

Topic 4: Descriptive statistics & visualization

Lectures 8 & 9

- Descriptive statistics
- Spurious correlations
- Visualization challenges

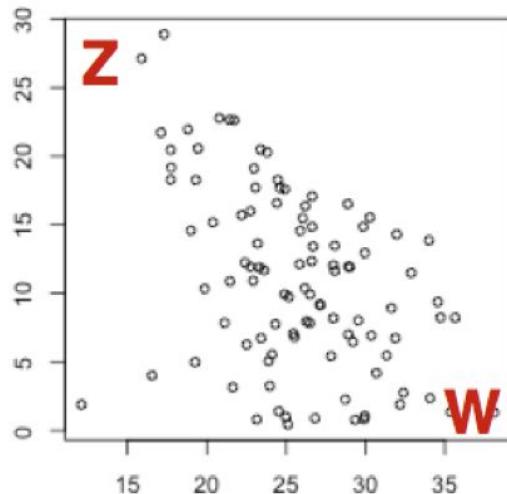
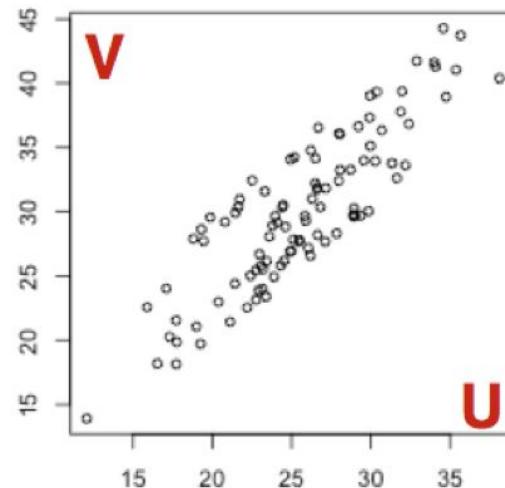
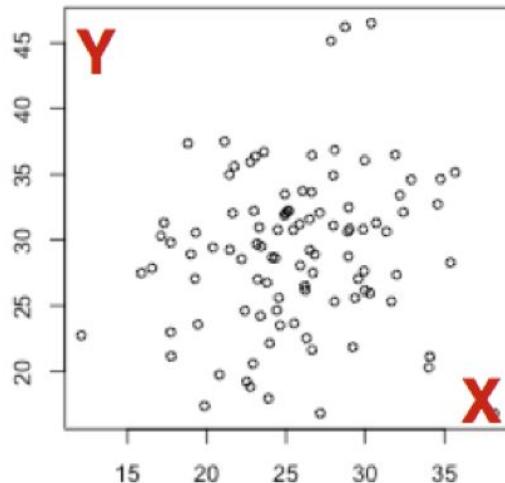
Calculating correlation

Variables

Attributes / Features



x	10	8	13	9	11	14	6	4	12	7	5
y	8.04	6.95	7.58	8.81	8.33	9.96	7.24	4.26	10.84	4.82	5.68



Correlation coefficient

Pearson Correlation Coefficient

- Measures 'linear' relationship between variables.

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

where:

- n is the sample size
- x_i, y_i are the single samples indexed with i
- $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ (the sample **mean**); and analogously for \bar{y}

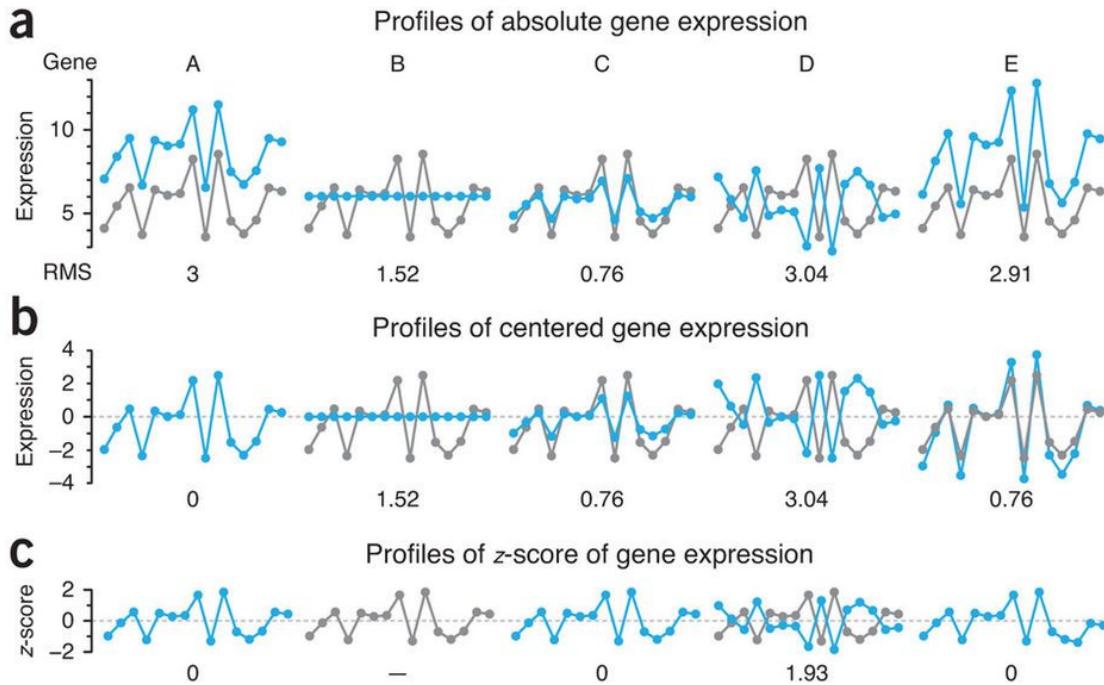
$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

Correlation coefficient

Pearson Correlation Coefficient

- Captures the relationship between 2 vectors after centering each vector by its mean and scaling by its standard deviation.
- The final quantities for each vector are called z-scores.

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

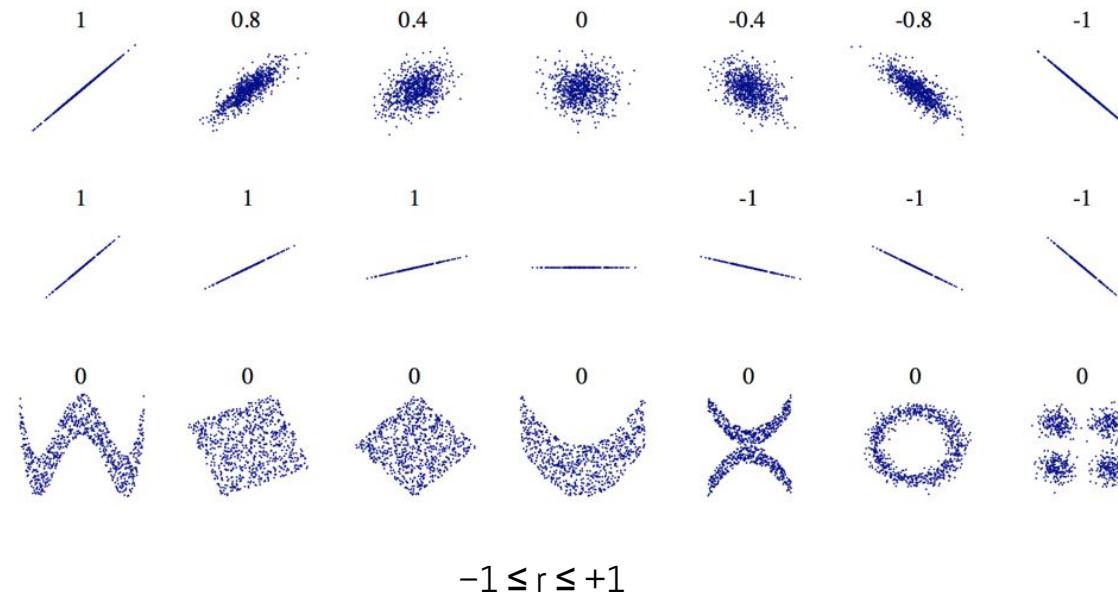


Correlation coefficient

Pearson Correlation Coefficient

- Measures 'linear' relationship between variables.

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

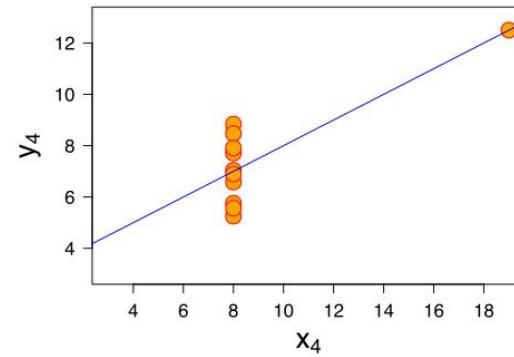
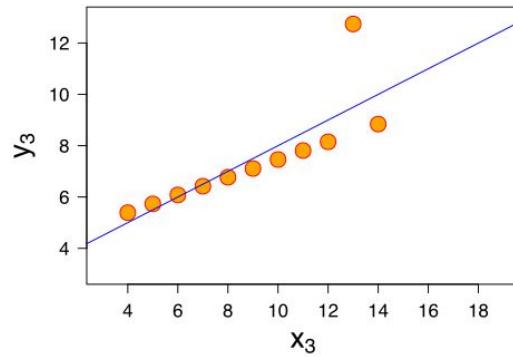
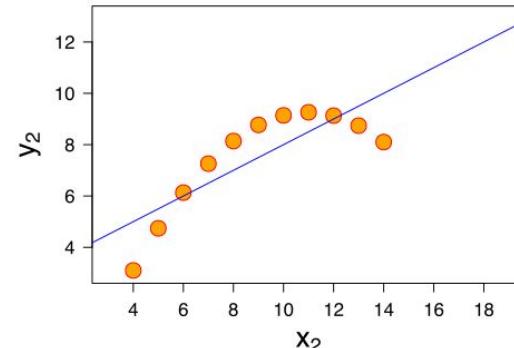
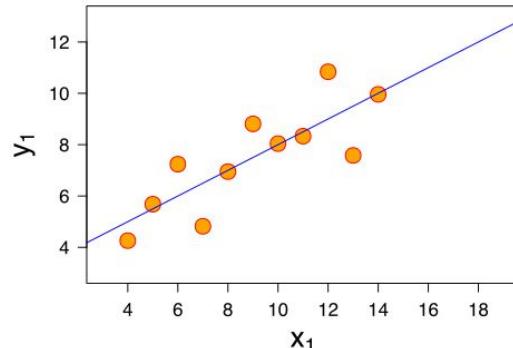


-1 is total -ve correlation | 0 is no correlation | +1 is total +ve correlation

Anscombe's quartet: “calculation are exact; graphs are rough!”

11 datapoints

- Mean (x) = 9
- Var (x) = 11
- Mean (y) = 7.50
- Var (y) ~ 4.12
- Cor (x, y) = 0.816
- Linear regression line:
 - $y = 3.00 + 0.500x$

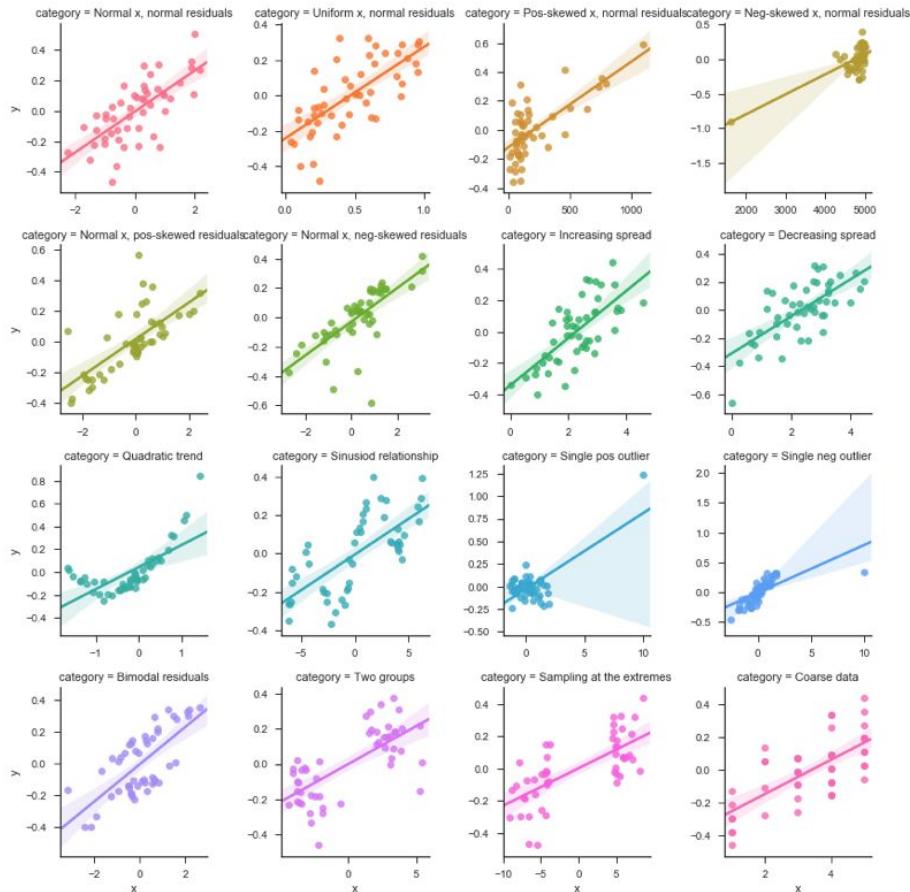
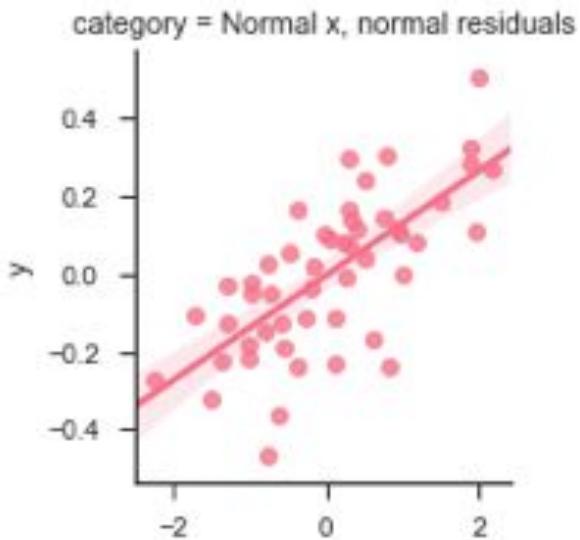


Anscombe, F. J. (1973). "Graphs in Statistical Analysis". American Statistician 27 (1): 17–21.

