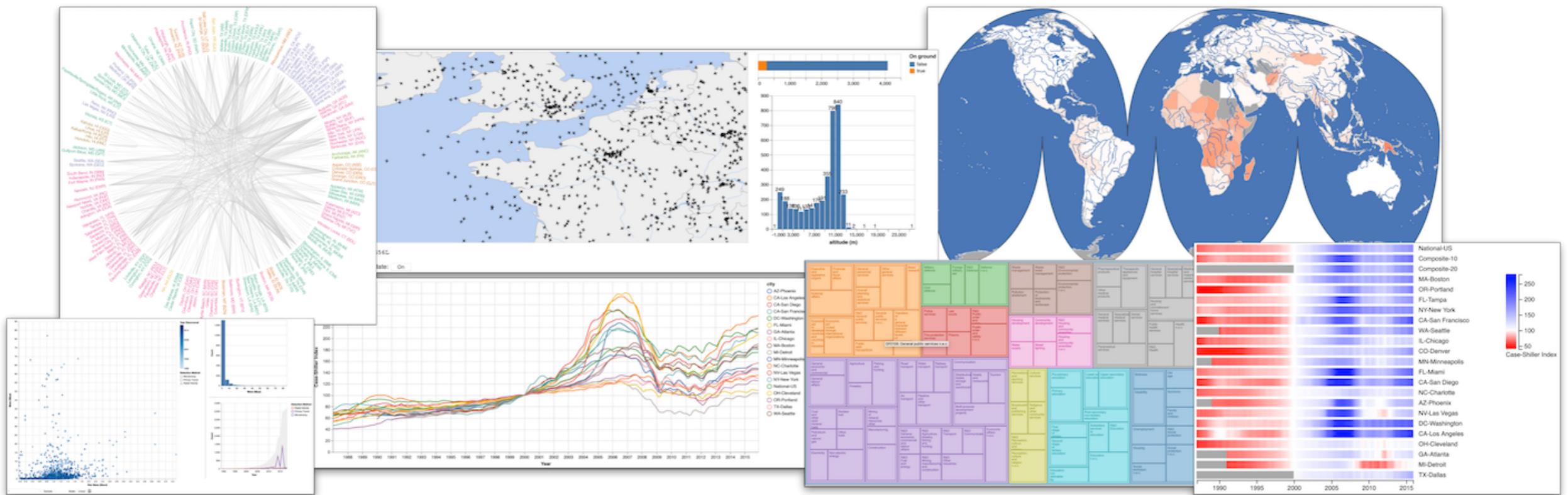


Data Visualization

INF552 (2023-2024)

