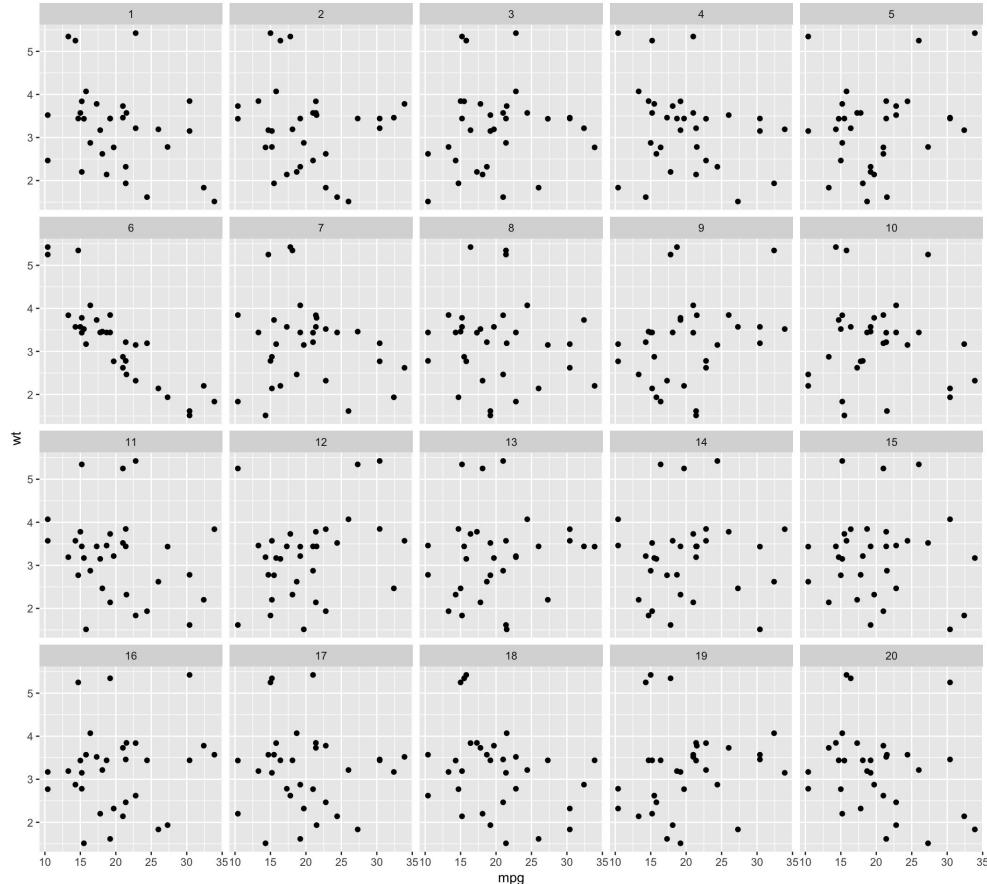


Visualization challenges

Day 08

- A survey of plots
- Issues with each type of plot

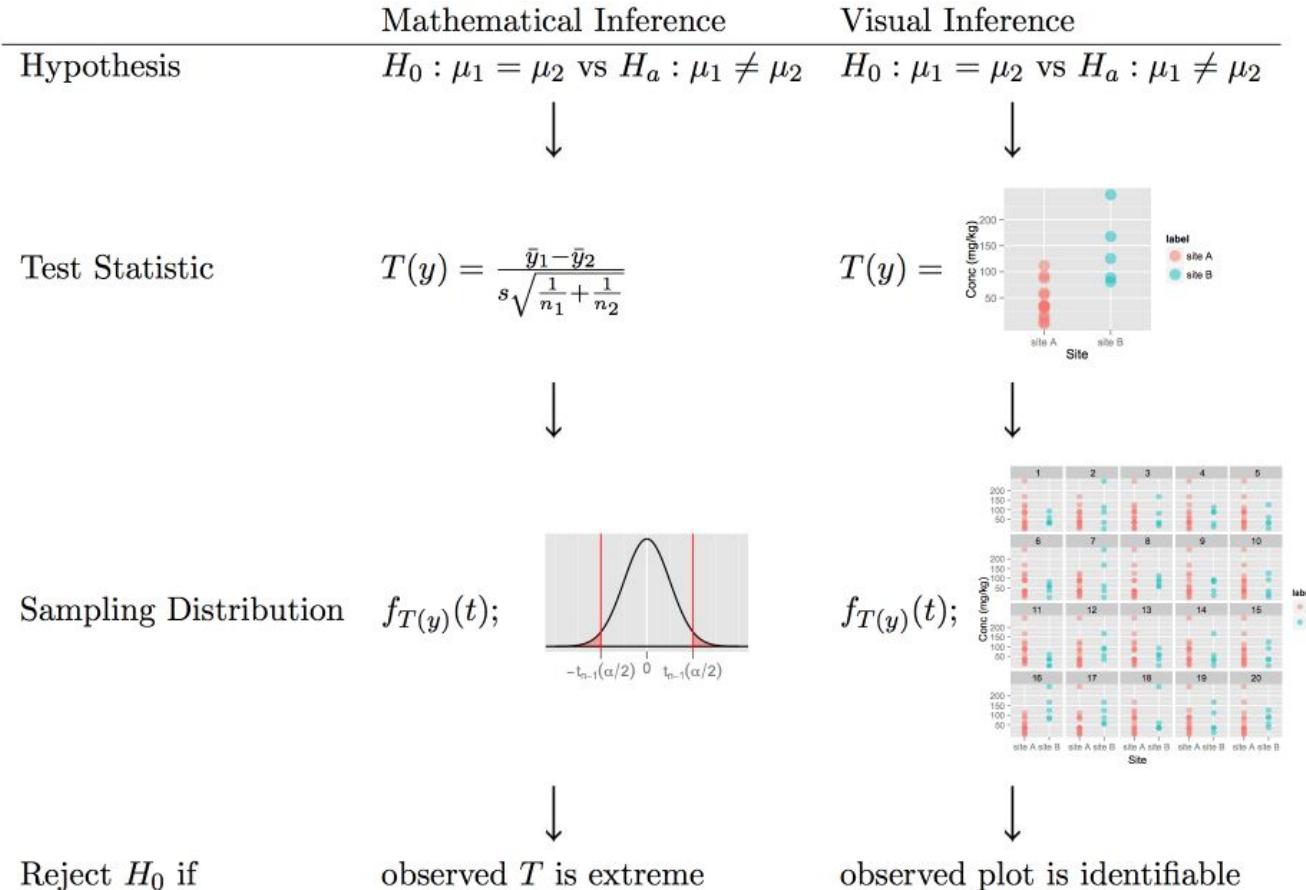
Spurious correlations – But it *looks* associated!



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.

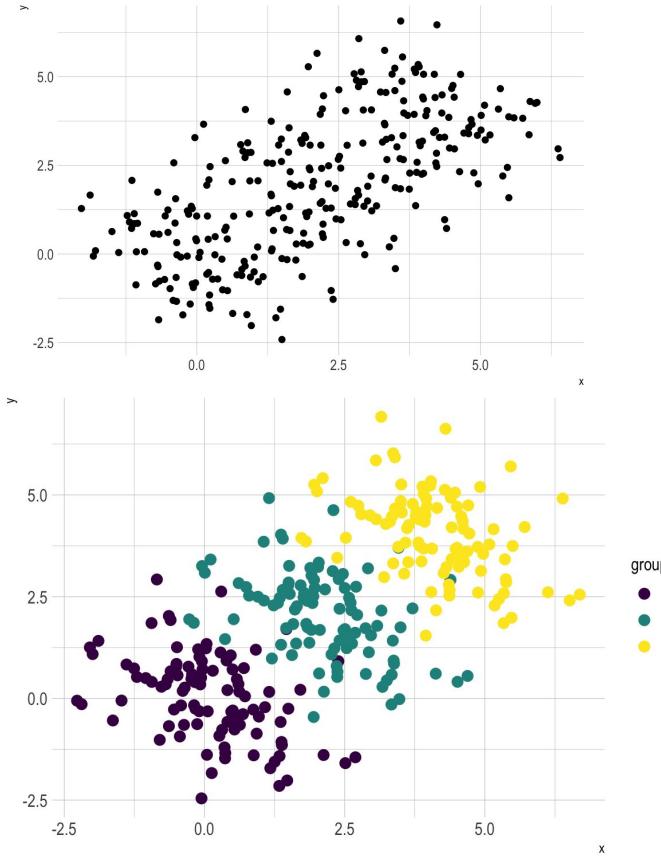
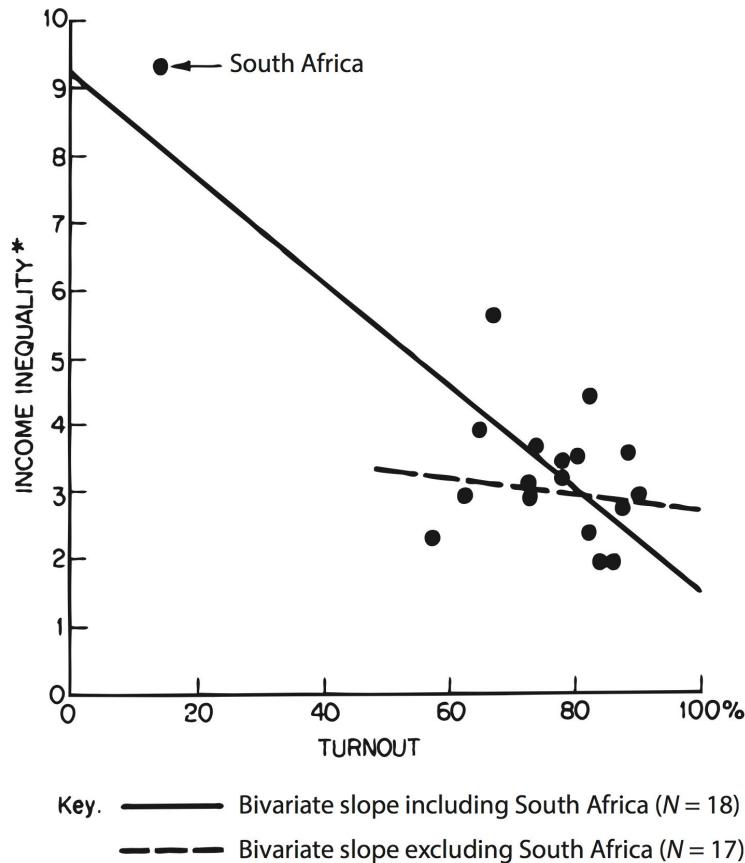
Spurious correlations – But it *looks* associated!