Wikipedia

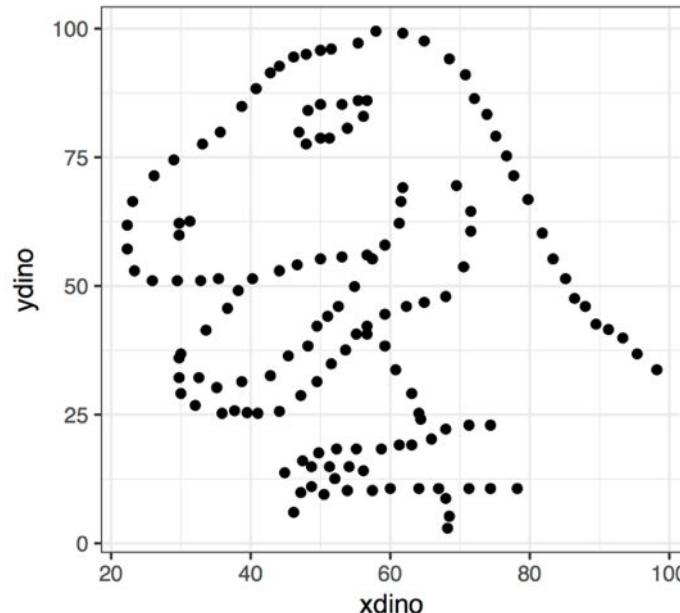
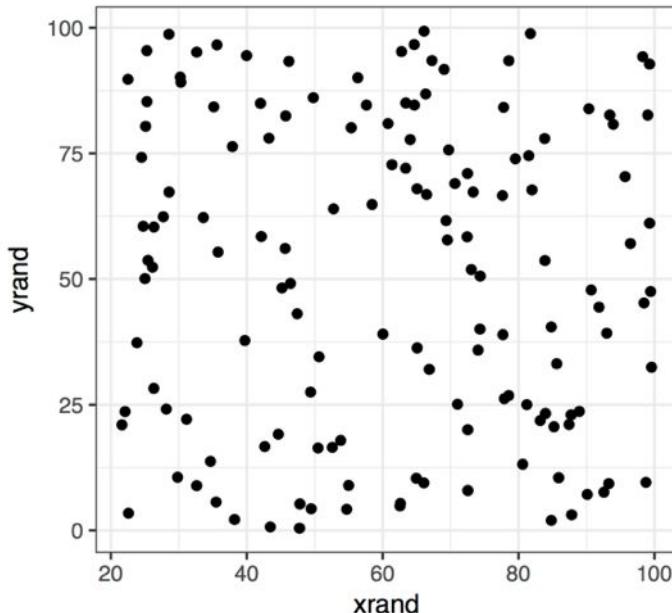
What does a correlation coefficient tell you about the data?

Correlation = 0.7



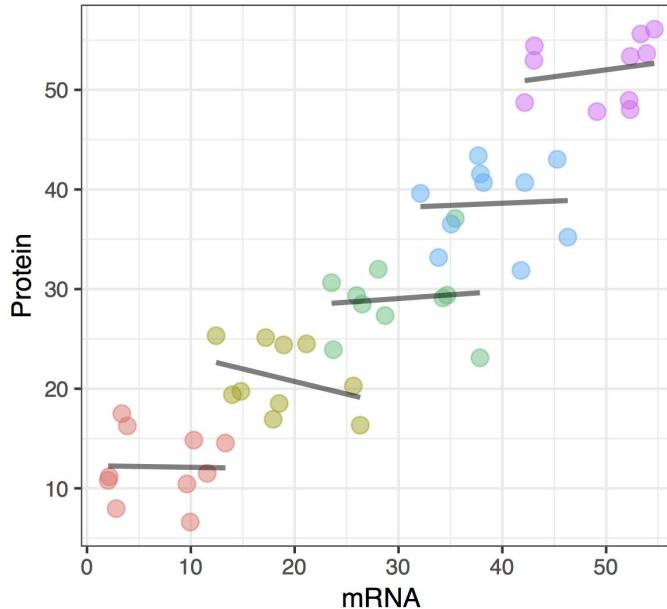
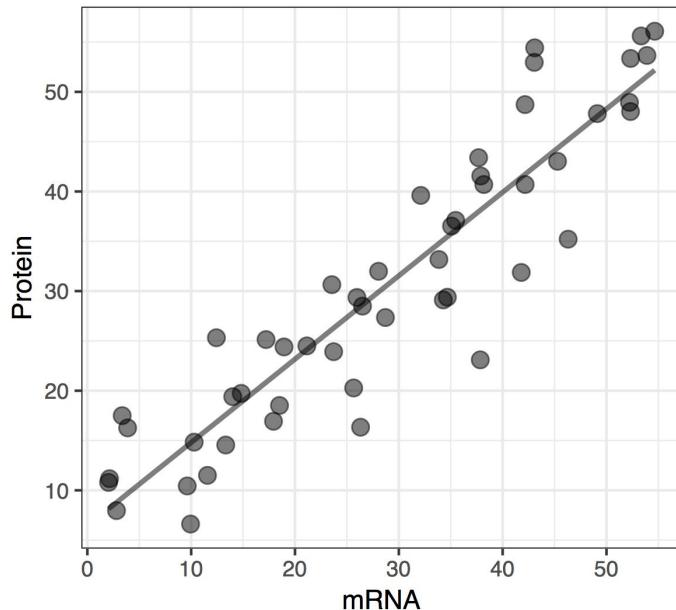
What does a correlation coefficient tell you about the data?

Correlation = -0.06



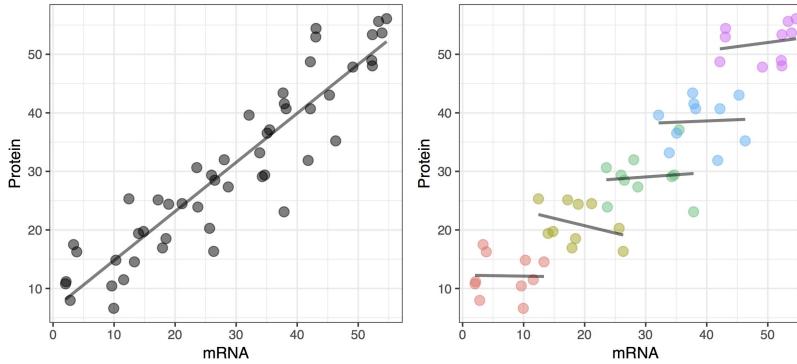
What does a correlation coefficient tell you about the data?

Simpson's Paradox



What does a correlation coefficient tell you about the data?

Simpson's Paradox

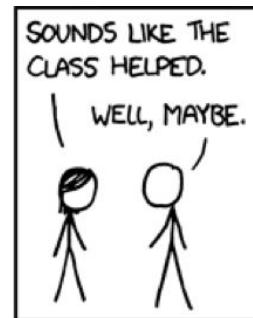
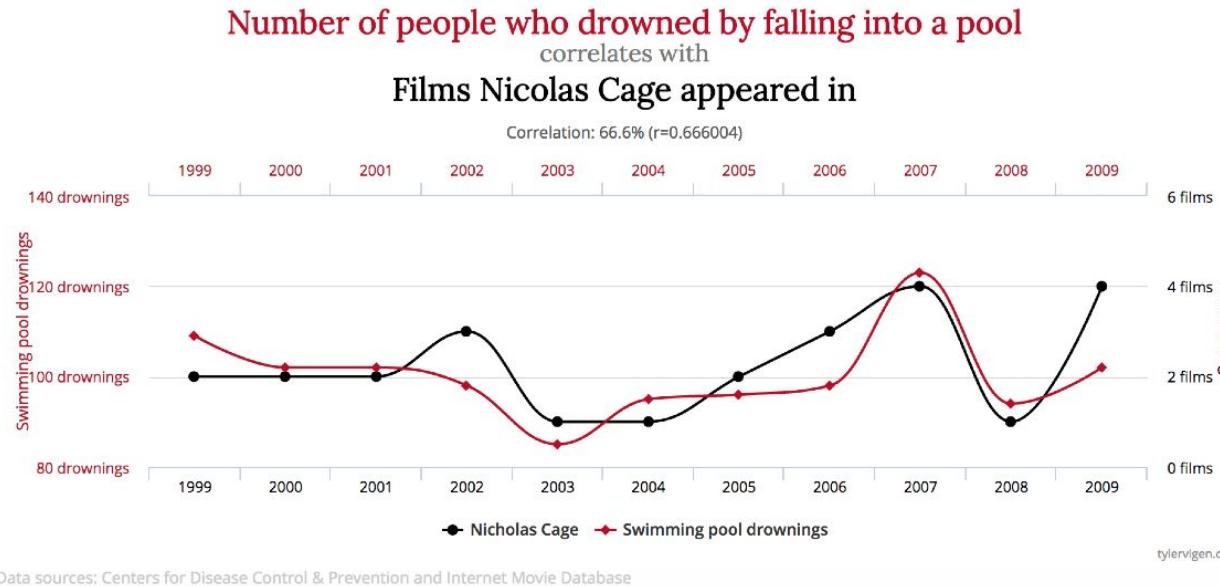
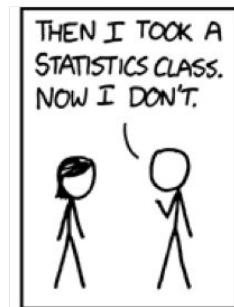
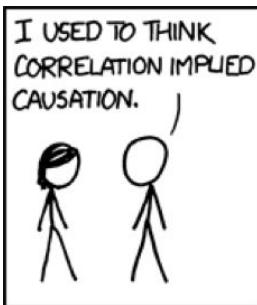


Success rates of kidney stone removal surgeries

Treatment	Diameter < 2 cm	Dia. \geq 2 cm	Overall
Open surgery	93%	73%	78%
Percutaneous nephrolithotomy	87%	69%	83%

Spurious correlations

What does Nicholas Cage have to do with people drowning in swimming pools?



Checkout <https://www.google.com/trends/correlate>

Spurious correlations

Simulate fluctuations in correlation coefficients

- Repeat 10,000: Calculate correlation coefficients of $n = 10$ samples of two independent normally distributed variables ($\mu = 0, \sigma = 1$). Plot a histogram.
- Mark statistically significant coefficients ($\alpha = 0.05$).
- Plot the samples with the three largest and smallest correlation coefficients (statistically significant).
- Vary $\sigma = \{0.1, 0.5, 1.0\}$ and vary sample size $n = \{5, 10, 50\}$.

