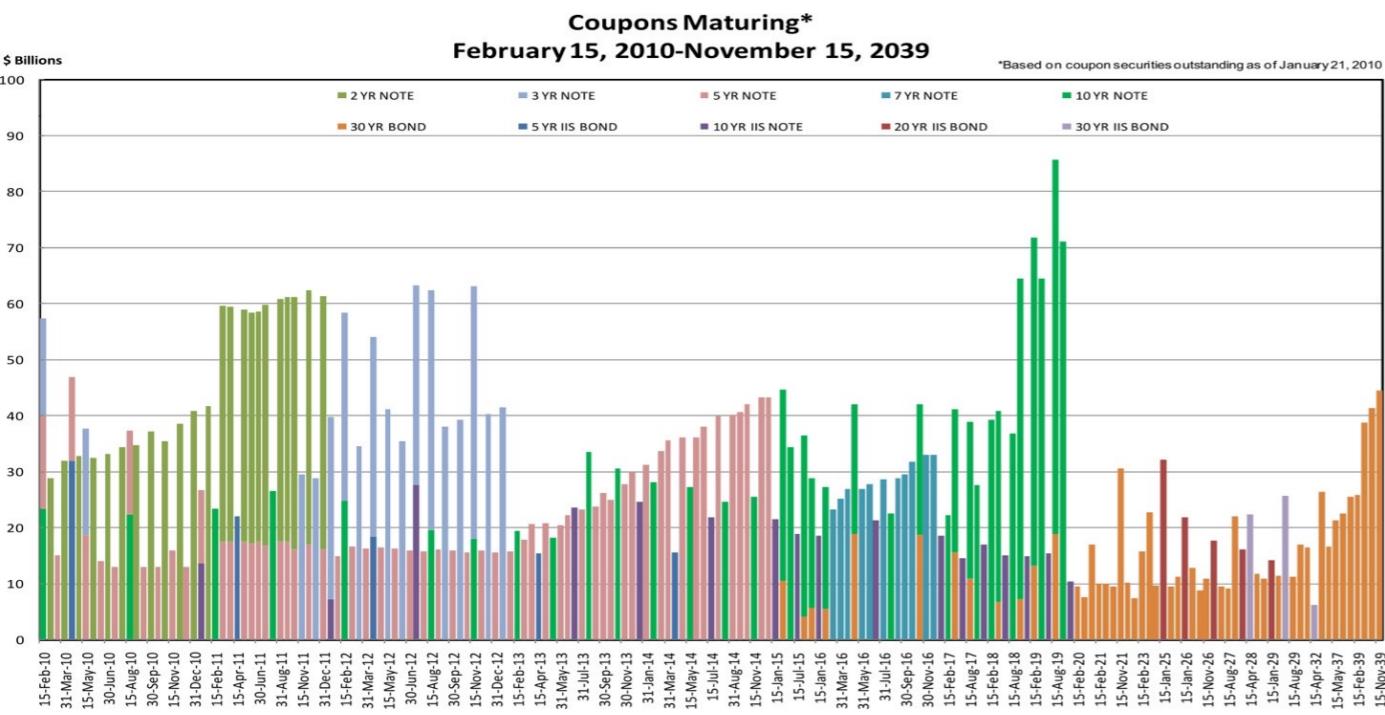
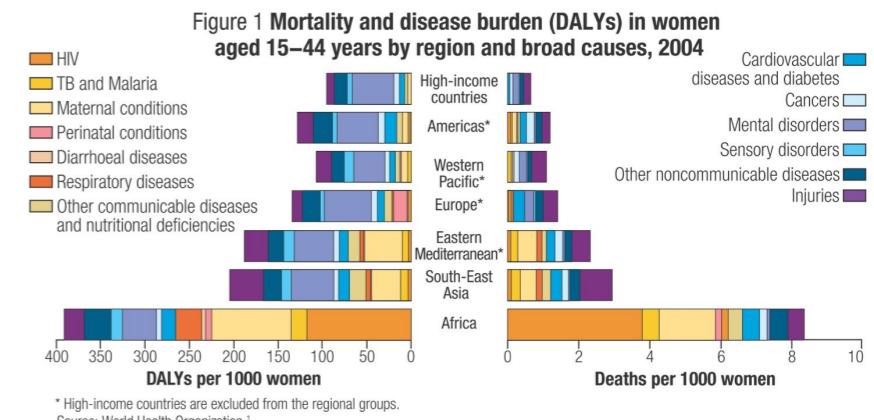


Figure 5.2 Mean prevalence rates of *Cryptosporidium* oocysts by animal species.

Sources: US Treasury and WHO reports

EFFECTIVE EDA VISUALIZATION

1. Have graphical integrity
2. Keep it simple
3. Use the right display
4. Use color sensibly

I. GRAPHICAL INTEGRITY



Same Veritas. More Lux.

Yale Summer Session

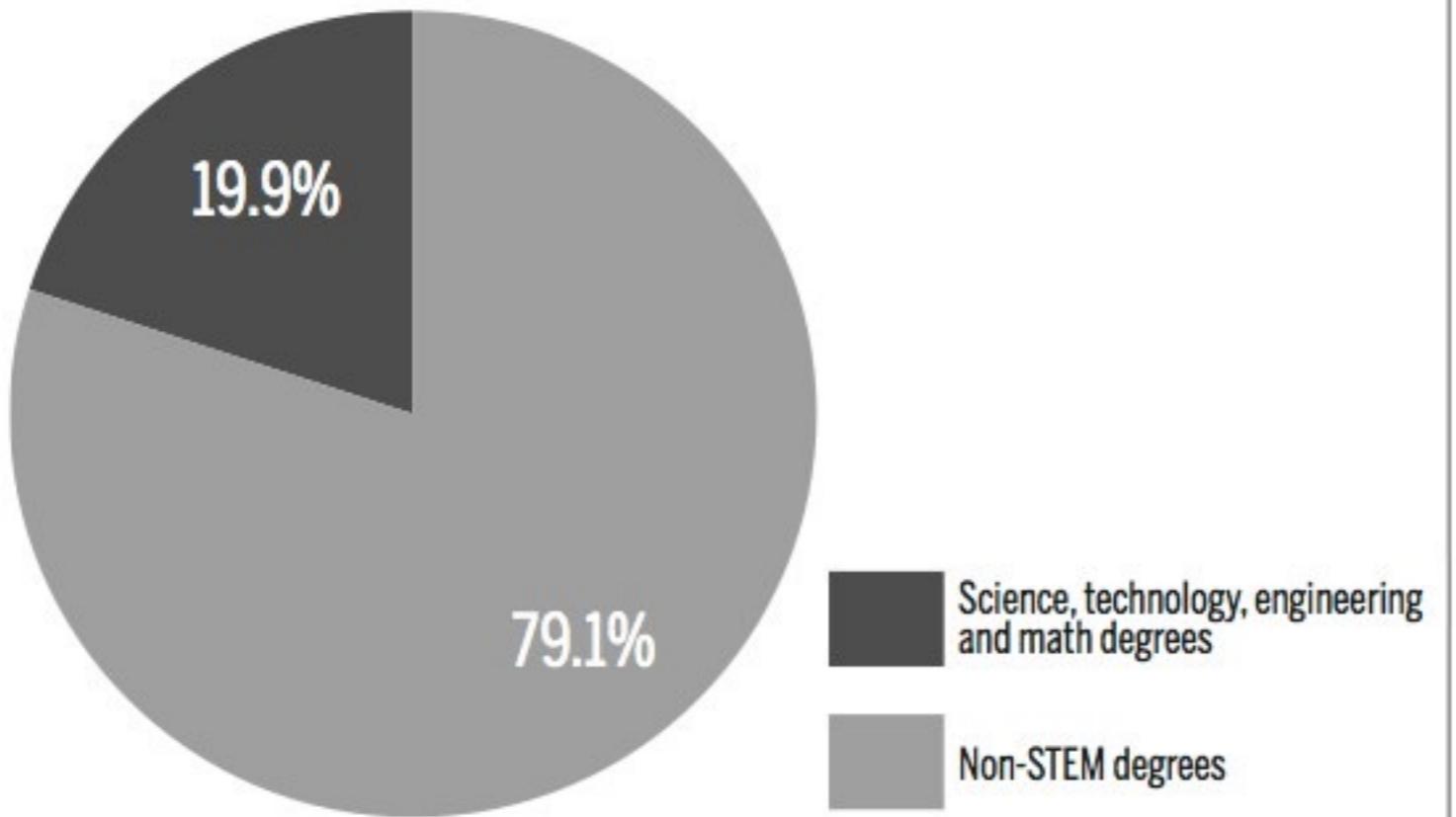
Over 200 full-credit courses.

June 4 – July 6 , July 9 – Aug 10



2012 *experience Yale*

CHART YALE GRADUATES' MAJORS, CLASS OF 2011



Facebook Recommendations



[Shake Shack to open in New Haven](#)
277 people recommend this.



[Popular anti-religion creates false dichotomy](#)
15 people recommend this.



[Friends remember Foucher LAW '14](#)
10 people recommend this.



[AIDS activist speaks about documentary film](#)
8 people recommend this.



[Panel outlines changes in hip-hop](#)
30 people recommend this.

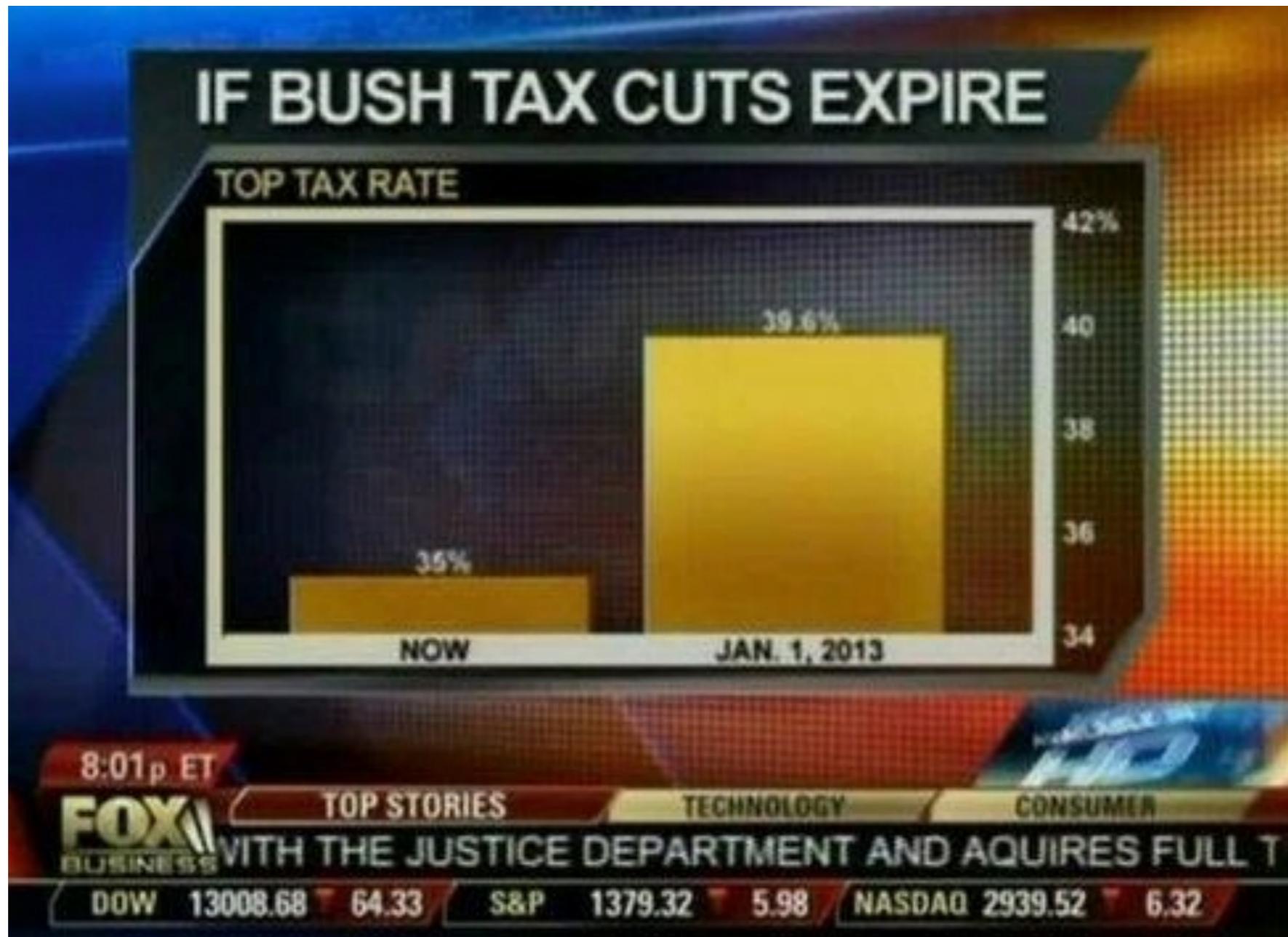


Facebook social plugin

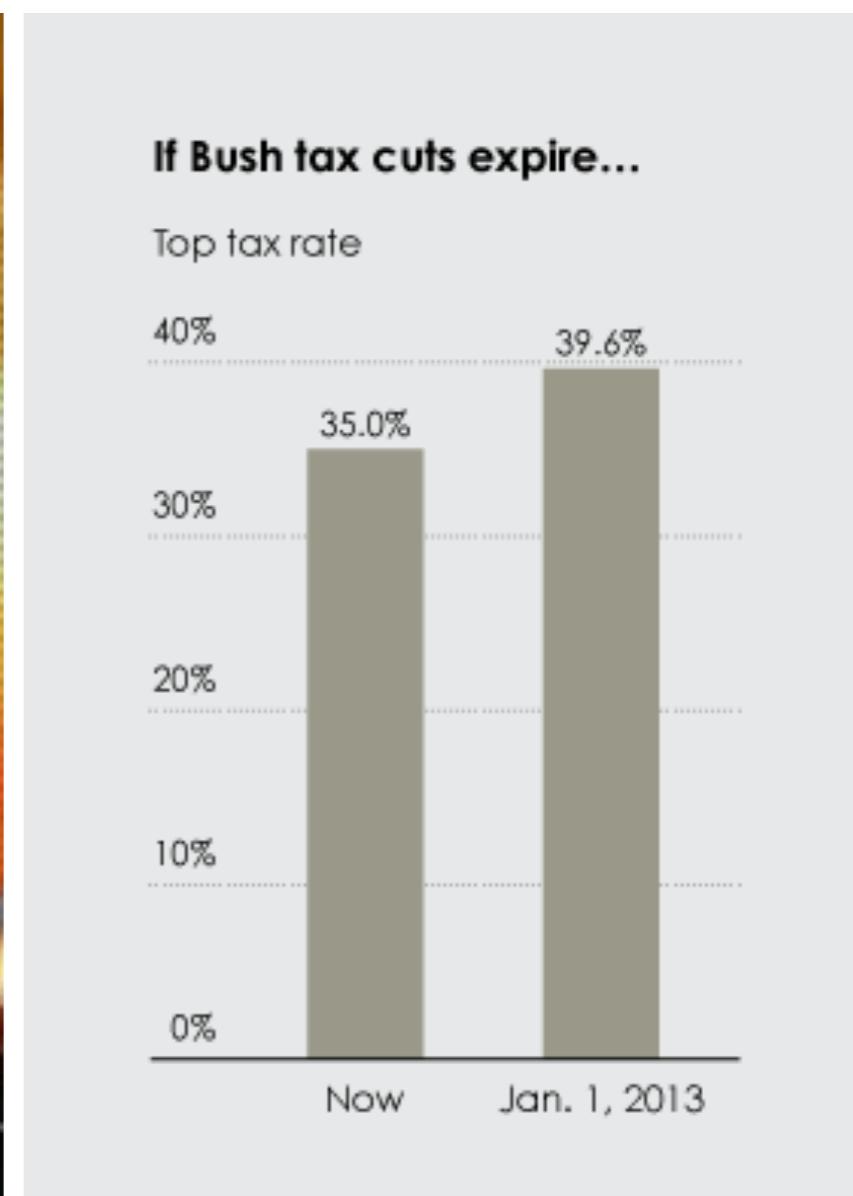
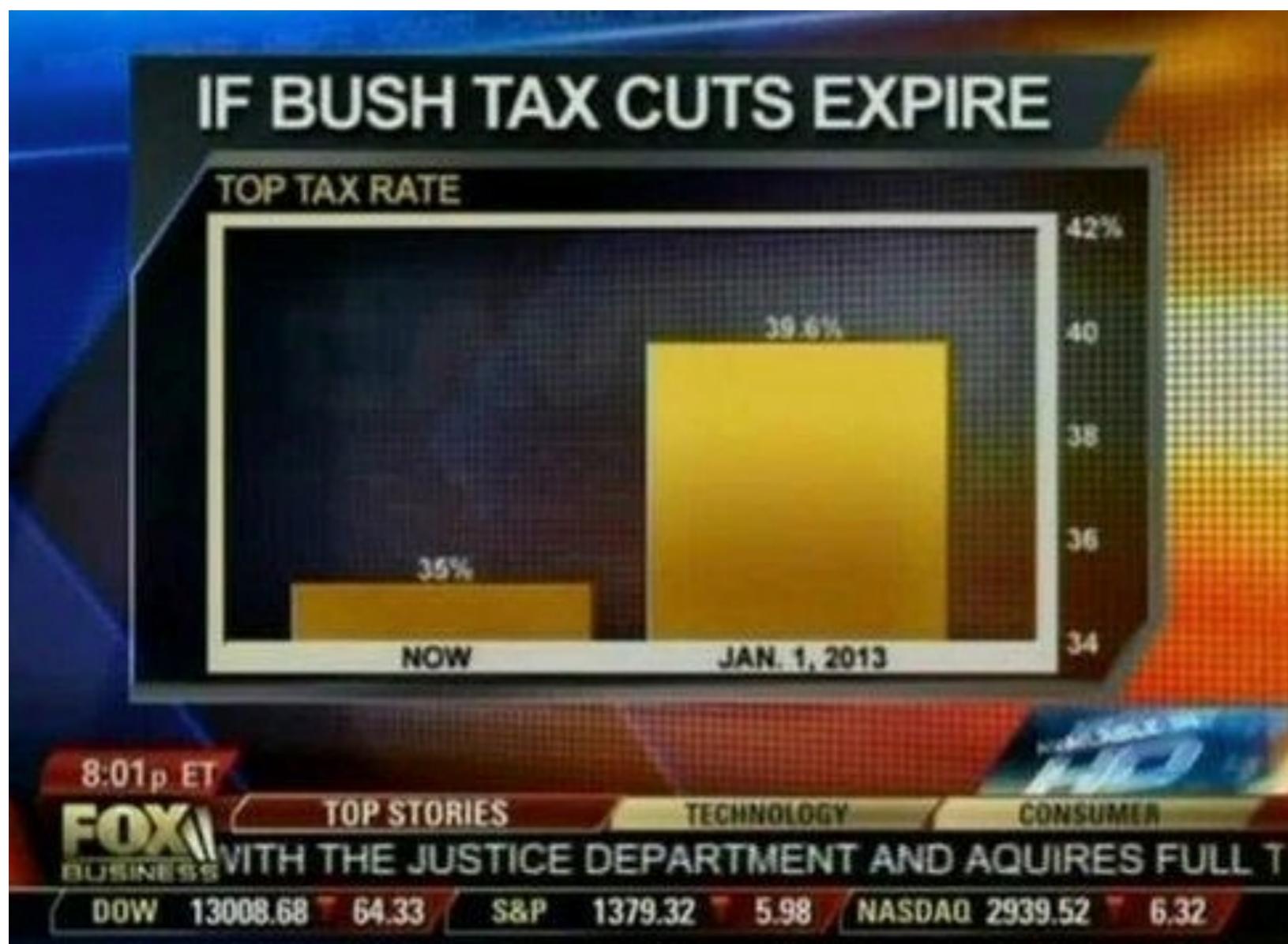
Advertisement

Featured
Jobs

GRAPHICAL INTEGRITY



SCALE DISTORTIONS

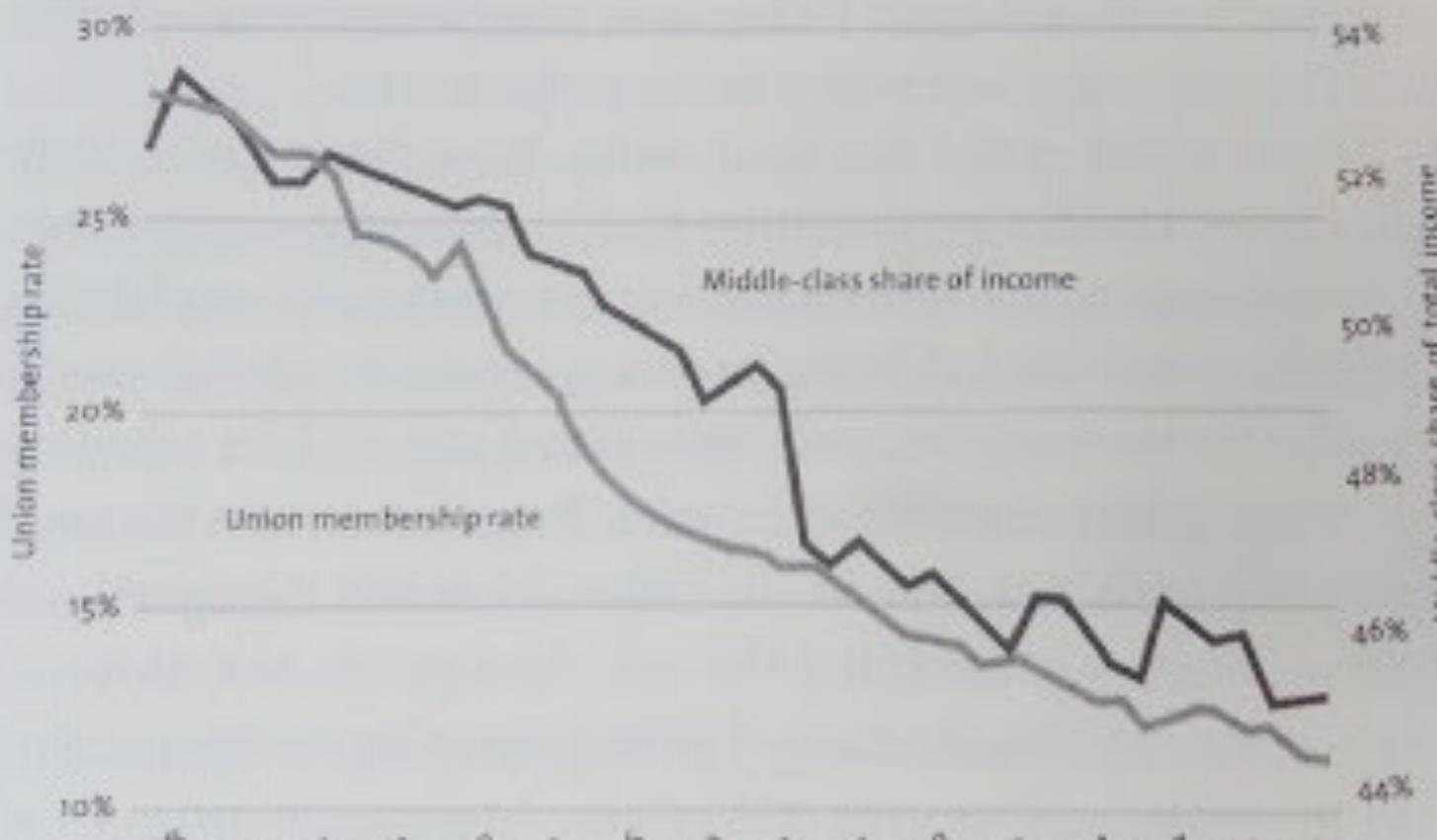


SCALE DISTORTIONS

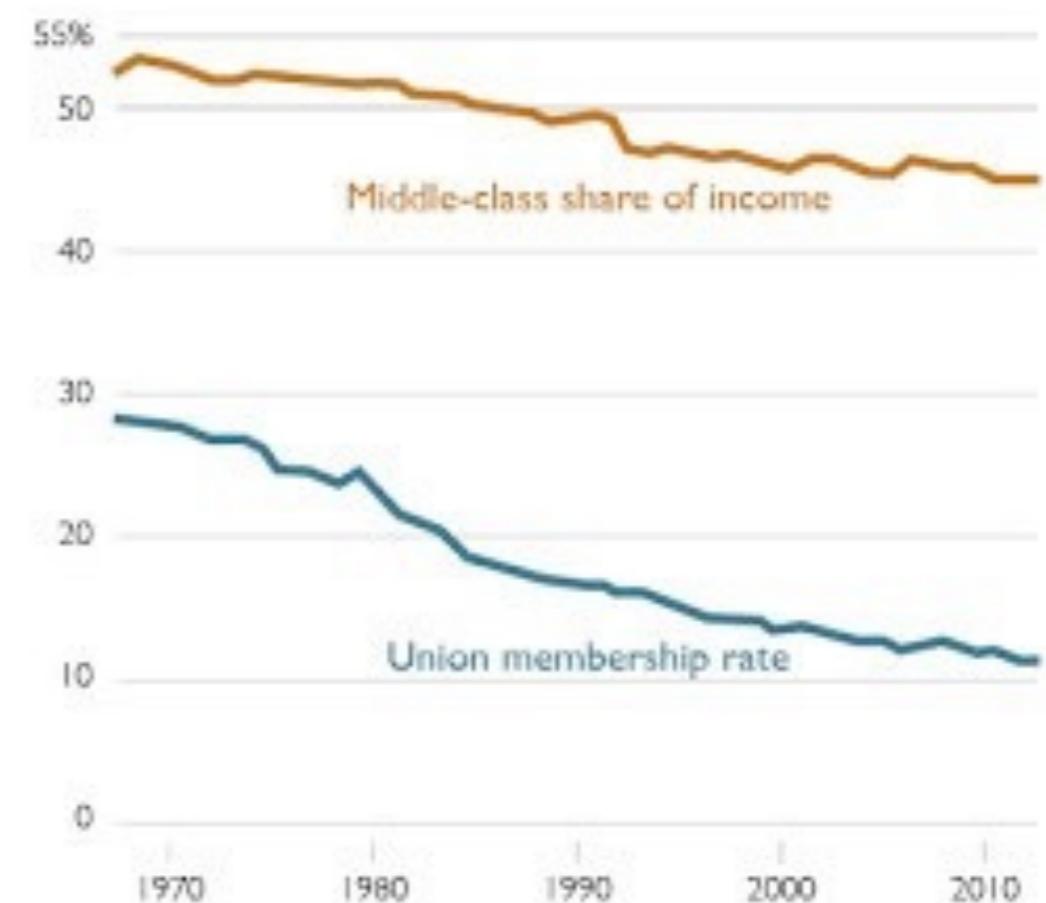


“DOUBLE THE AXES, DOUBLE THE MISCHIEF” (QUOTE FROM GARY SMITH’S STANDARD DEVIATIONS)

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



NEW VERSION



Graphic from Robert Reich's Saving Capitalism

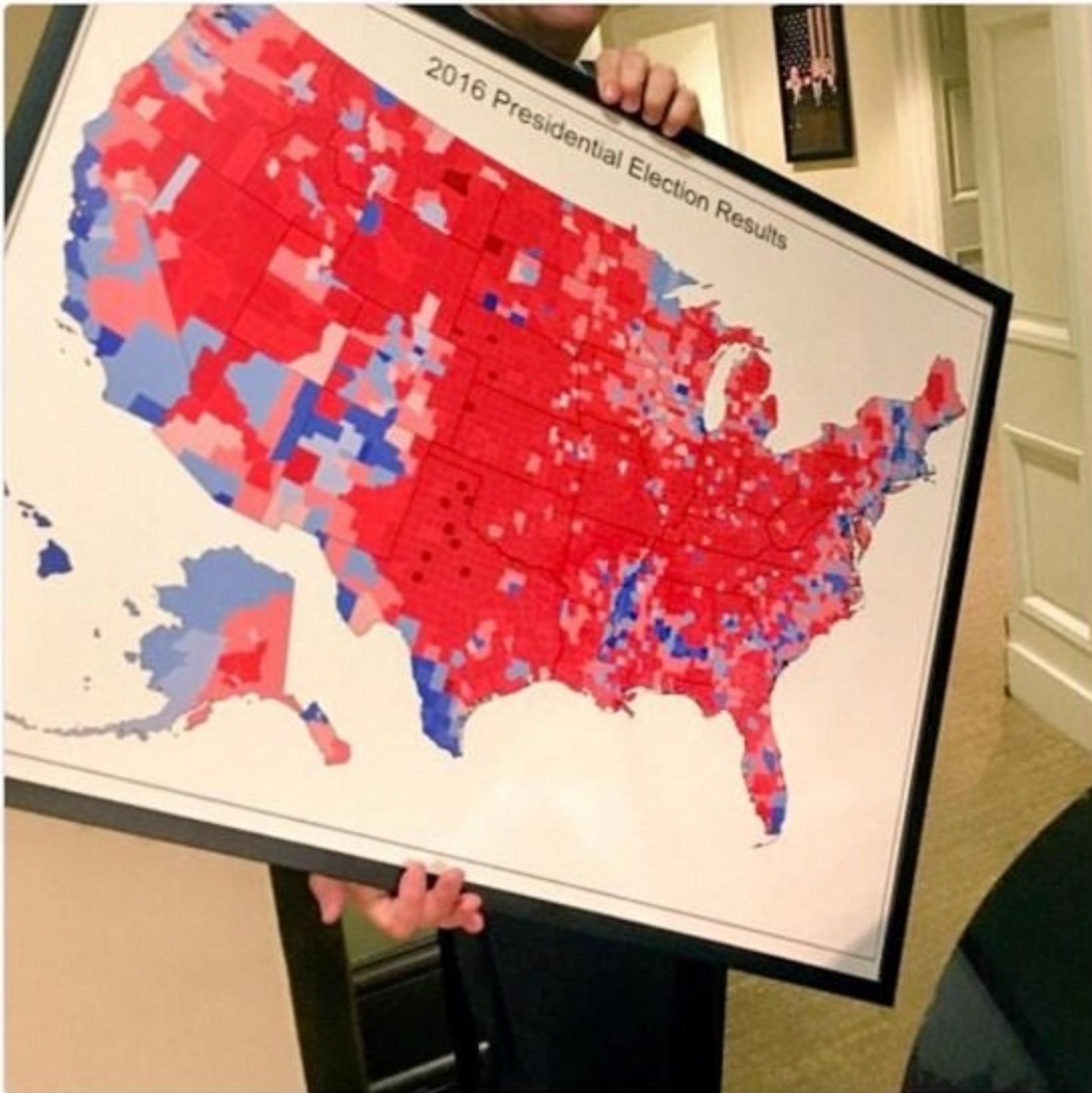
<http://www.thefunctionalart.com/2015/10/double-axes-double-mischief.html>