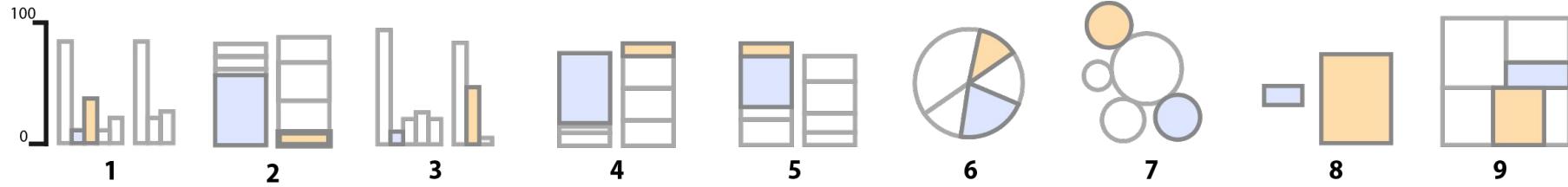
Purpose of data visualization

- Immediately convey information about study design.
- Illustrate important findings.
- Allow the reader to confirm that the statistical analysis is appropriate for the study design.
- Allow the reader to critically evaluate the data.

Visualization challenges – Don't do statistics without visualization



Visualization challenges



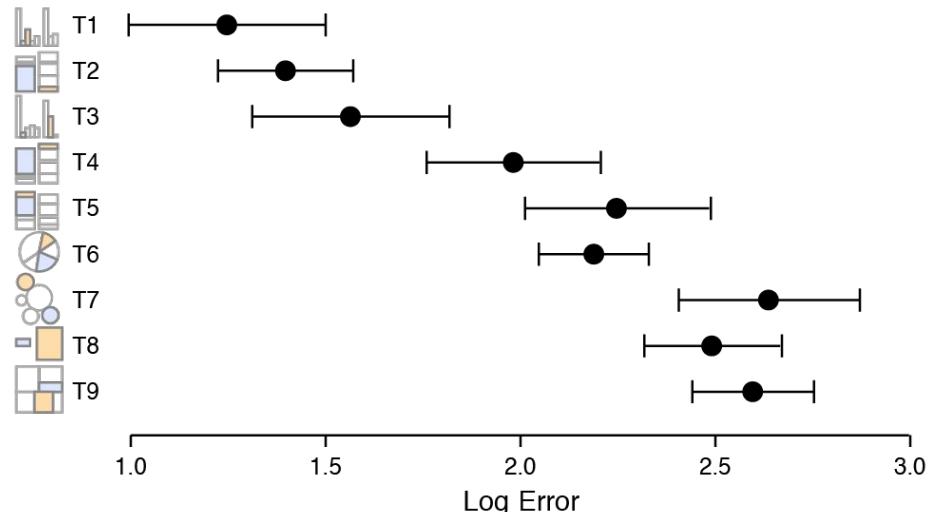
Position

Length

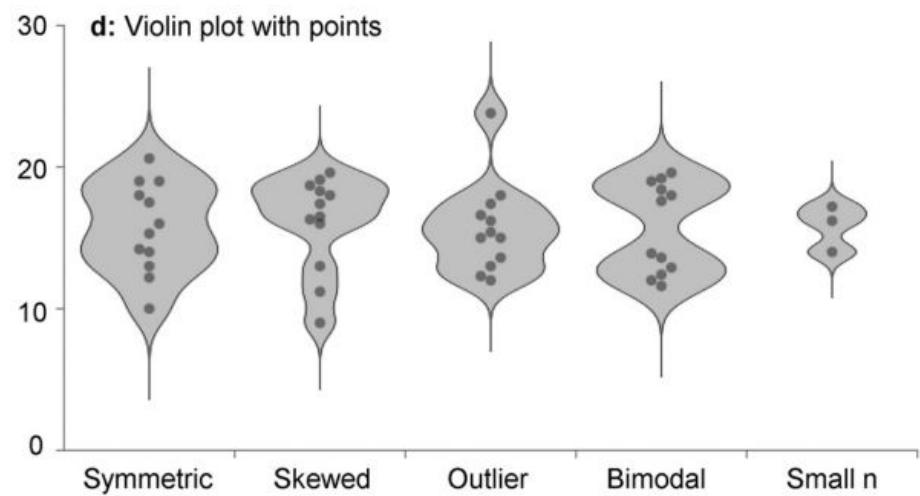
Angle

Circular Area

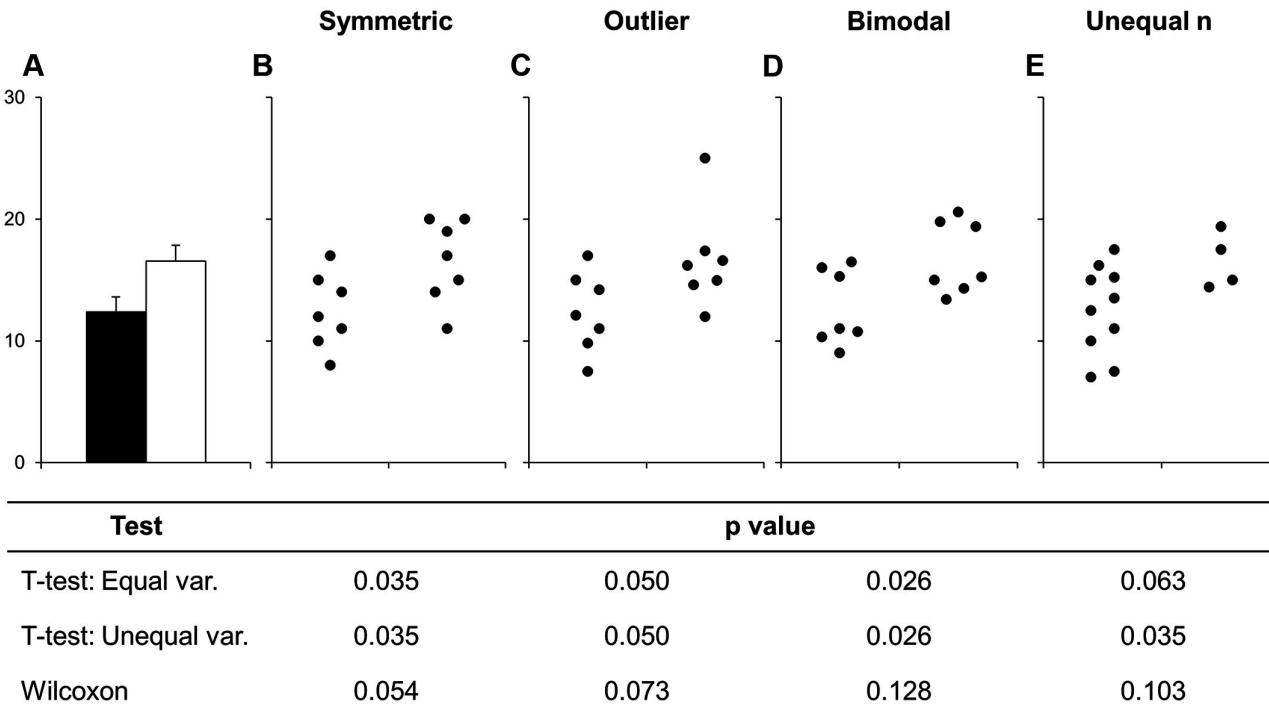
Rectangular Areas



Visualization challenges – Bar plots, error bars