Many correlation/distance measures

Pearson Correlation Coefficient

$$d = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Spearman Rank Correlation

Euclidean Distance

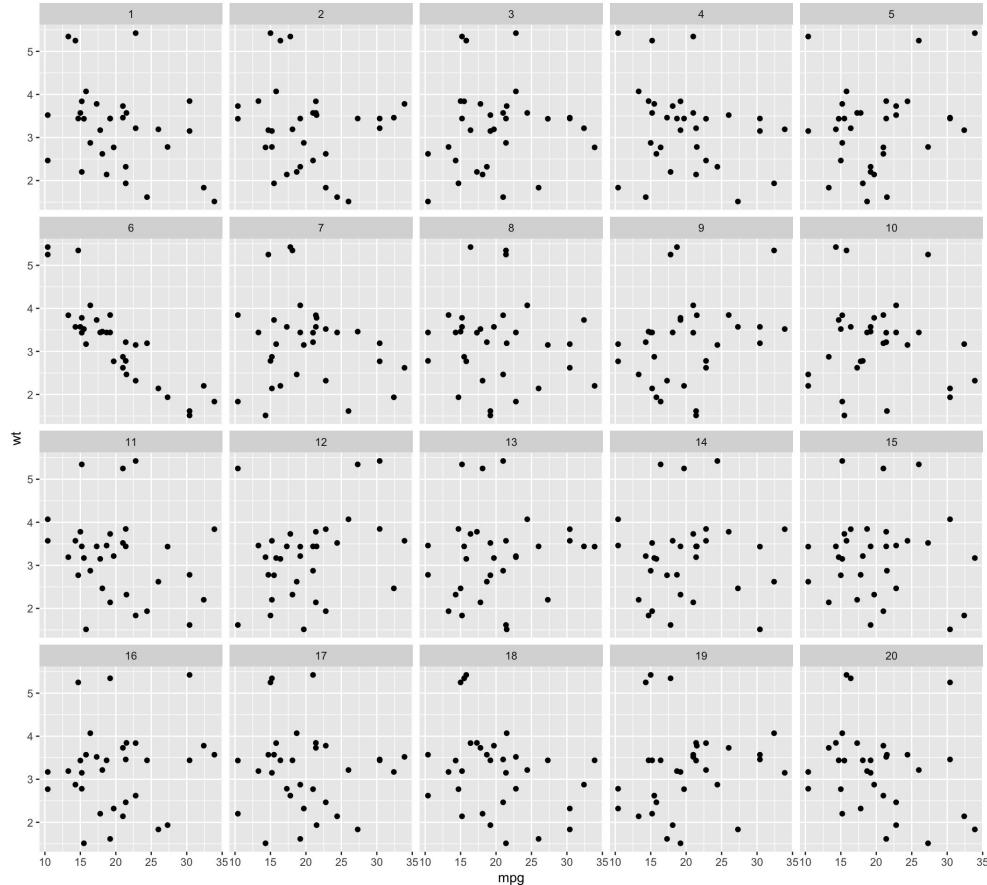
$$r = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma_x} \right) \left(\frac{y_i - \bar{y}}{\sigma_y} \right)$$

Mutual Information

...

$$\rho = 1 - \frac{6 \sum_{i=1}^n [rank(x_i) - rank(y_i)]}{n(n^2 - 1)}$$

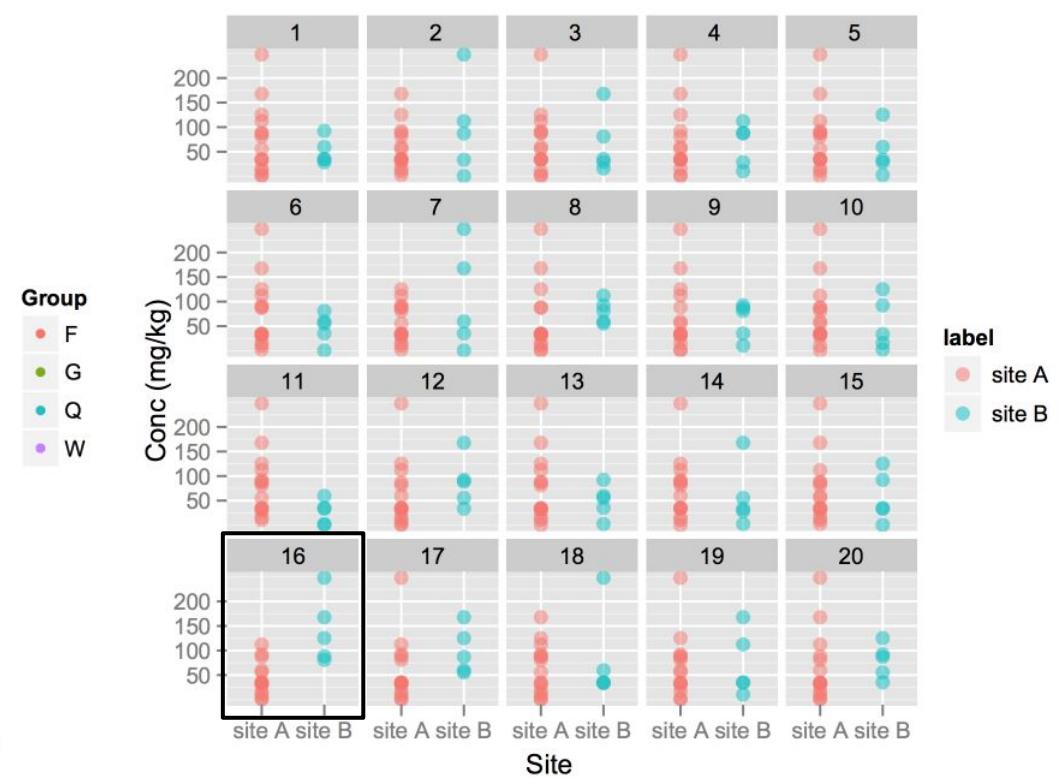
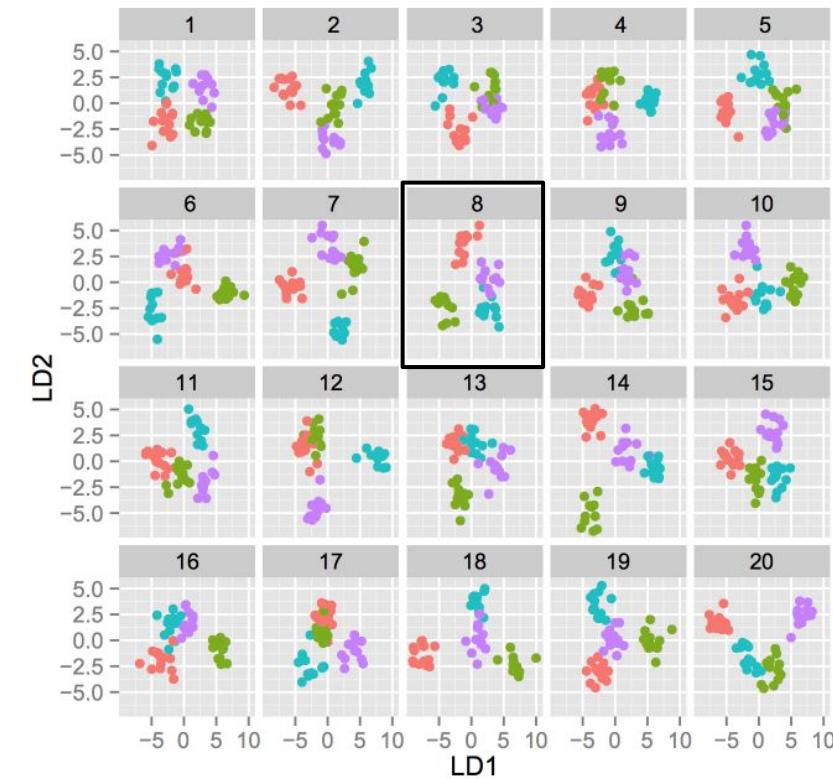
Spurious correlations – But it *looks* associated!



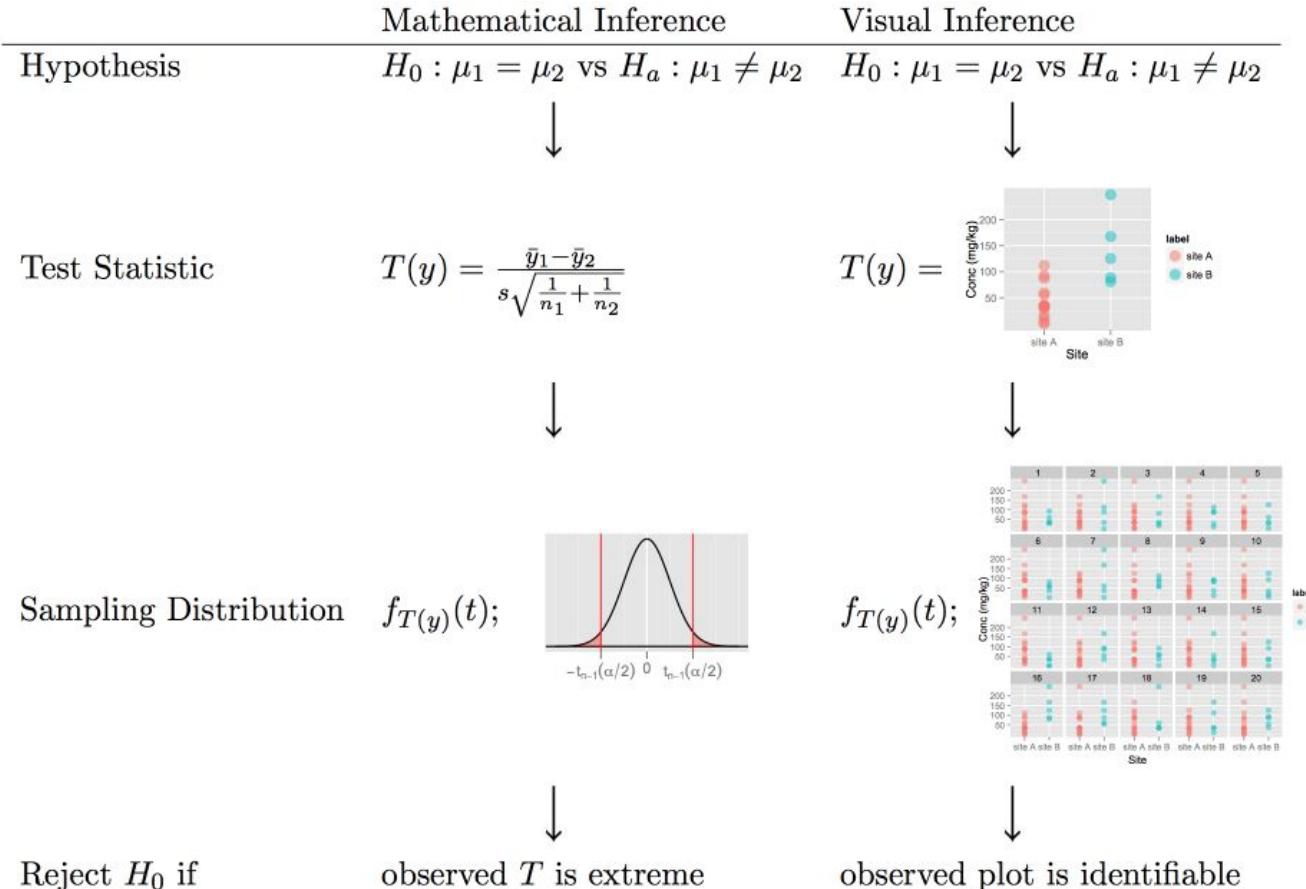
Create a lineup for visual inference

- Place the plot of the real data amongst a set of null plots to create a lineup; Null plots are generated in a way consistent with the null hypothesis.
- If the observer can pick the real data as different from the others, this puts weight on the statistical significance of the structure in the plot.