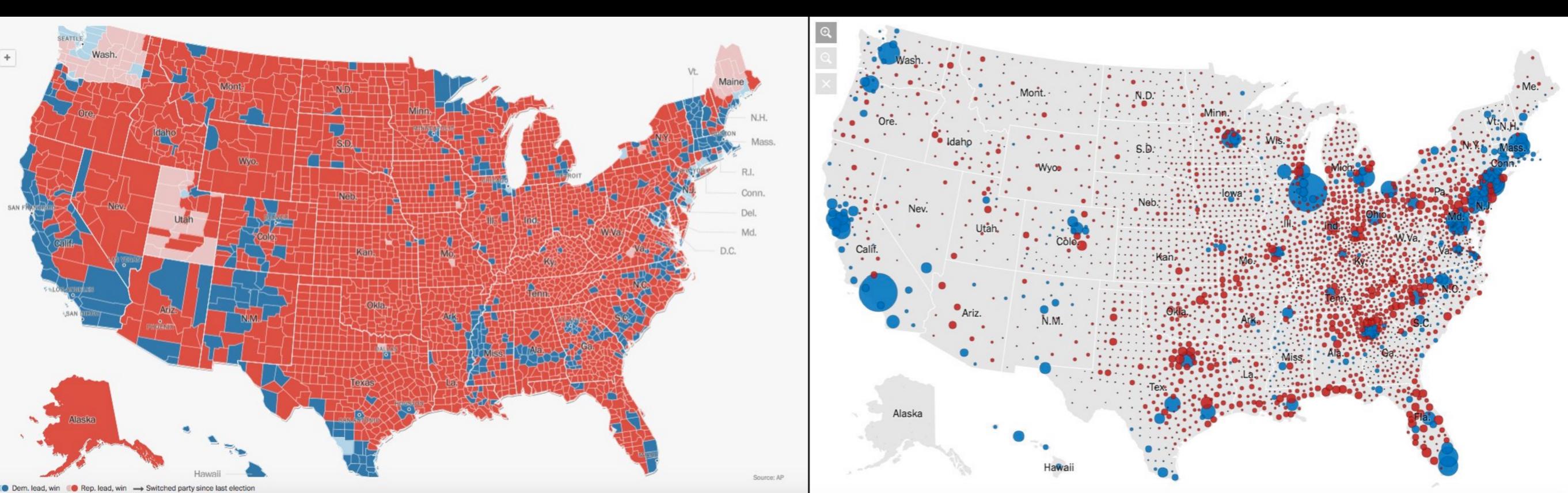
BE PROPORTIONAL



Trey Yingst • @TreyYingst • May 11

Spotted: A map to be hung somewhere in the West Wing





US Presidential Election 2016

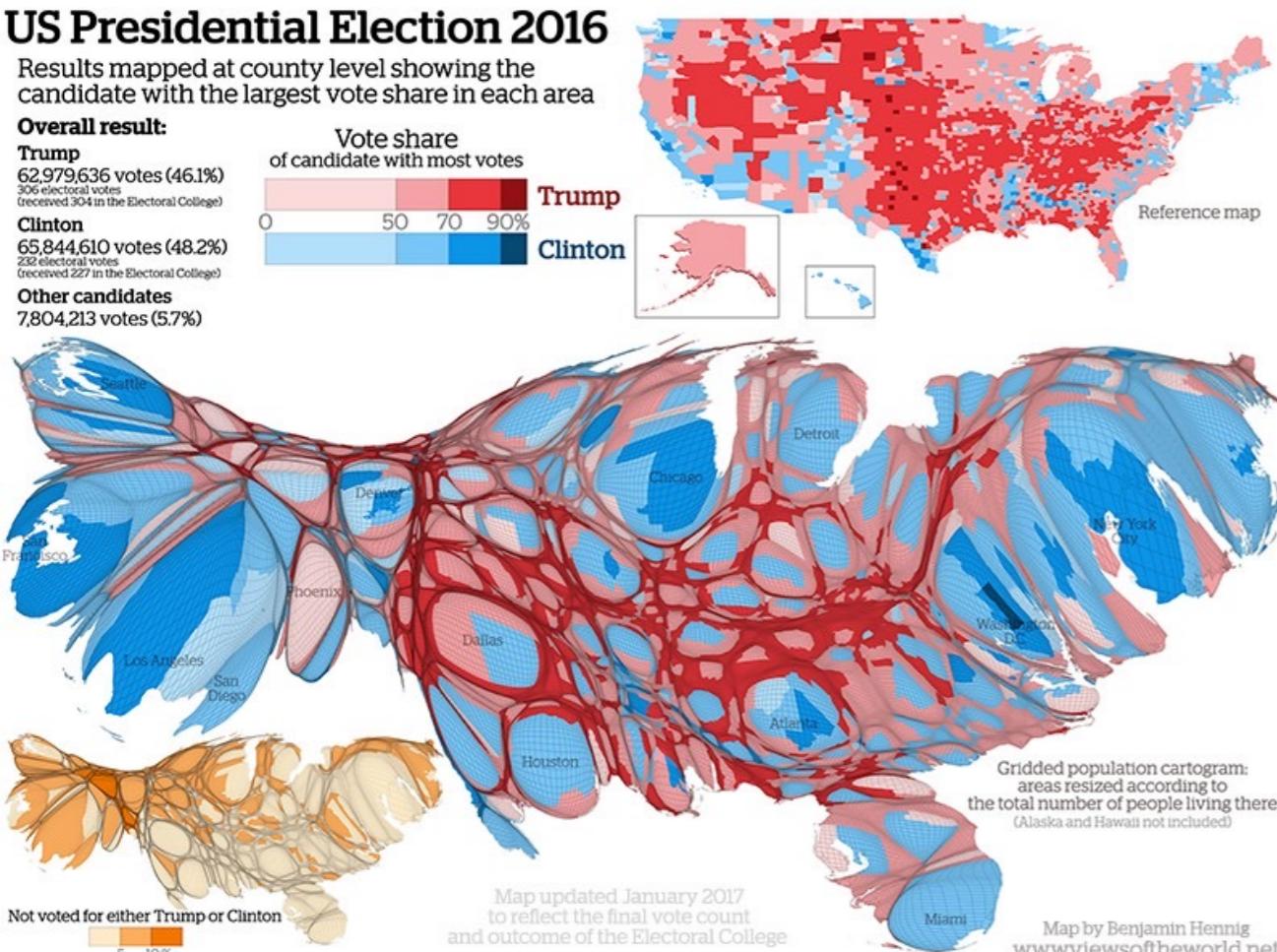
Results mapped at county level showing the candidate with the largest vote share in each area

Overall result:

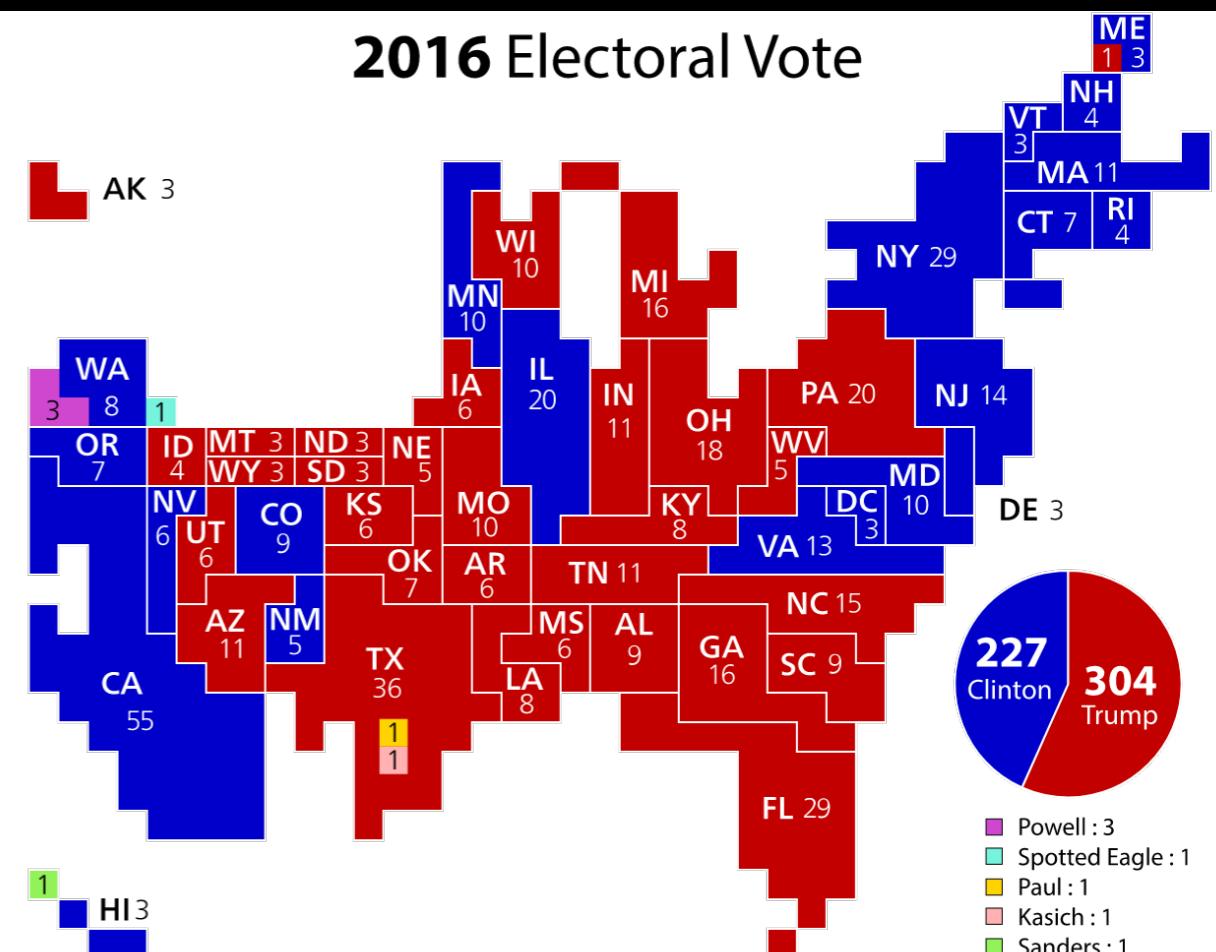
Trump
62,979,636 votes (46.1%)
306 electoral votes
(received 304 in the Electoral College)

Clinton
65,844,610 votes (48.2%)
232 electoral votes
(received 227 in the Electoral College)

Other candidates
7,804,213 votes (5.7%)



2016 Electoral Vote



INCLUDE UNCERTAINTY

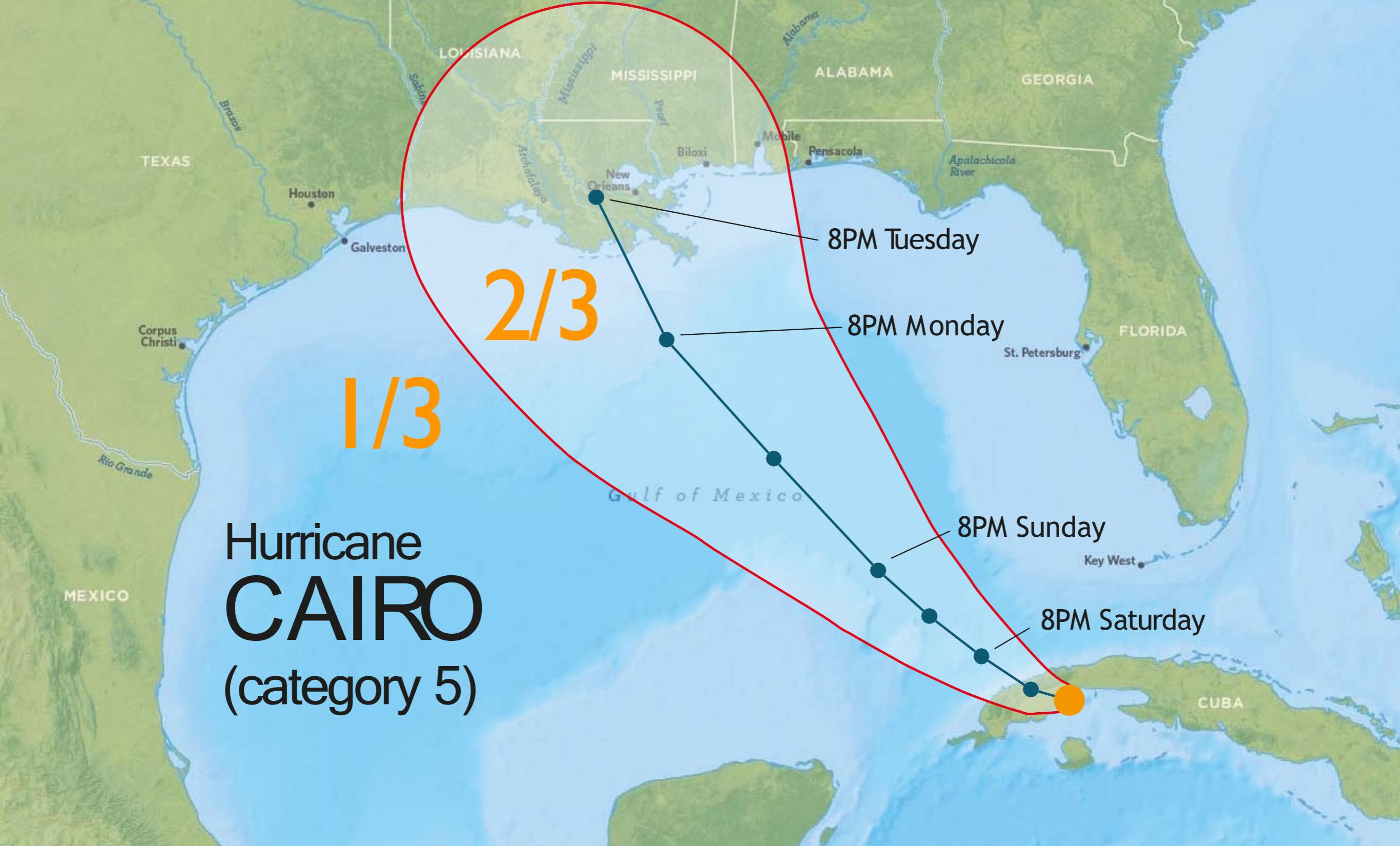
Note: The cone contains the probable path of the storm center but does not show the size of the storm. Hazardous conditions can occur outside of the cone.



Hurricane **CAIRO** (category 5)



What you show

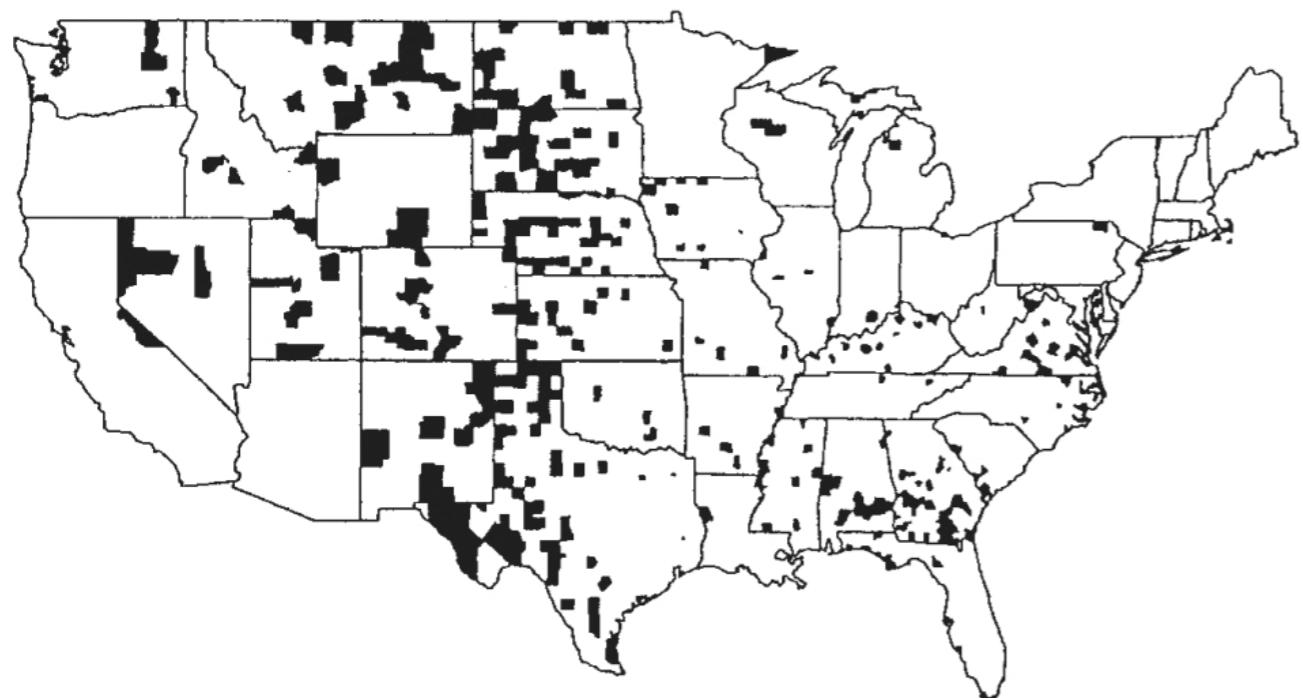


What non-scientists are not aware of (cone is just 66% probability)

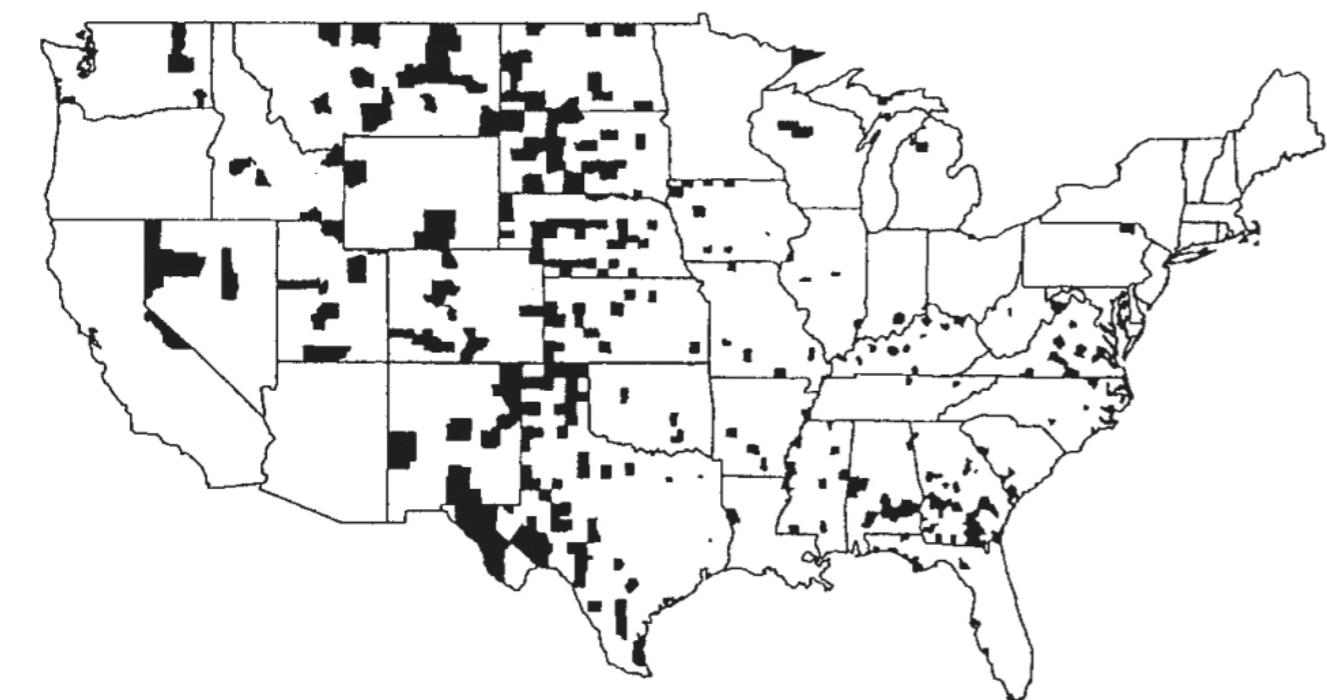
Hurricane **CAIRO** (category 5)

What we could be showing instead

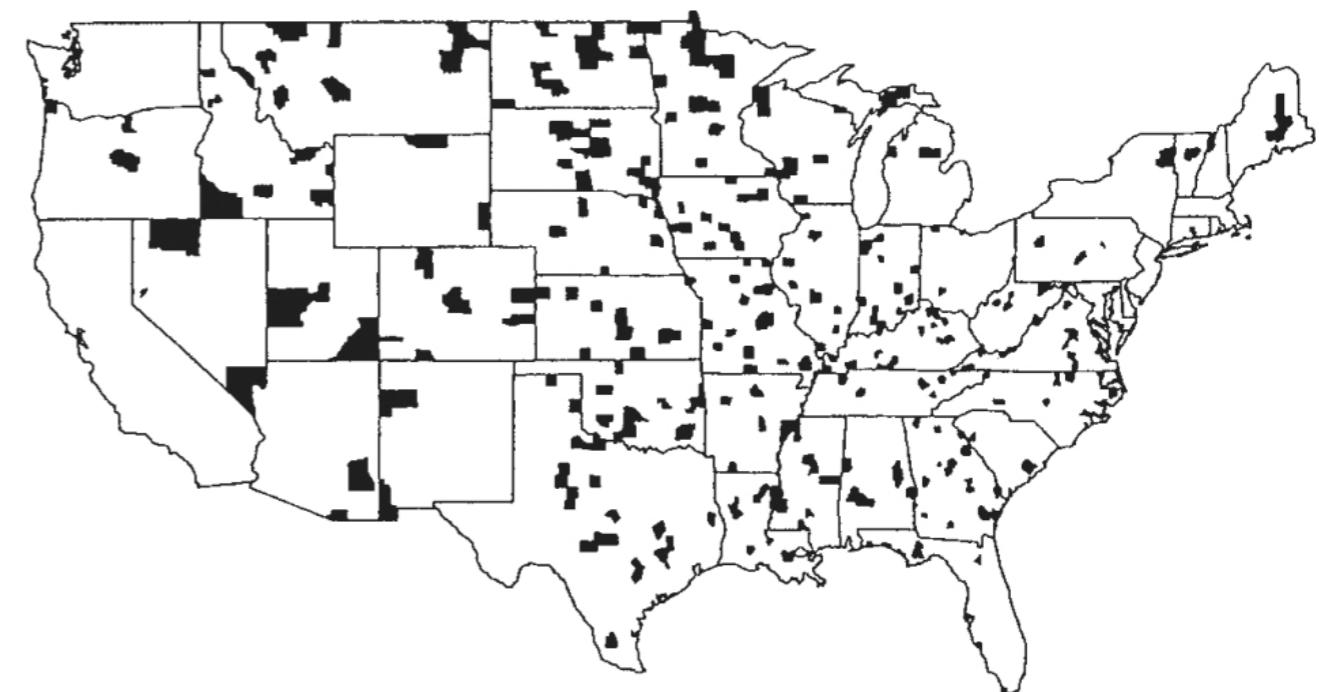
PLOT ALL YOUR DATA



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



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

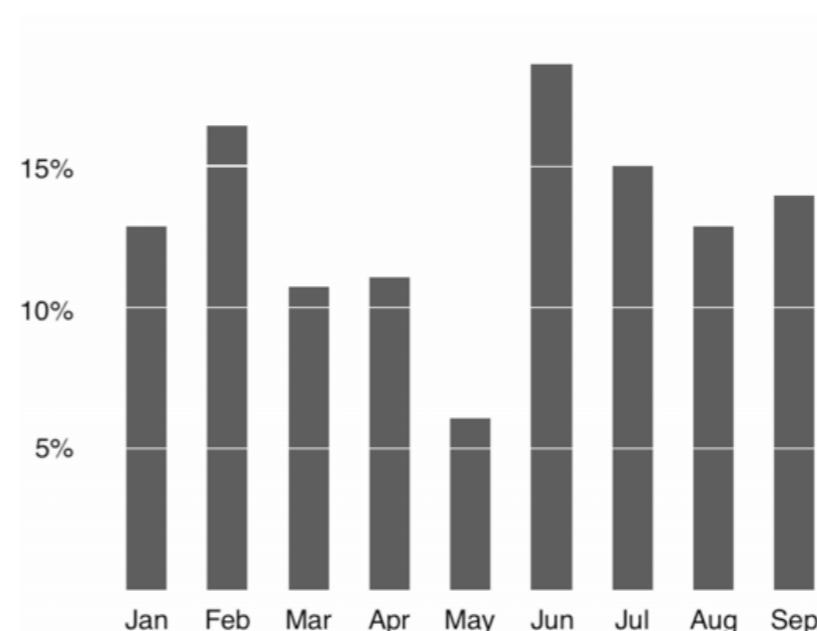
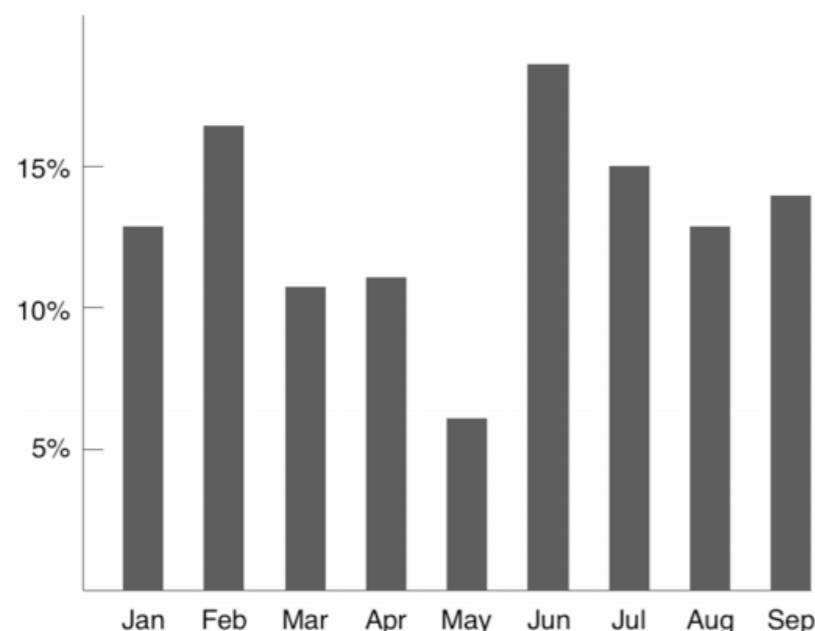
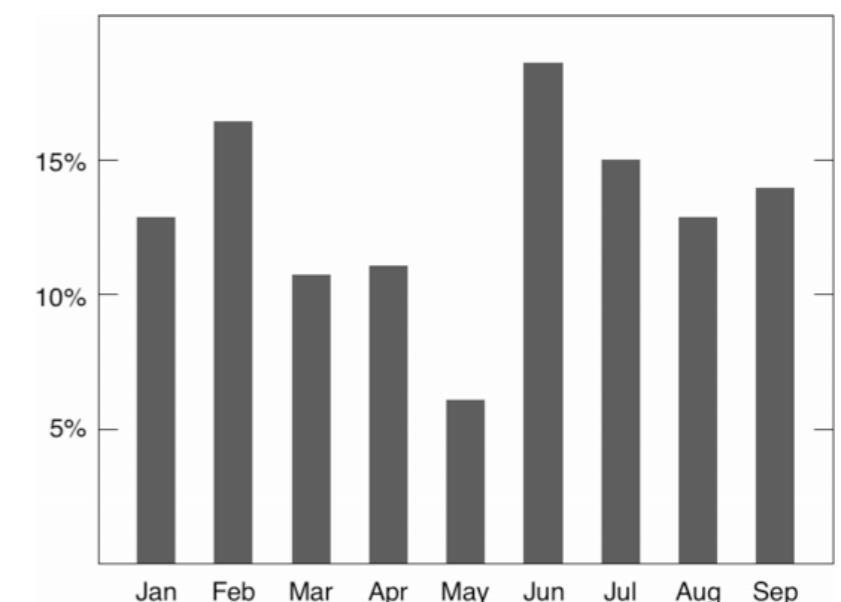
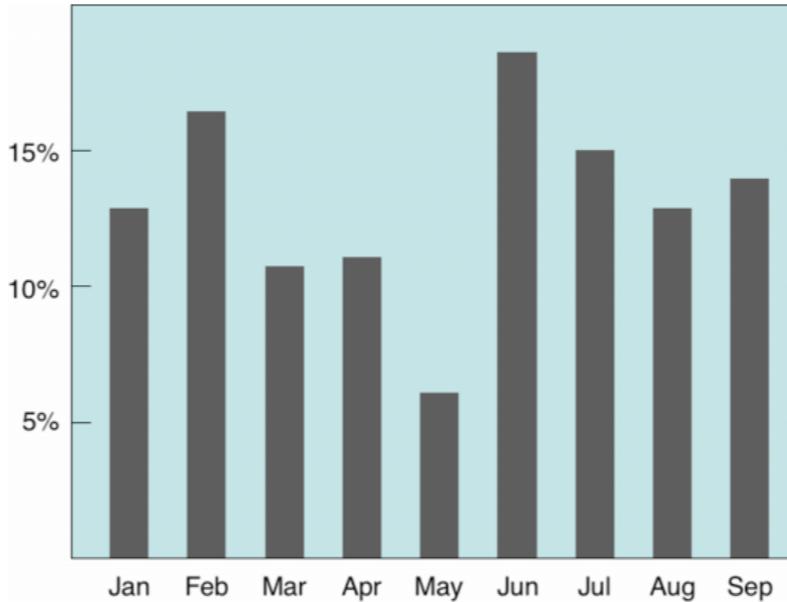
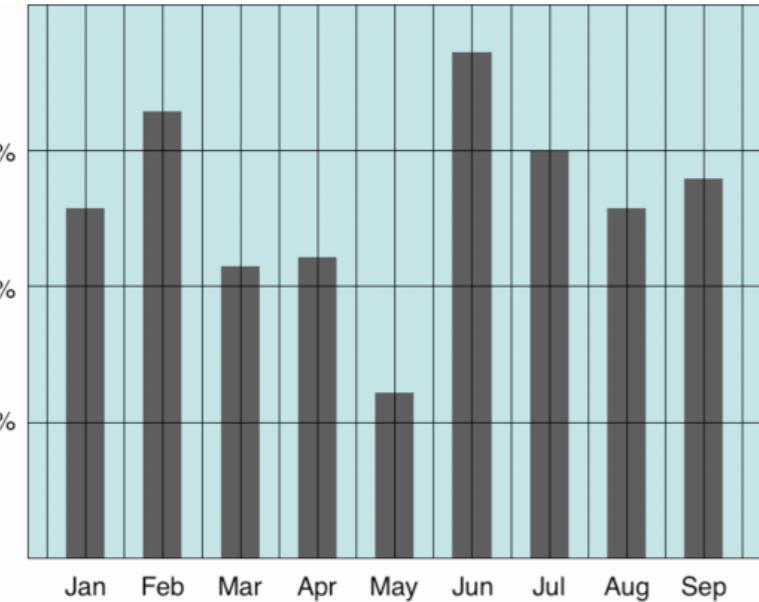