Visualization challenges – Bar plots, error bars



Parametric vs.
Non-parametric test

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

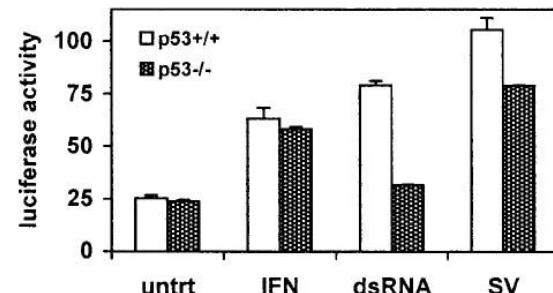
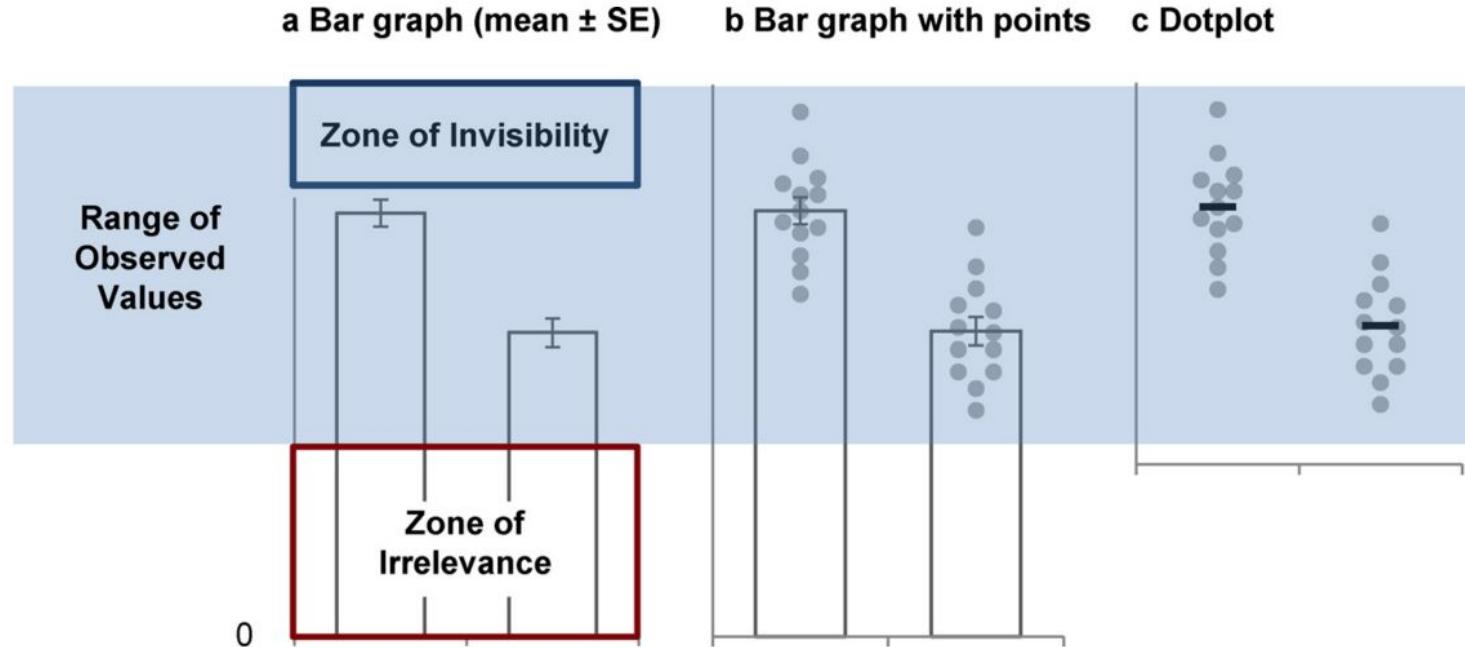
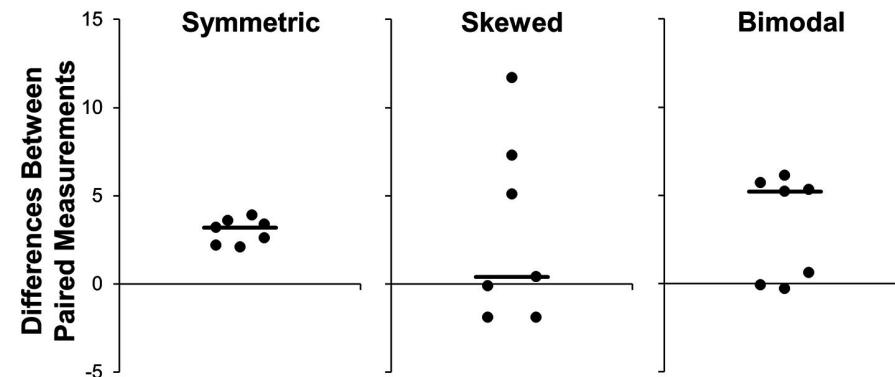
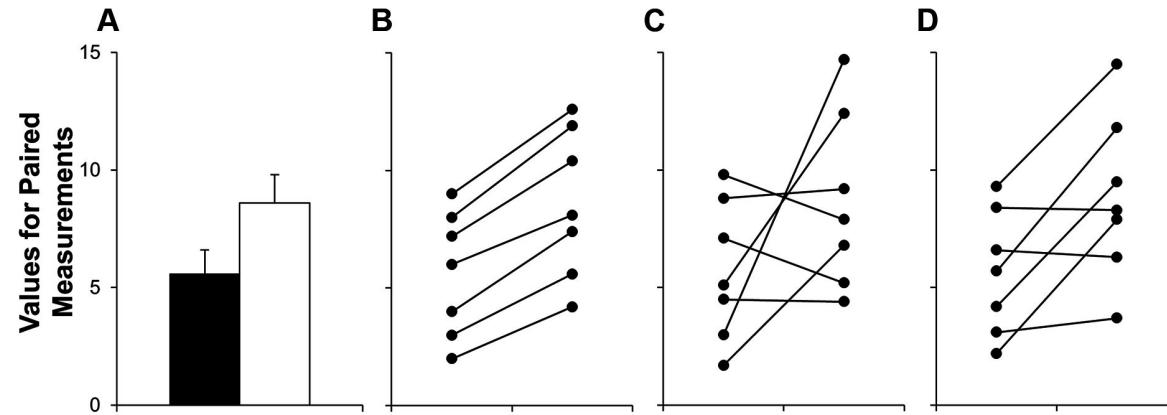


FIG. 4. ISG15 promoter activity mimics endogenous ISG15 mRNA regulation by p53, dsRNA, and virus. Cells (6×10^5 HCT 116) were

Visualization challenges – Bar plots, error bars

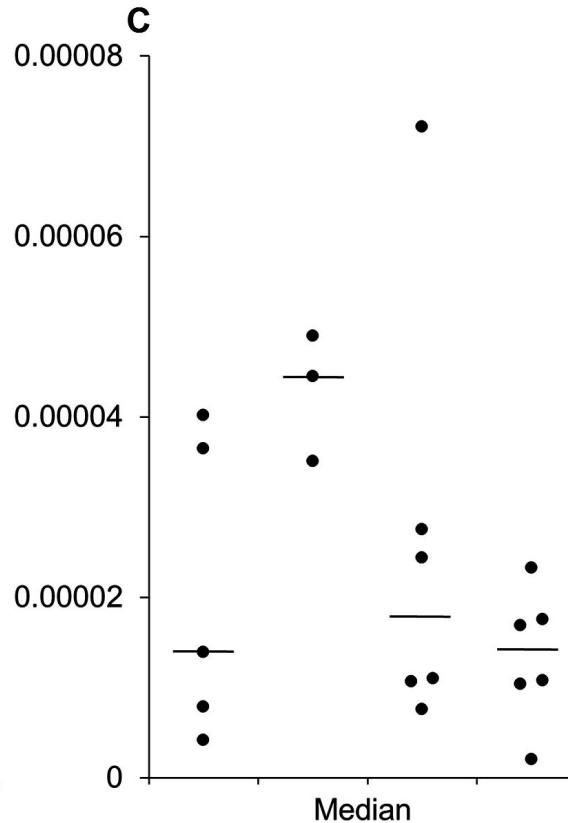
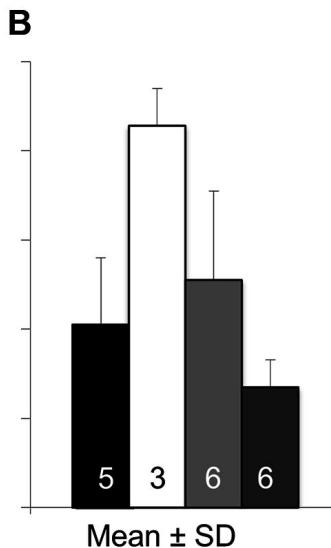
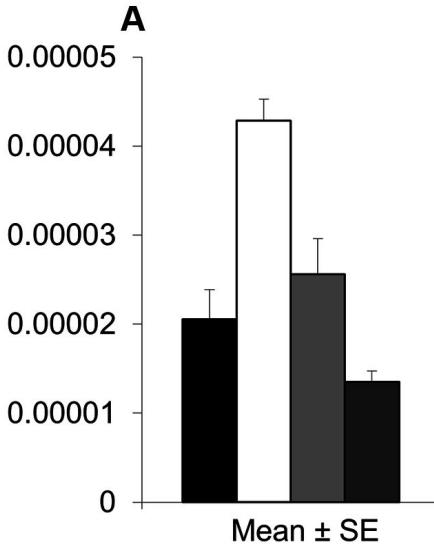