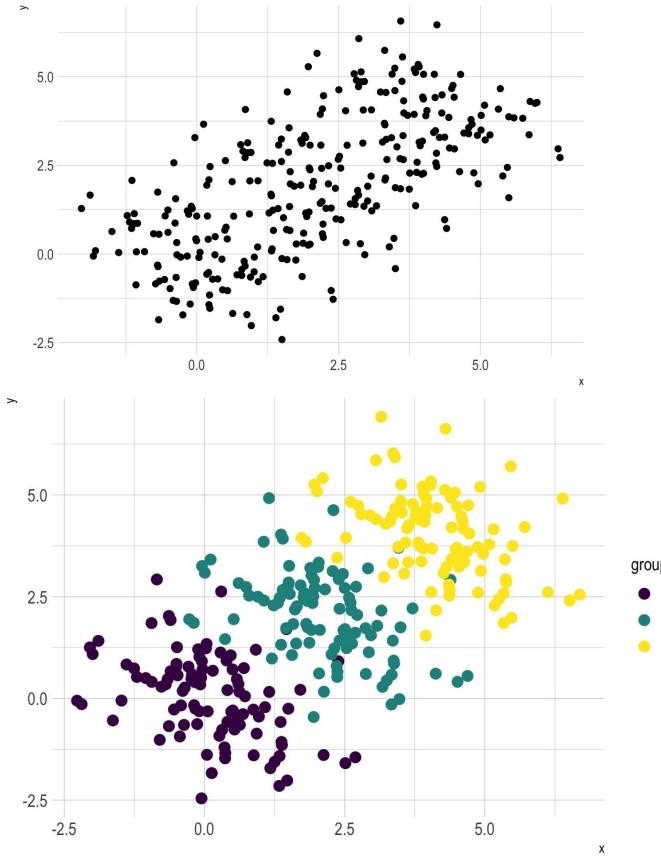
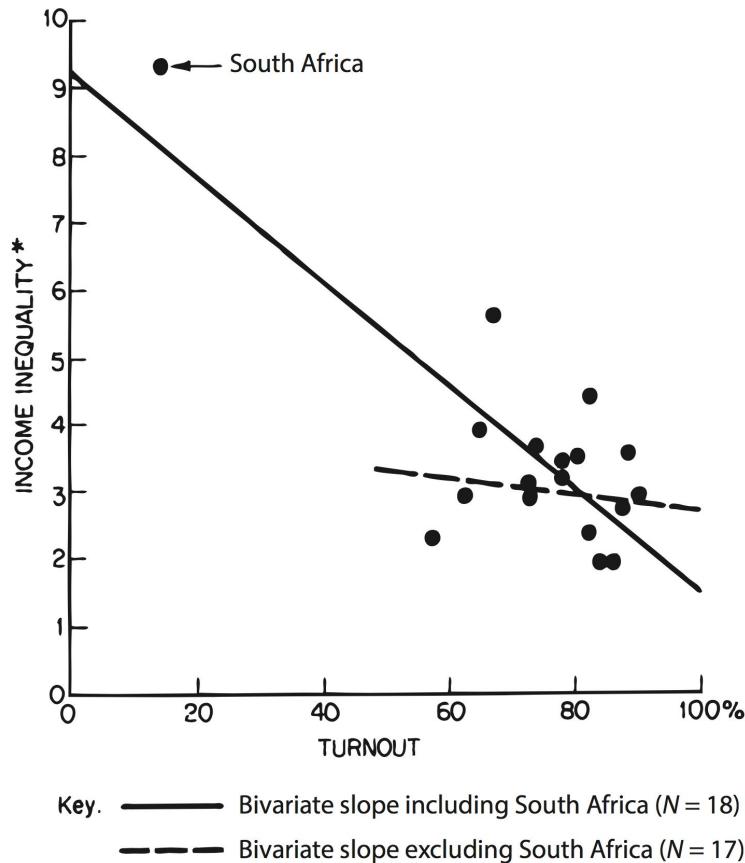
Spurious correlations – But it *looks* associated!



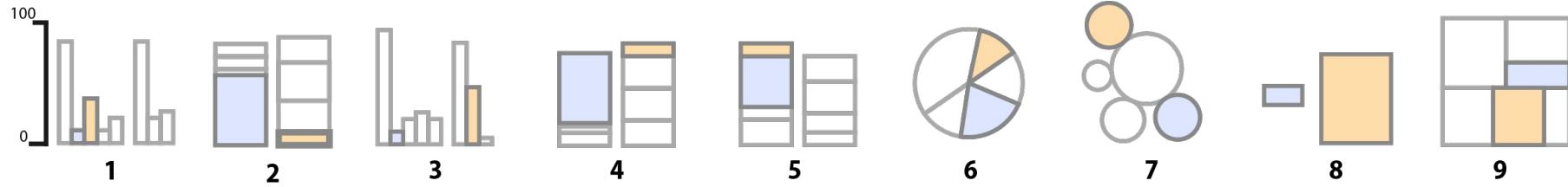
Spurious correlations – But it *looks* associated!



Visualization challenges – Don't do statistics without visualization



Visualization challenges



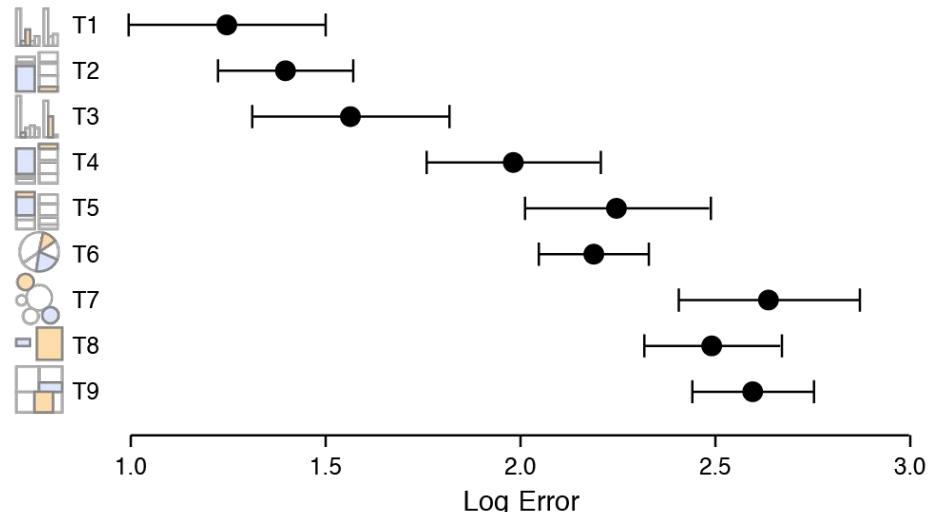
Position

Length

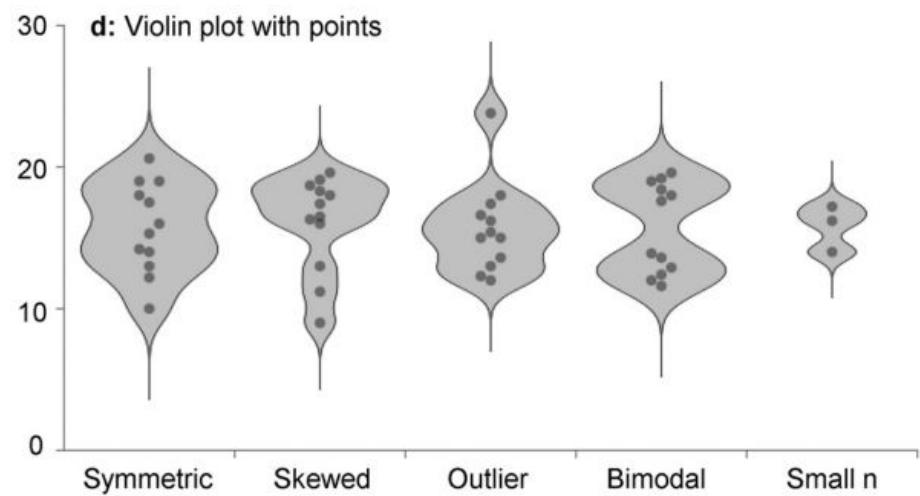
Angle

Circular Area

Rectangular Areas



Visualization challenges – Bar plots, error bars



Visualization challenges – Bar plots, error bars

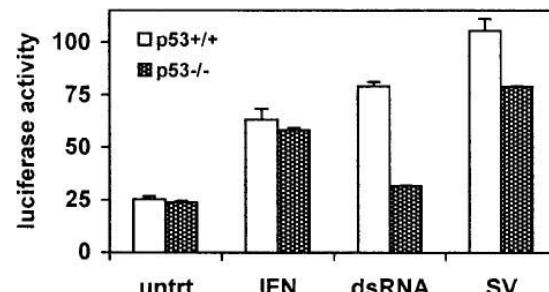
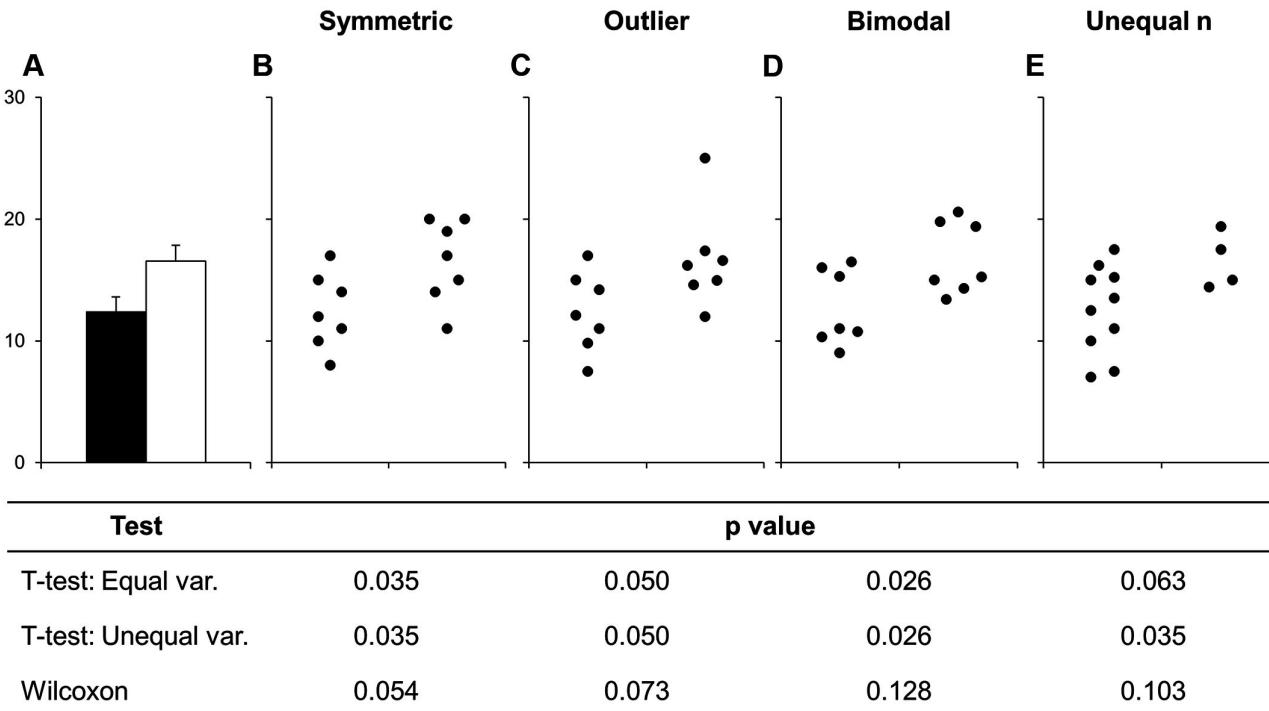
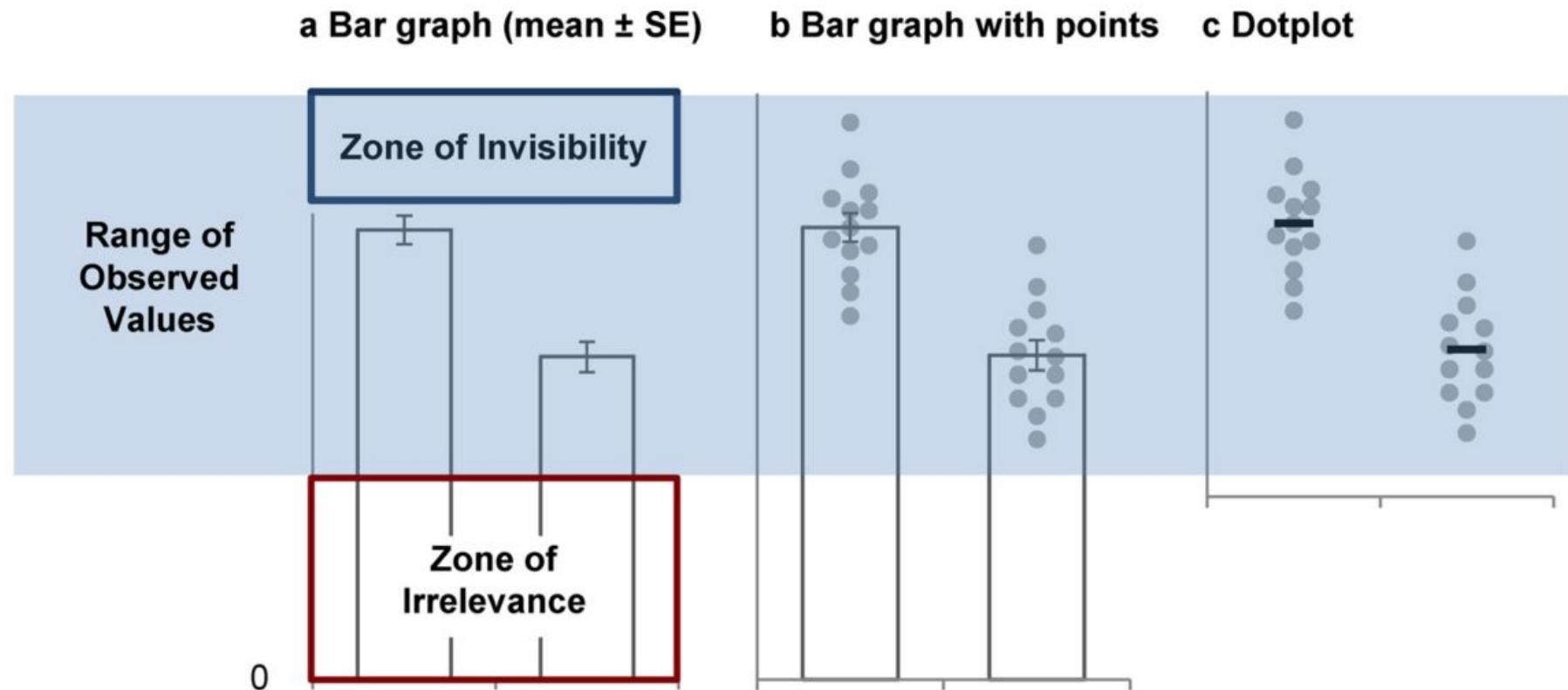