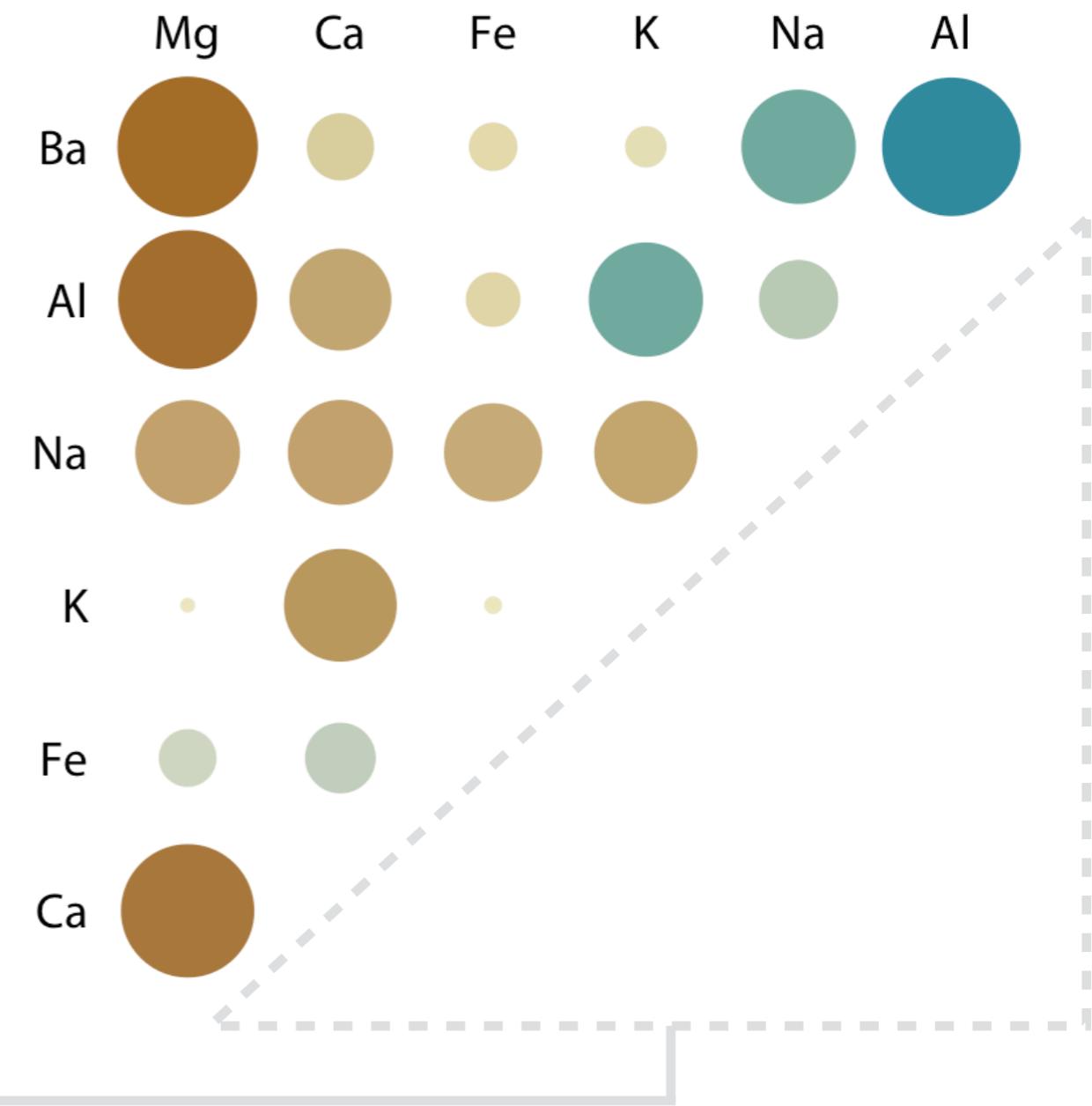
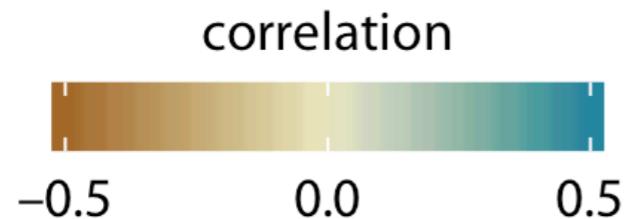
Session 04 Multi-variate Data Visualisation (Part II)



Correlograms

When interested in identifying correlations, if there are too many variables, it can be more efficient to compute those coefficients and visualize them directly:

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}}$$

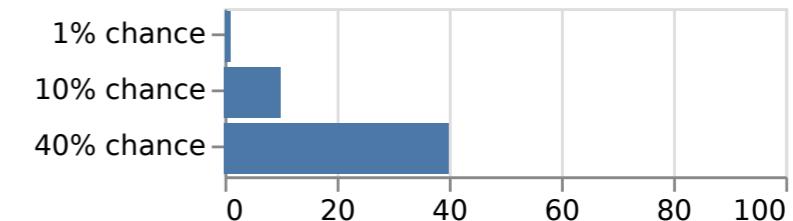


Symmetric formula in x_i and y_i

This is a relatively abstract representation though...