Visualization challenges – Bar plots, error bars

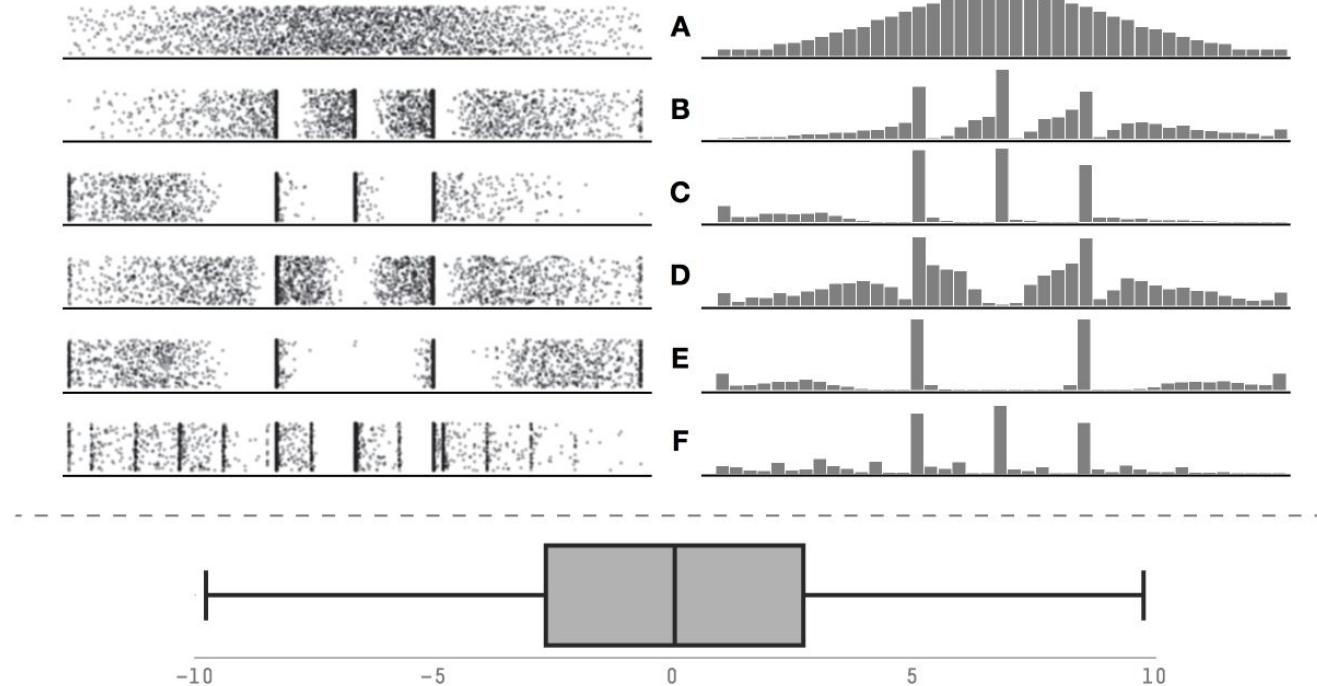


Visualization challenges – Bar plots, error bars

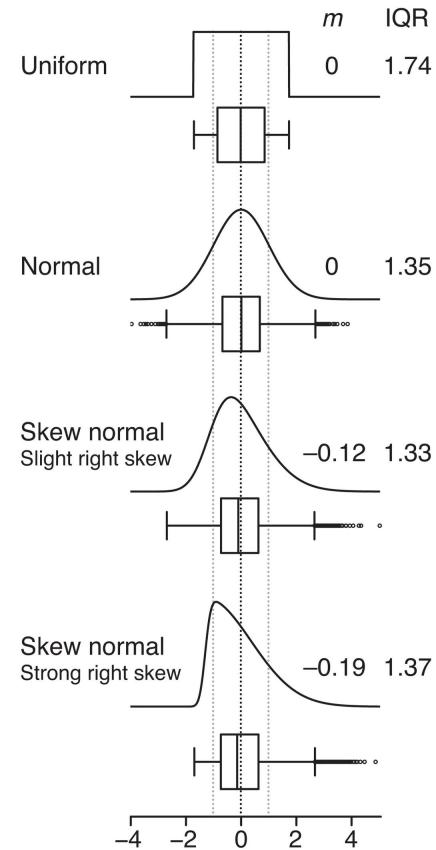
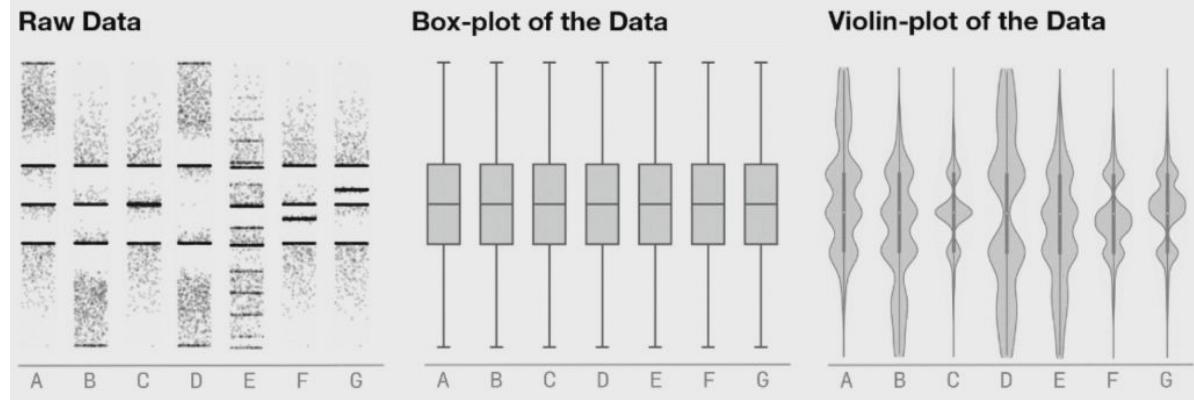
Underlying data is
inscrutable!



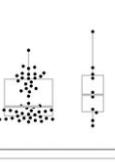
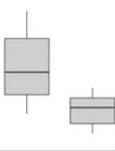
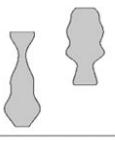
Visualization challenges – Different distributions



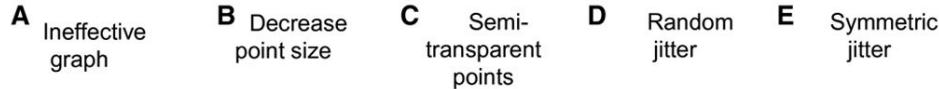
Visualization challenges – Different distributions



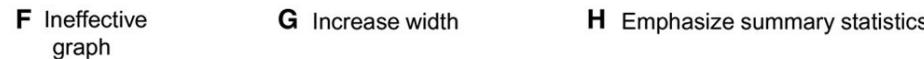
Visualization challenges – Different distributions

Figure Types	Example	Type of Variable	What the Plot Shows	Sample Size	Data Distribution	Best Practices
Dot plot		Continuous	Individual data points & mean or median line Other summary statistics (i.e. error bars) can be added for larger samples	Very small OR small; can also be useful with medium samples	Sample size is too small to determine data distribution OR Any data distribution	<ul style="list-style-type: none"> Make all data points visible - use symmetric jittering Many groups: Increase white space between groups, emphasize summary statistics & de-emphasize points Only add error bars if the sample size is large enough to avoid creating a false sense of certainty Avoid "histograms with dots"
Dot plot with box plot or violin plot		Continuous	Combination of dot plot & box plot or violin plot (see descriptions above and below)	Medium	Any	<ul style="list-style-type: none"> Make all data points visible (symmetric jittering) Smaller n: Emphasize data points and de-emphasize box plot, delete box plot and show only median line for groups with very small n Larger n: Emphasize box plot and de-emphasize points
Box plot		Continuous	Horizontal lines on box: 75 th , 50 th (median) and 25 th percentile Whiskers: varies; often most extreme data points that are not outliers Dots above or below whiskers: outliers	Large	Do not use for bimodal data	<ul style="list-style-type: none"> List sample size below group name on x-axis Specify what whiskers represent in legend
Violin plot		Continuous	Gives an estimated outline of the data distribution. The precision of the outline increases with increasing sample size.	Large	Any	<ul style="list-style-type: none"> List sample size below group name on x-axis The violin plot should not include biologically impossible values
Bar graph		Counts or proportions	Bar height shows the value of the count or proportion	Any	Any	<ul style="list-style-type: none"> Do not use for continuous data

Visualization challenges – Different distributions



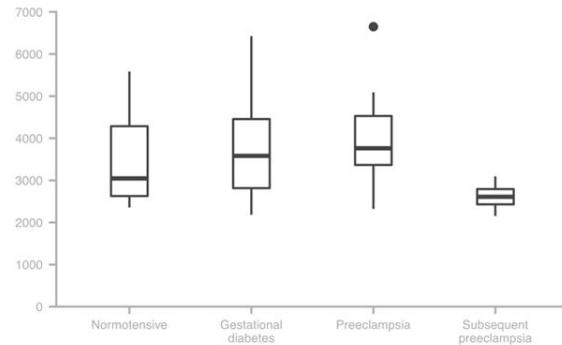
Step 1:
Make all
points
visible



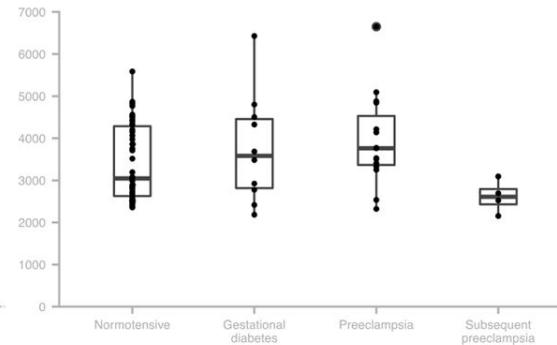
Step 2:
Emphasize
summary
statistics

Visualization challenges – Different distributions

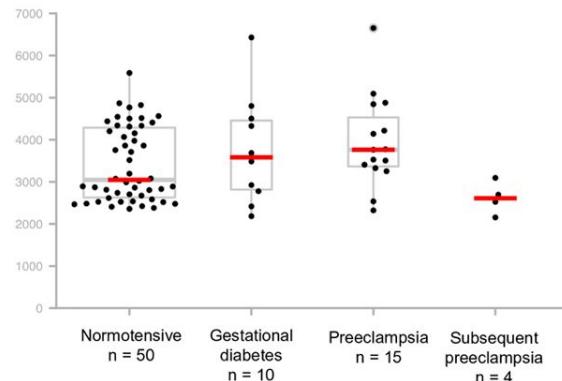
A Box plot



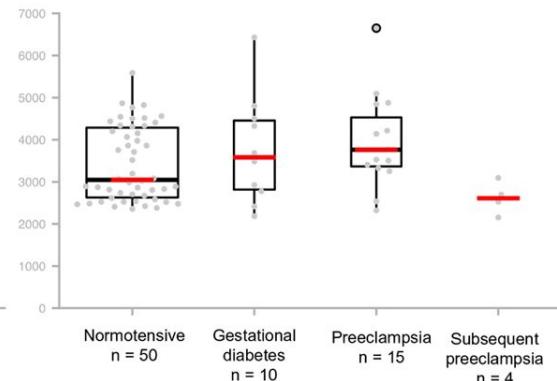
B Box plot with unjittered dot plot (strip plot)



C Emphasizing the dot plot



D Emphasizing the box plot



Visualization challenges – Reflecting study design

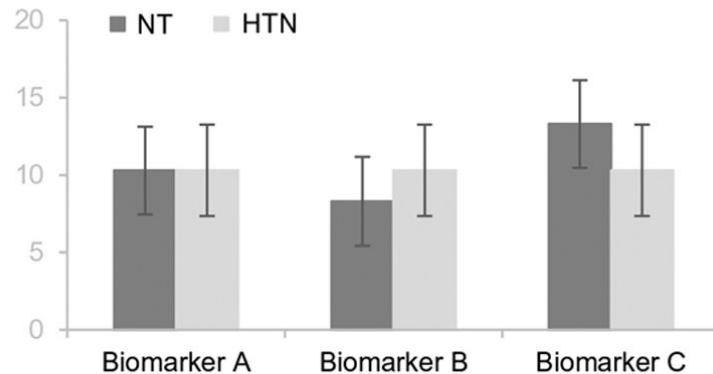
Experimental goal: Compare normotensive (NT) vs. hypertensive (HTN) patients

Statistical analysis: t-tests were used to compare values for each dependent variable (biomarker A, B and C)

A

Sending mixed messages

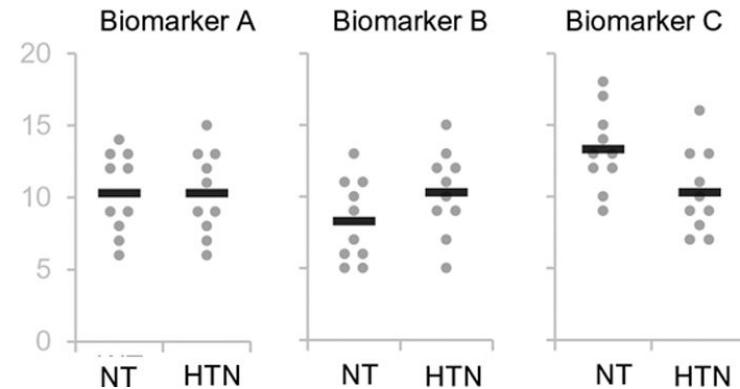
Figure structure erroneously suggests that authors also intended to compare biomarkers A, B and C