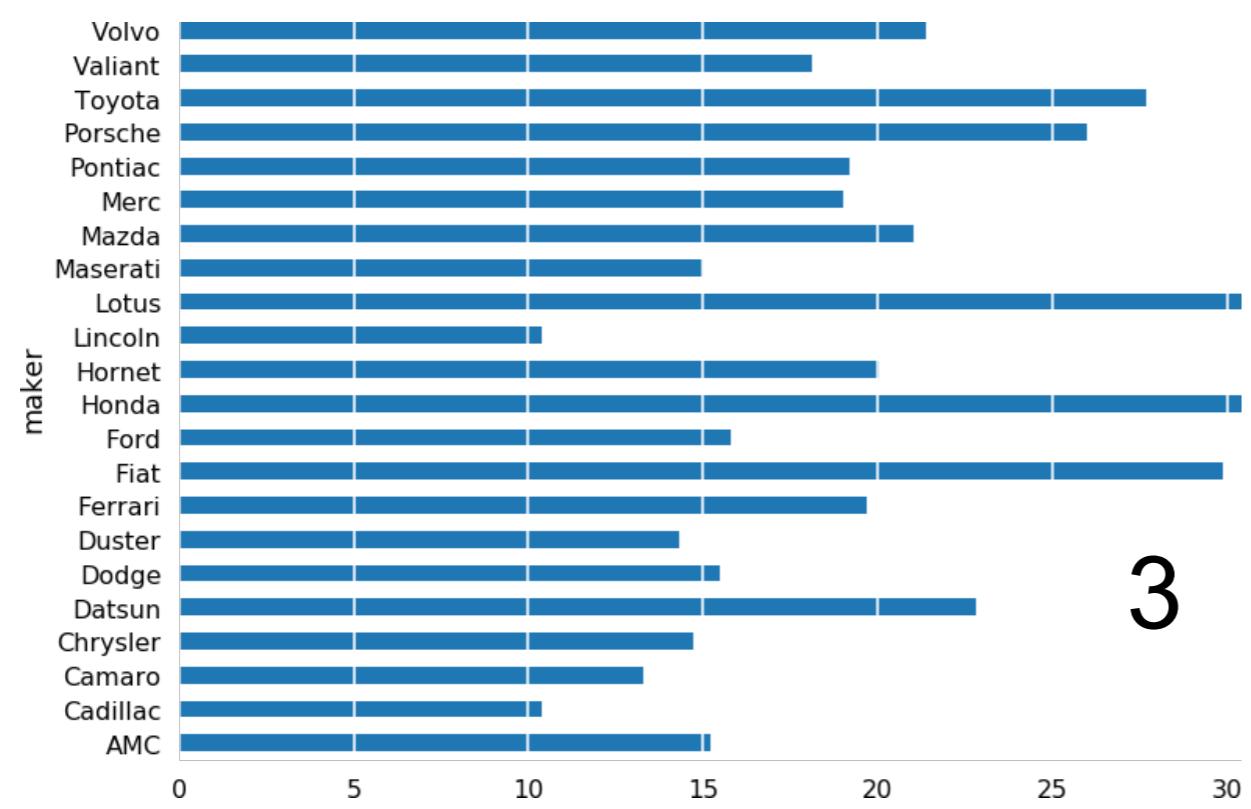
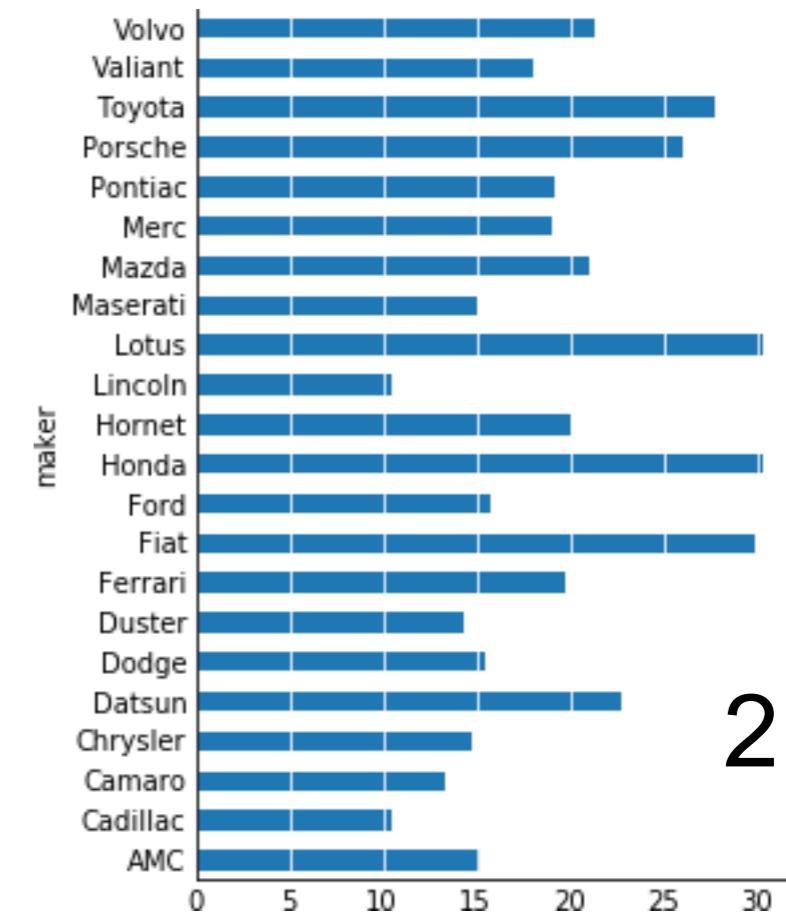
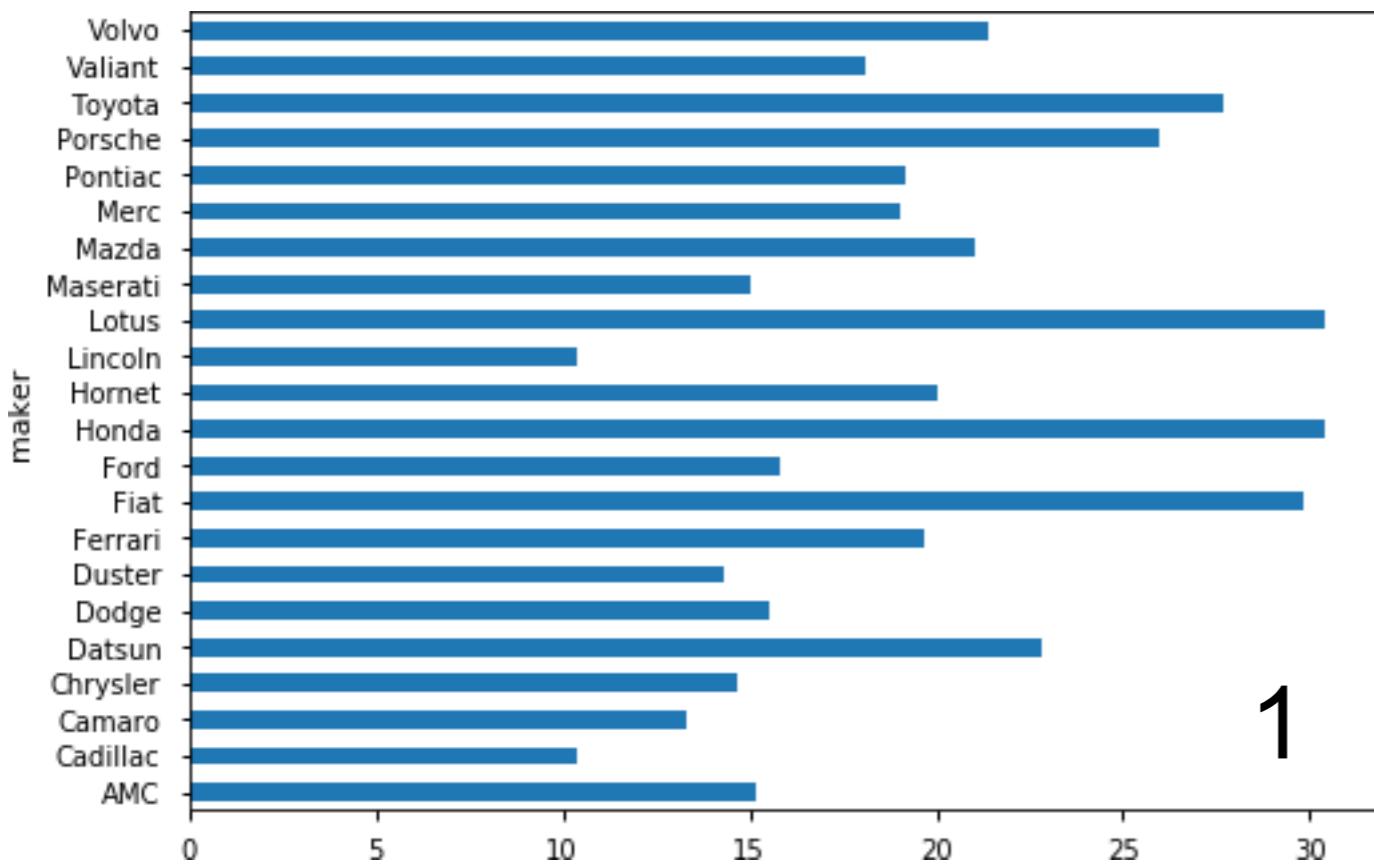
Counties with the HIGHEST
kidney cancer death rates
(1980-1989)

2. KEEP IT SIMPLE

AVOID CHARTJUNK

EXTRANEOUS VISUAL ELEMENTS THAT DISTRACT FROM THE MESSAGE



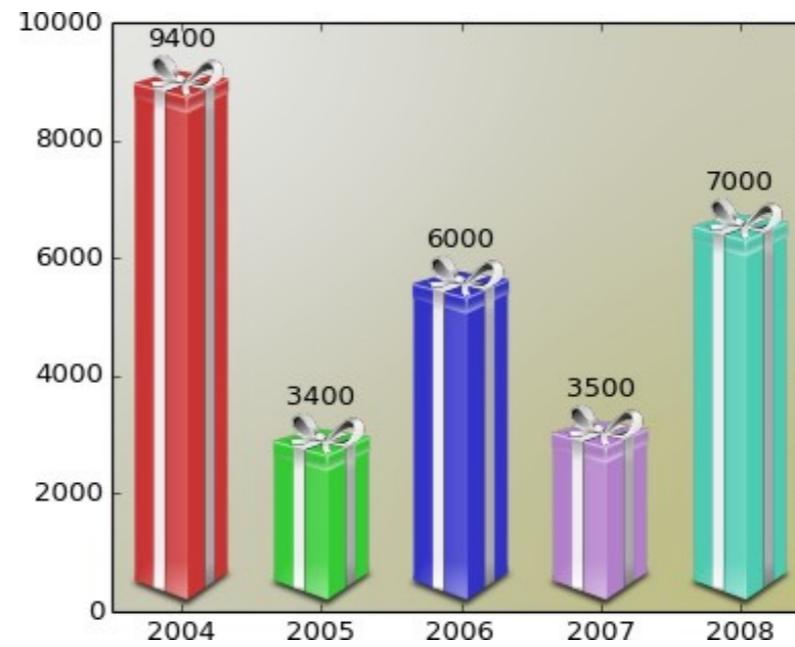
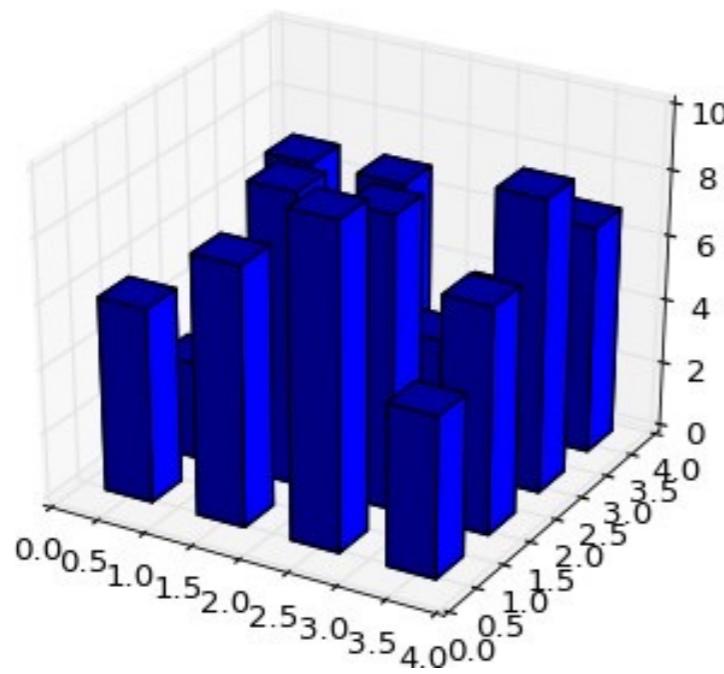


Horizontal bar chart 4 showing car makers and their values. The y-axis is labeled "maker" and lists various car brands. The x-axis ranges from 0 to 30 with increments of 5. The bars are blue. Numerical values are displayed at the end of each bar.

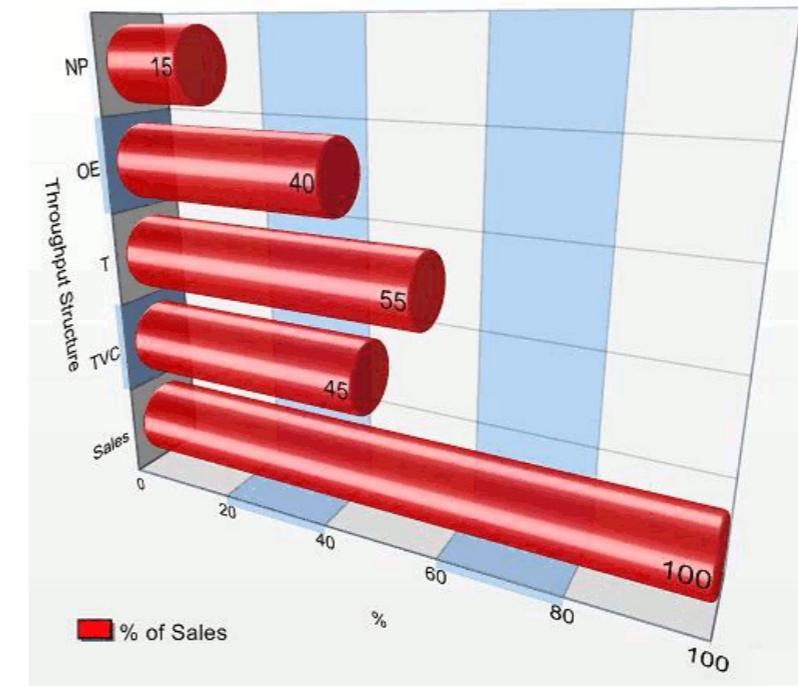
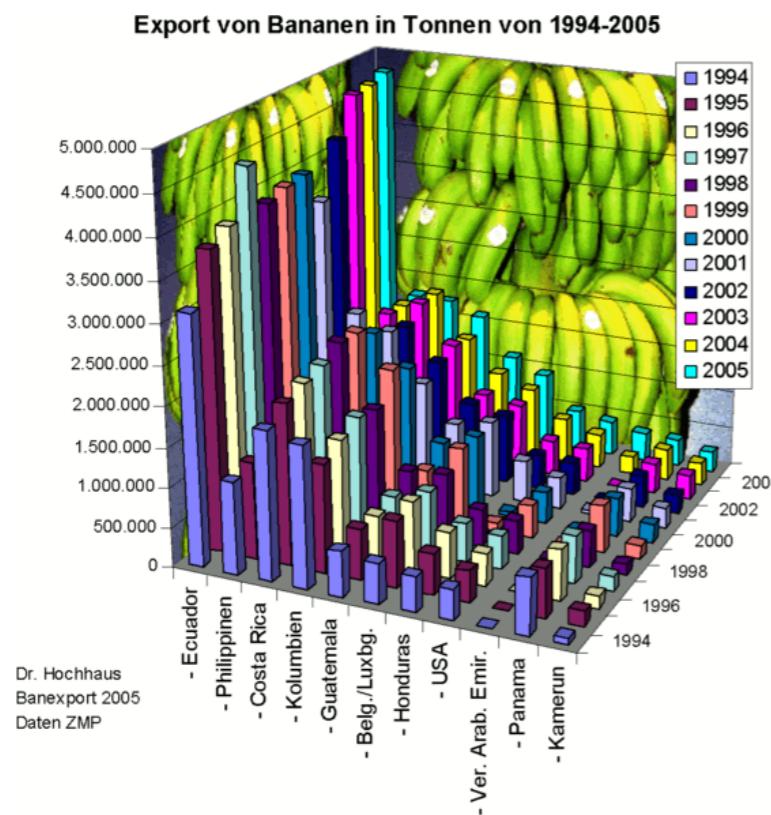
maker	Value
Honda	30.4
Lotus	30.4
Fiat	29.85
Toyota	27.7
Porsche	26.0
Datsun	22.8
Volvo	21.4
Mazda	21.0
Hornet	20.05
Ferrari	19.7
Pontiac	19.2
Merc	19.0142857143
Valiant	18.1
Ford	15.8
Dodge	15.5
AMC	15.2
Maserati	15.0
Chrysler	14.7
Duster	14.3
Camaro	13.3
Lincoln	10.4
Cadillac	10.4

4

DON'T!



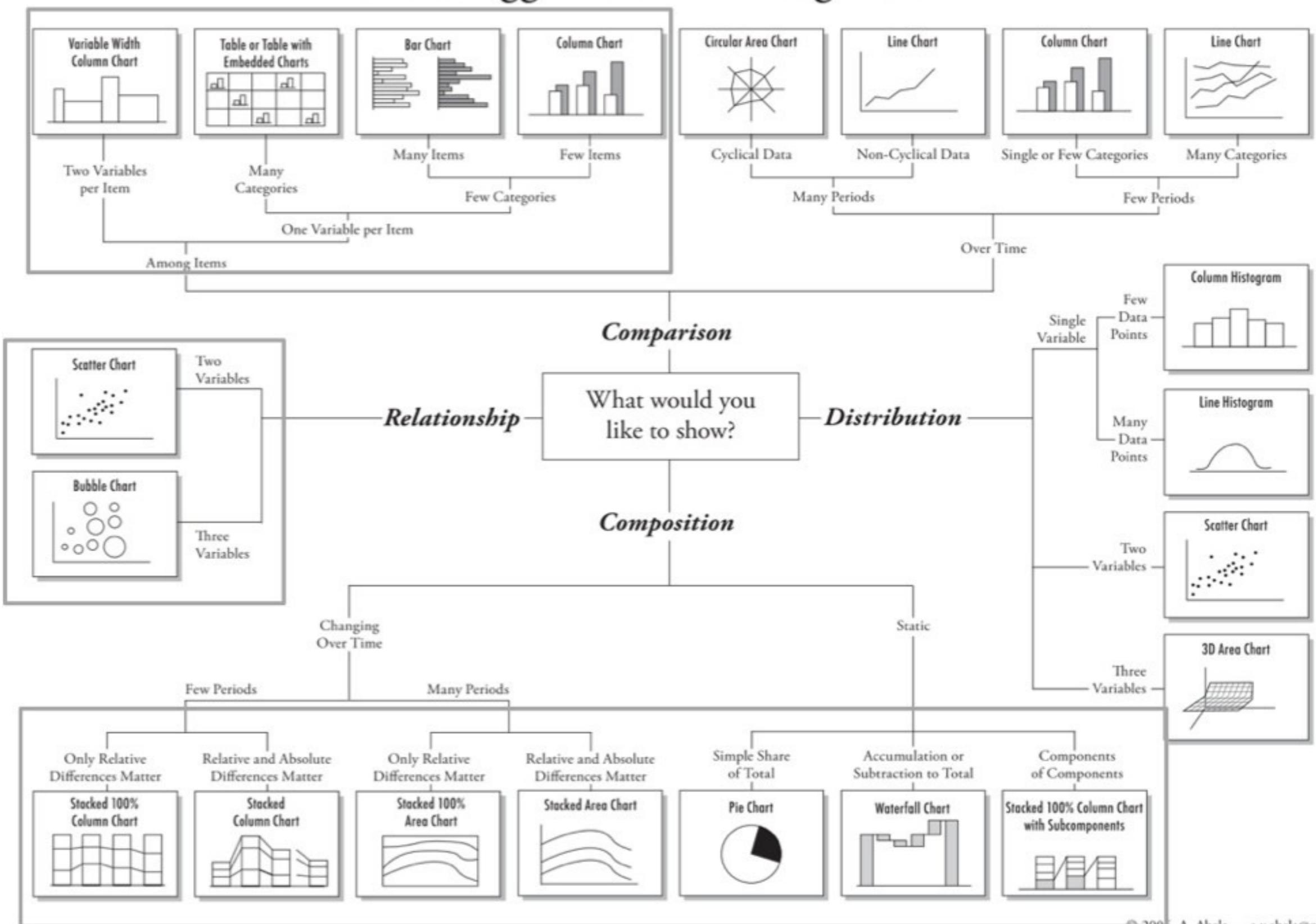
matplotlib gallery



Excel Charts Blog

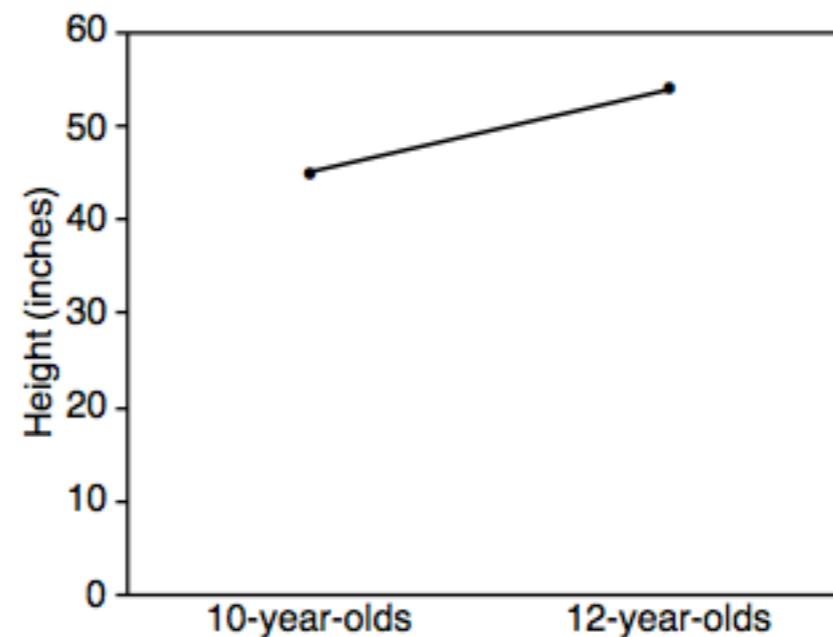
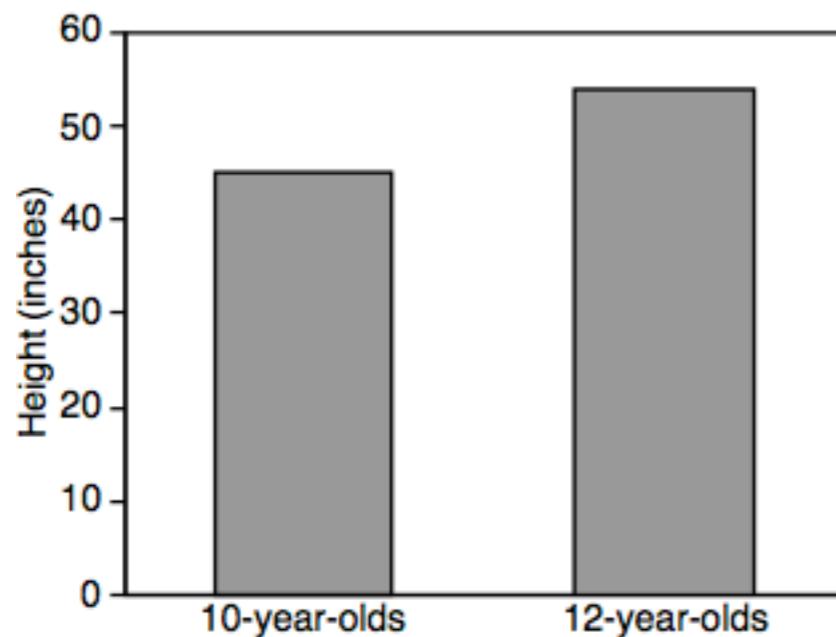
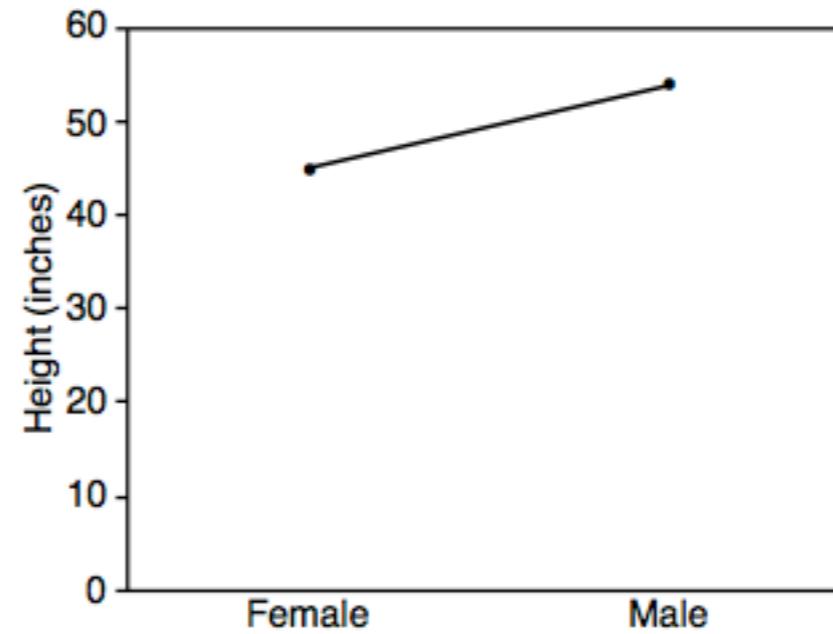
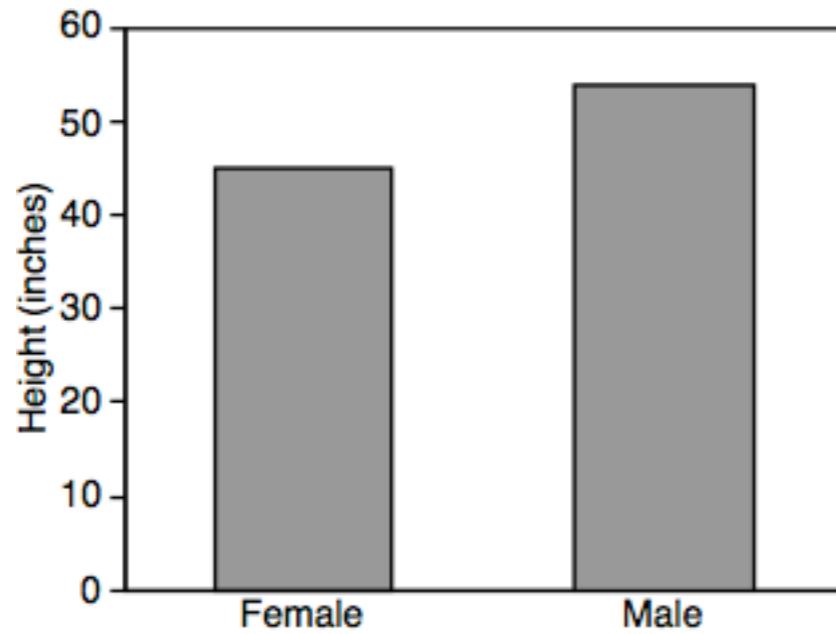
3. USE THE RIGHT DISPLAY