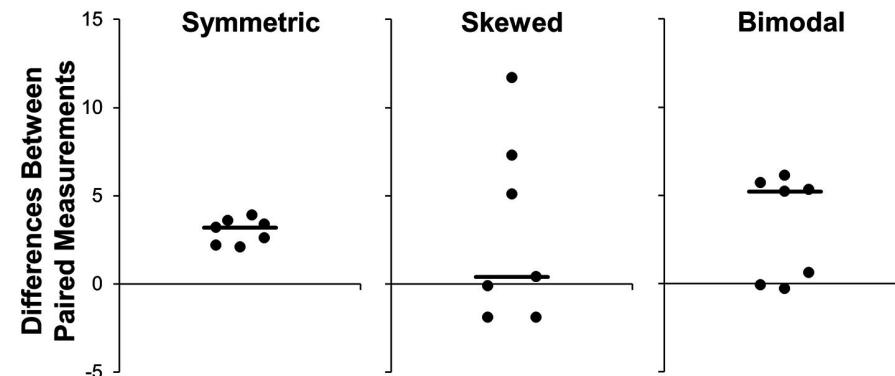
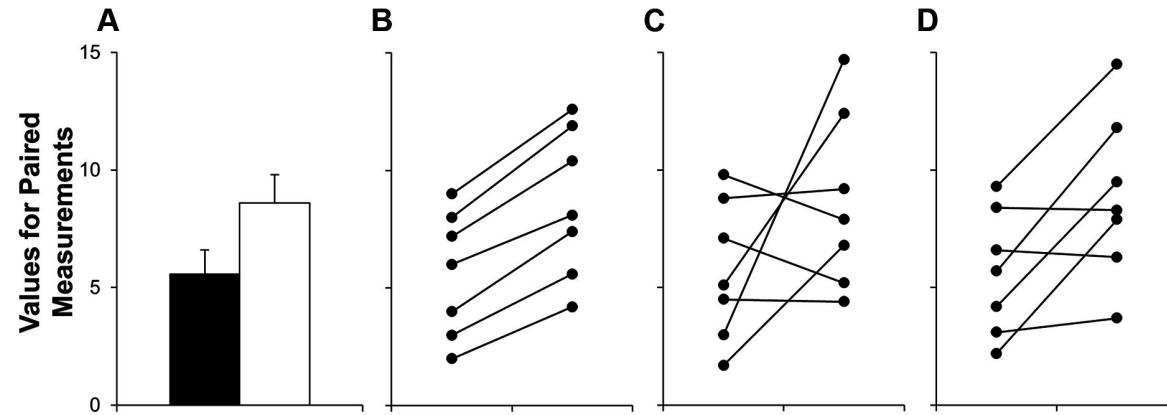
FIG. 4. ISG15 promoter activity mimics endogenous ISG15 mRNA regulation by p53, dsRNA, and virus. Cells (6×10^5 HCT 116) were

differences in transfection efficiency. Each data point is the mean of triplicate samples \pm the standard error; the data presented are repre-

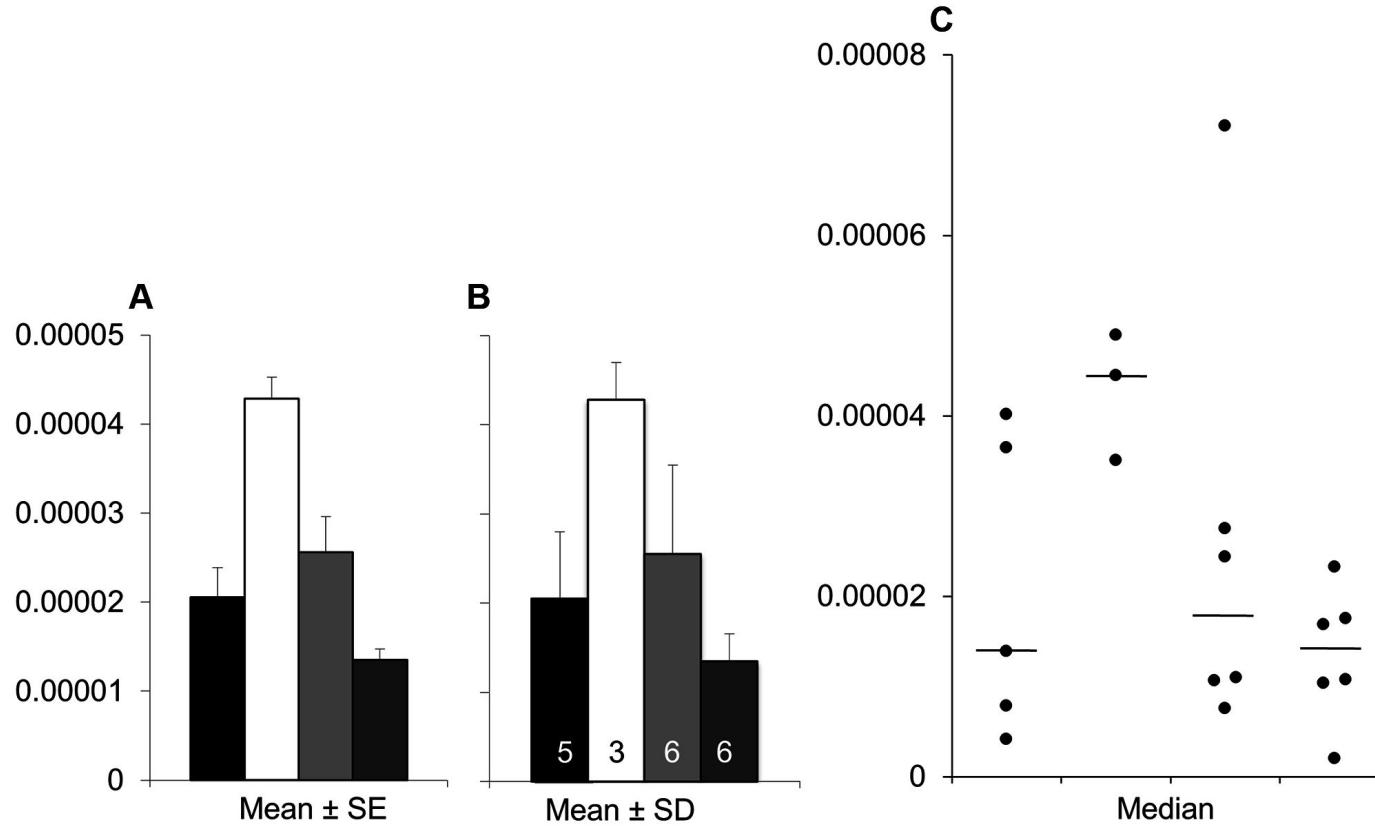
Visualization challenges – Bar plots, error bars



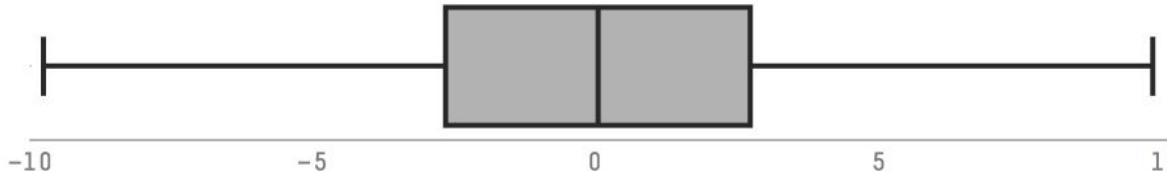
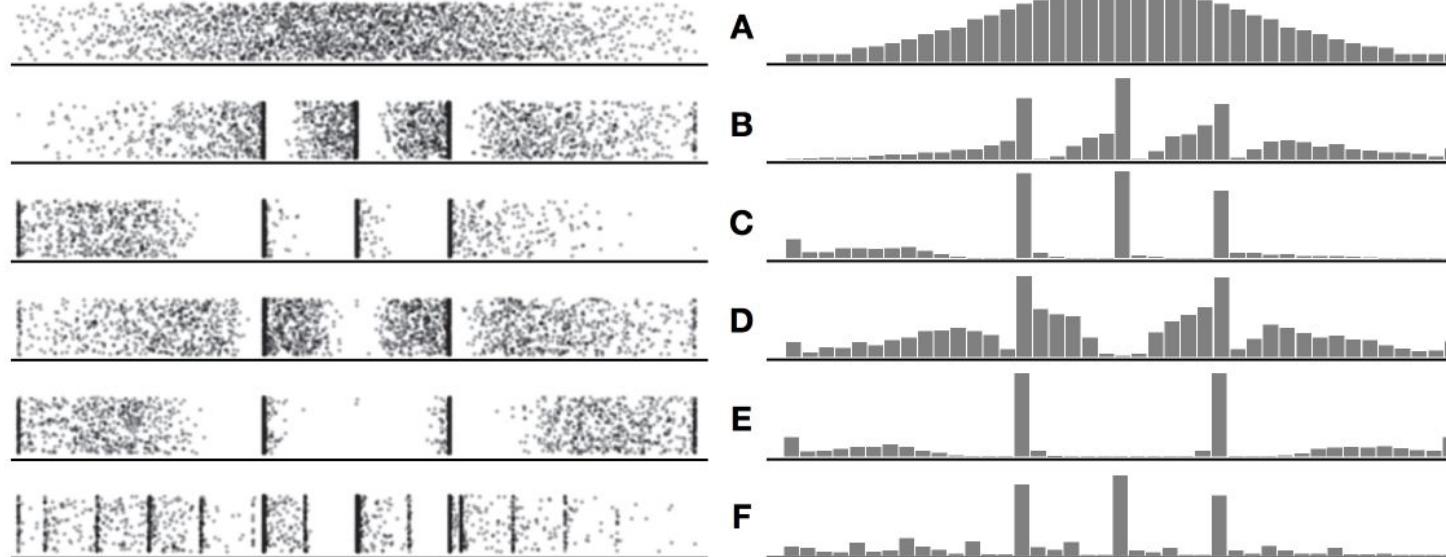
Visualization challenges – Bar plots, error bars