B

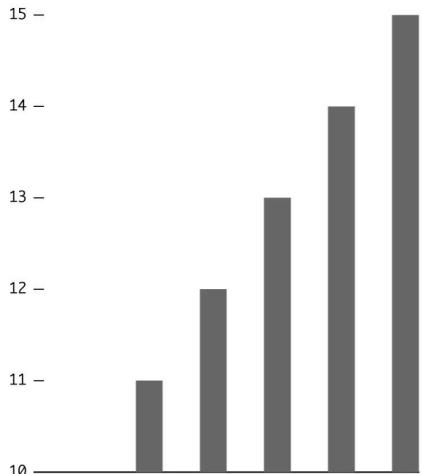
Clear communication

Figure structure matches study design & analysis, shows that the authors did not intend to compare biomarkers

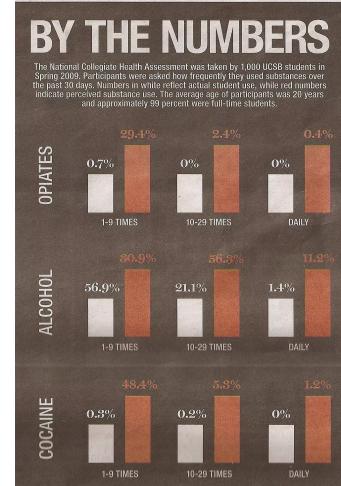
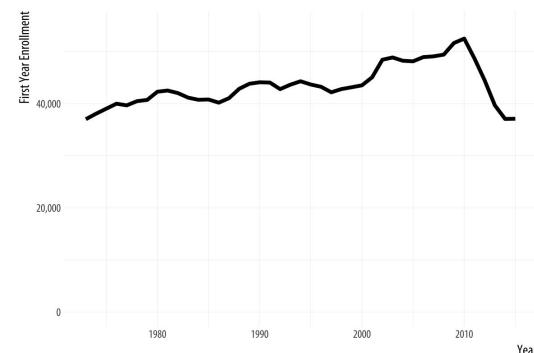
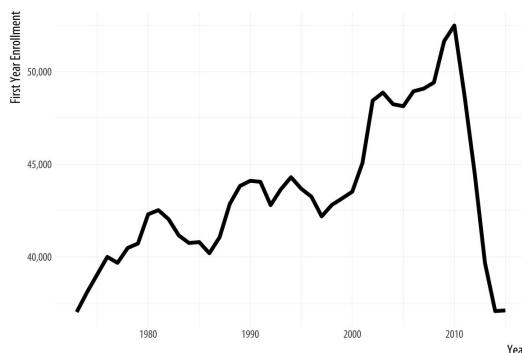
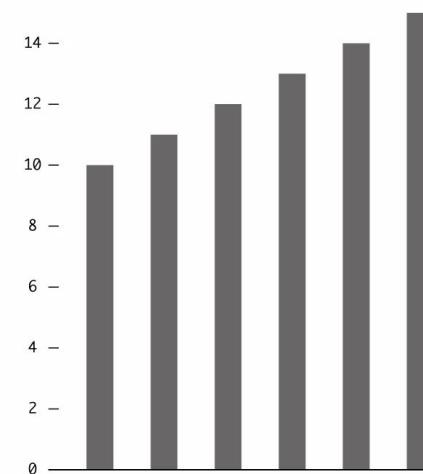


Visualization challenges – Truncated axis

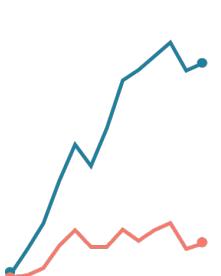
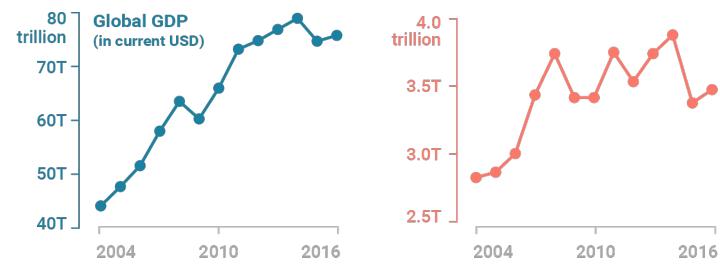
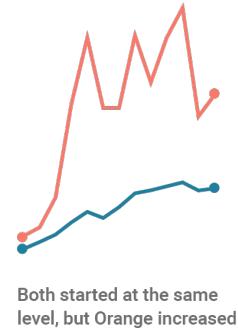
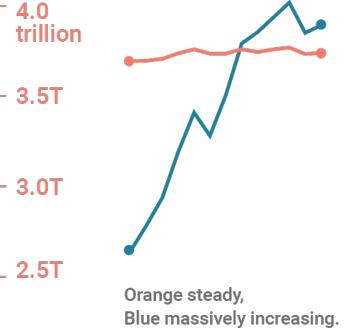
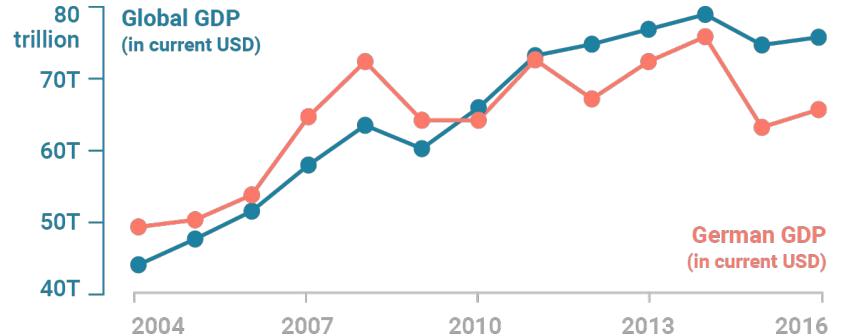
The value axis starts at ten. Liar, liar, pants on fire.



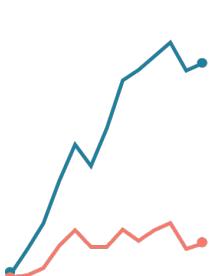
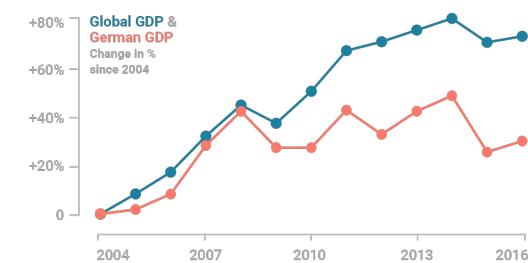
The value axis starts at zero. Good.



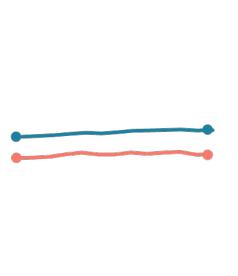
Visualization challenges – Dual axes



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.



Both steady.