Chart Suggestions—A Thought-Starter



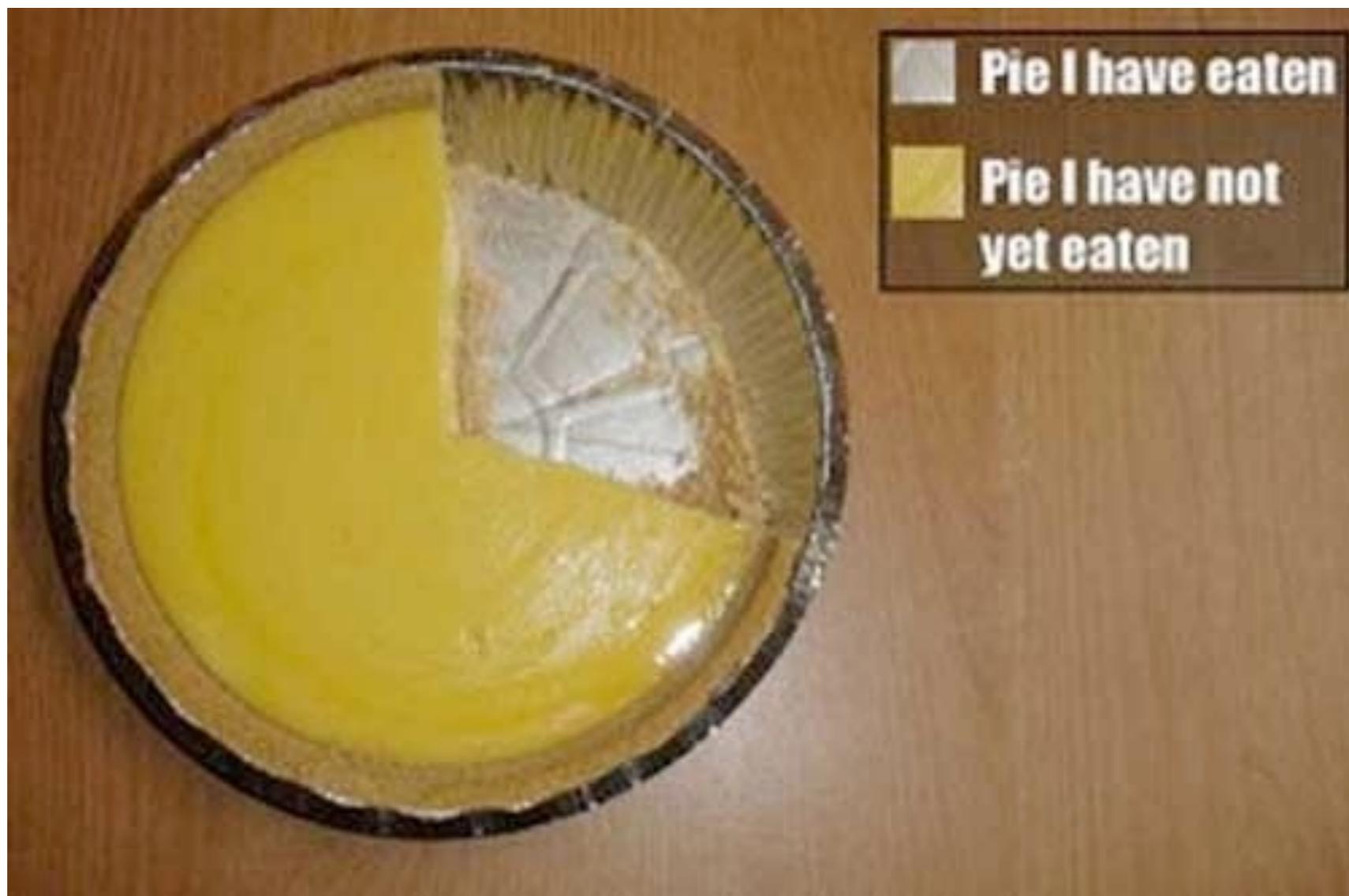
COMPARISONS

BARS vs LINES

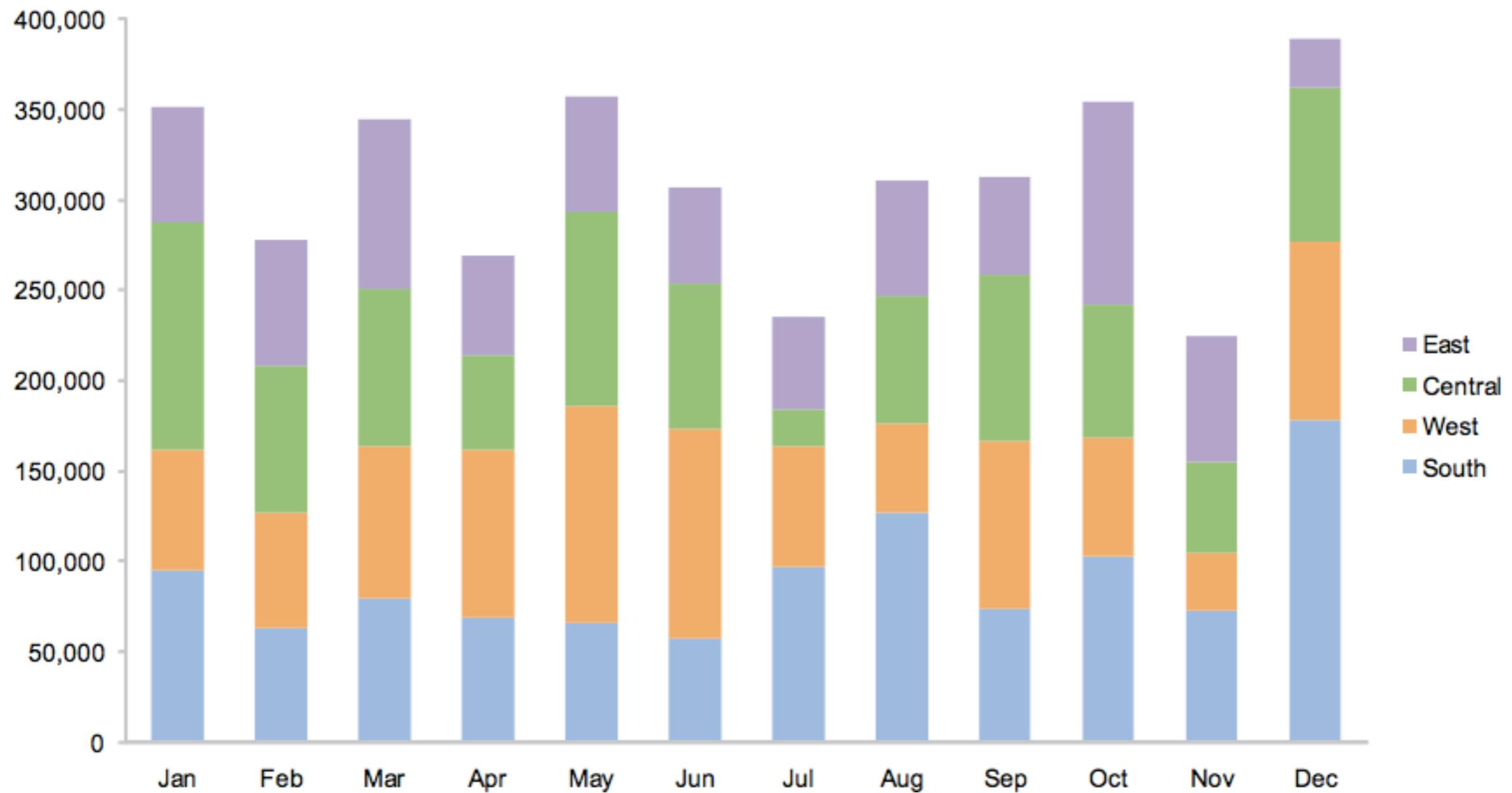


PROPORTIONS

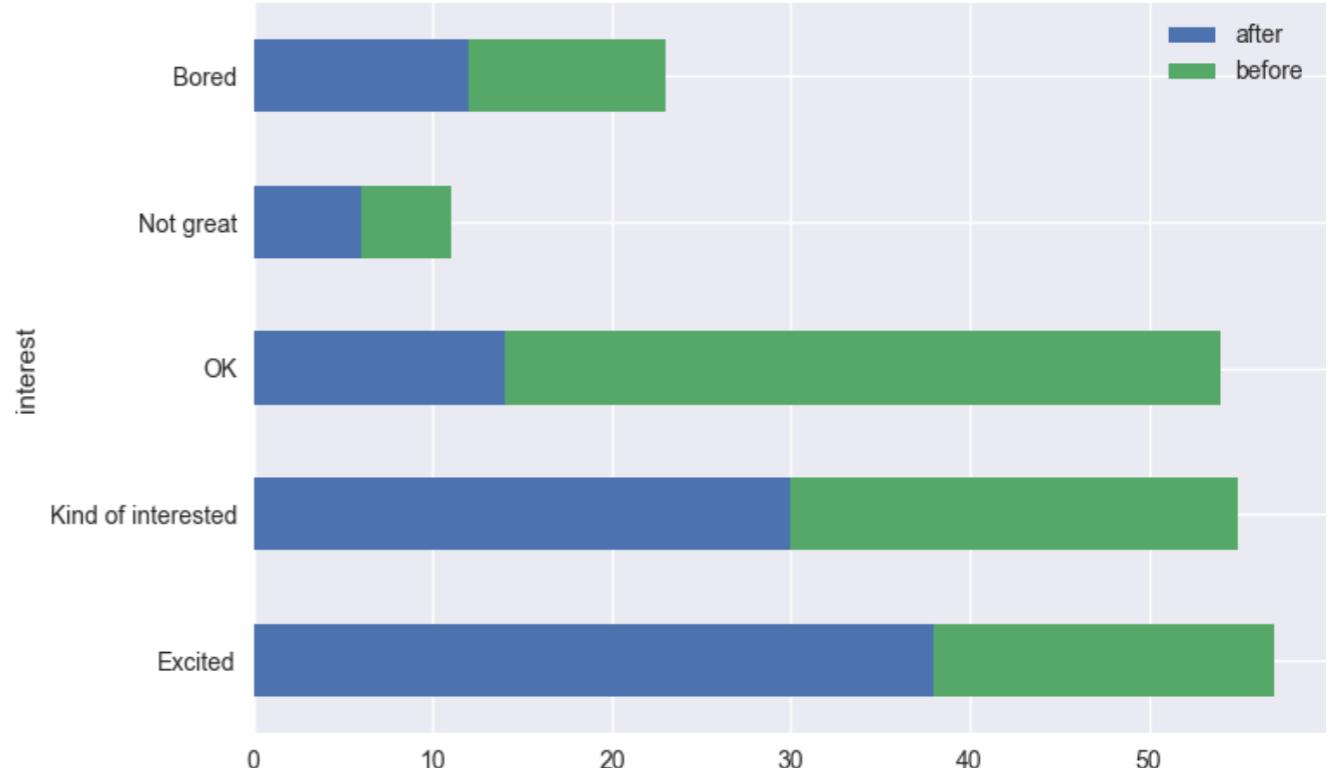
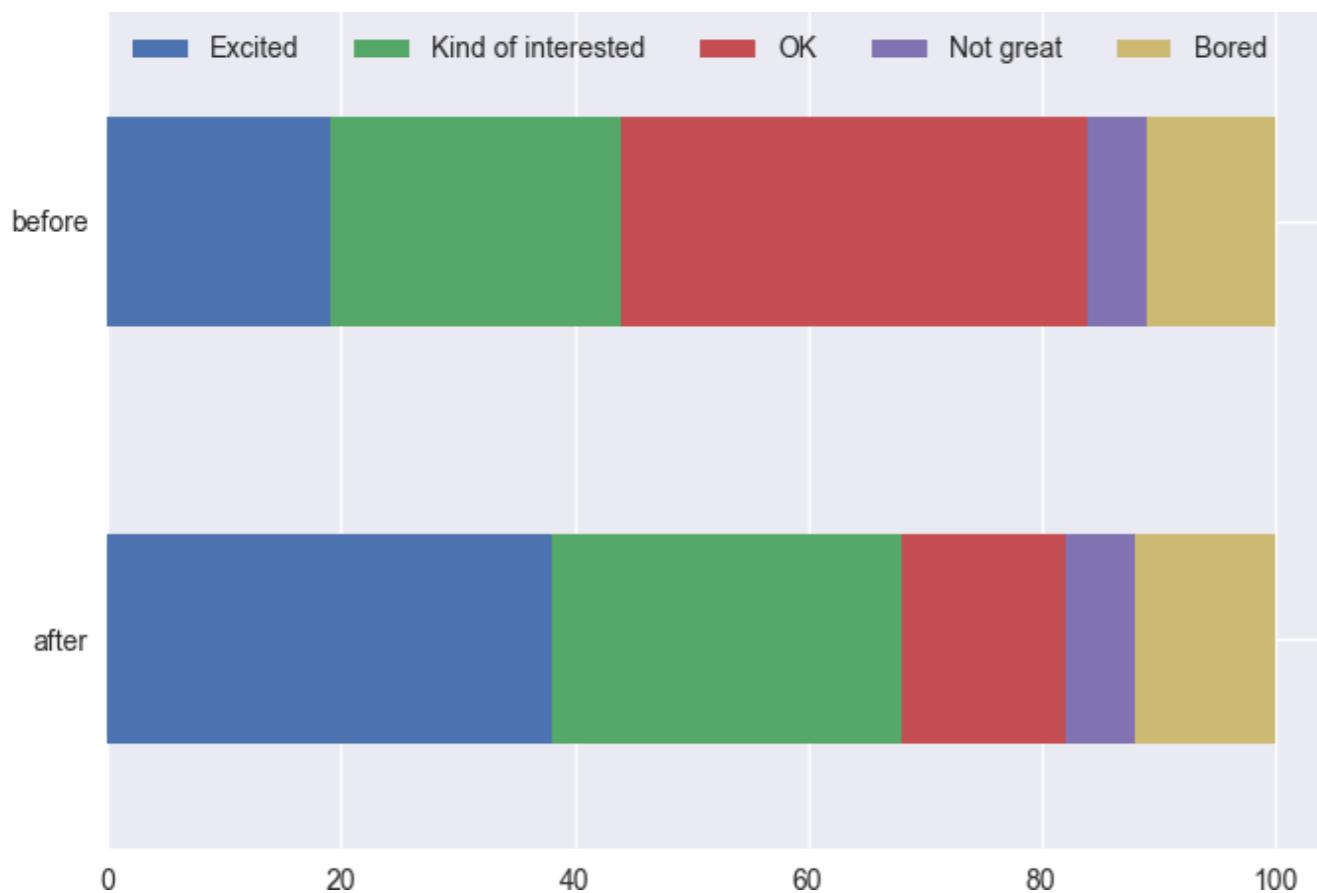
PIE CHARTS



STACKED BAR CHART

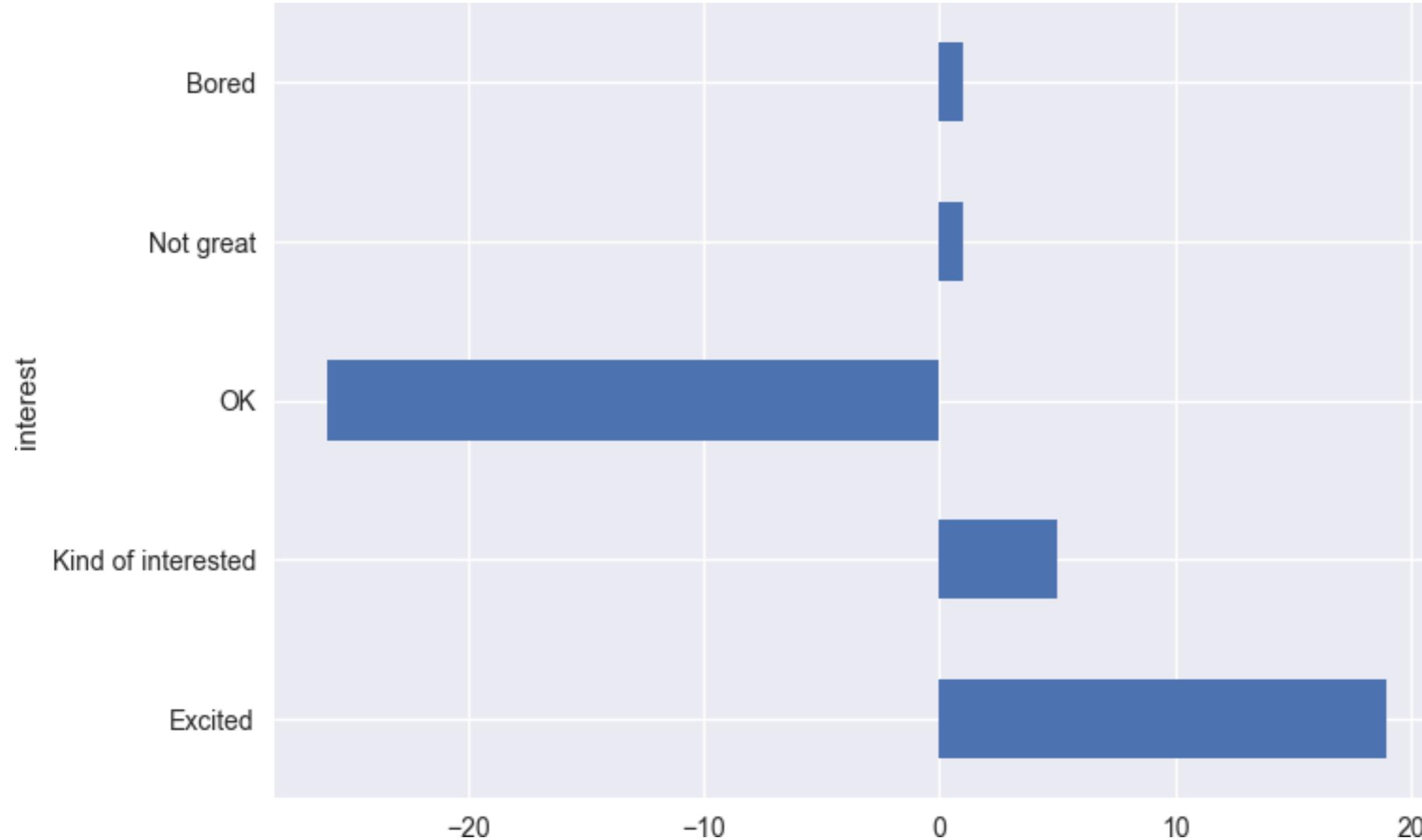


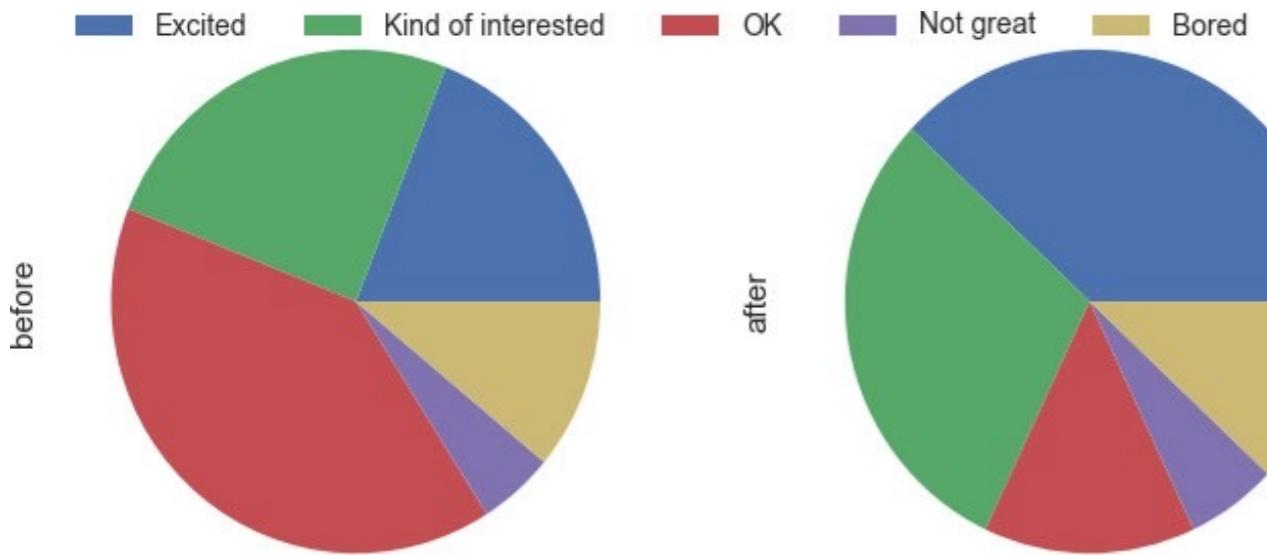
Stacked bar, not very useful



Data Transposed Bar Chart

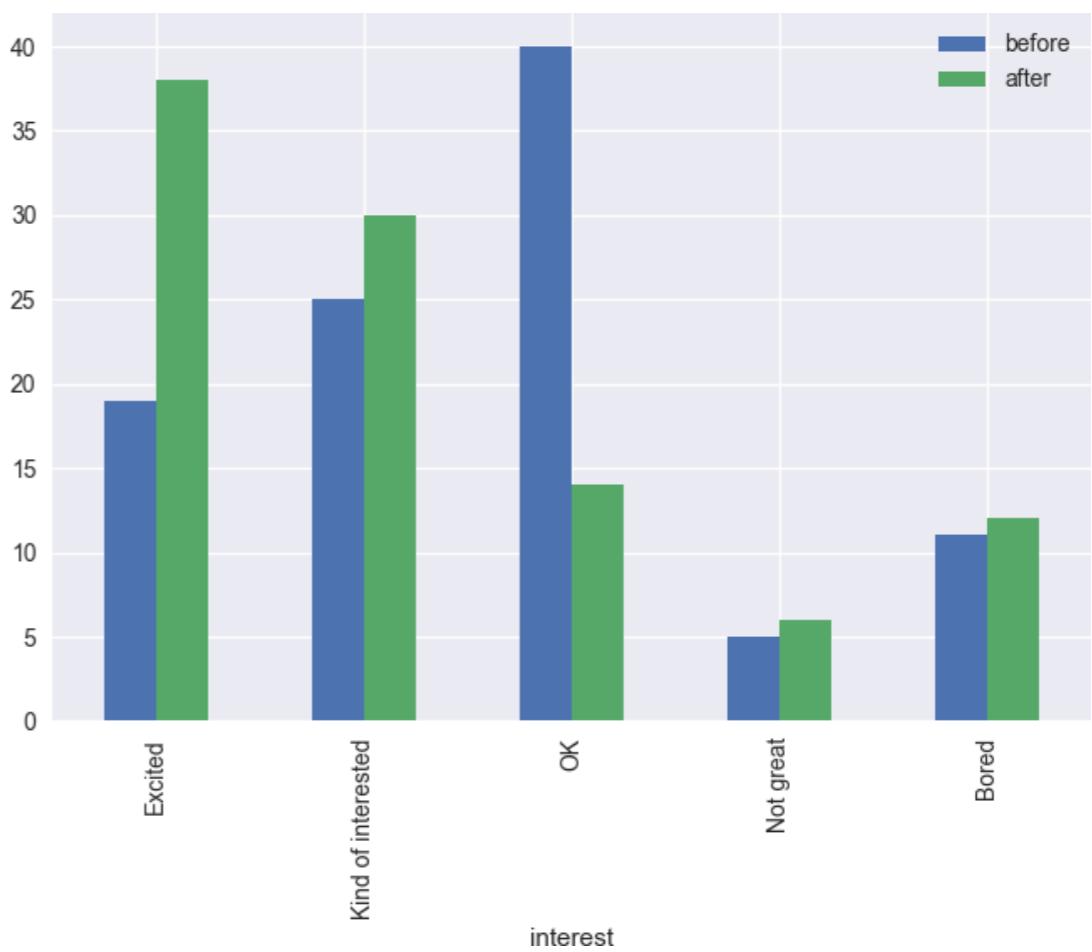
DIFFERENCE BAR CHART



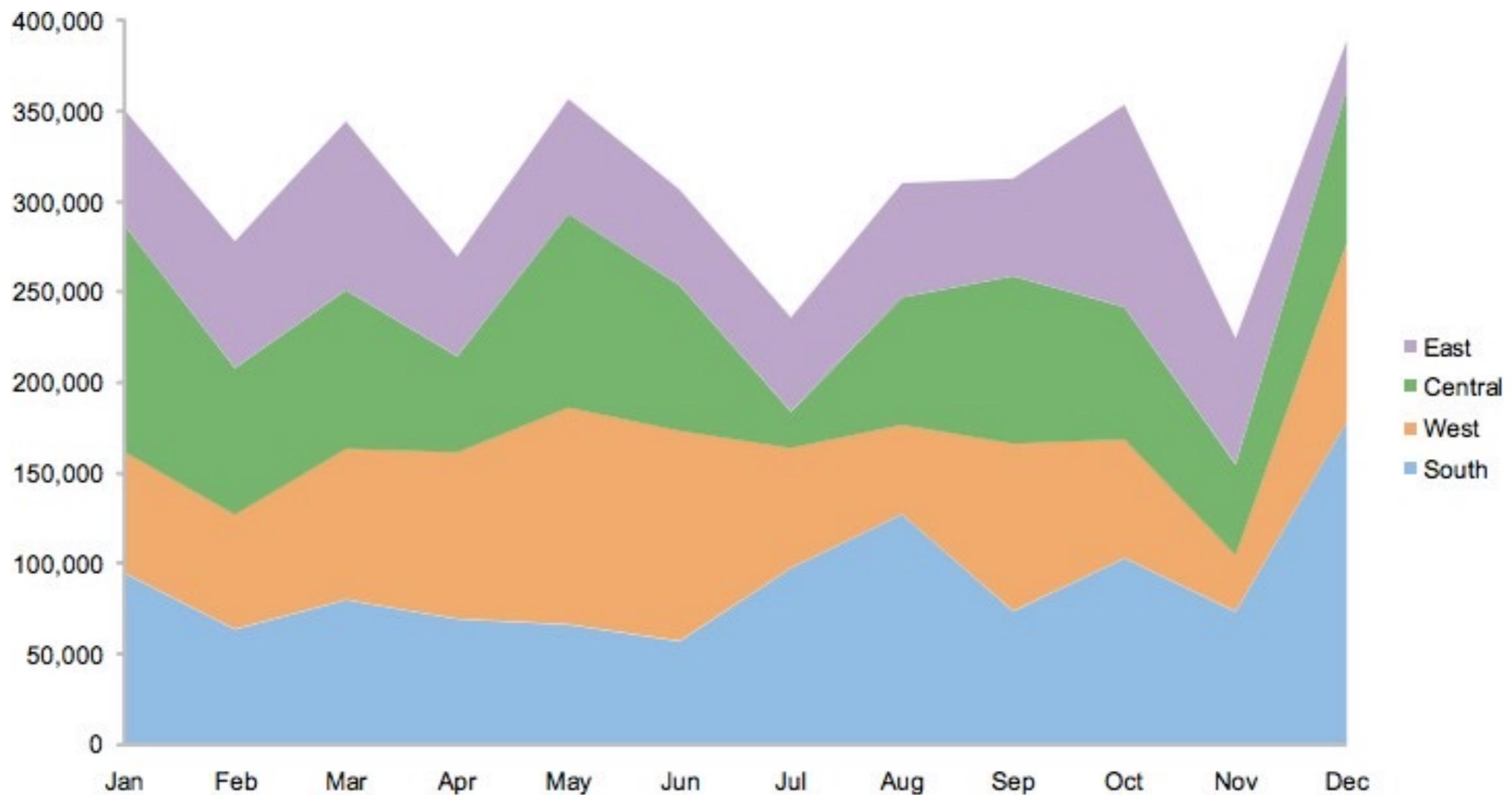


Pie

Side by side bar

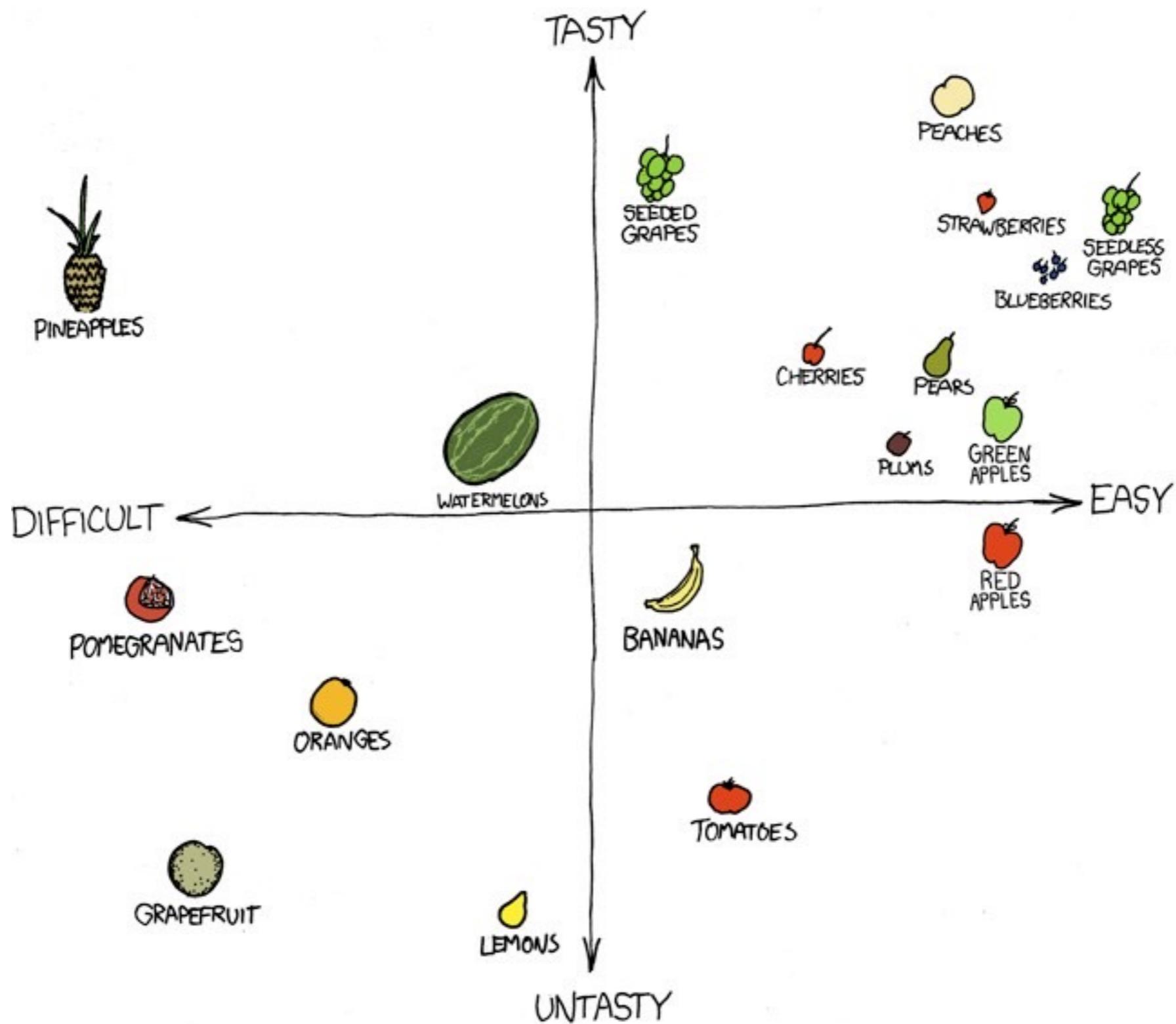


STACKED AREA CHART

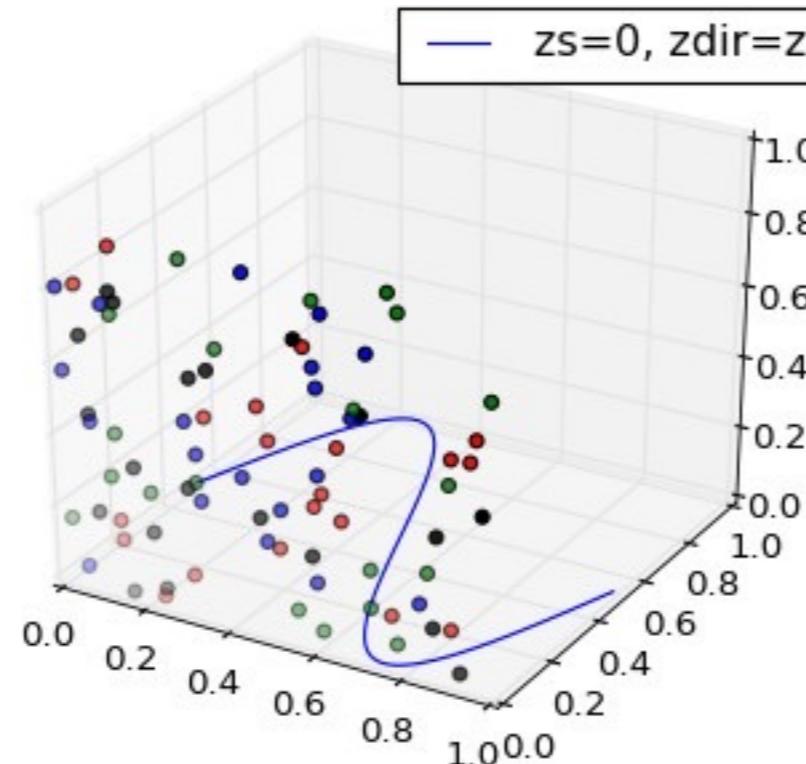
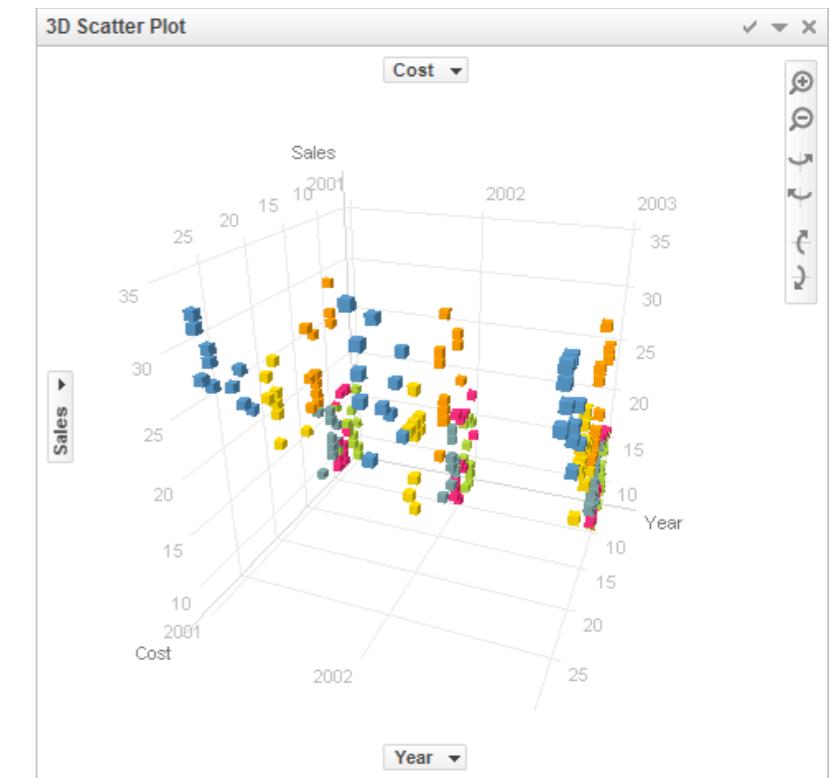
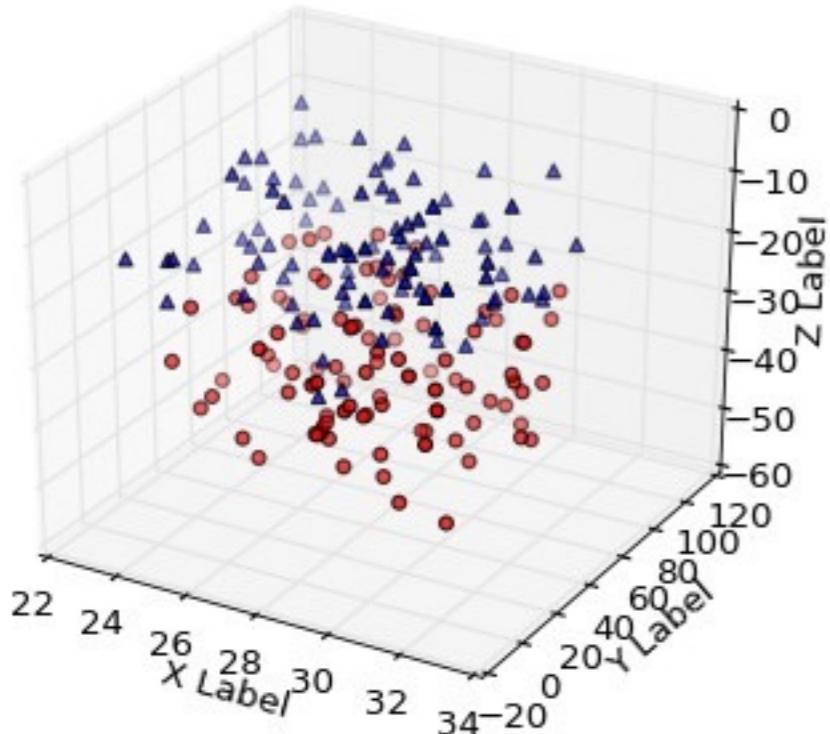