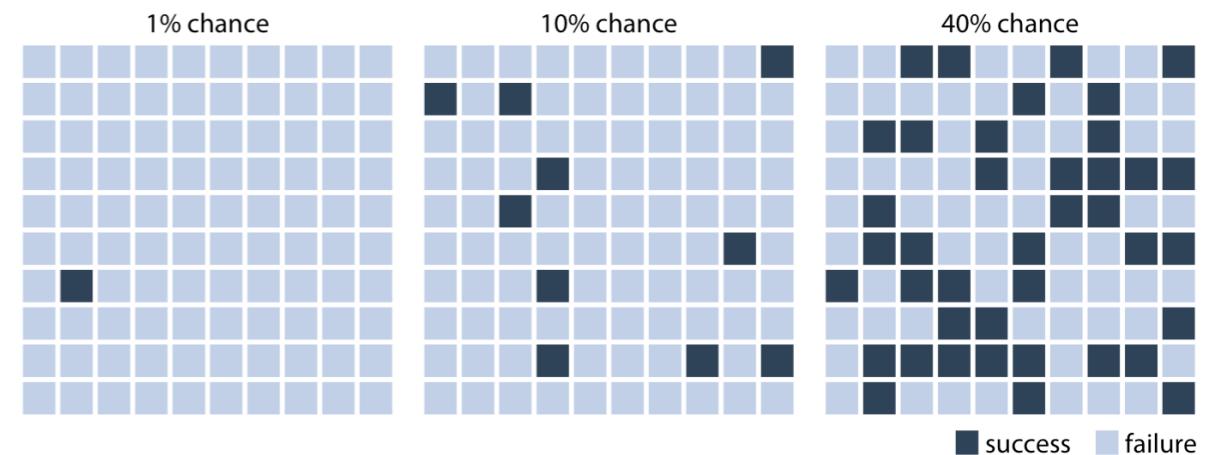
Discrete outcome visualization / Frequency framing

Showing the probability as a single number:



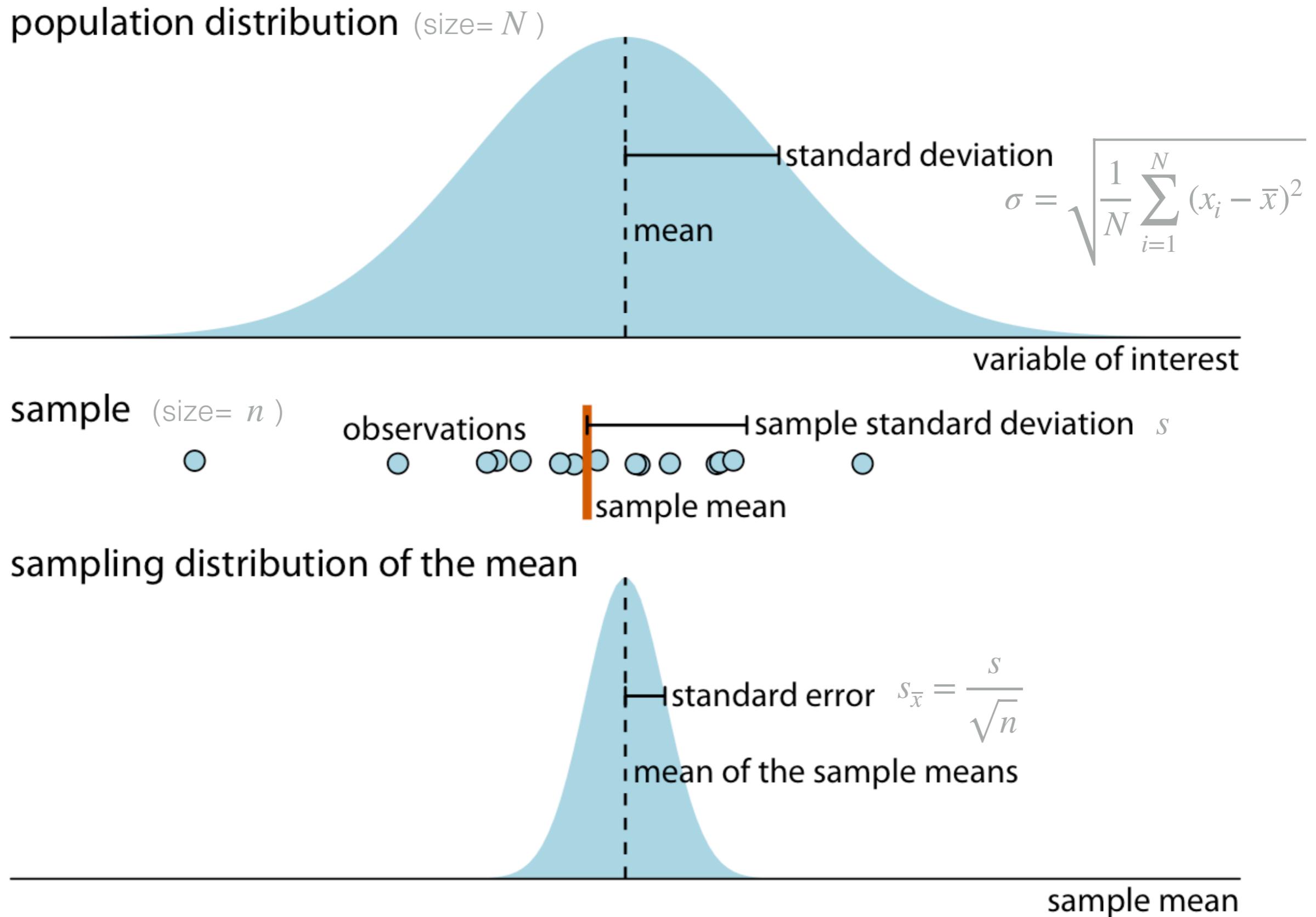
VS.