Visualization challenges – Inappropriate axes

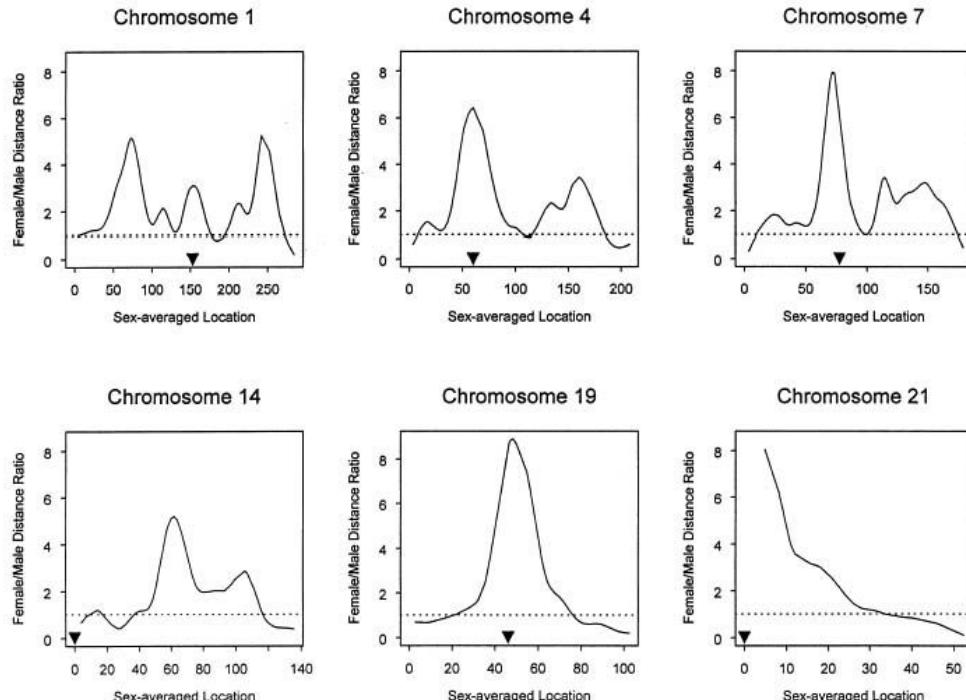
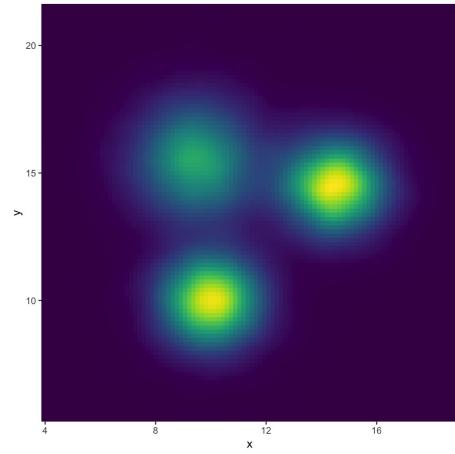
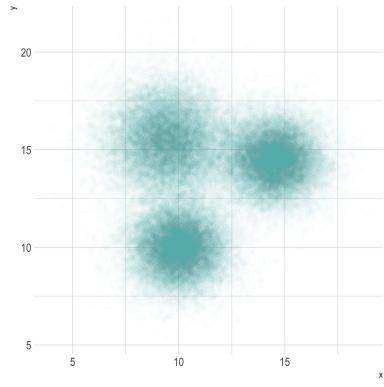
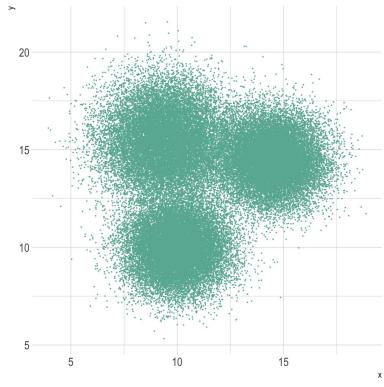
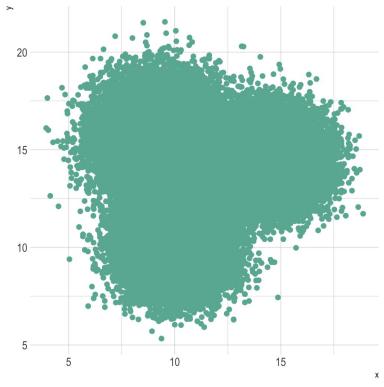
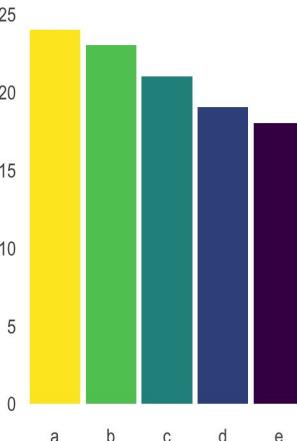
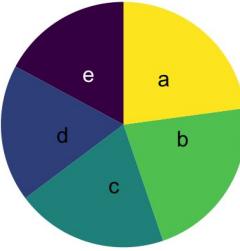
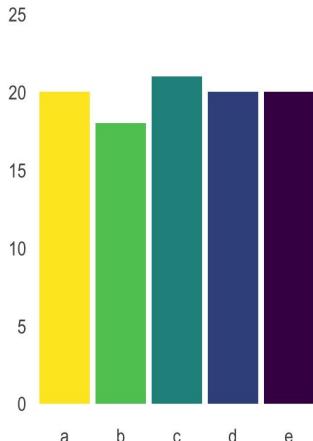
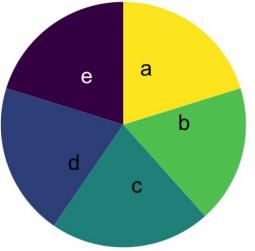
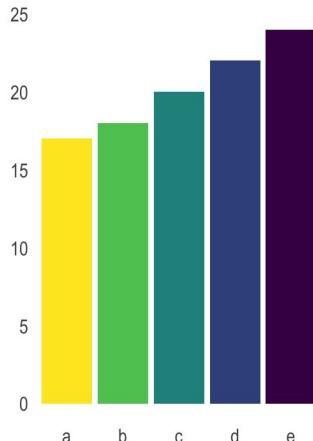
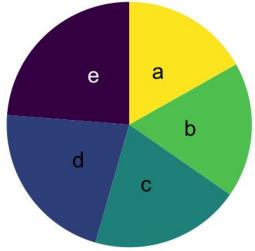


Figure 1 Plots of the female:male genetic-distance ratio against sex-averaged genetic location (in cM) along six selected chromosomes. Approximate locations of the centromeres are indicated by the triangles. The dashed lines correspond to equal female and male distances.

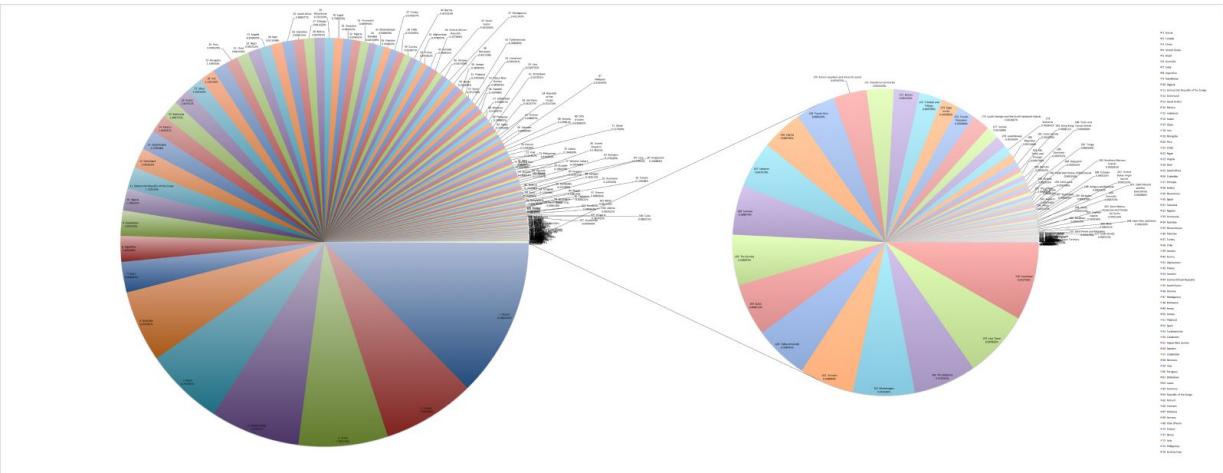
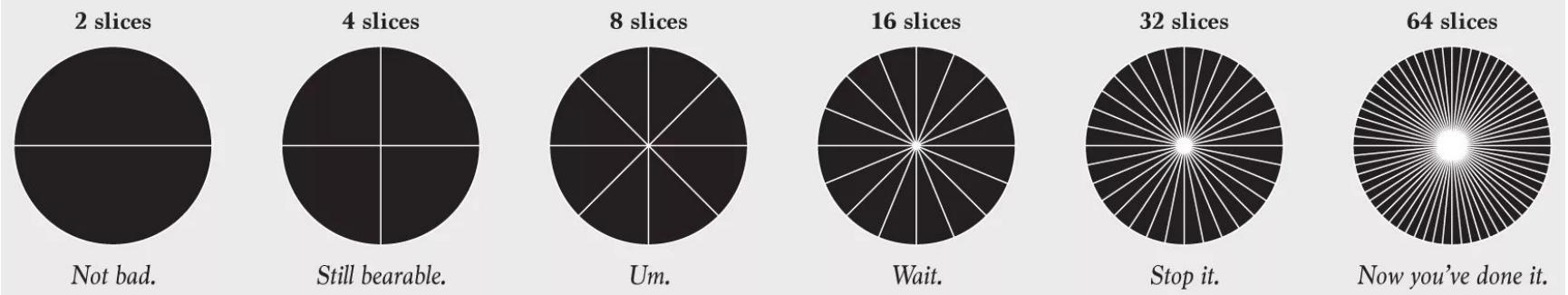
Visualization challenges – Overplotting



Visualization challenges – Pie charts - almost never a good idea



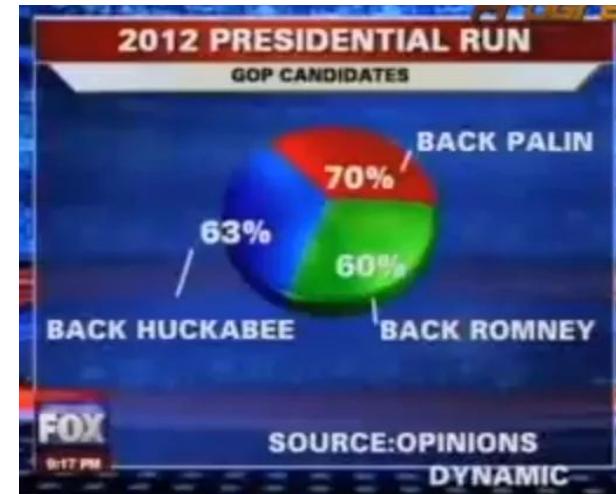
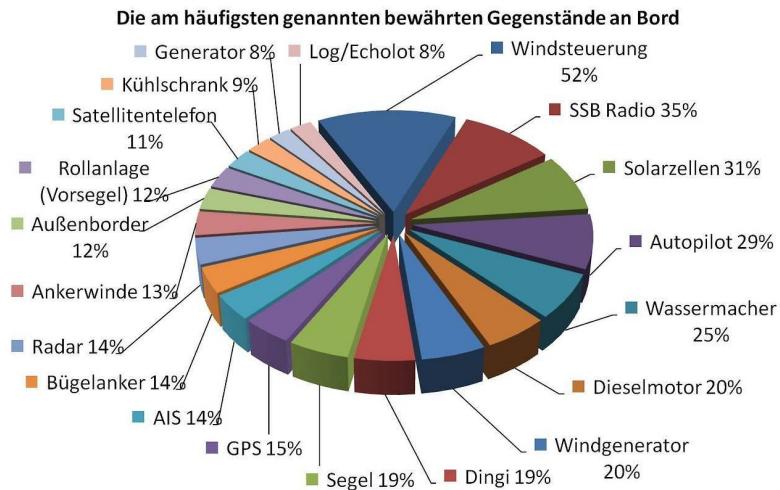
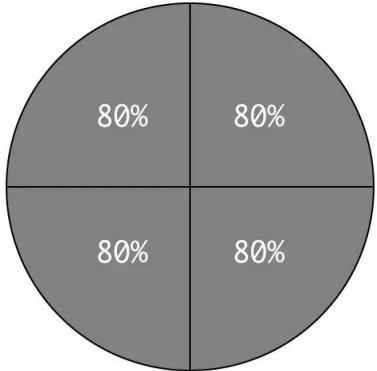
Visualization challenges – Pie charts - almost never a good idea



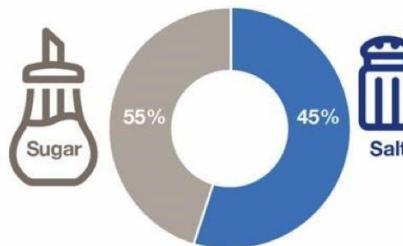
<https://flowingdata.com/2015/08/11/real-chart-rules-to-follow/>