CORRELATIONS

SCATTERPLOTS

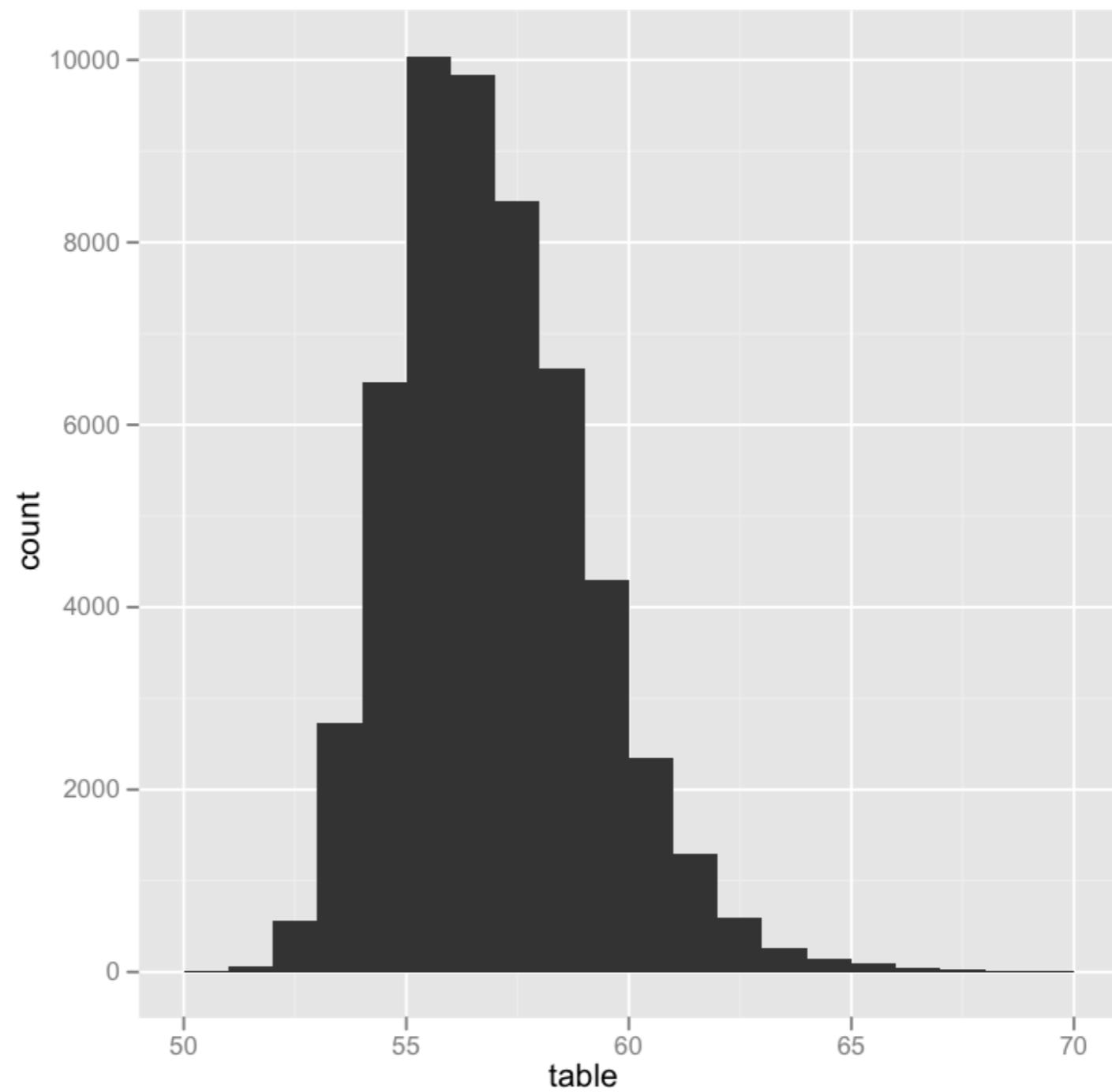


DON'T!

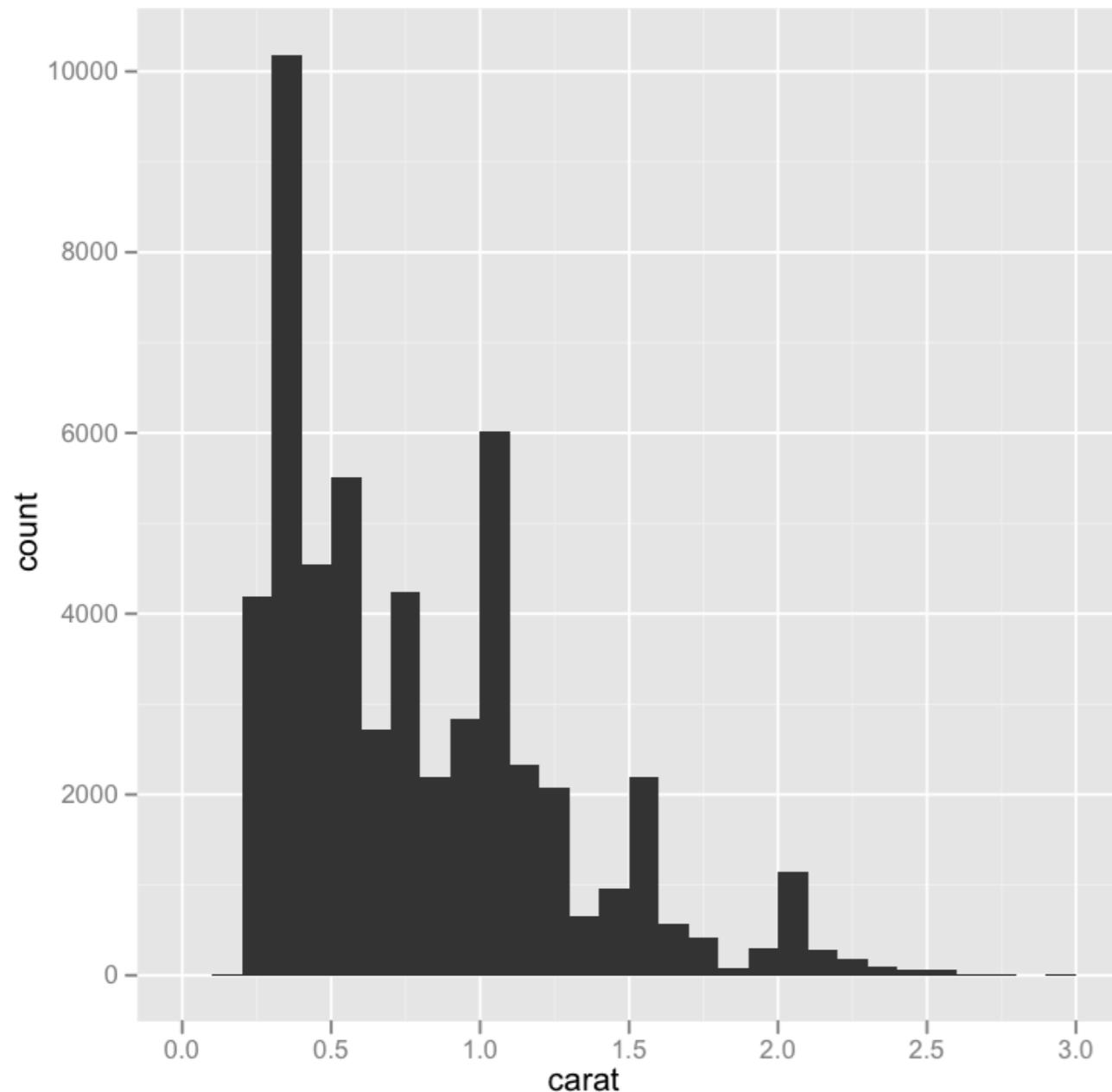


DISTRIBUTIONS

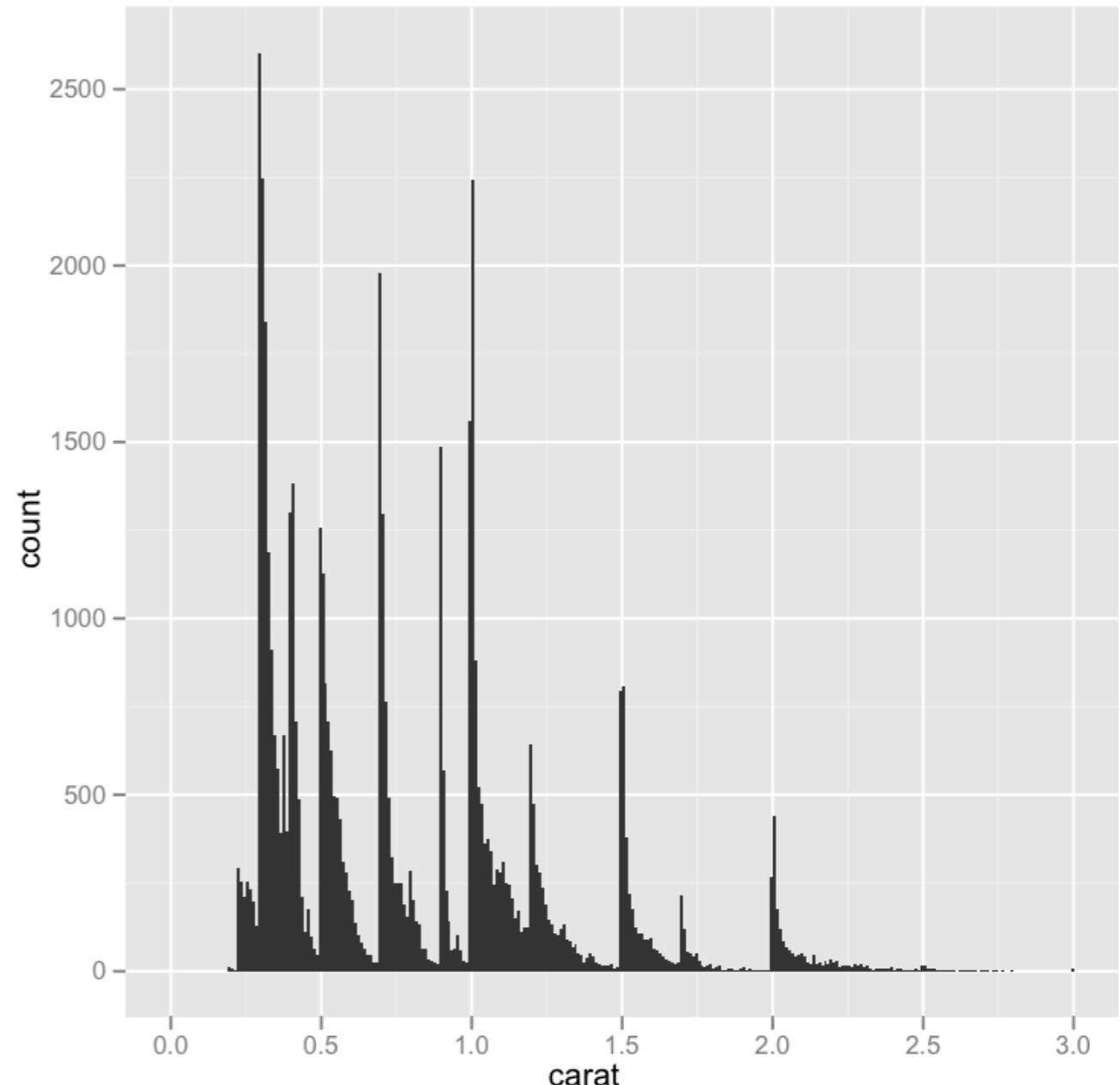
HISTOGRAM



BIN WIDTH

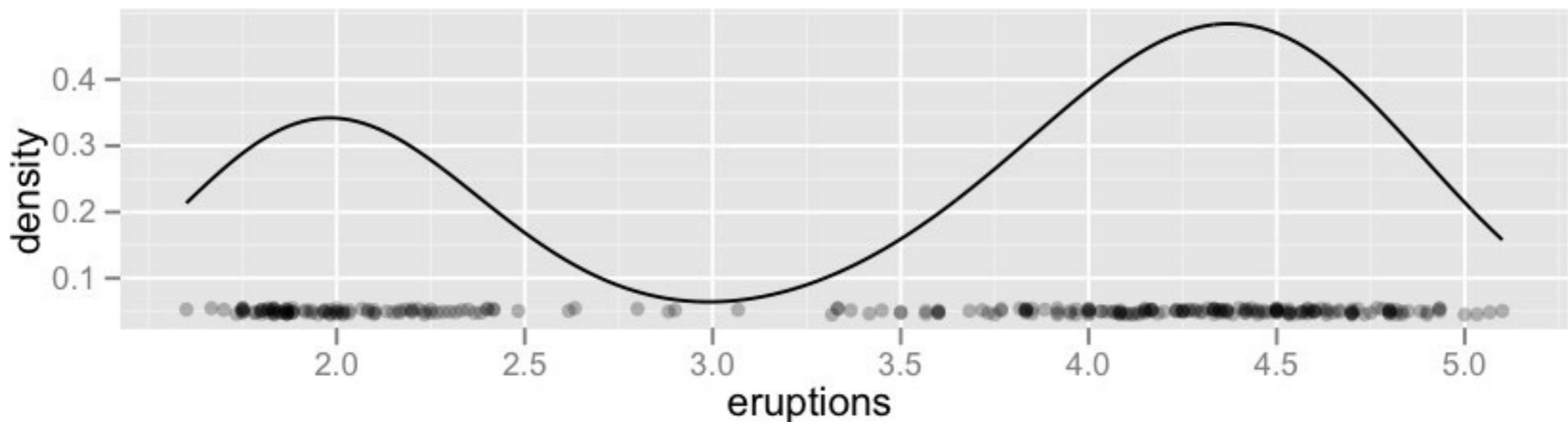


binwidth = 0.1

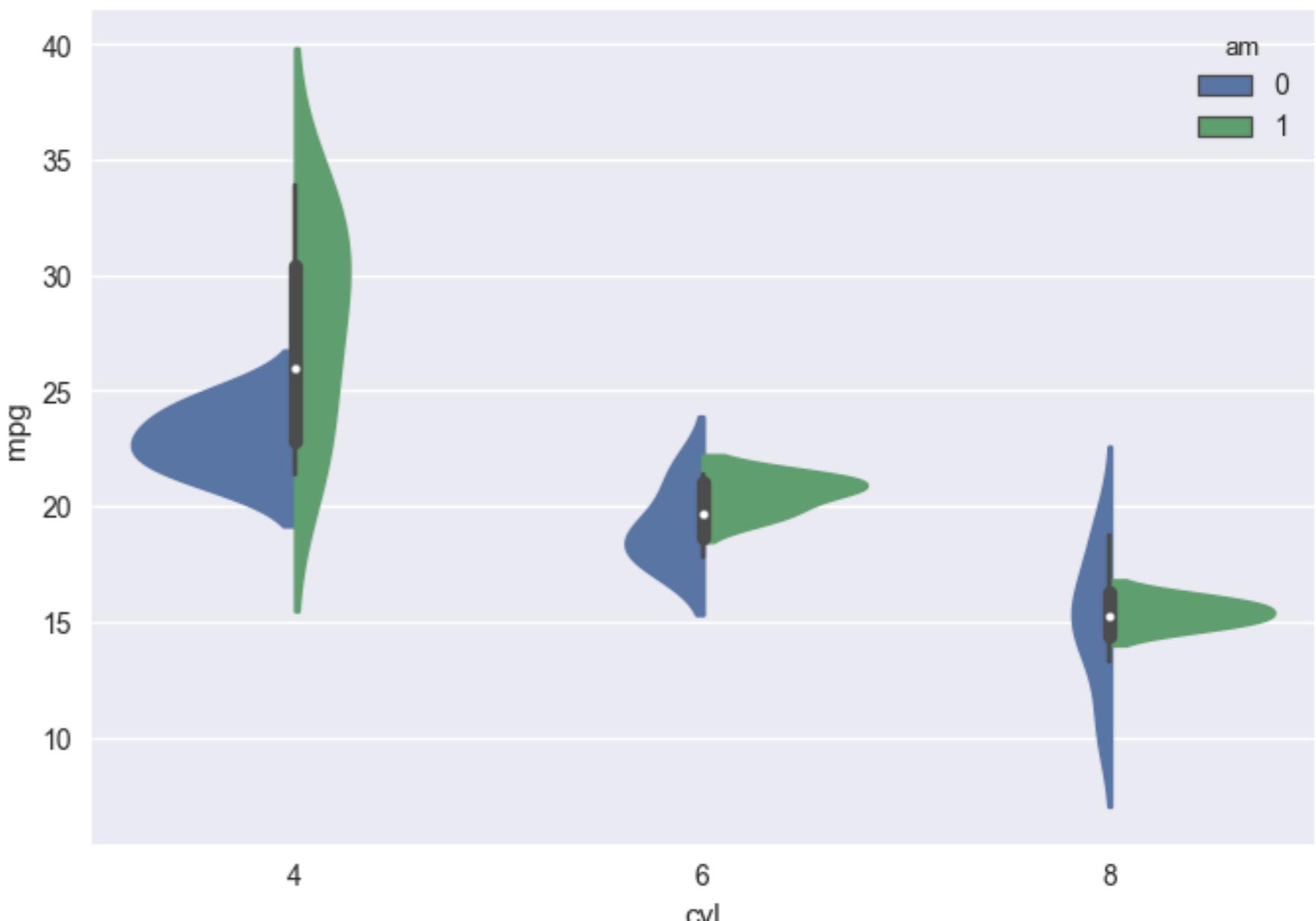
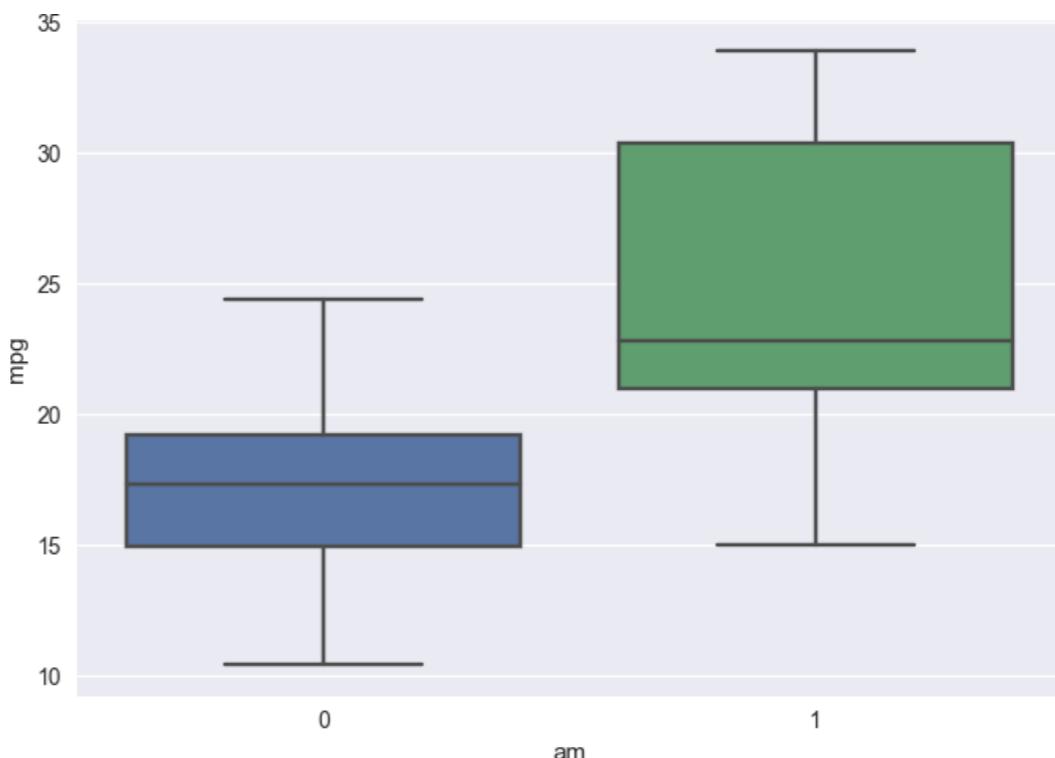
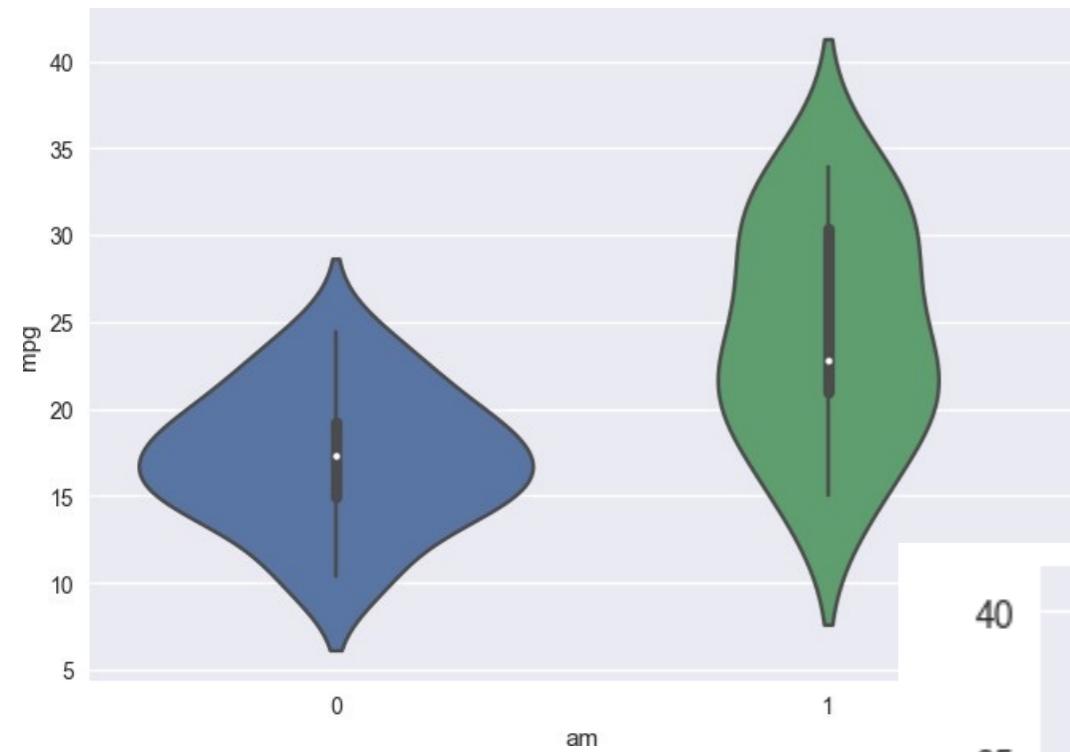


binwidth = 0.01

DENSITY PLOTS

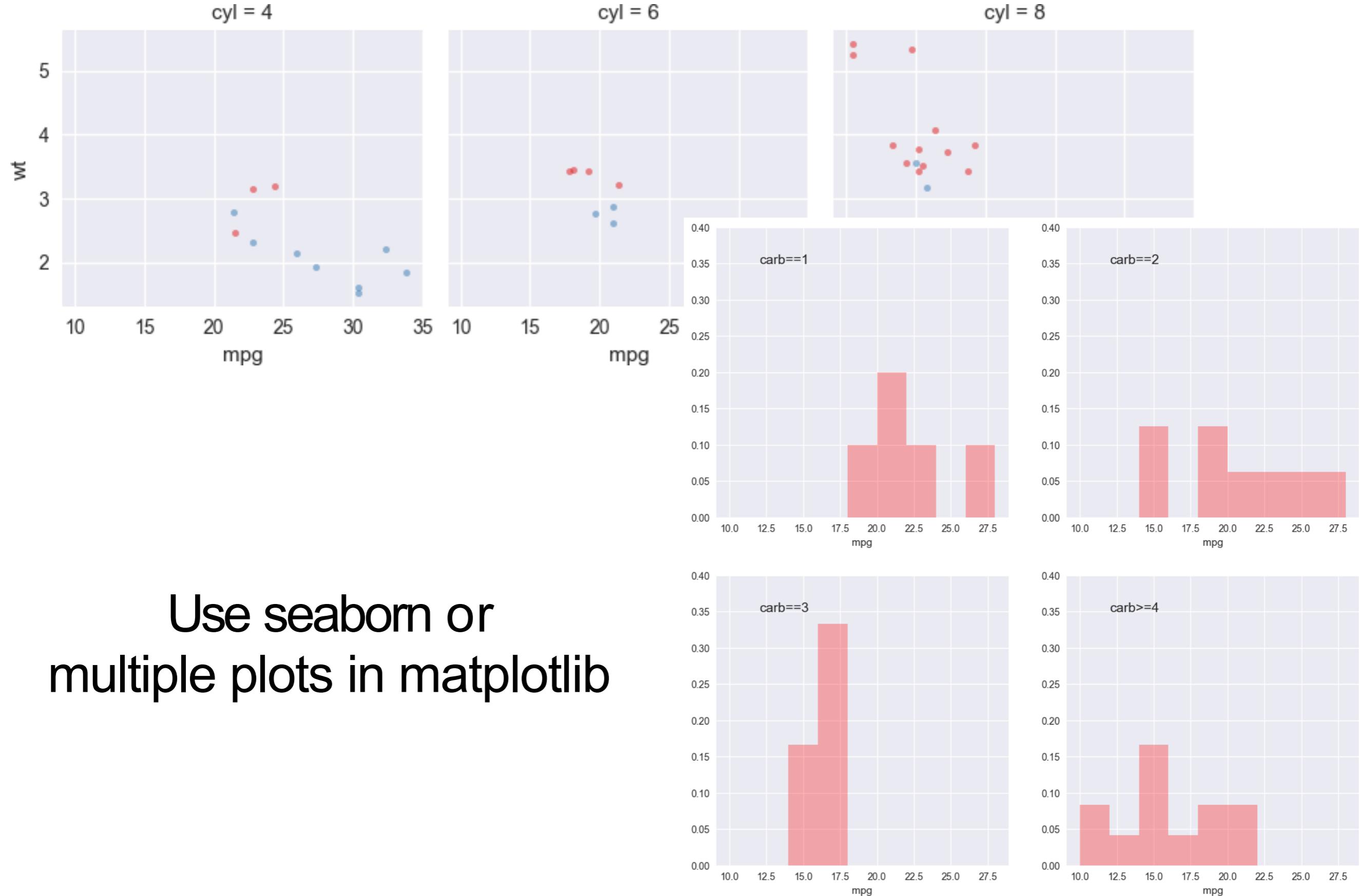


GROUP

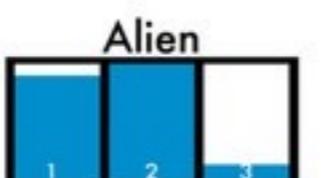
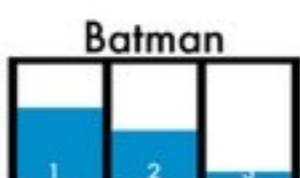
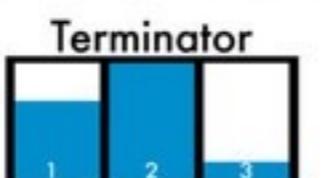
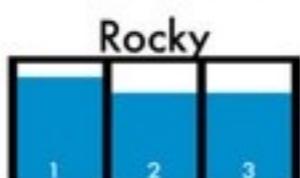
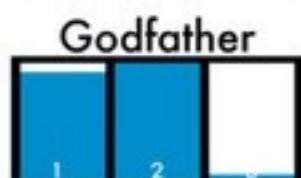
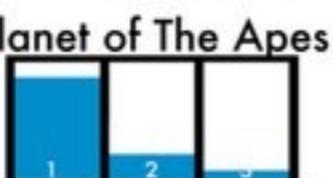
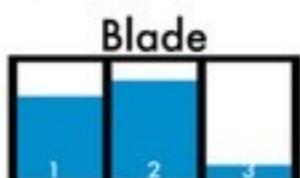
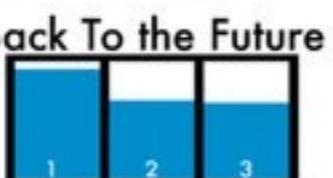
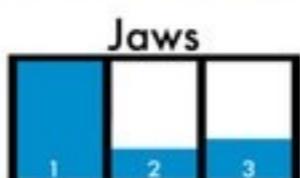
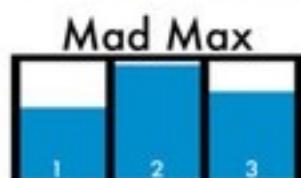
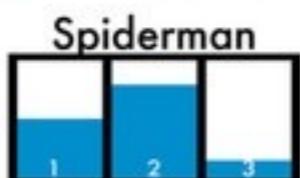
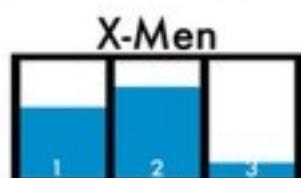
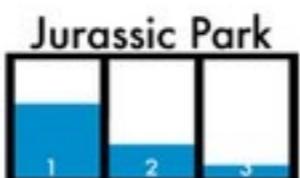
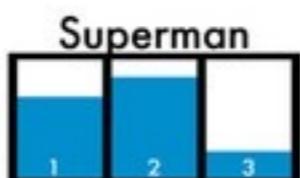
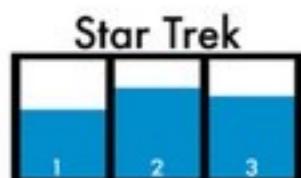
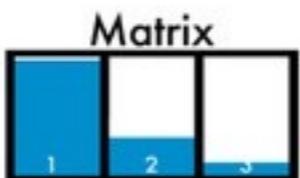


getting complex...