Showing discrete outcomes:

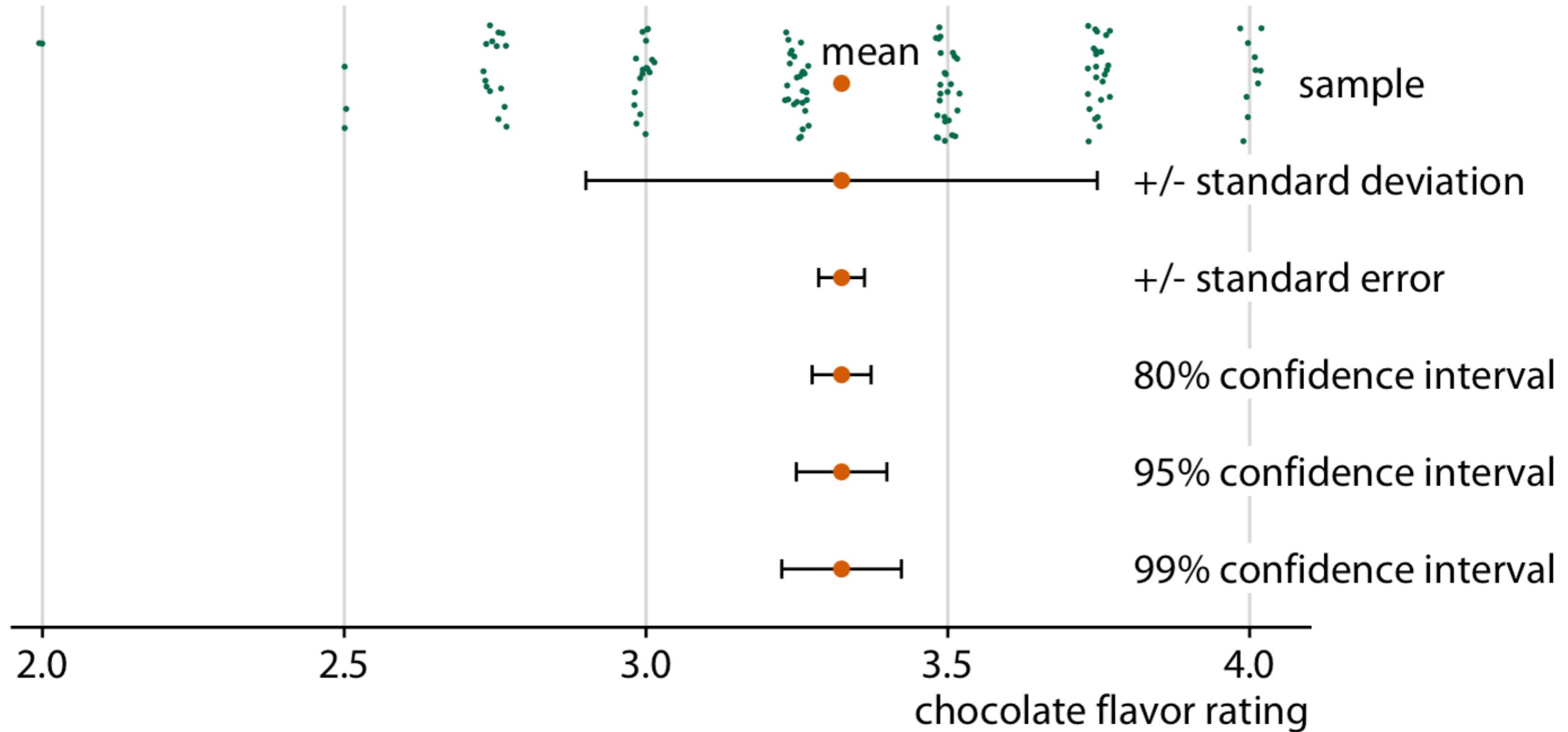


Key idea: emphasize the frequency aspect and the unpredictability of a random trial