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

Visualization challenges – When things don't add up

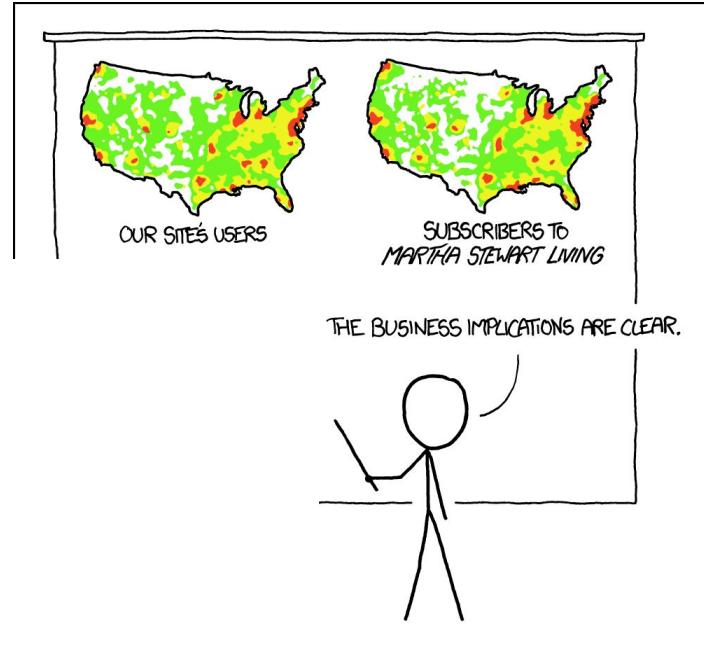
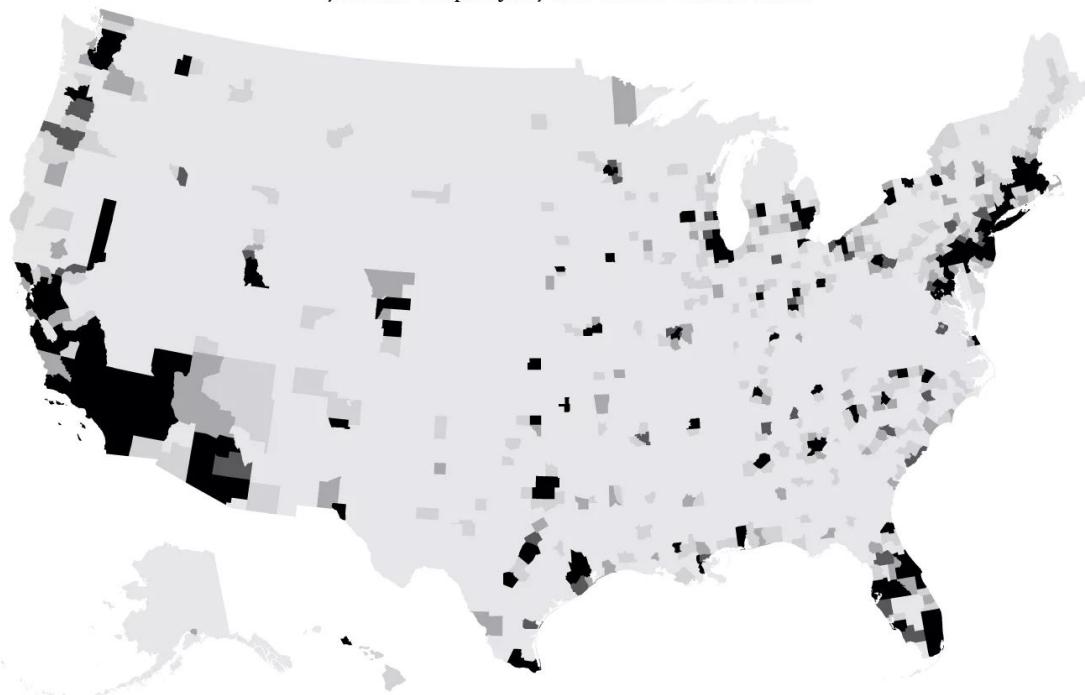


Last Week's Results
Which of these would you have a harder time giving up, salt or sugar?



Visualization challenges – Absolutes vs. relative values

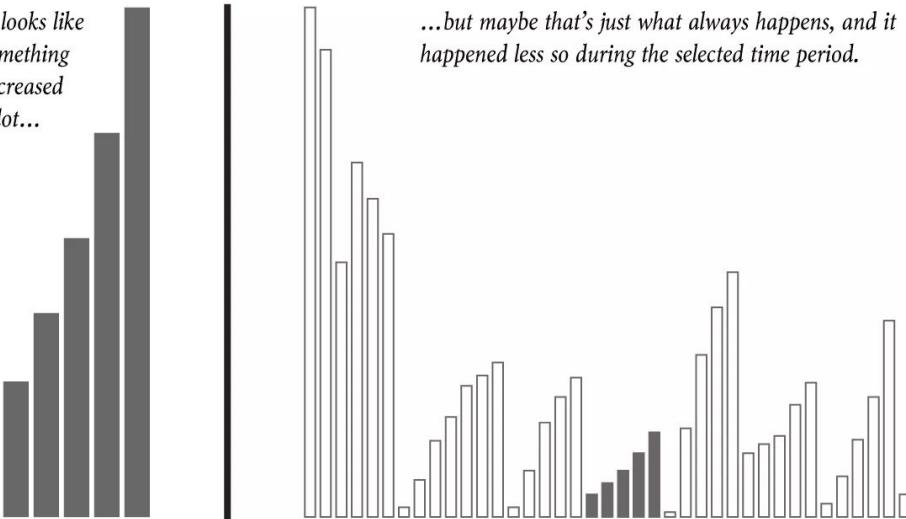
This is just population. When comparing across places, categories, or groups, you must compare fairly and consider relative values.



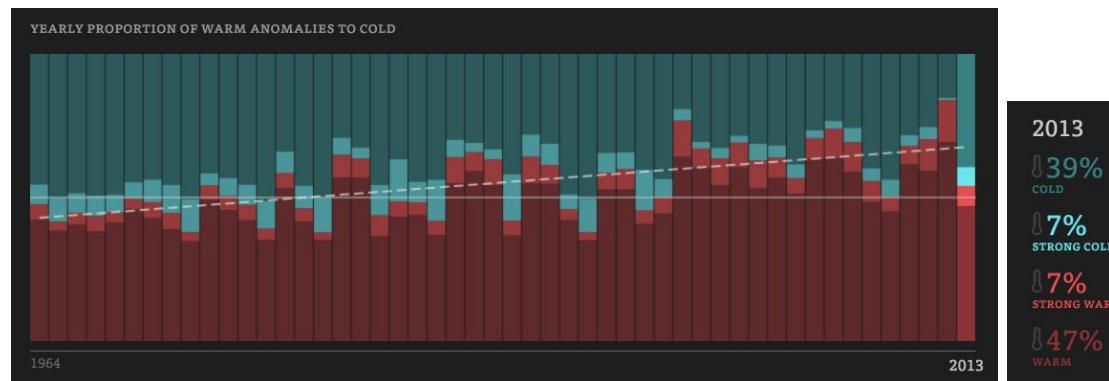
PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

Visualization challenges – Limited scope

It looks like something increased a lot...

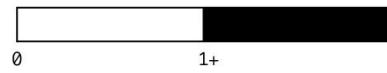


...but maybe that's just what always happens, and it happened less so during the selected time period.

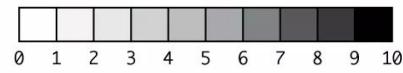


Visualization challenges – Choice of plot & data binning

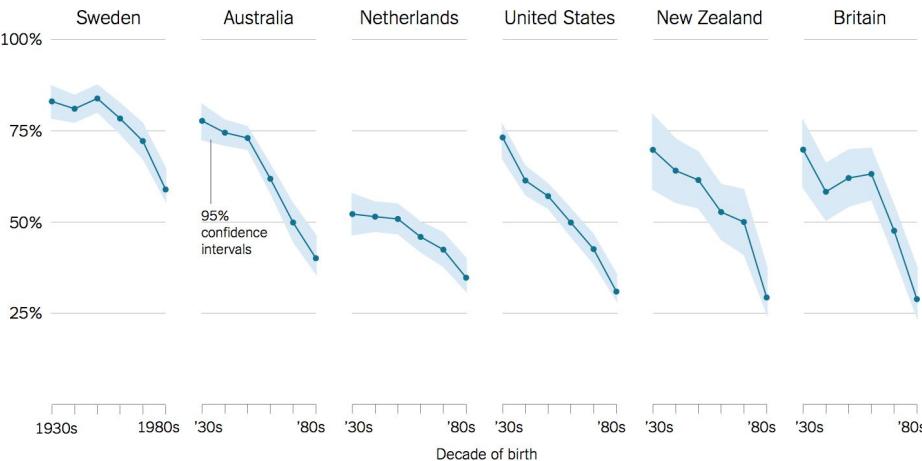
Two bins. What's really in the 1+ category?
Might be hiding something.



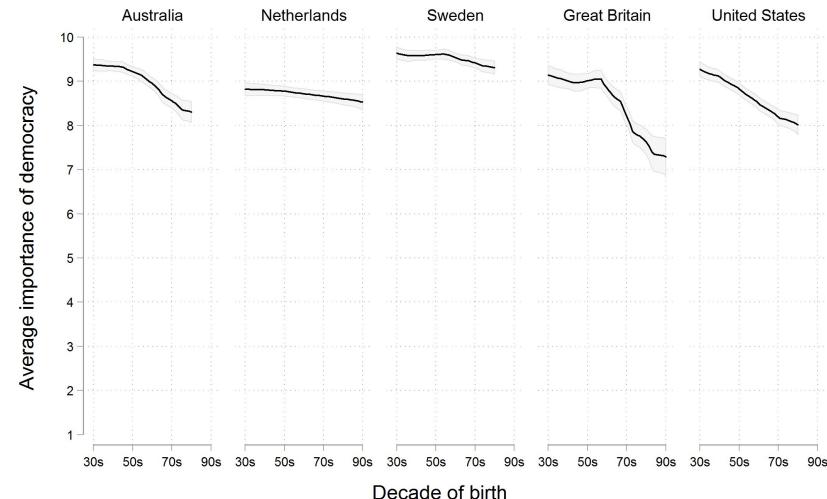
That's better. It can show more variation.