FACETING AND SMALL MULTIPLES



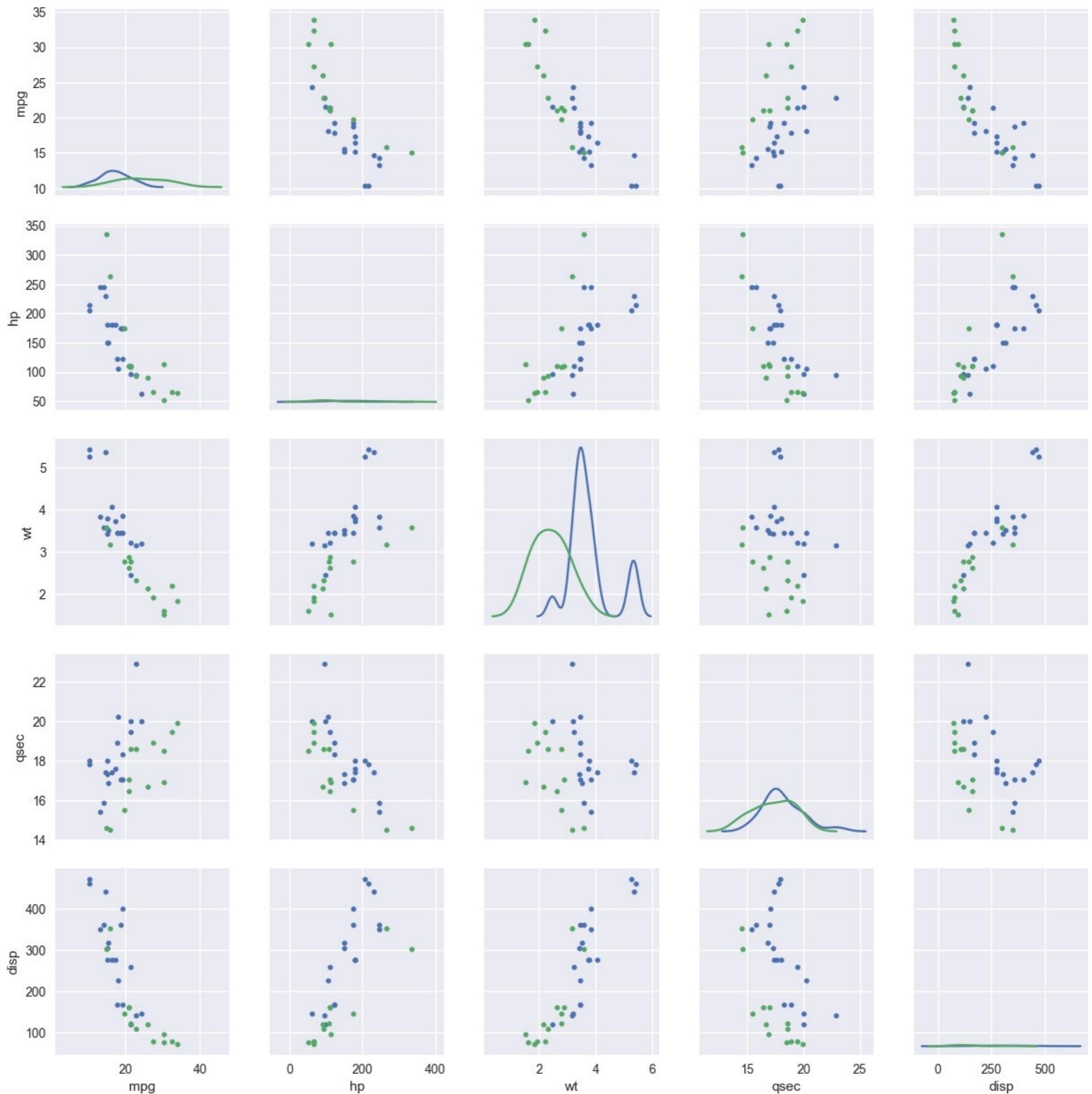
THE TRILOGY METER



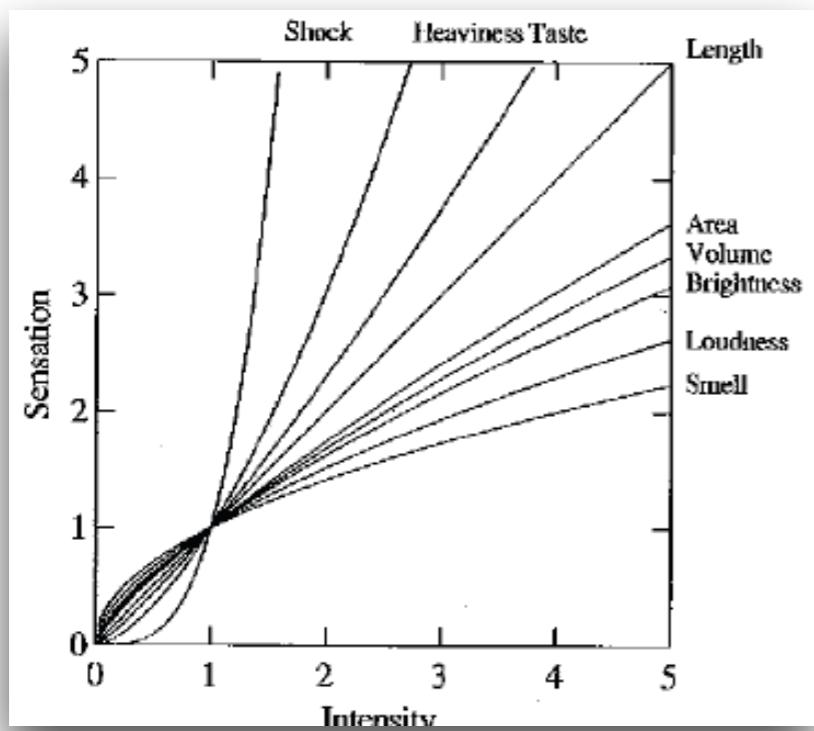
SMALL
MULTIPLES

SPLOM

The
scatterplot
matrix



PERCEPTUAL EFFECTIVENESS



Stephen's Power Law, 1961

	Nominal	Ordinal	Quantitative
Position	✓	✓	✓
Size	✓	✓	~
(Grey)Value	✓	✓	~
Texture	✓	~	✗
Color	✓	✗	✗
Orientation	✓	✗	✗
Shape	✓	✗	✗

✓ = Good

~ = OK

✗ = Bad

J. Bertin, 1967

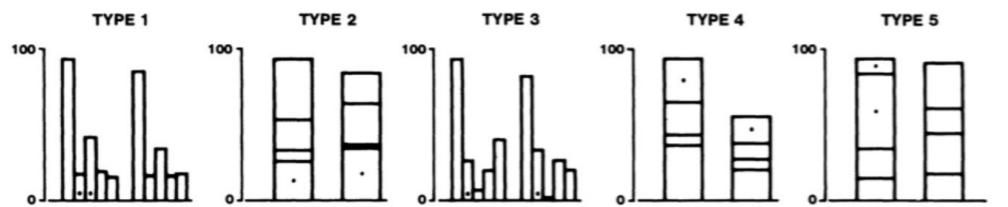


Figure 4. Graphs from position-length experiment.

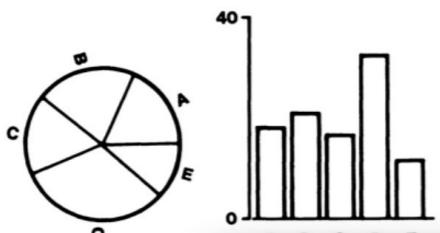
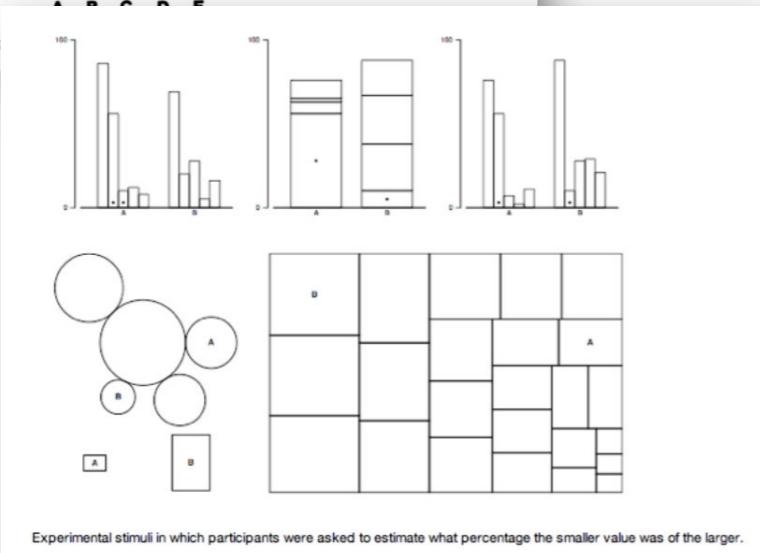
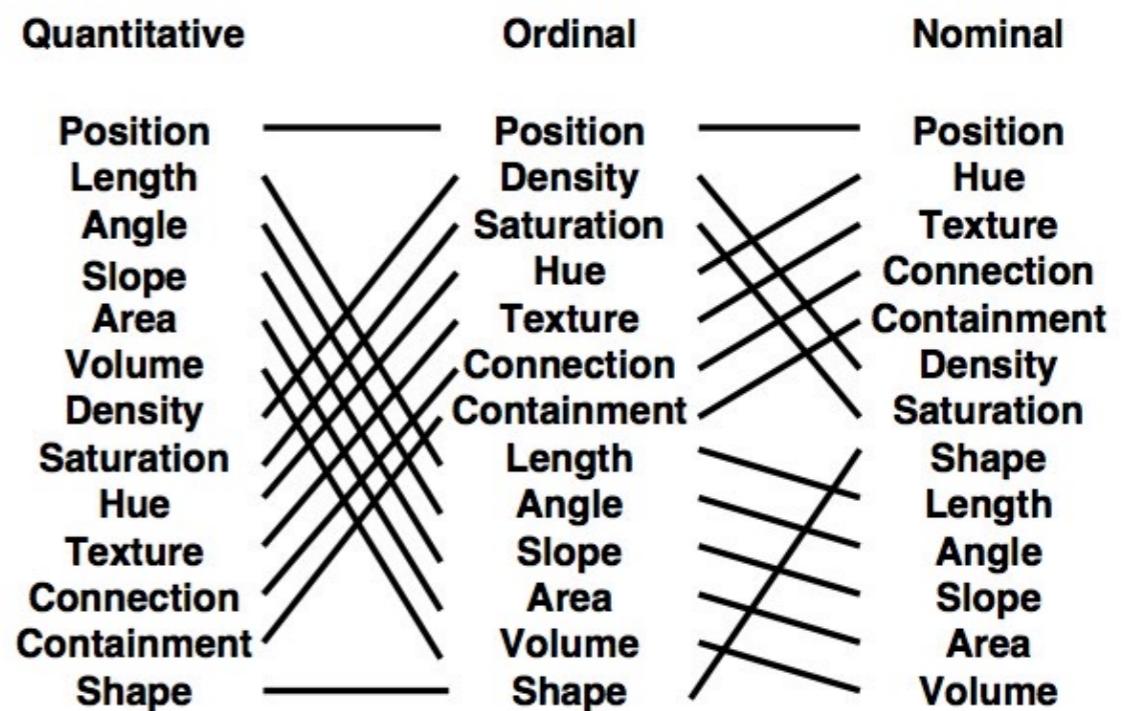


Figure 3. Graphs from position-length experiment.

Cleveland / McGill, 1984

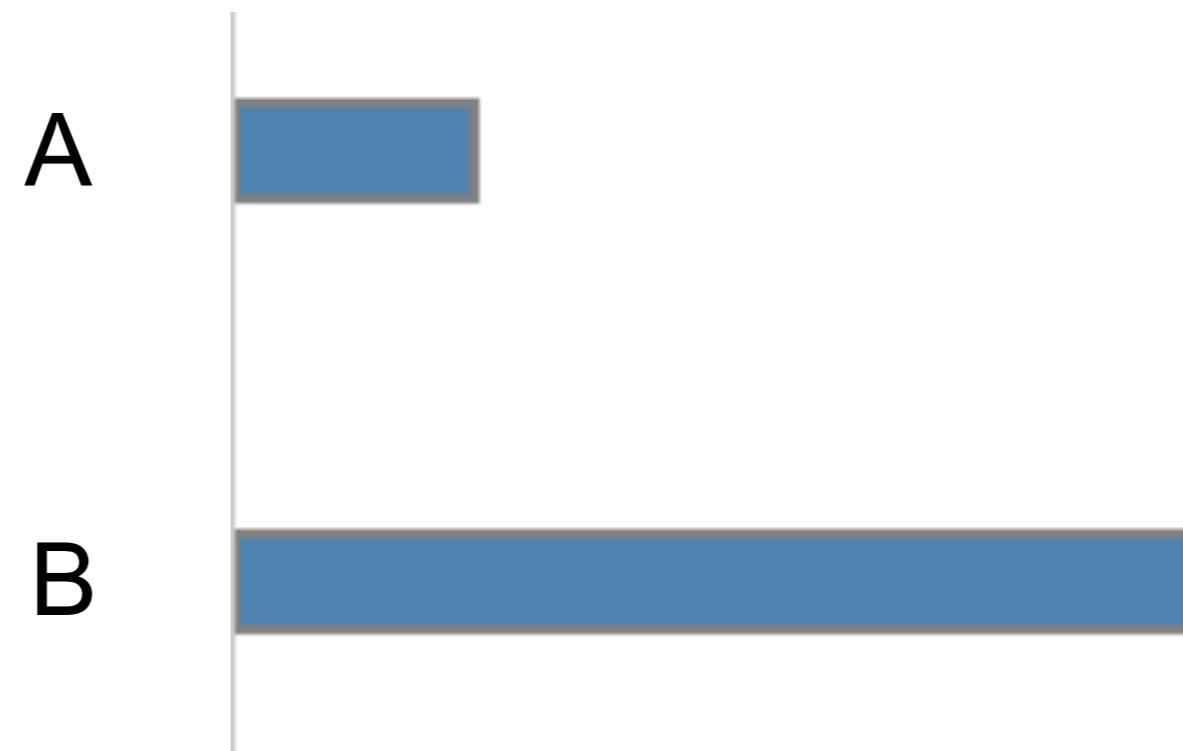


Heer / Bostock, 2010

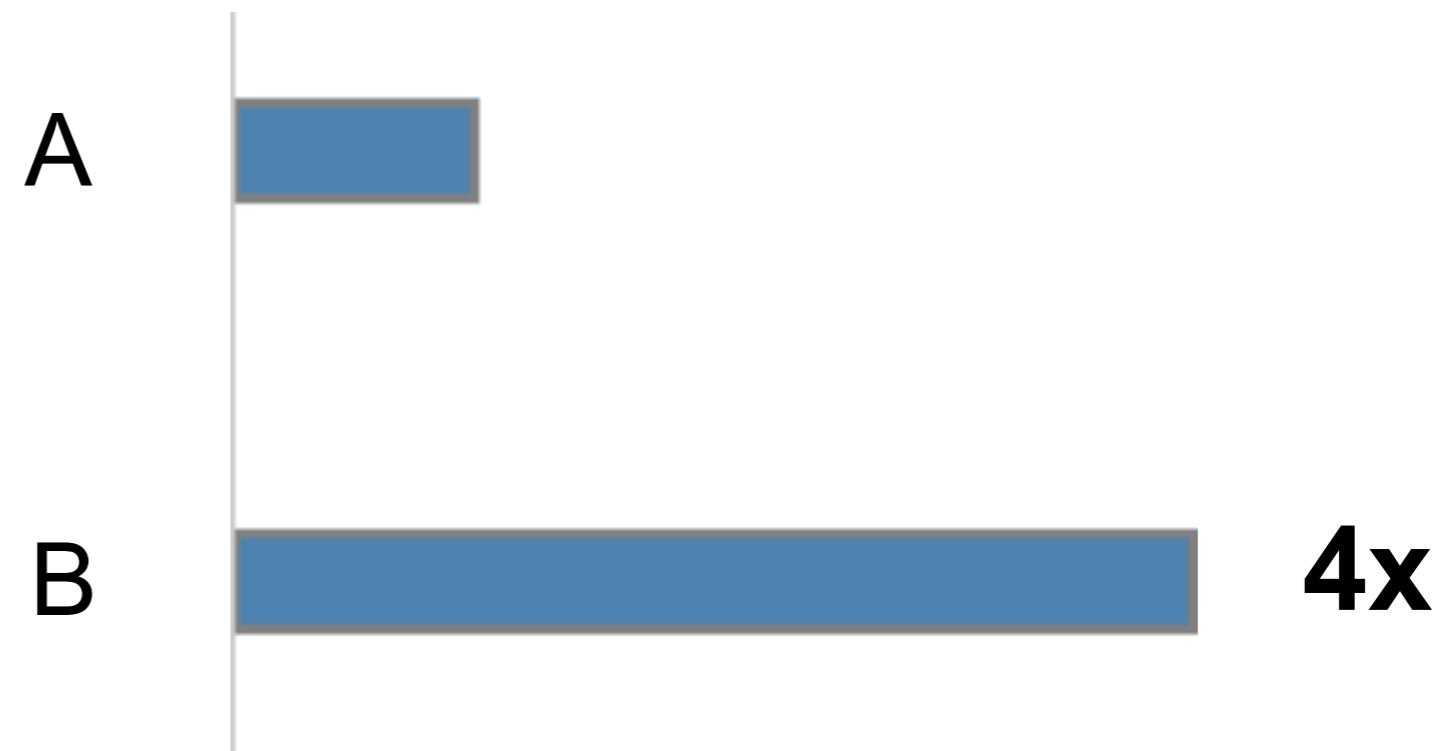


J. Mackinlay, 1986

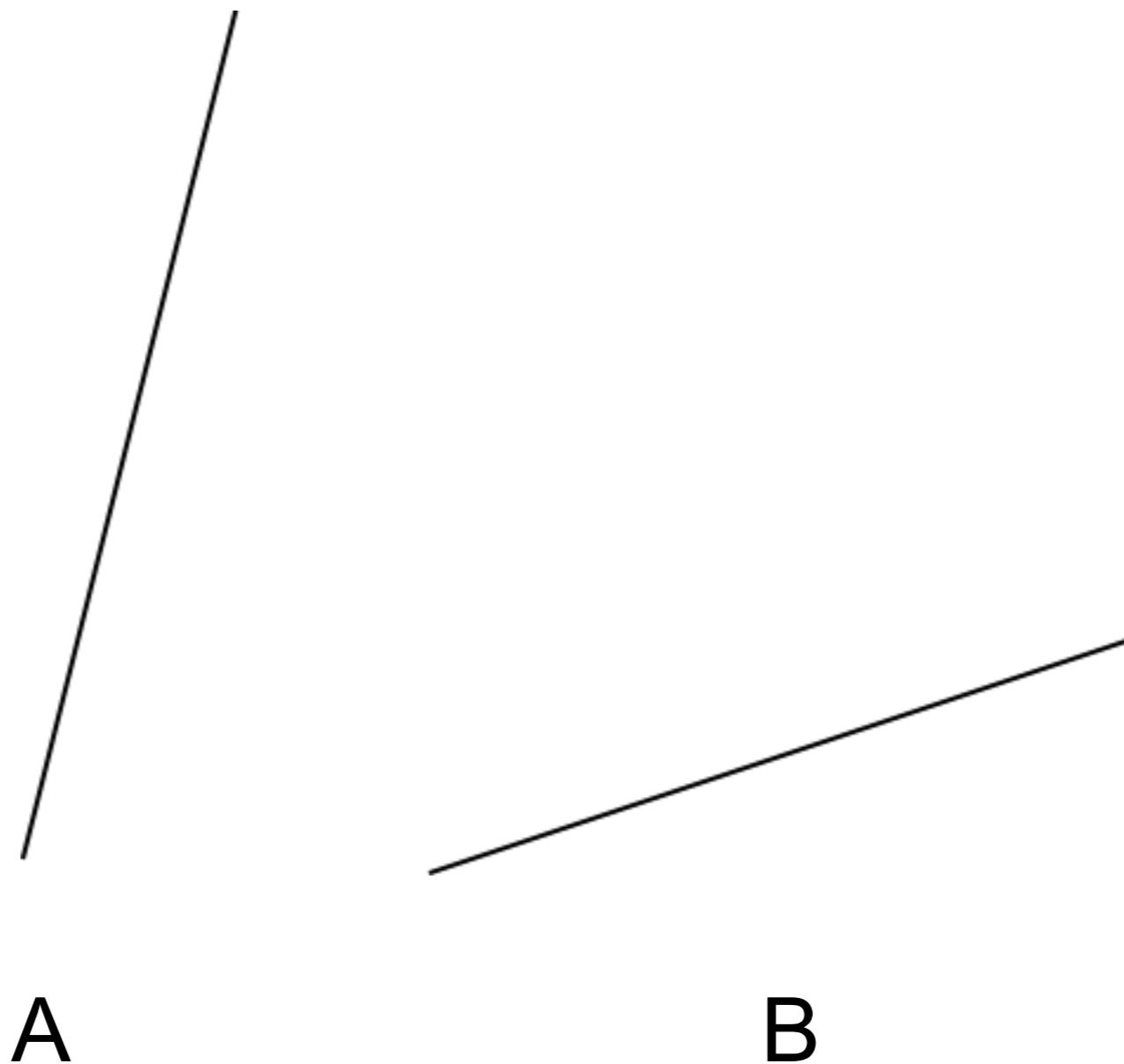
HOW MUCH LONGER?



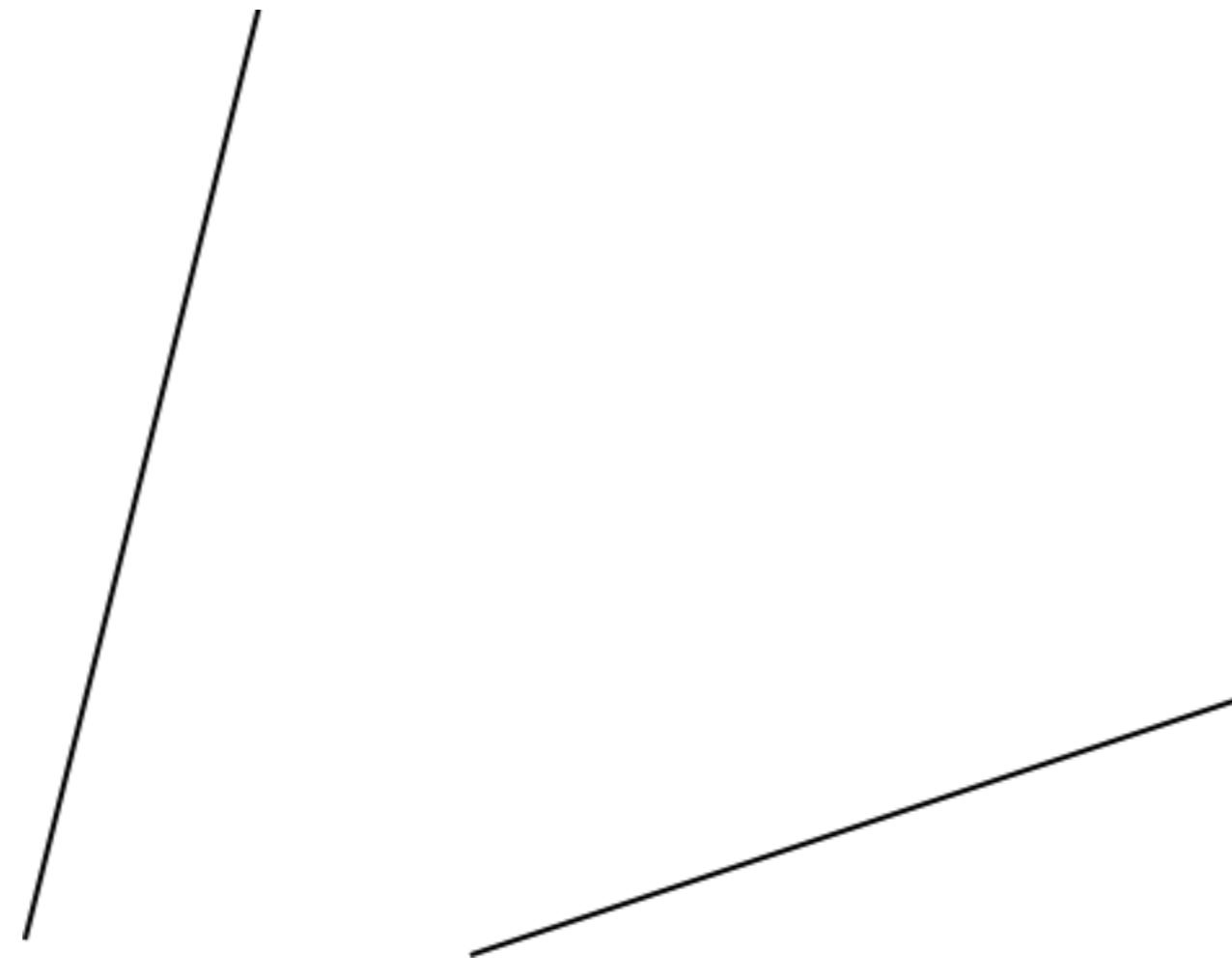
HOW MUCH LONGER?



HOW MUCH STEEPER SLOPE?



HOW MUCH STEEPER SLOPE?

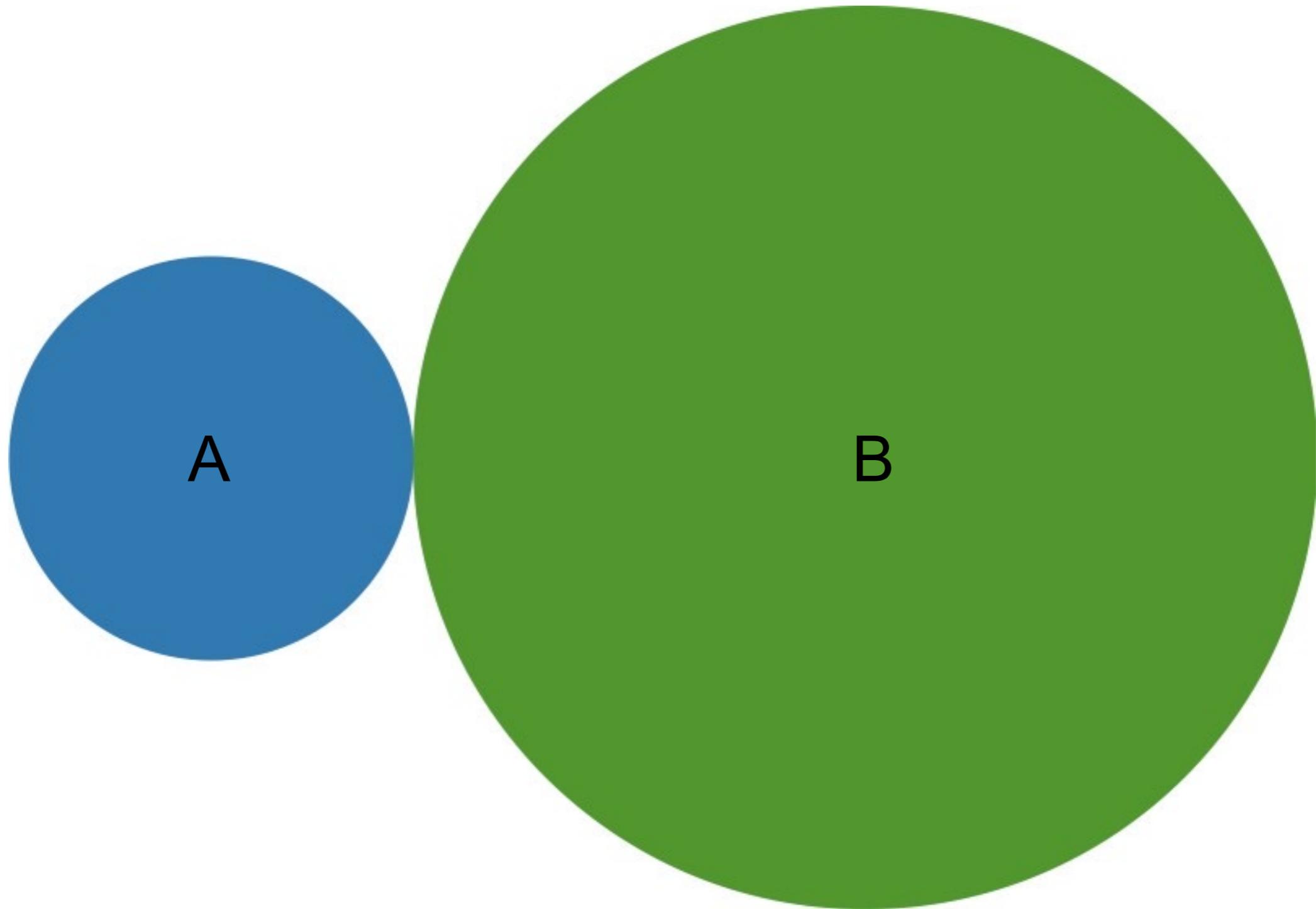


A

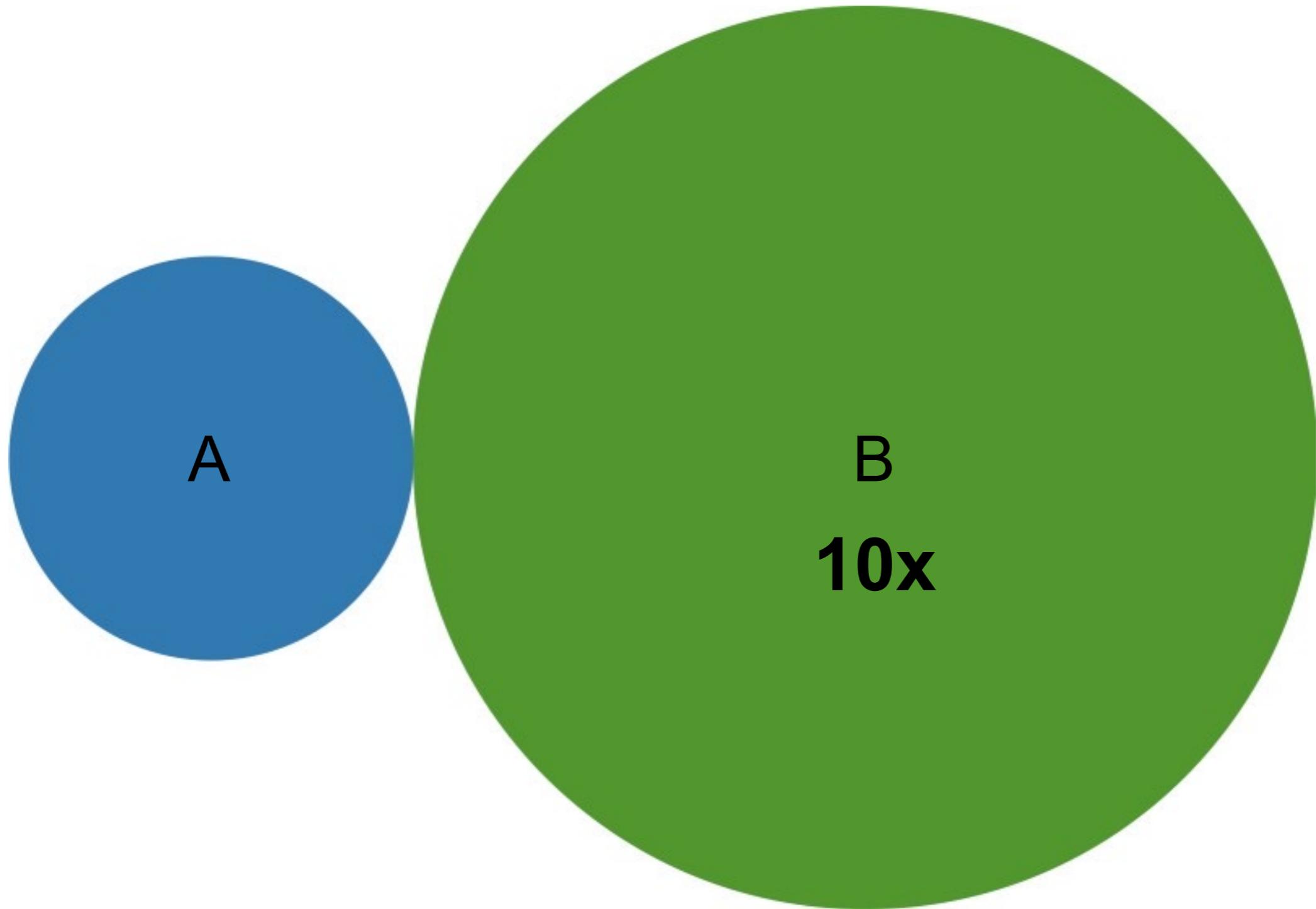
B

4x

HOW MUCH LARGER AREA?