Reminder - Descriptive Statistics



Error Bars

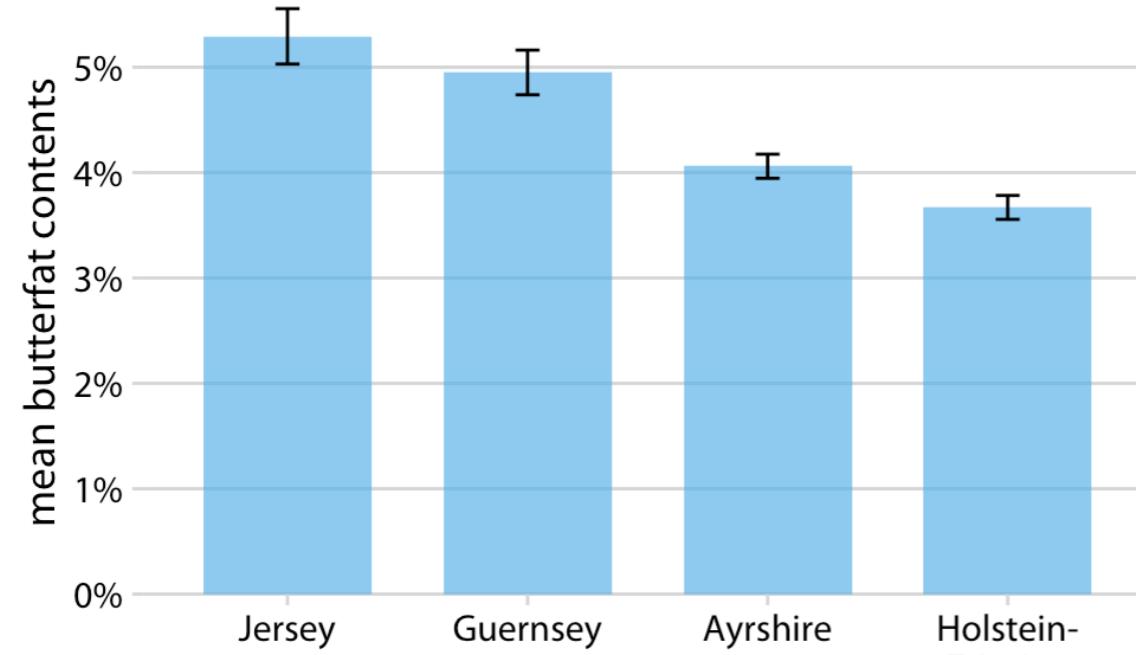


Always indicate what the error bar represents (quantity, confidence level)