Visualization challenges – Bar plots, error bars

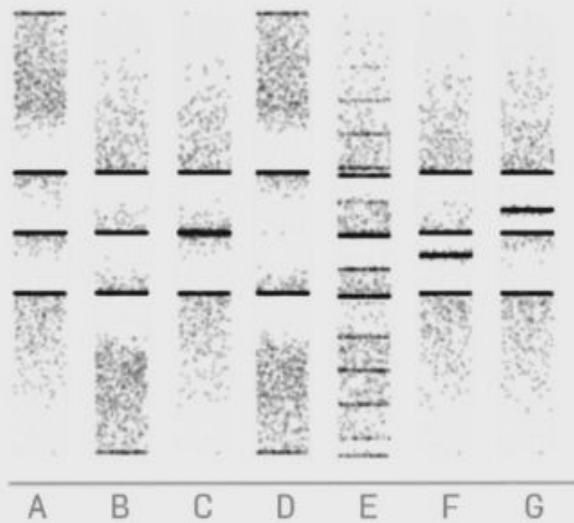


Visualization challenges – Different distributions

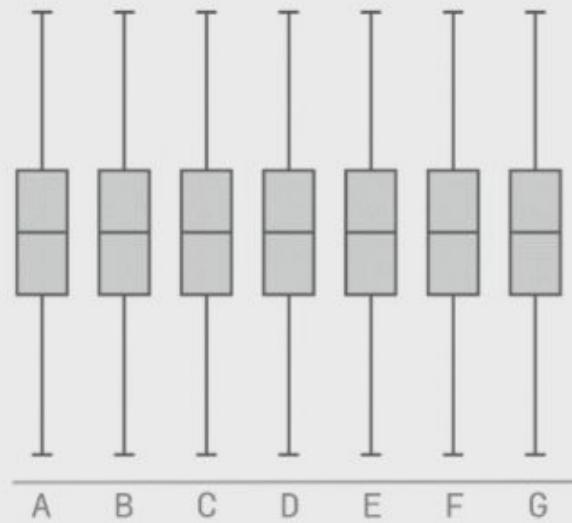


Visualization challenges – Different distributions

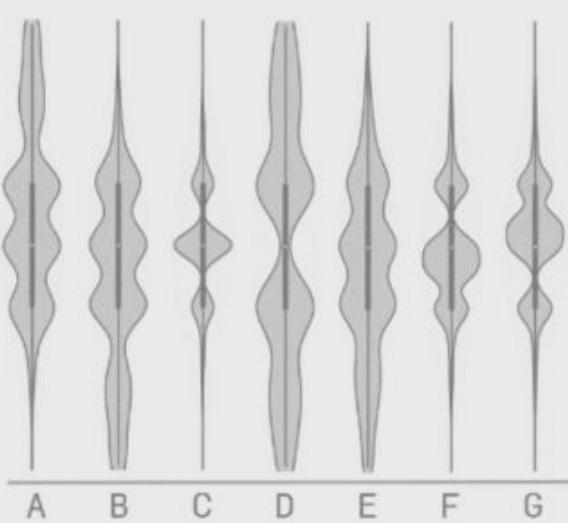
Raw Data



Box-plot of the Data

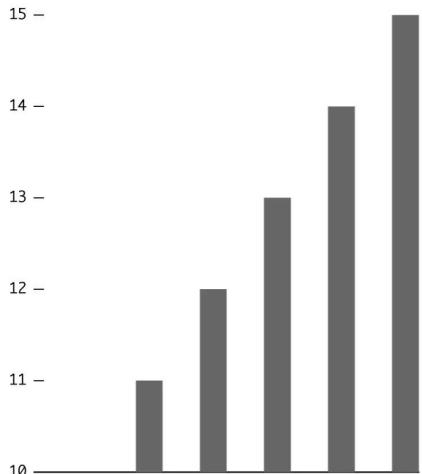


Violin-plot of the Data

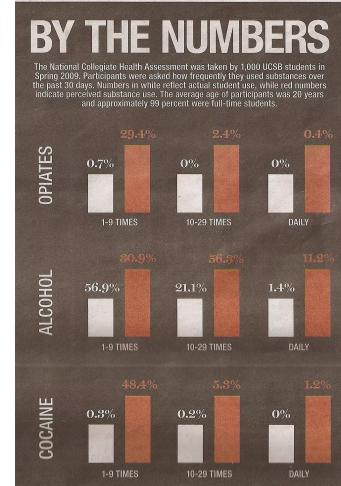
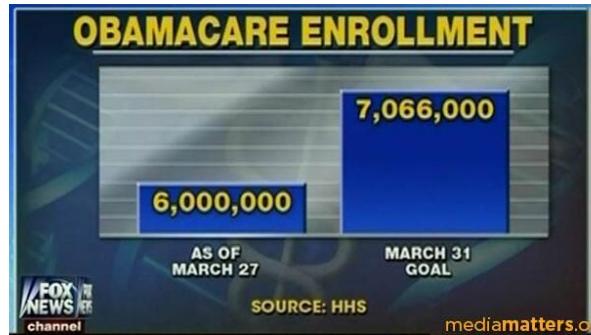
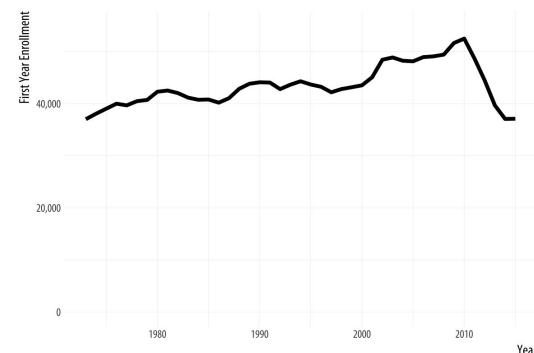
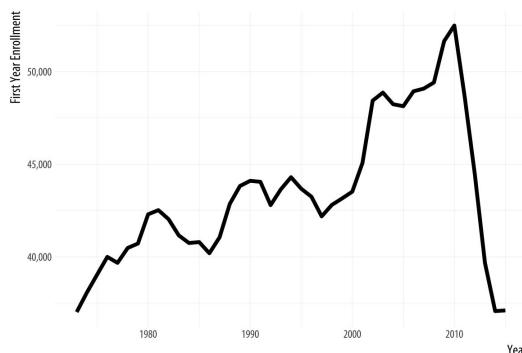
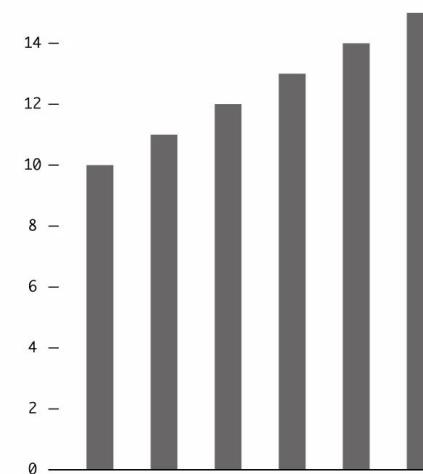


Visualization challenges – Truncated axis

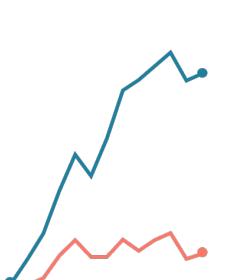
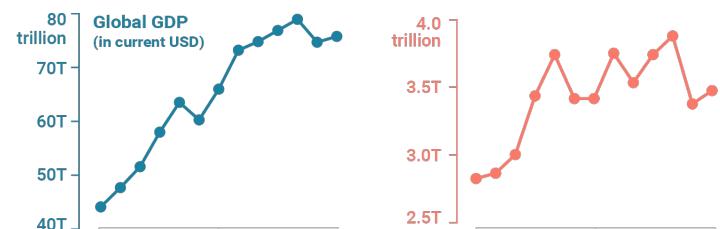
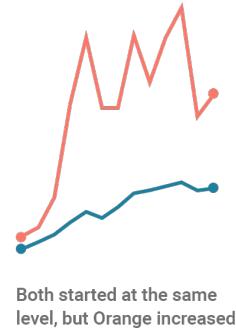
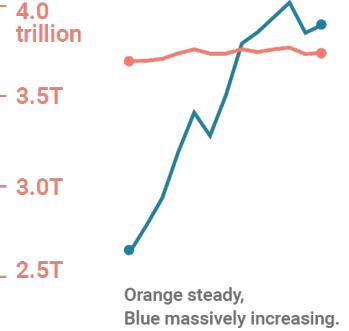
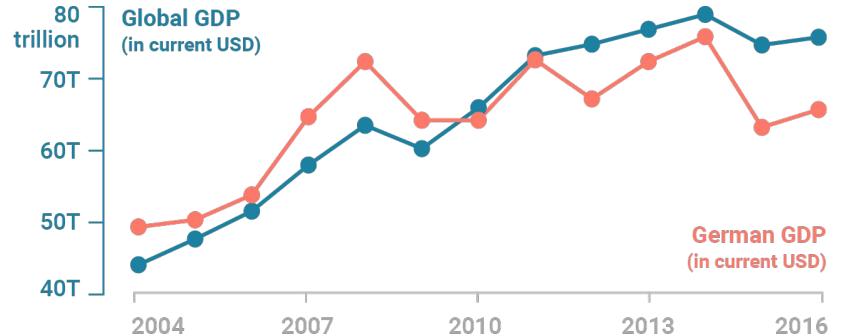
The value axis starts at ten. Liar, liar, pants on fire.



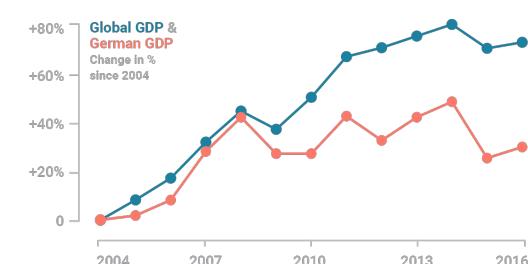
The value axis starts at zero. Good.



Visualization challenges – Dual axes



Both started at the same level, but Blue increased far more than Orange.



Both started with the same increase, then Blue raced to the top.

Visualization challenges – Inappropriate axes

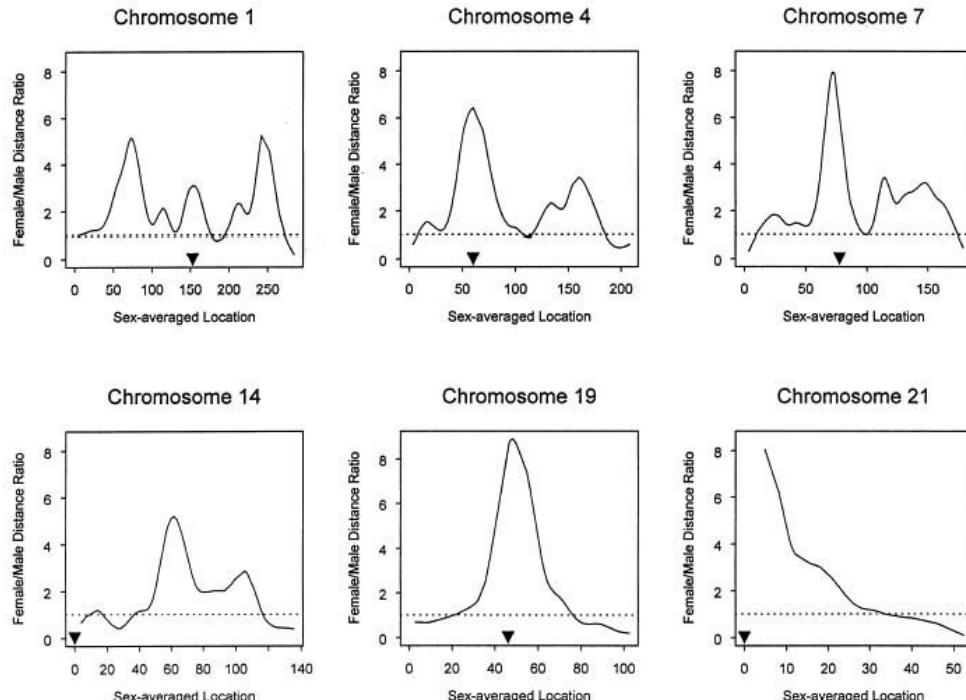
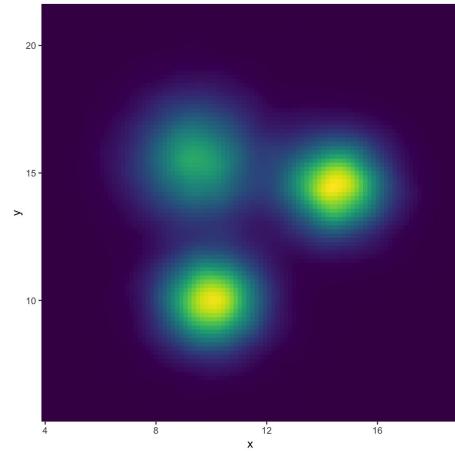
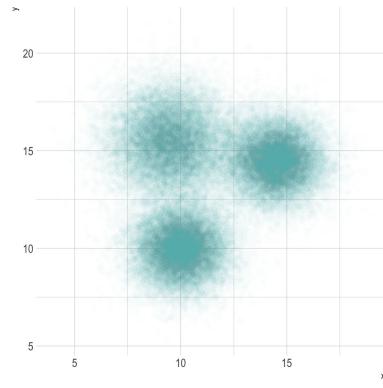
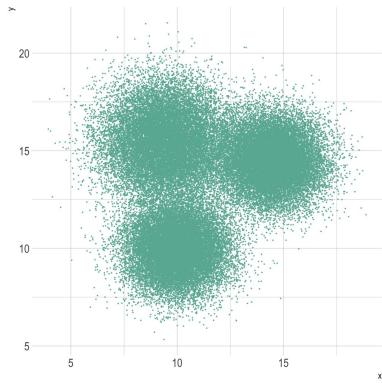
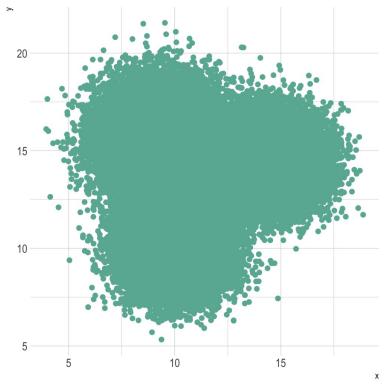
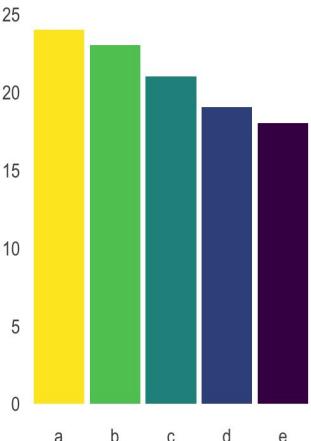
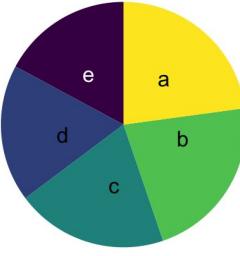
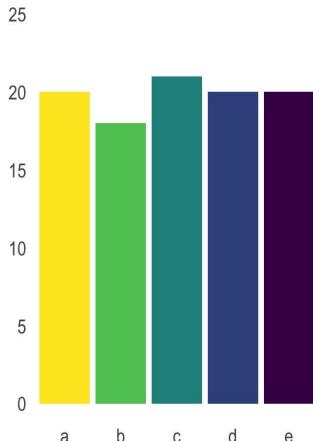
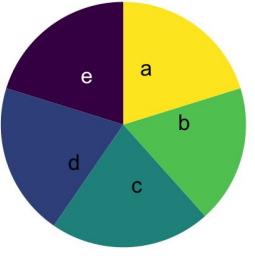
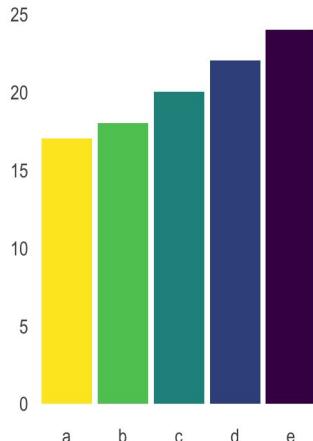
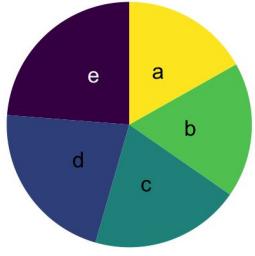


Figure 1 Plots of the female:male genetic-distance ratio against sex-averaged genetic location (in cM) along six selected chromosomes. Approximate locations of the centromeres are indicated by the triangles. The dashed lines correspond to equal female and male distances.