HOW MUCH LARGER AREA?



HOW MUCH DARKER?



A



B

HOW MUCH DARKER?



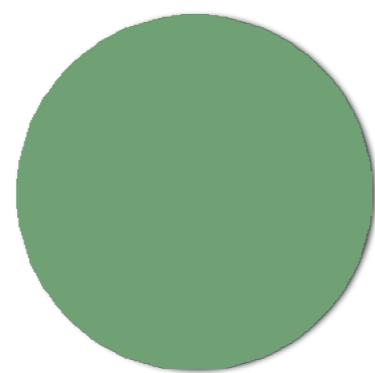
A



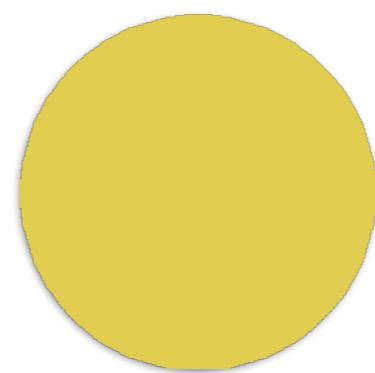
B

2x

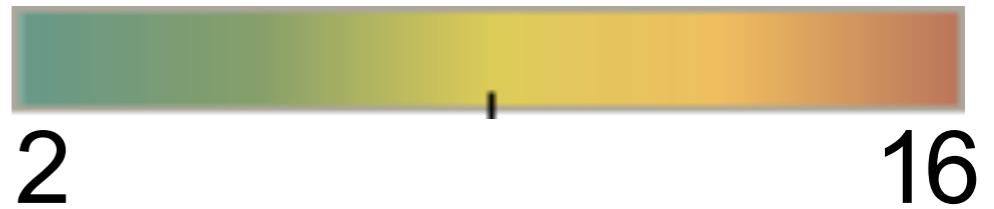
HOW MUCH BIGGER VALUE?



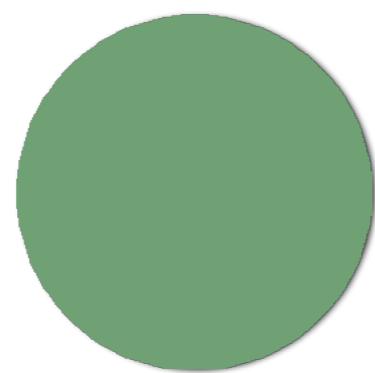
A



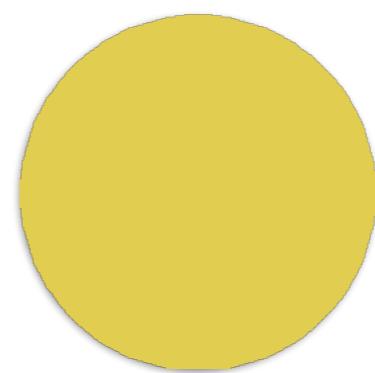
B



HOW MUCH BIGGER VALUE?

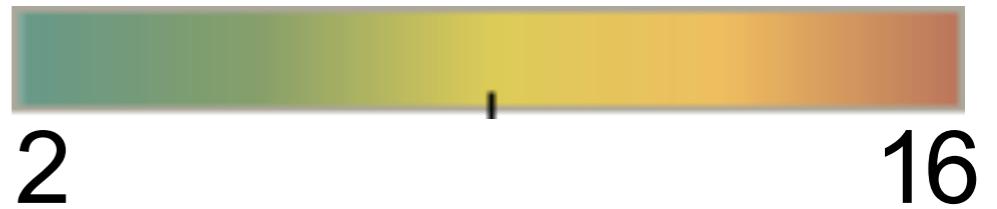


A



B

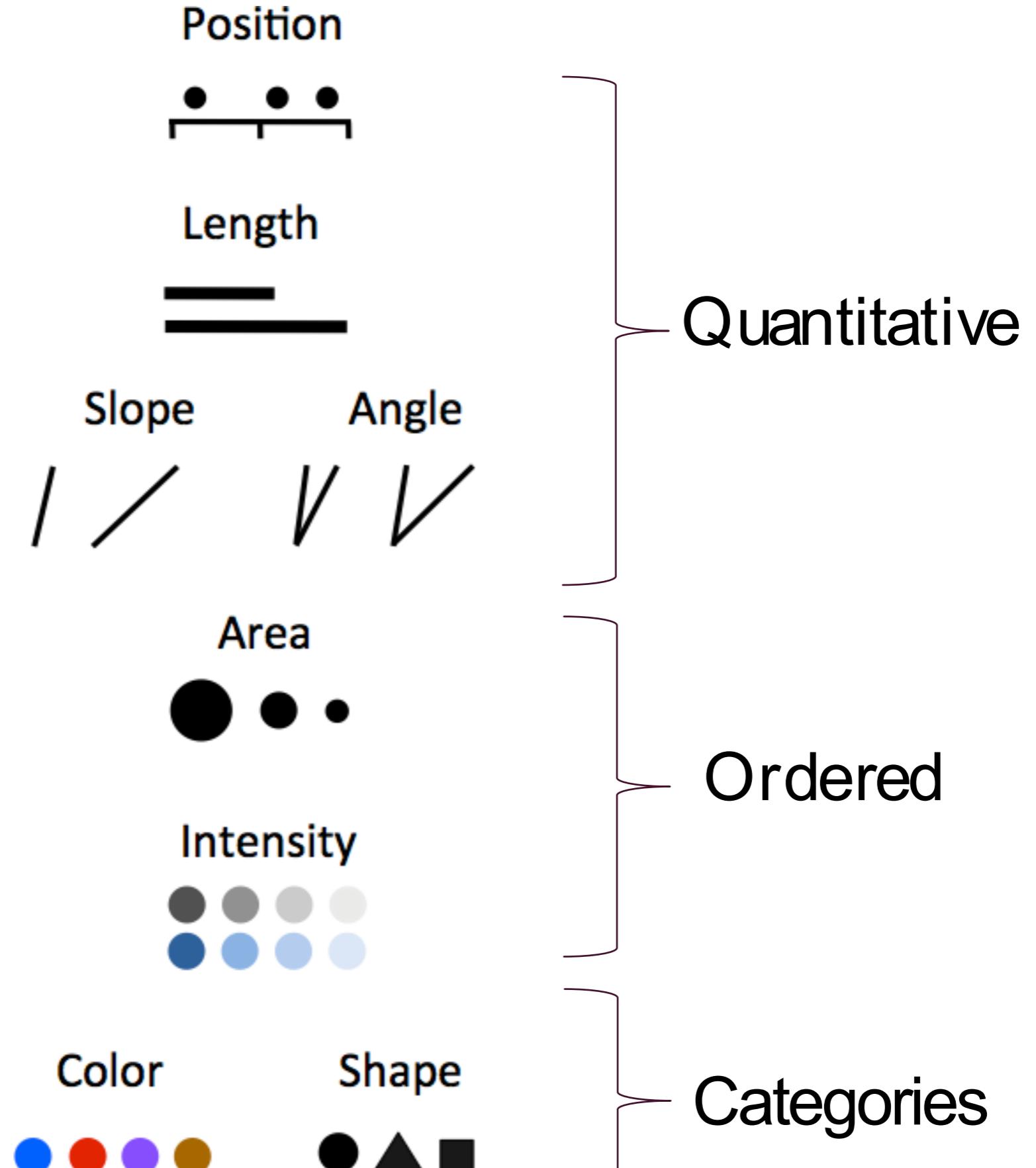
4x



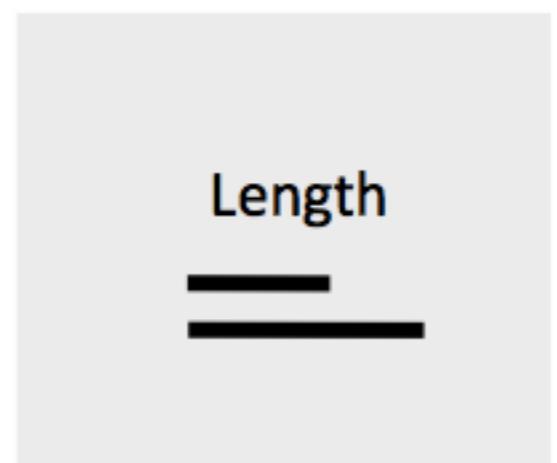
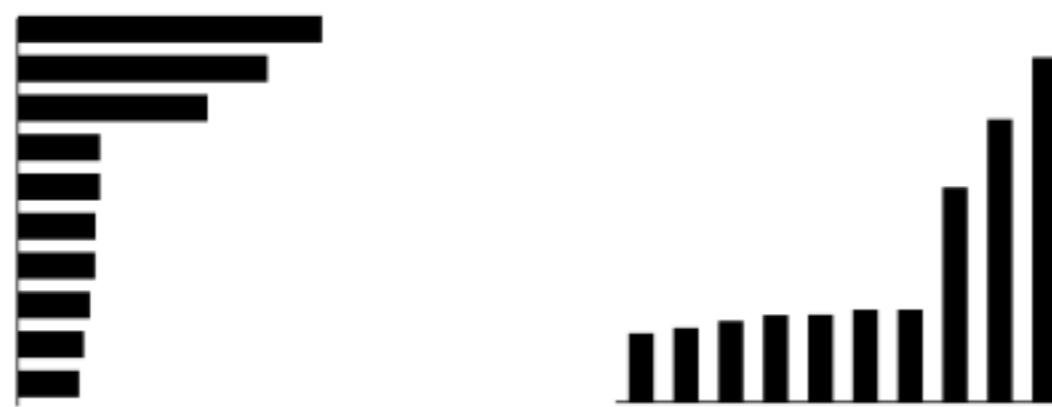
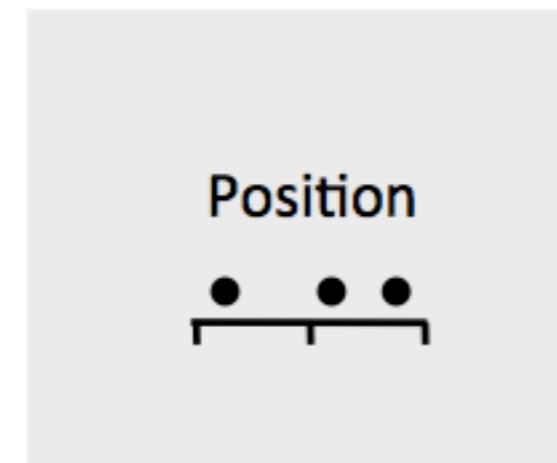
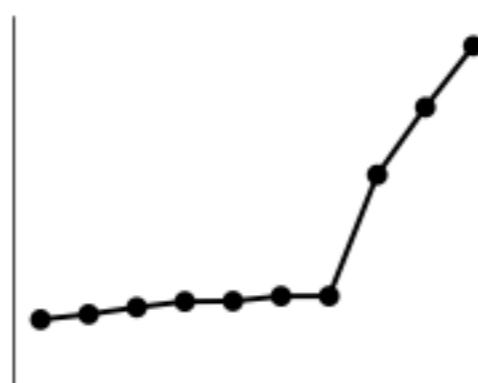
Most
Efficient



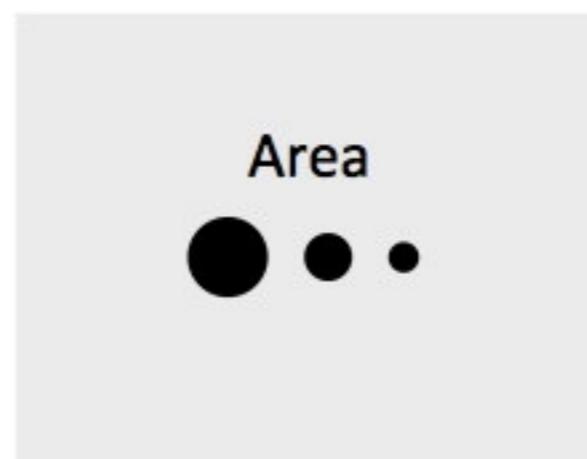
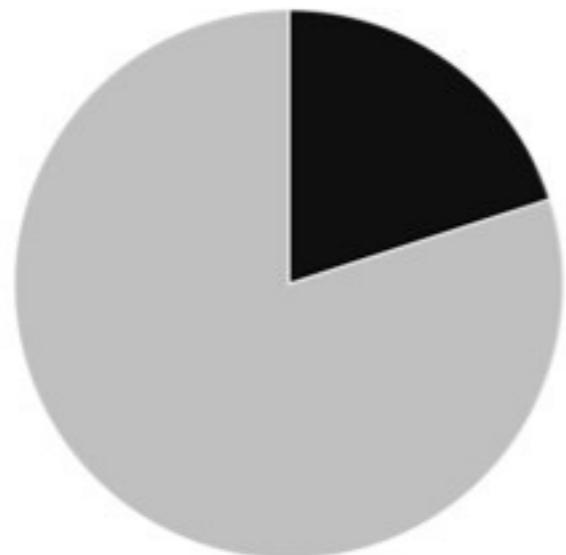
Least
Efficient



MOST EFFECTIVE

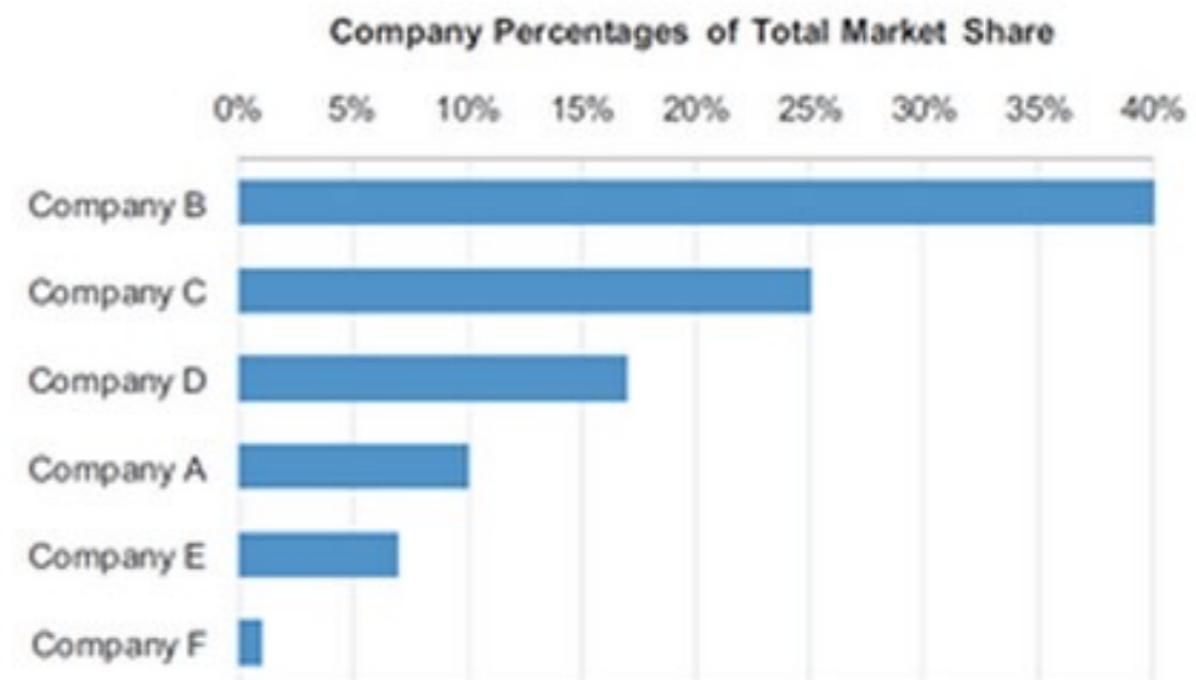
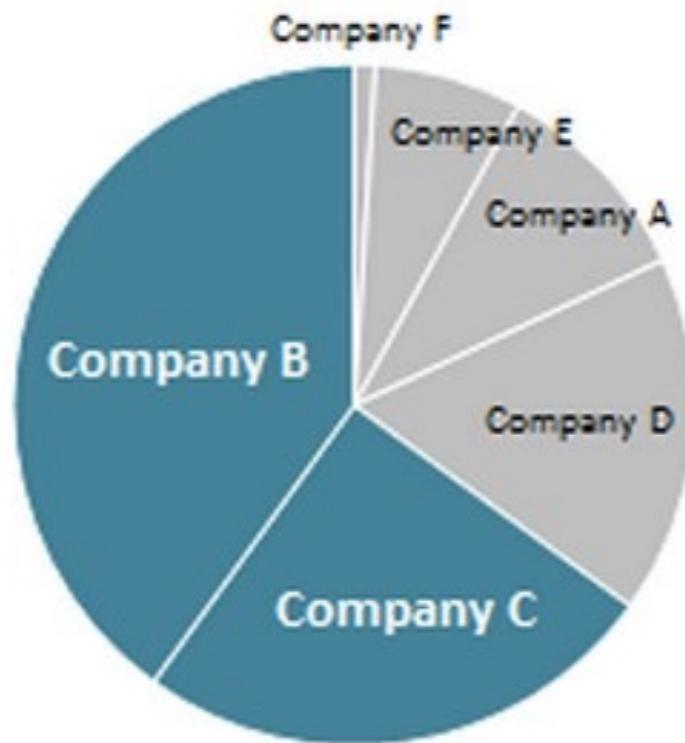


LESS EFFECTIVE

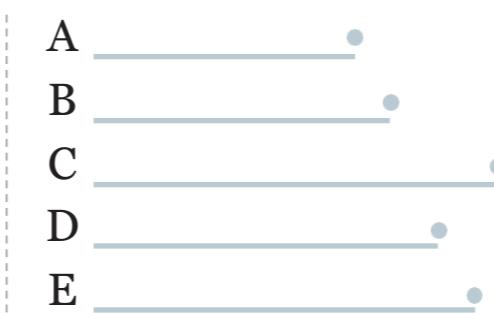
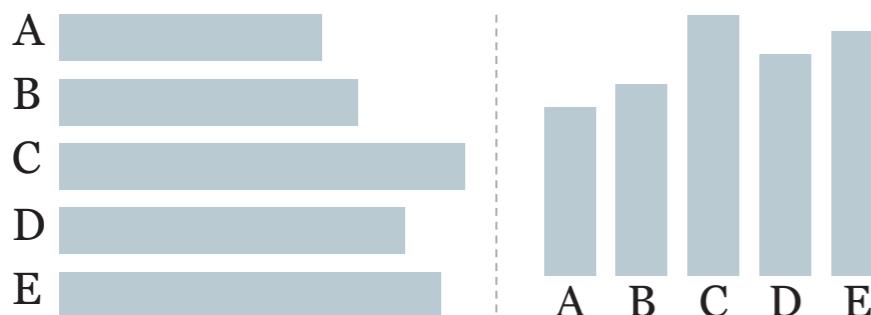


PIE vs BAR CHARTS

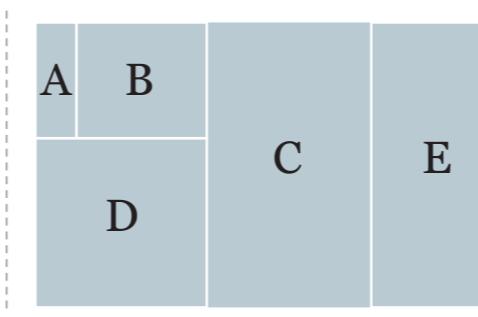
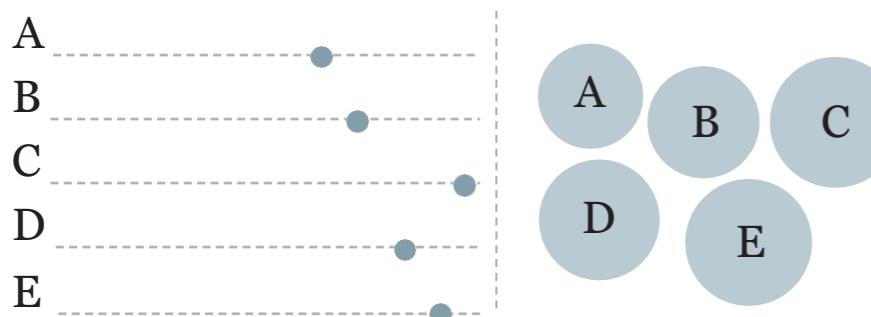
65% of the market is controlled by companies B and C



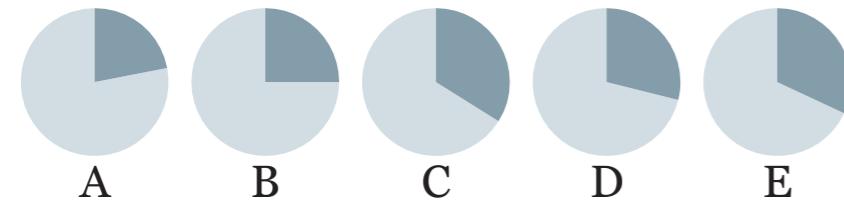
LENGTH OR HEIGHT



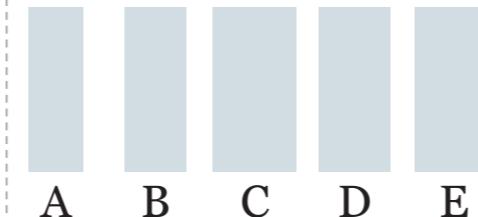
Position



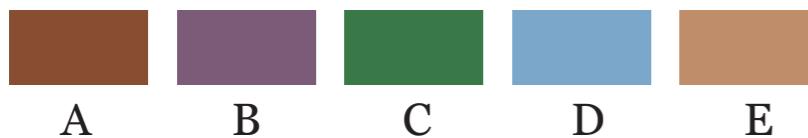
Angle/area



Line weight



Hue and shade



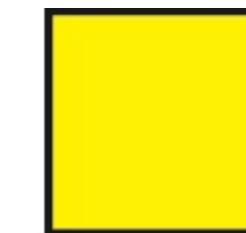
Figures represented
in all these graphics:
22%, 25%, 34%, 29%, 32%

Data visualization
and visual encoding

4. USE COLOR STRATEGICALLY

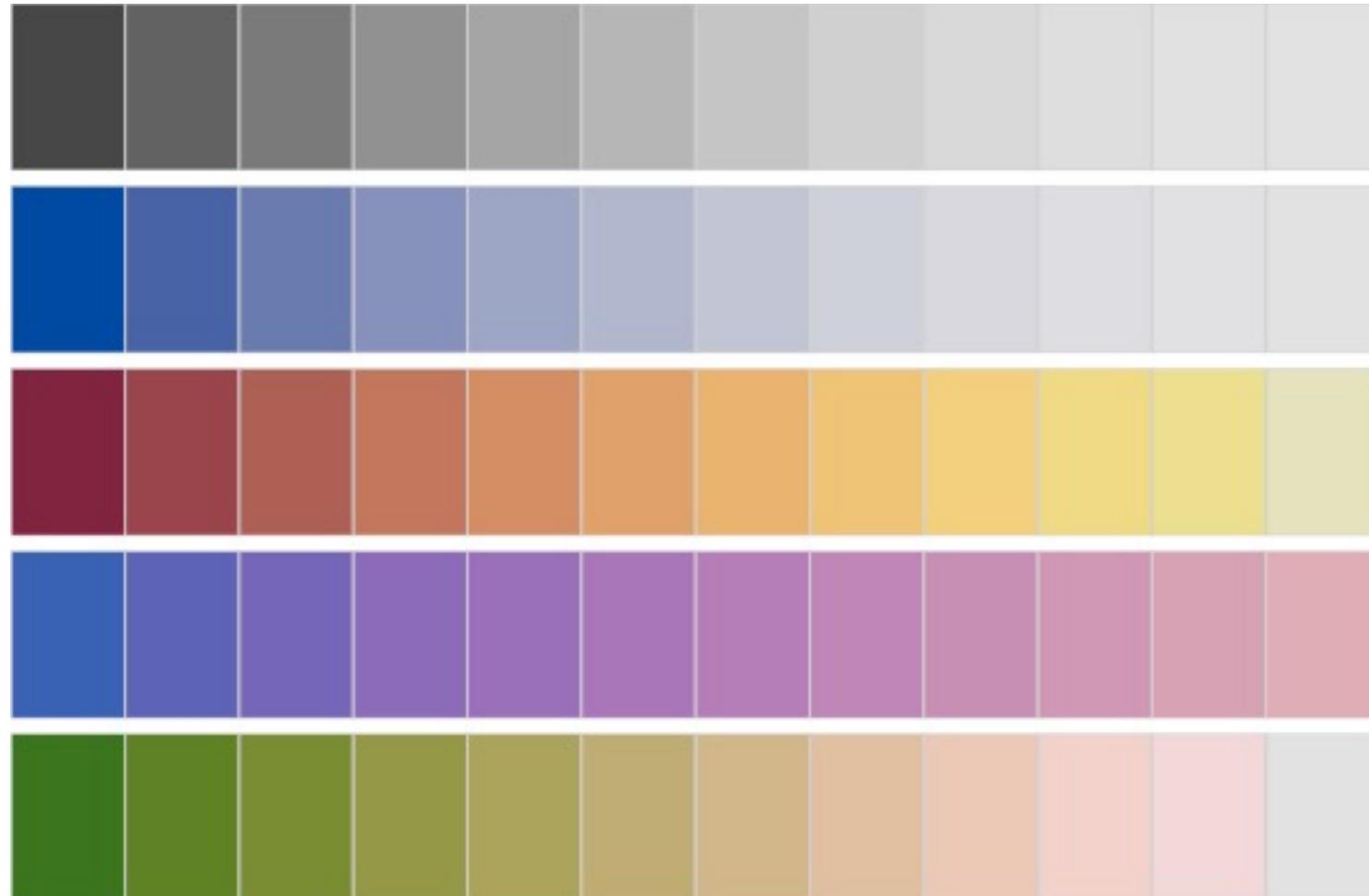
Colors for Categories

Do not use more than 5-8 colors at once!