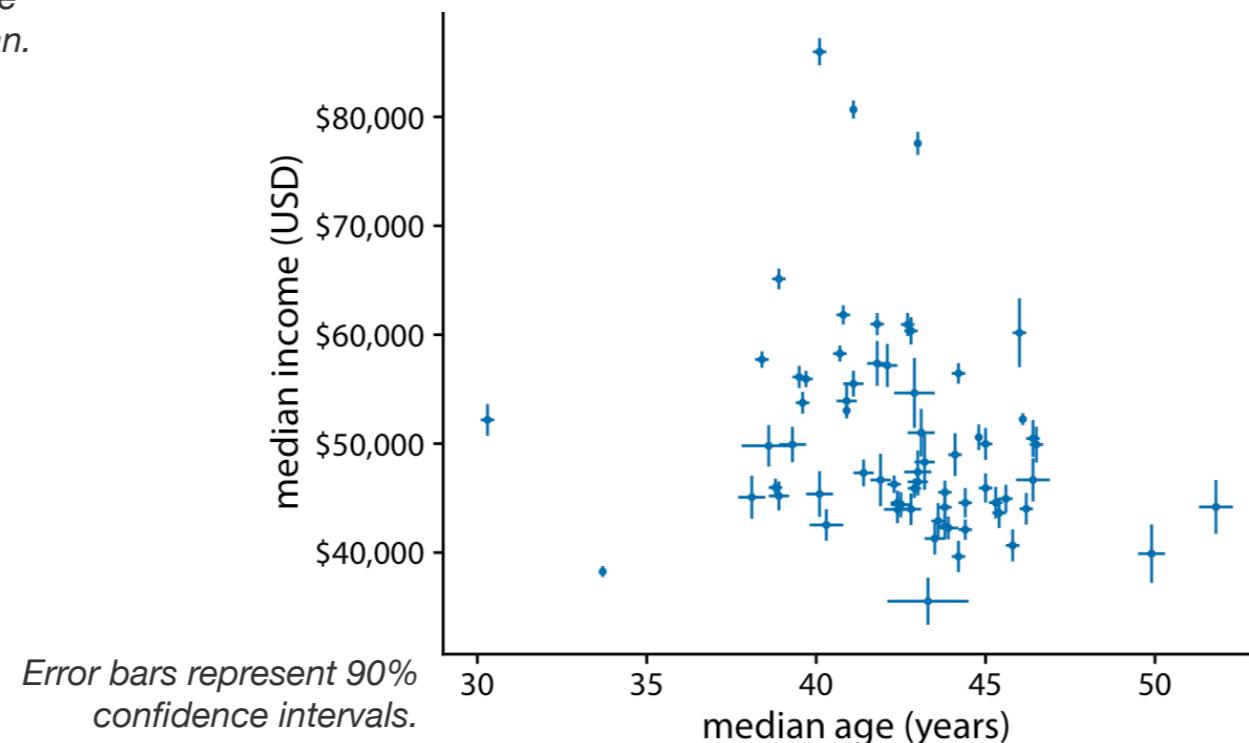
Advantages of error bars

Small visual footprint (usually cause little visual interference)

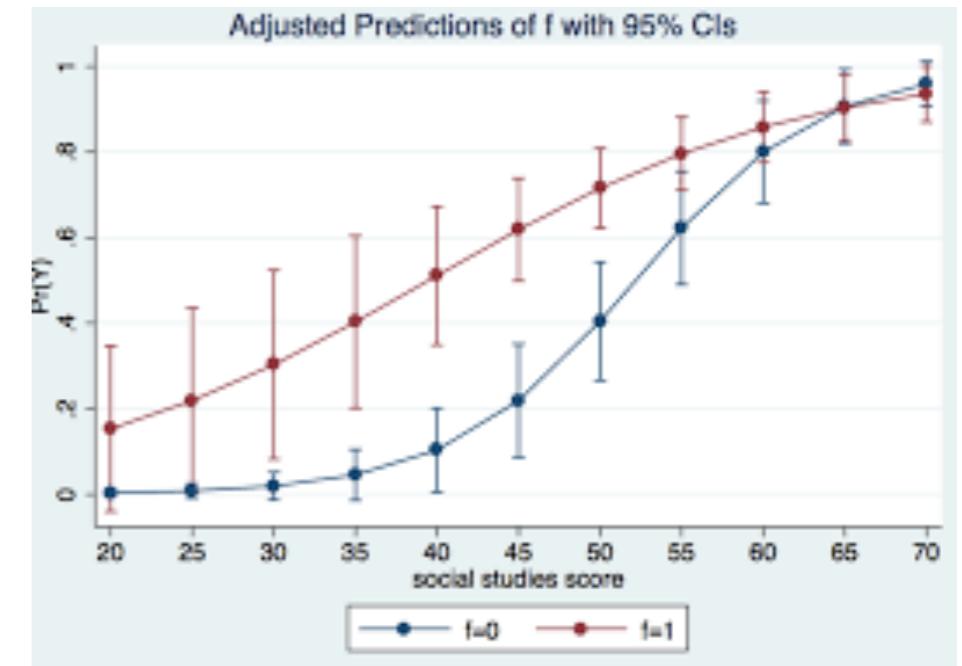
Apply to other types of plots than bar charts:



Error bars indicate +/- one standard error of the mean.