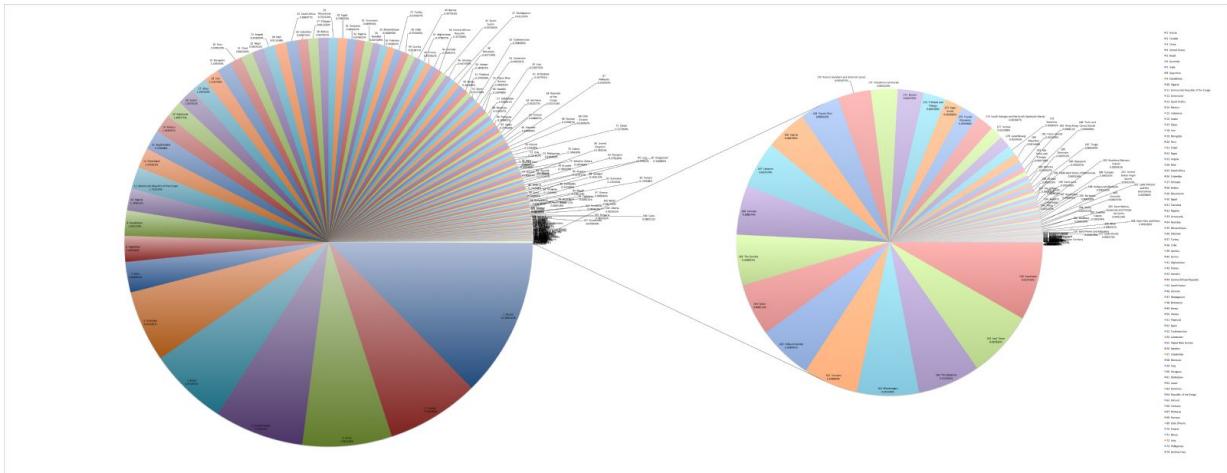
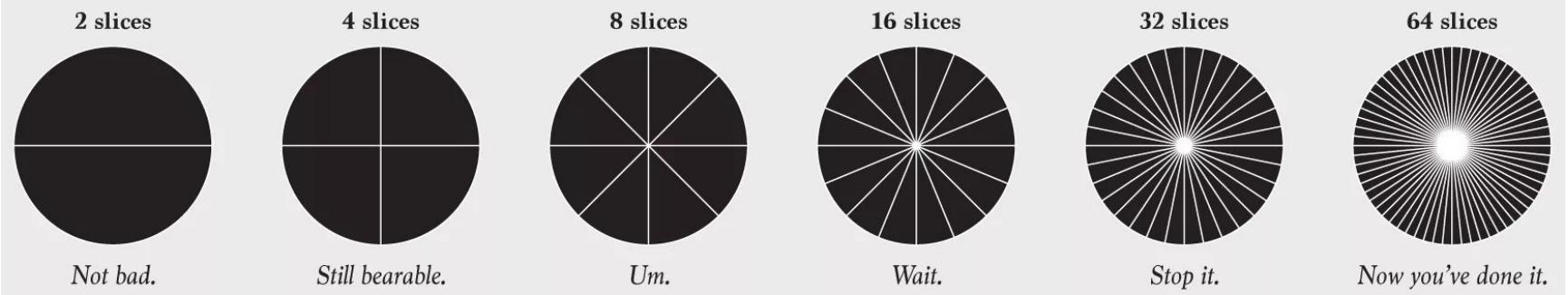
Visualization challenges – Overplotting



Visualization challenges – Pie charts - almost never a good idea



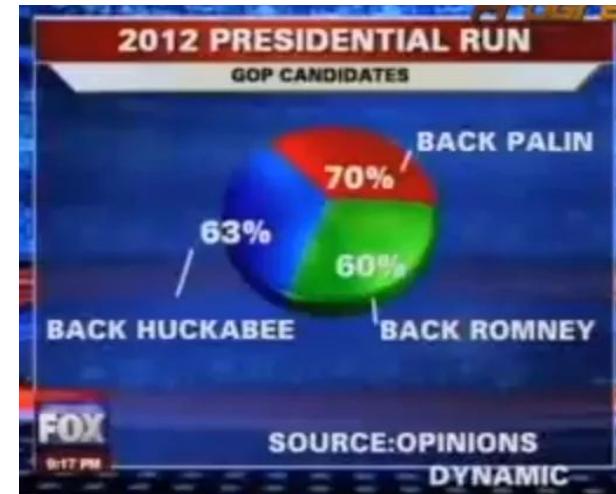
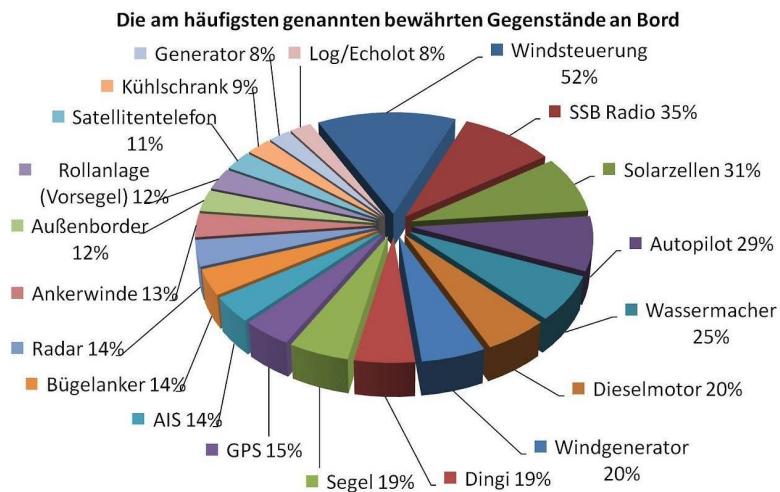
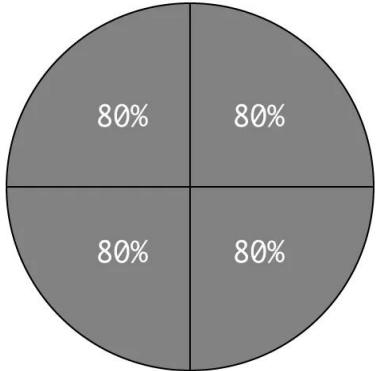
Visualization challenges – Pie charts - almost never a good idea



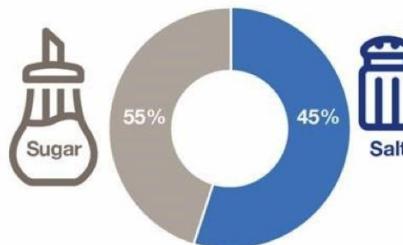
<https://flowingdata.com/2015/08/11/real-chart-rules-to-follow/>

https://commons.wikimedia.org/wiki/File:Pie_chart_of_countries_by_area.png

Visualization challenges – When things don't add up

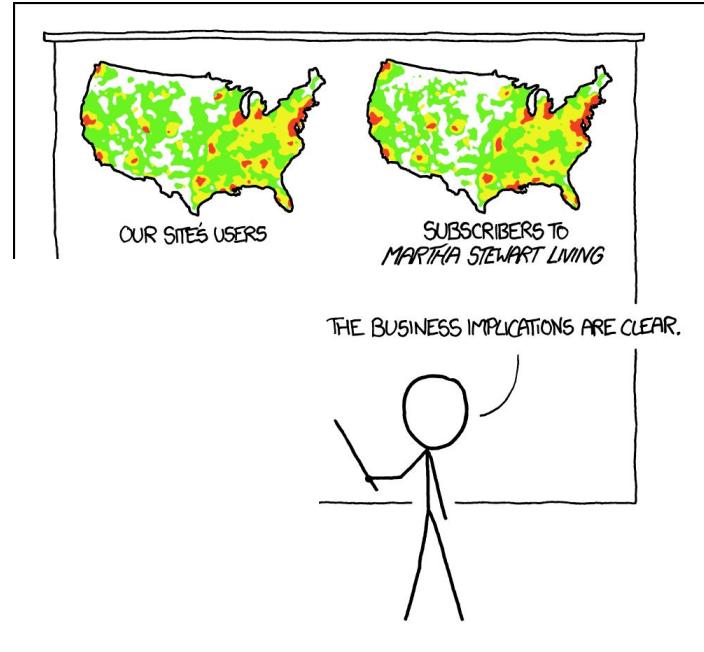
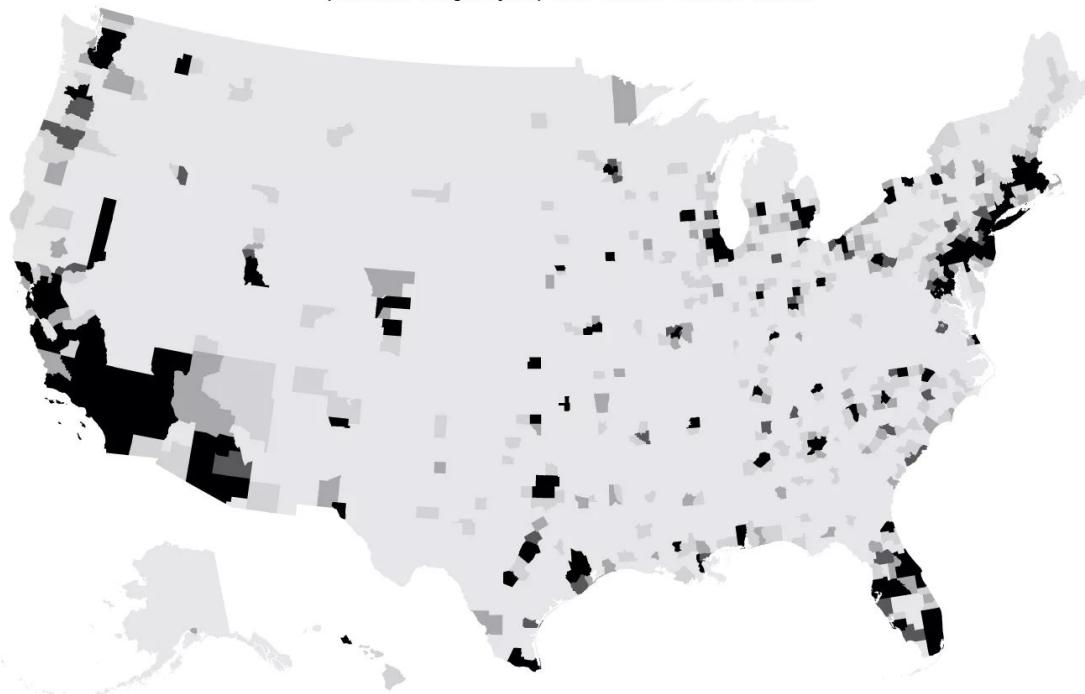


Last Week's Results
Which of these would you have a harder time giving up, salt or sugar?



Visualization challenges – Absolutes vs. relative values

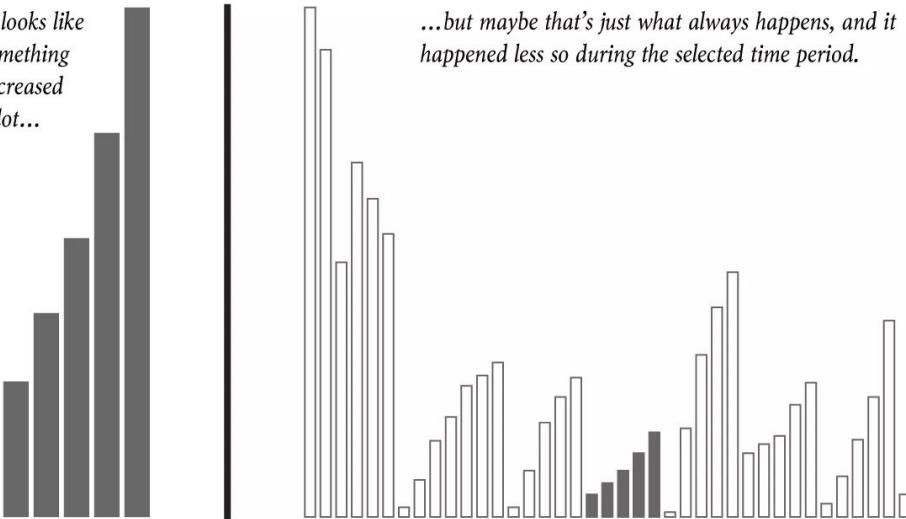
This is just population. When comparing across places, categories, or groups, you must compare fairly and consider relative values.