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



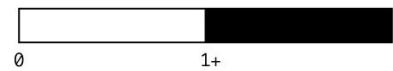
Source: Yascha Mounk and Roberto Stefan Foa, “The Signs of Democratic Deconsolidation,” Journal of Democracy | By The New York Times



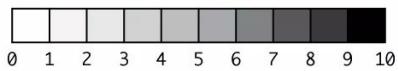
Graph by Erik Voeten, based on WVS 5

Visualization challenges – Choice of plot & data binning

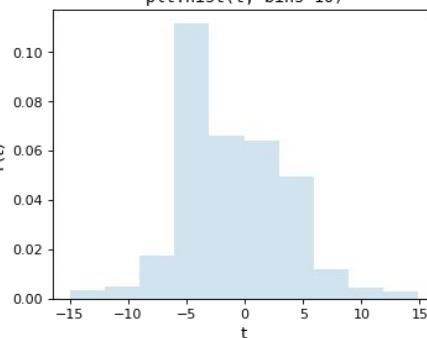
Two bins. What's really in the 1+ category?
Might be hiding something.



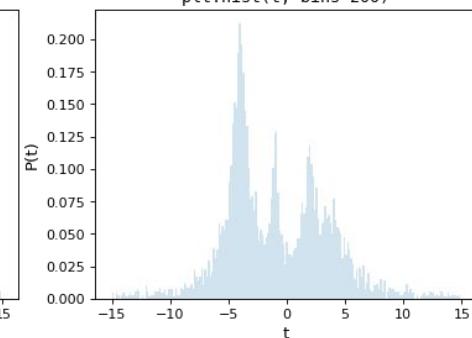
That's better. It can show more variation.



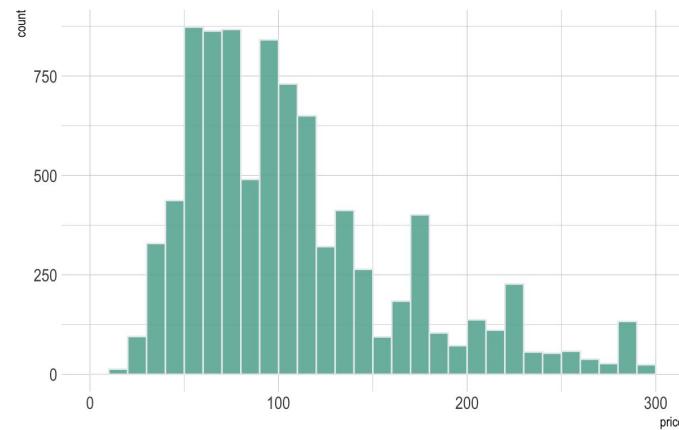
`plt.hist(t, bins=10)`



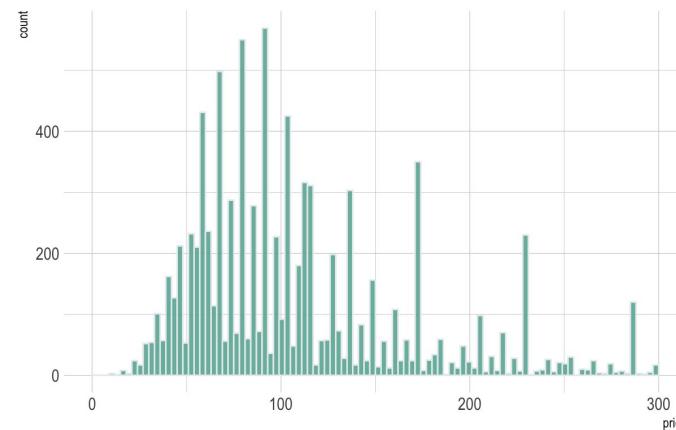
`plt.hist(t, bins=200)`



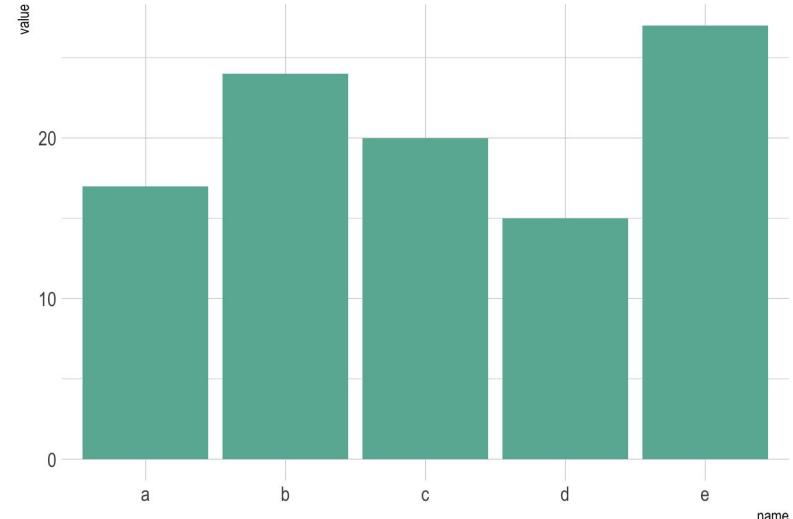
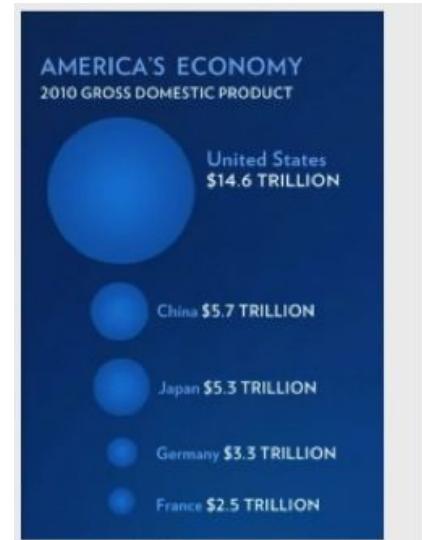
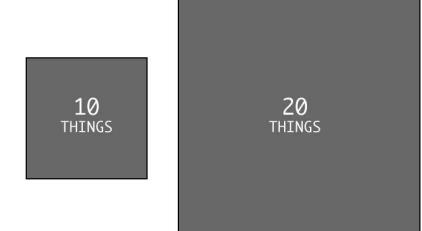
Night price distribution of Airbnb appartements



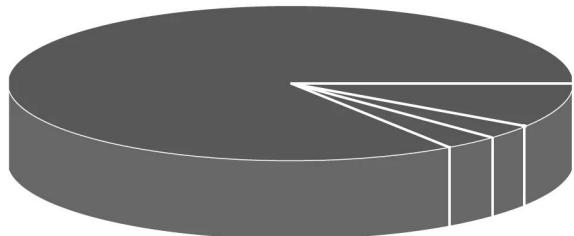
Night price distribution of Airbnb appartements



Visualization challenges – Area sized by dimension



Visualization challenges – Unwarranted extra dimensions



Distribution of All TFBS Regions

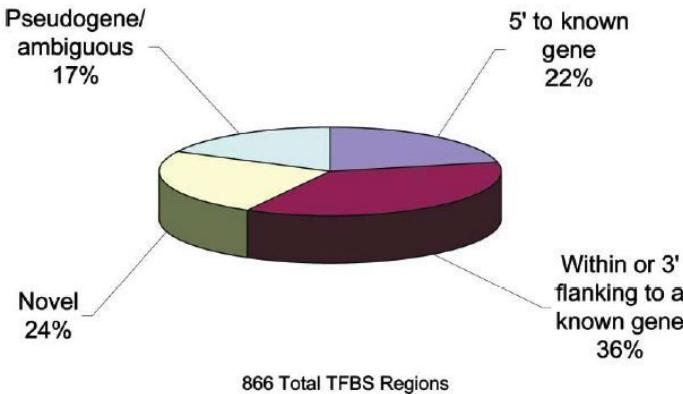
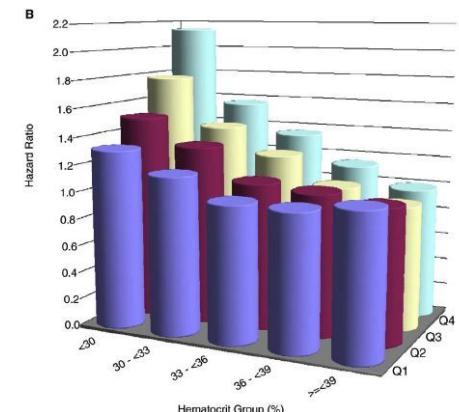
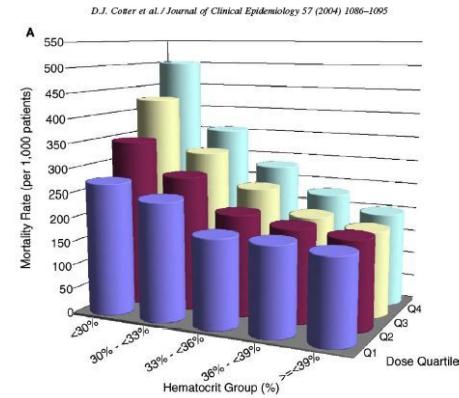


Figure 1. Classification of TFBS Regions

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).



Final exam

Select a primary research article published by you/your-group in the past 3–5 years that:

- Has a specific question (exploratory or hypothesis-driven).
- Contains experimental data (newly generated or previously published).
- Contains a “Methods” section.
- Involves at least **3 different results** – figures/tables in the main paper – obtained based on statistical data analysis. [This is critical.]

[Check with me to make sure the paper has enough statistical stuff to dig in for a final exam!]

Final exam

Carefully read the paper to dissect the statistical analyses and visualization choices that went into each of these results.

Remember that details about a single analysis/result might (unfortunately) be strewn across the entire paper:

- Results section
- Methods section
- Figure/table legend
- Supplementary materials/files, and
- Even older papers cited here!

For each result, take stock of all this information somewhere at least as rough notes, and bring them to the Final exam.

Complete online evaluation

<https://sirsonline.msu.edu>

I will also send you a survey.