Colors for Ordinal Data

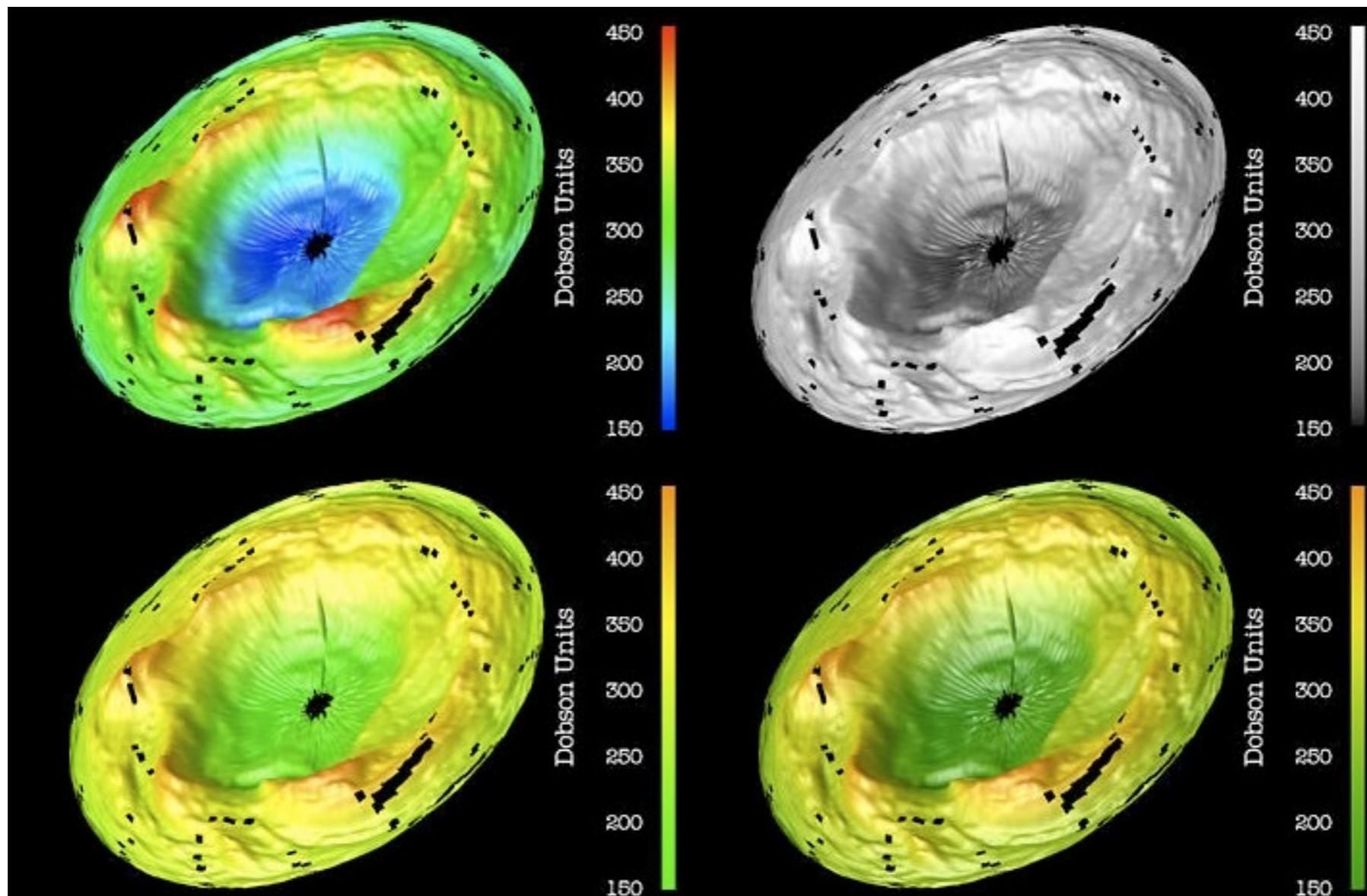
Vary luminance and saturation



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

COLORS FOR QUANTITATIVE DATA

Hue
(Rainbow)

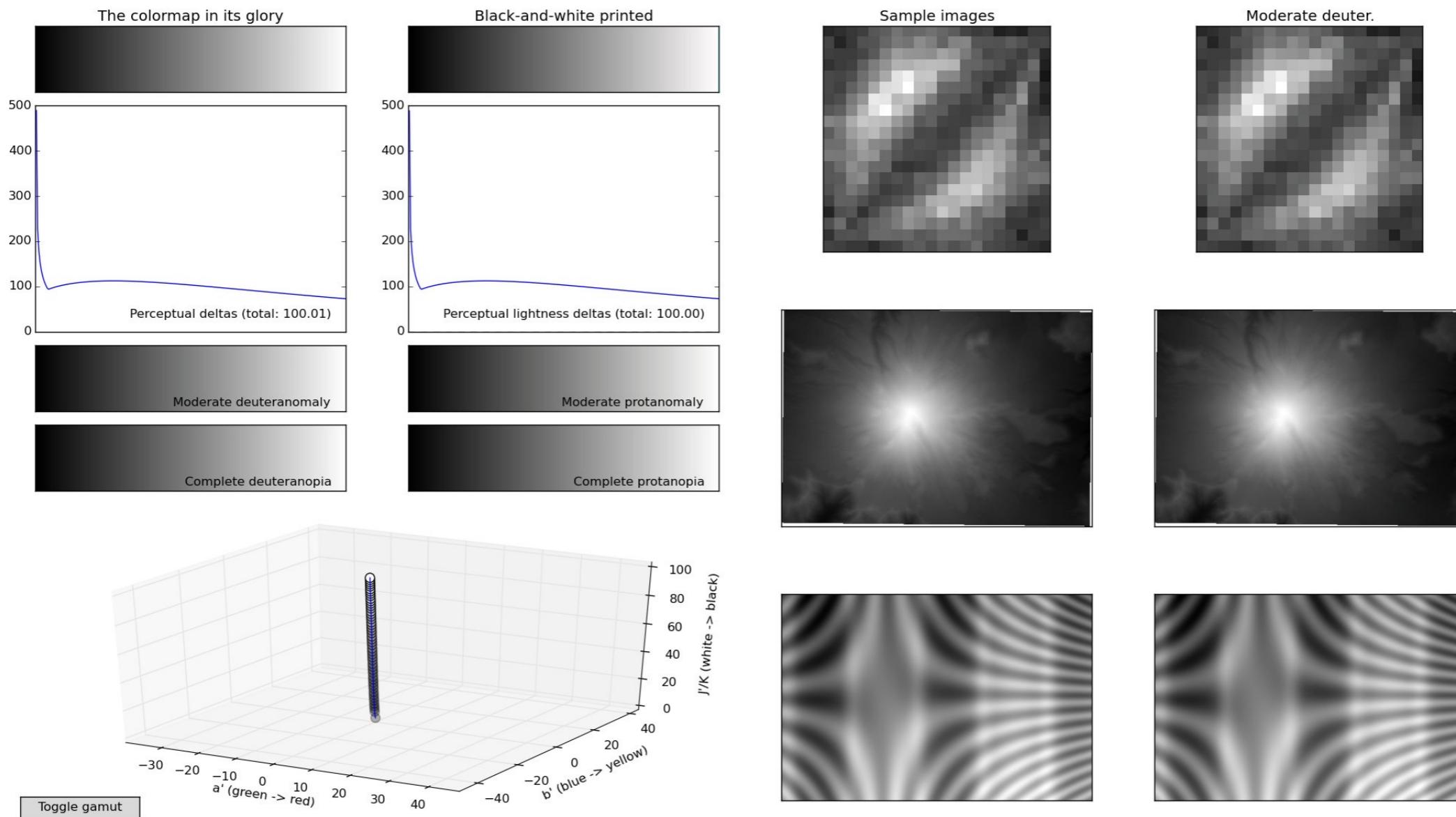


Luminance

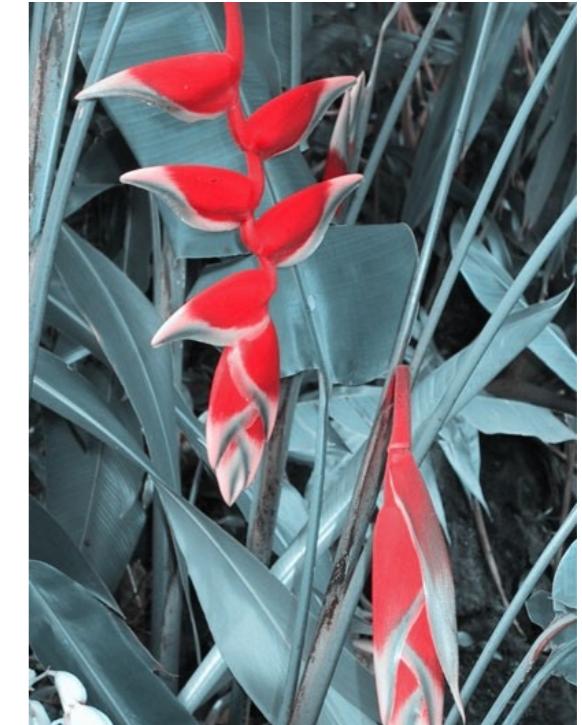
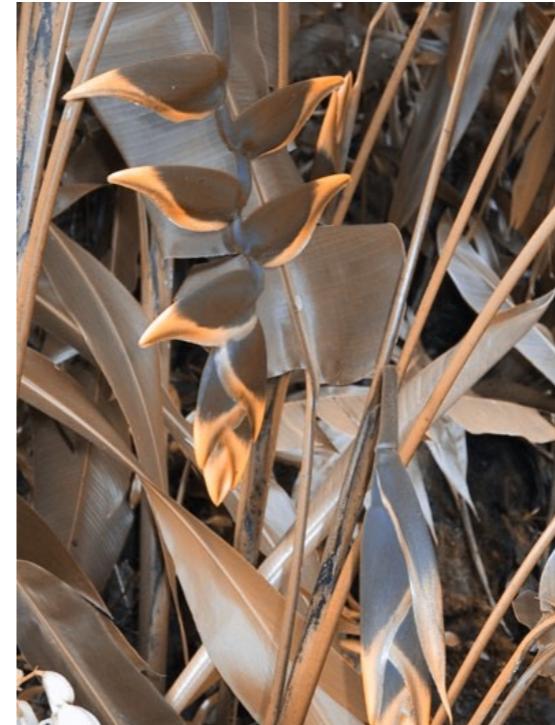
Luminance
& Hue

GRAY

Colormap evaluation: gray



COLOR BLINDNESS



Protanope

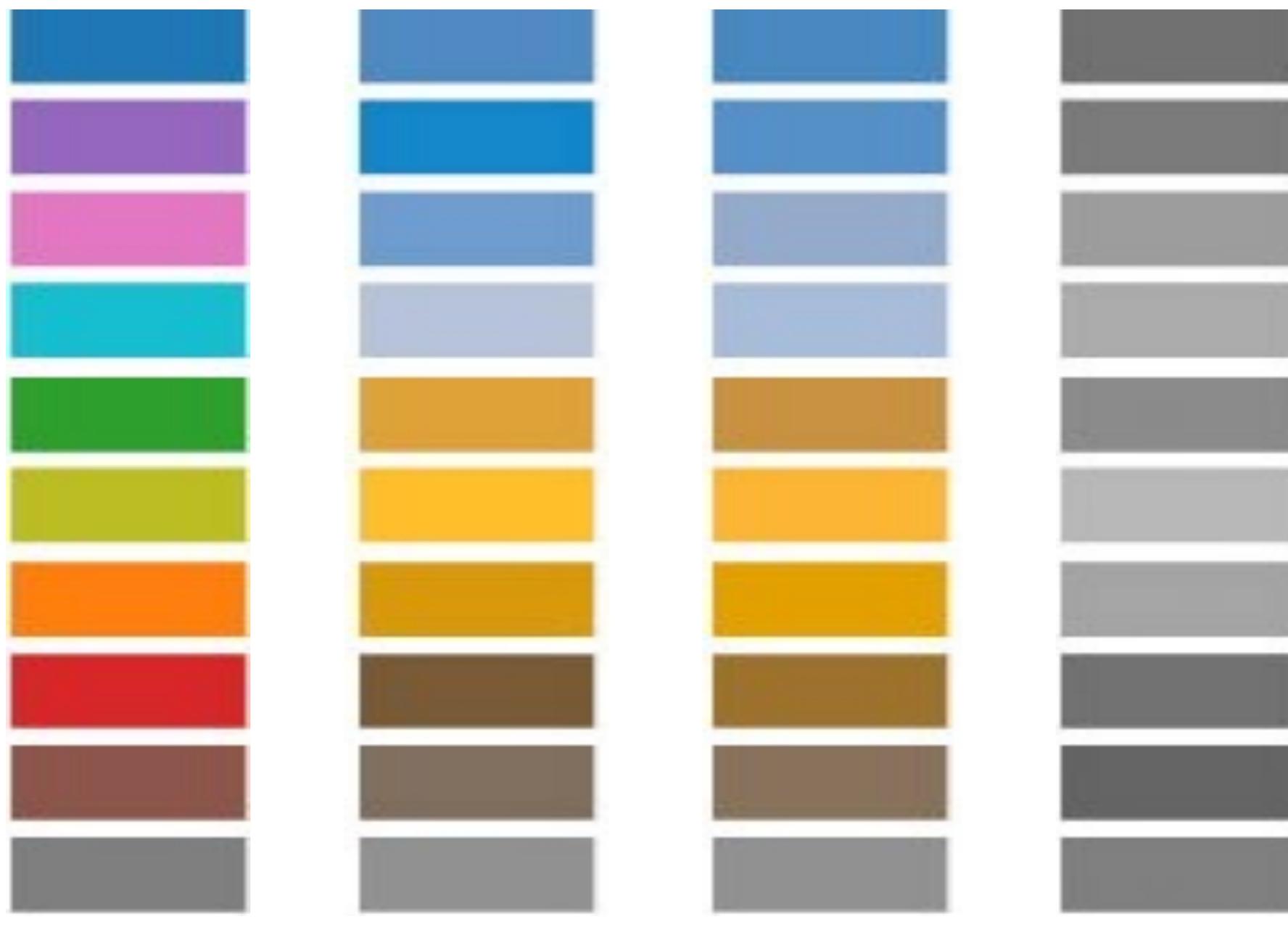
Deutanope

Tritanope

Red /green
deficiencies

Blue /Yellow
deficiency

COLOR BLINDNESS



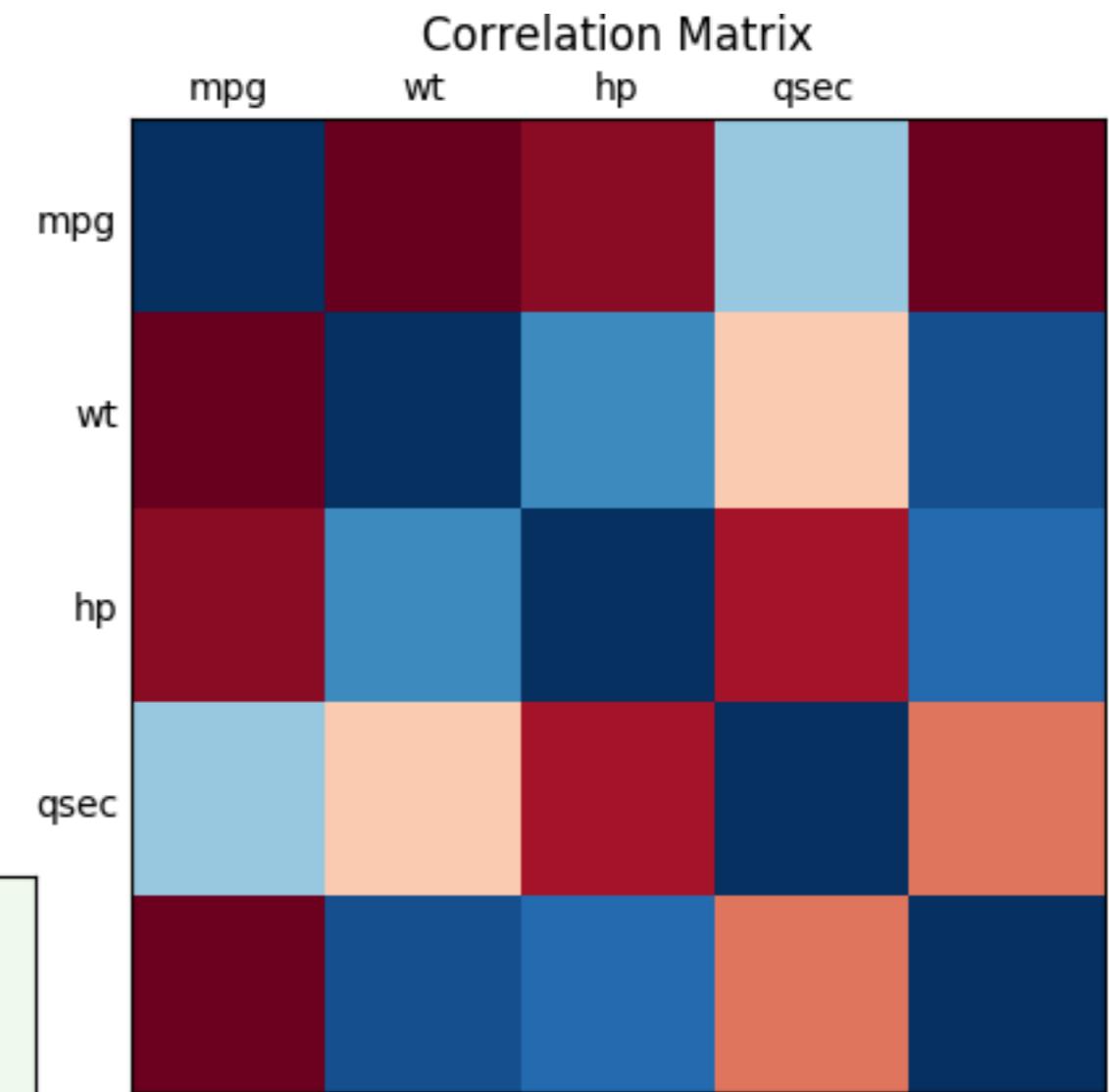
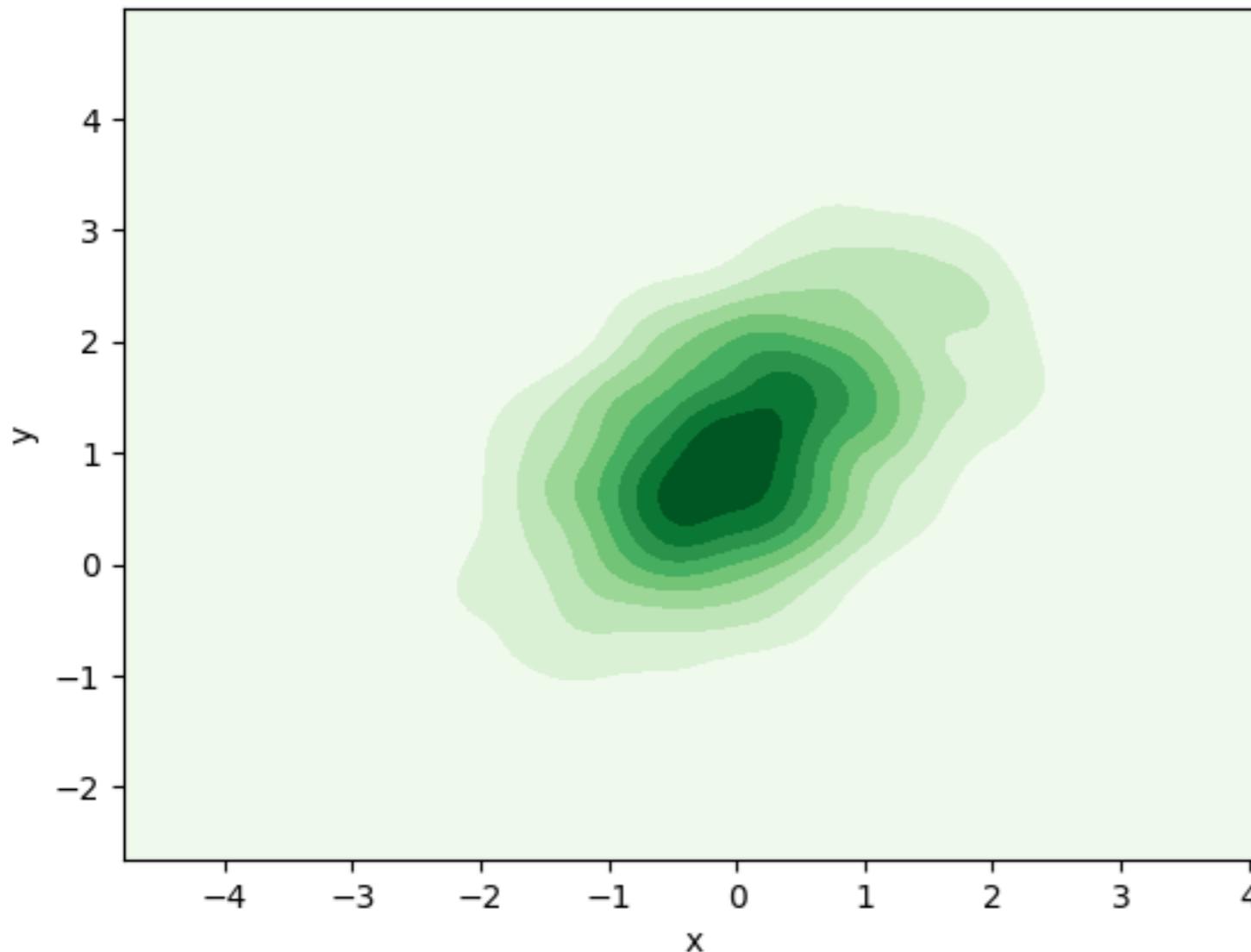
Normal

Protanope

Deutanope

Lightness

Diverging Palette for Quantitative or Ordinal



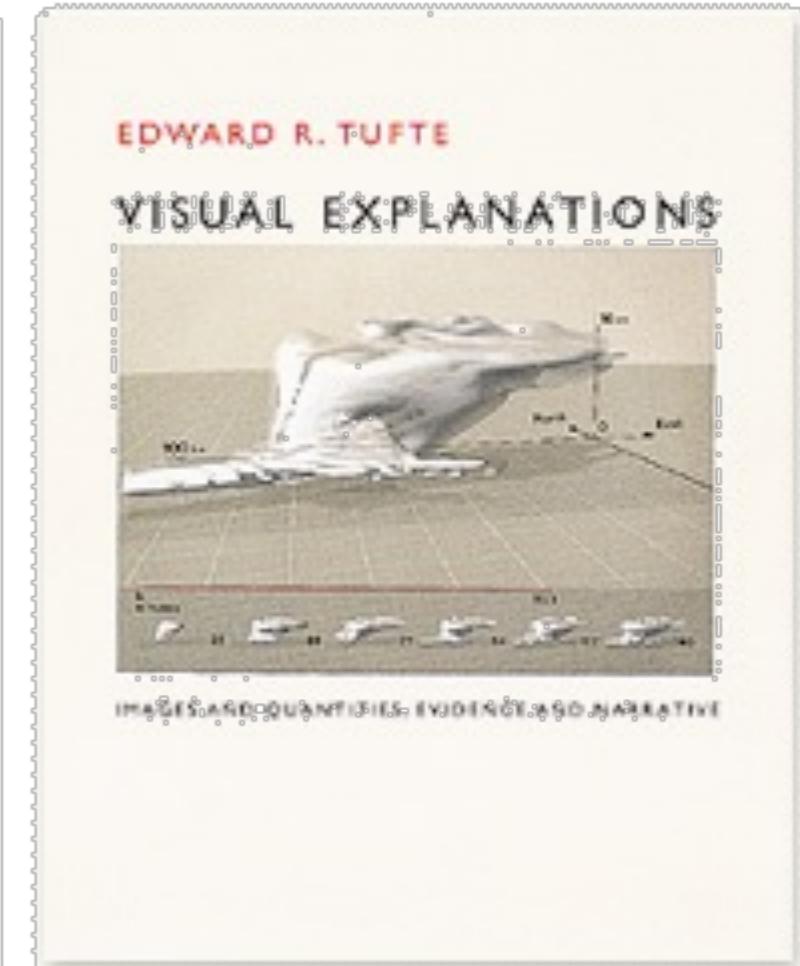
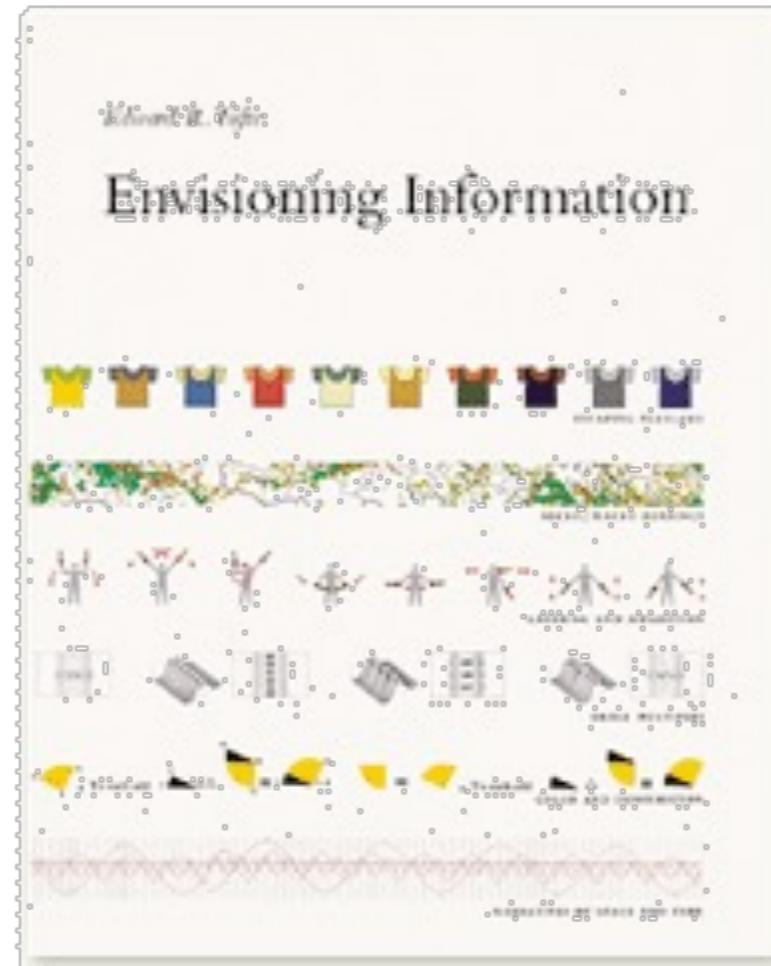
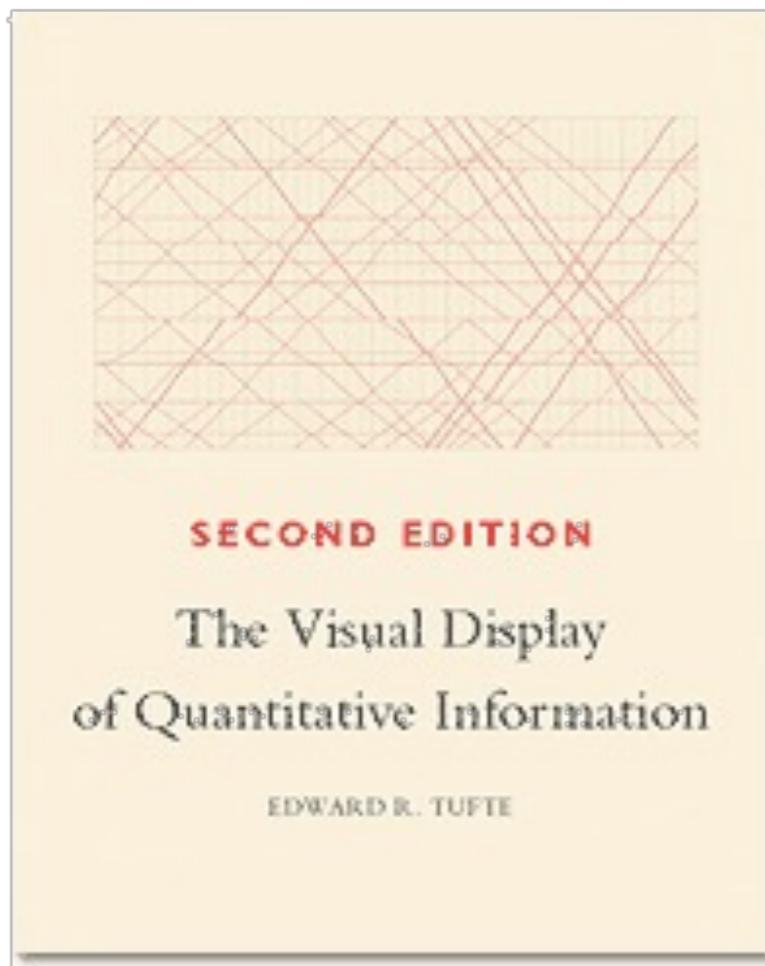
Sequential Palette for Densities

SUMMARY: EFFECTIVE VISUALIZATIONS

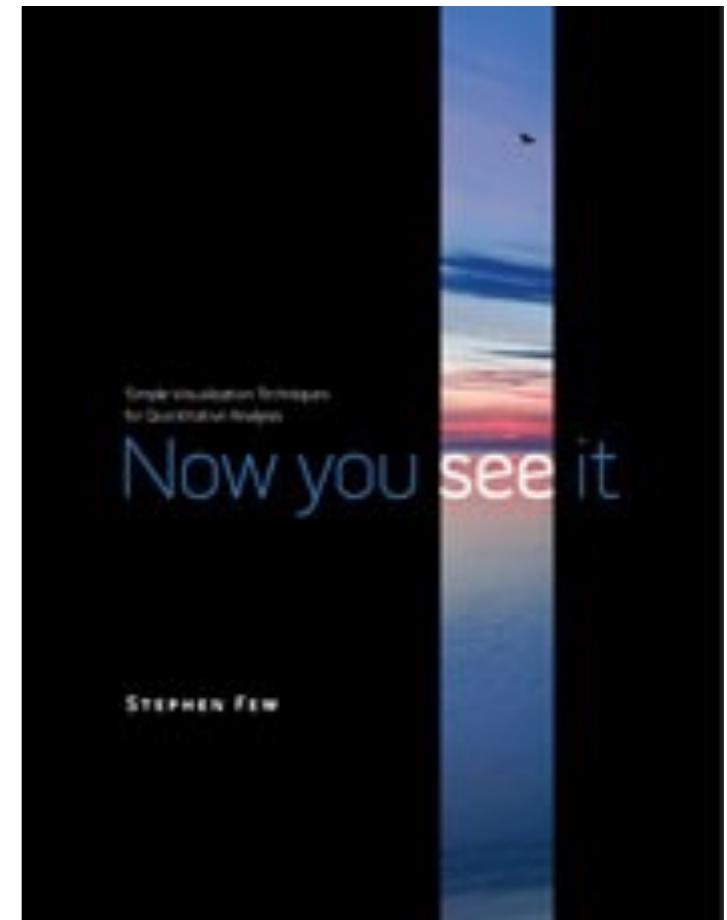
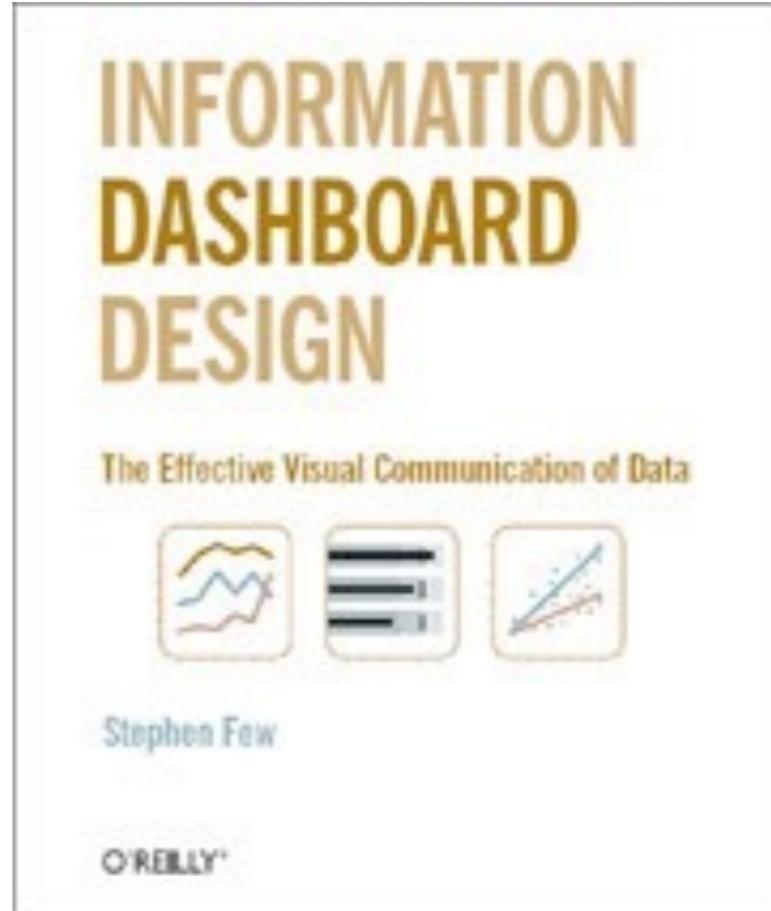
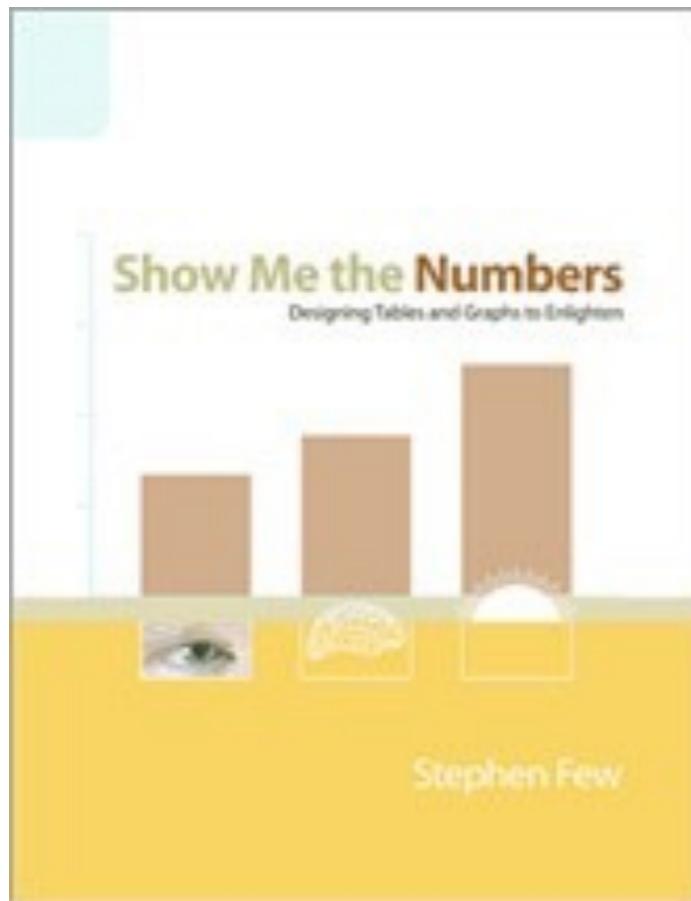
1. Have graphical integrity
2. Keep it simple
3. Use the right display
4. Use color strategically

FURTHER READING

EDWARD TUFTE



STEPHEN FEW



Thank you for your
attention





QUESTIONS?