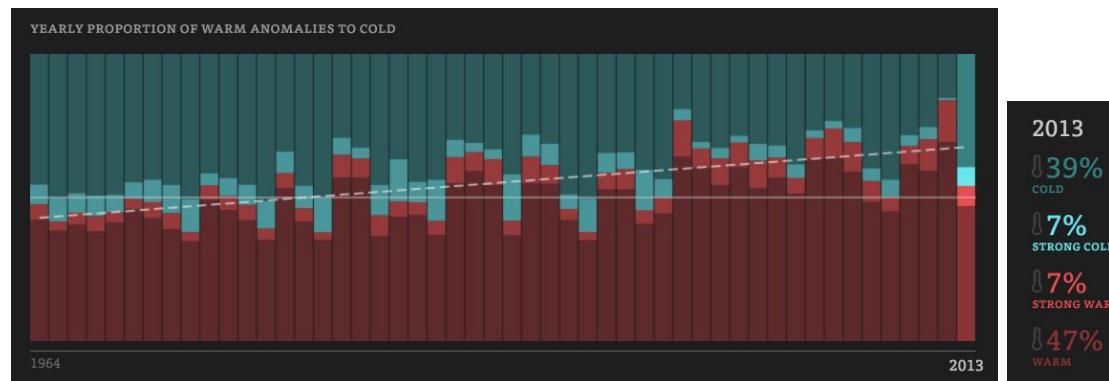
PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

Visualization challenges – Limited scope

It looks like something increased a lot...

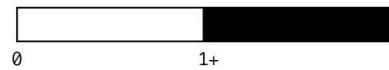


...but maybe that's just what always happens, and it happened less so during the selected time period.

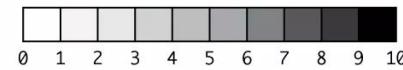


Visualization challenges – Choice of plot & data binning

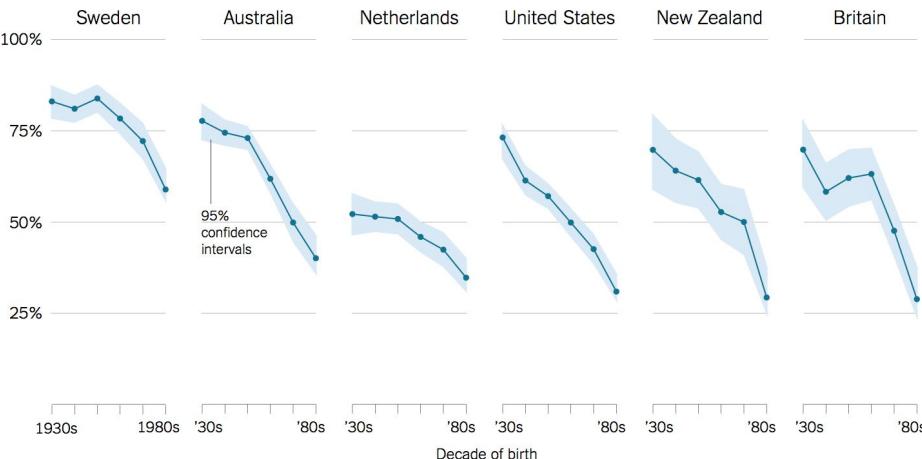
Two bins. What's really in the 1+ category?
Might be hiding something.



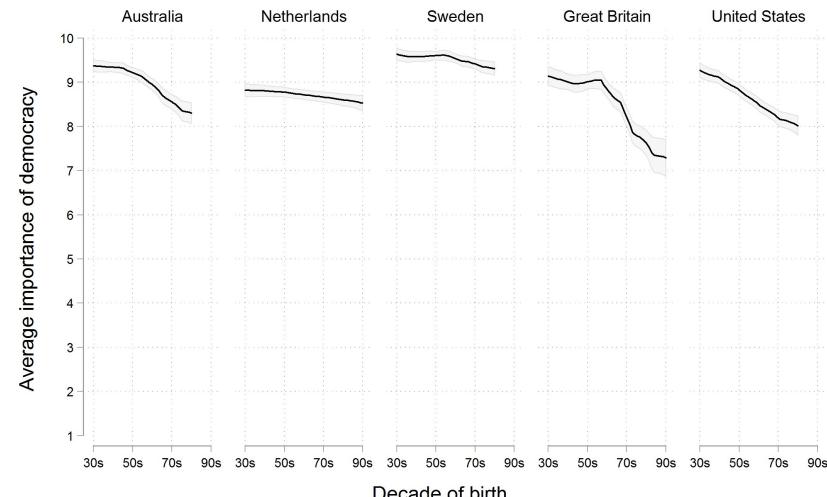
That's better. It can show more variation.



Percentage of people who say it is “essential” to live in a democracy



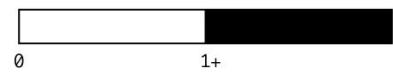
Source: Yascha Mounk and Roberto Stefan Foa, “The Signs of Democratic Deconsolidation,” Journal of Democracy | By The New York Times



Graph by Erik Voeten, based on WVS 5

Visualization challenges – Choice of plot & data binning

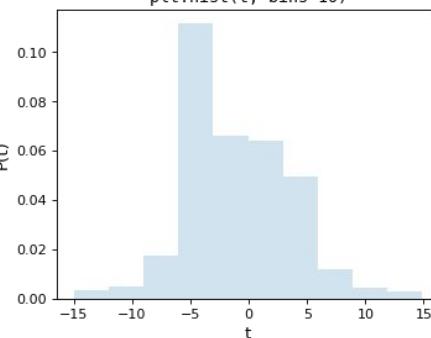
Two bins. What's really in the 1+ category?
Might be hiding something.



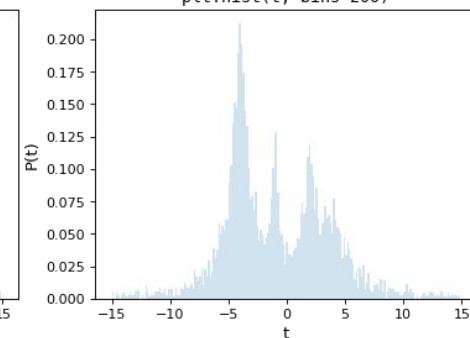
That's better. It can show more variation.



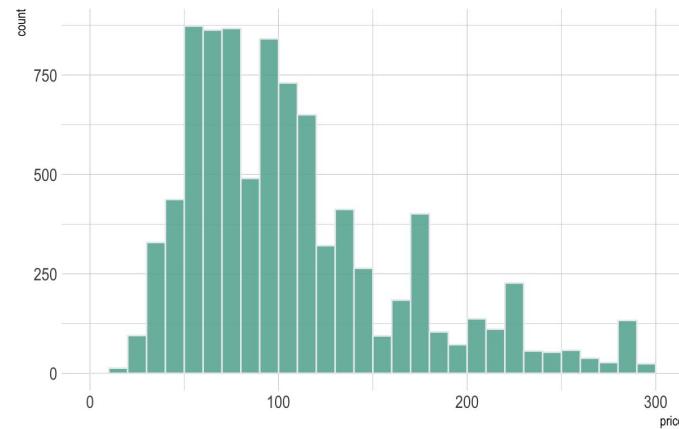
`plt.hist(t, bins=10)`



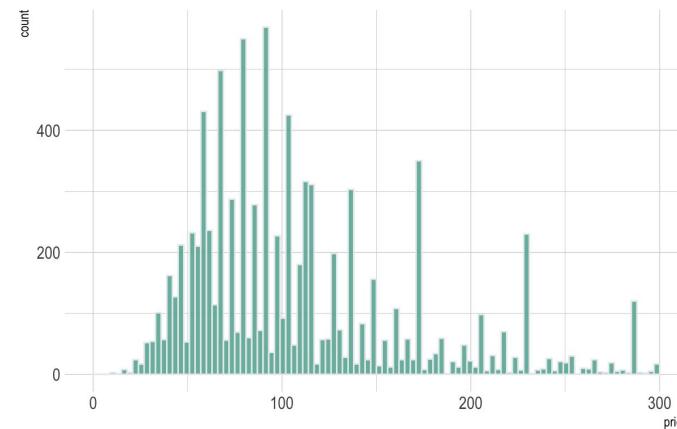
`plt.hist(t, bins=200)`



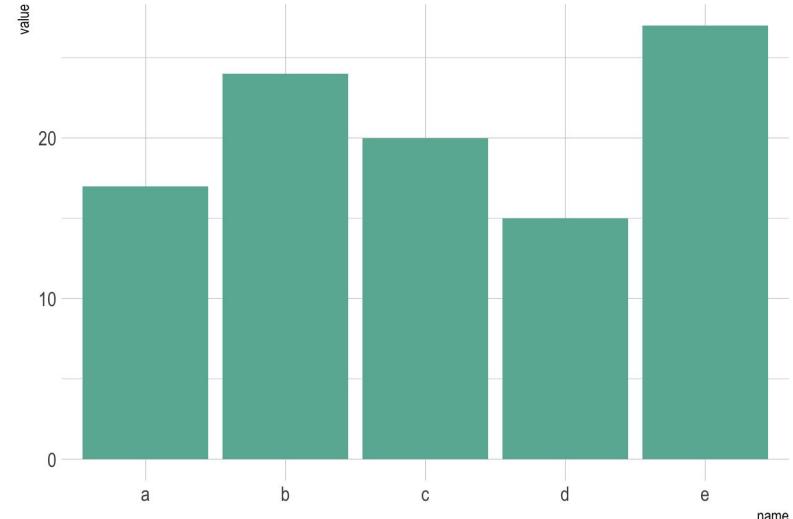
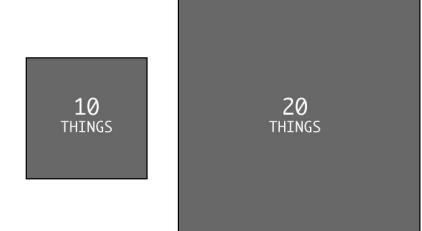
Night price distribution of Airbnb appartements



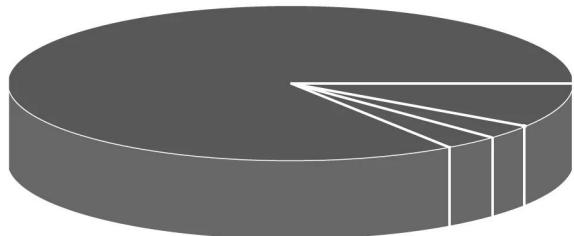
Night price distribution of Airbnb appartements



Visualization challenges – Area sized by dimension



Visualization challenges – Unwarranted extra dimensions



Distribution of All TFBS Regions

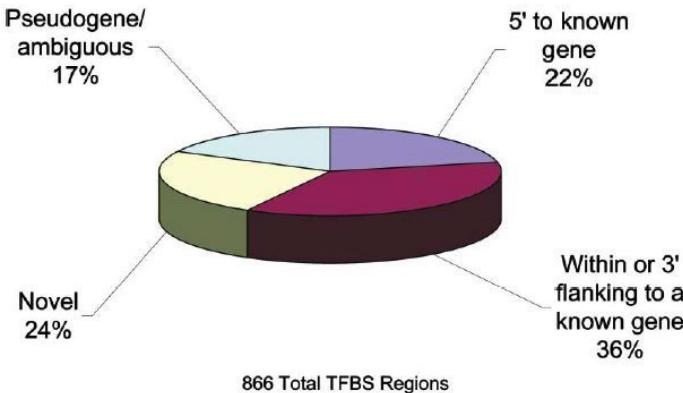


Figure 1. Classification of TFBS Regions

TFBS regions for Sp1, cMyc, and p53 were classified based upon proximity to annotations (RefSeq, Sanger hand-curated annotations, GenBank full-length mRNAs, and Ensembl predicted genes). The proximity was calculated from the center of each TFBS region. TFBS regions were classified as follows: within 5 kb of the 5' most exon of a gene, within 5 kb of the 3' terminal exon, or within a gene, novel or outside of any annotation, and pseudogene/ambiguous (TFBS overlapping or flanking pseudogene annotations, limited to chromosome 22, or TFBS regions falling into more than one of the above categories).