Error bars represent 90% confidence intervals.

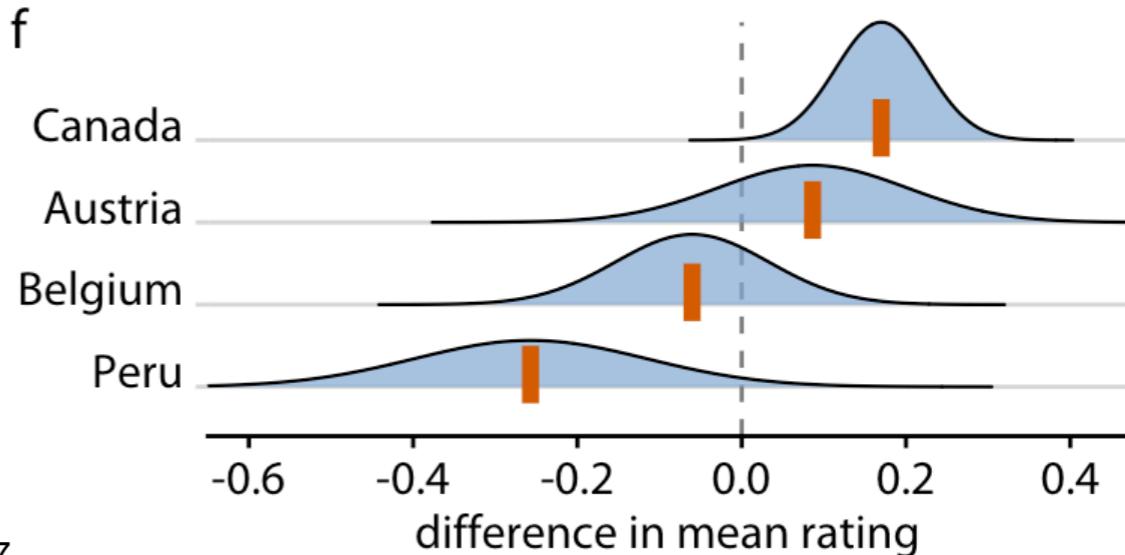
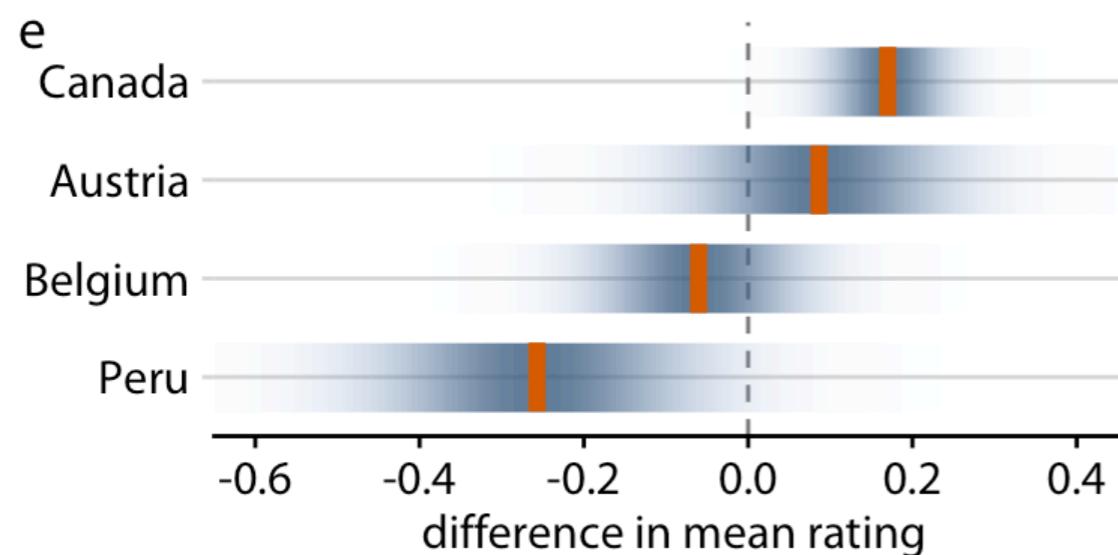
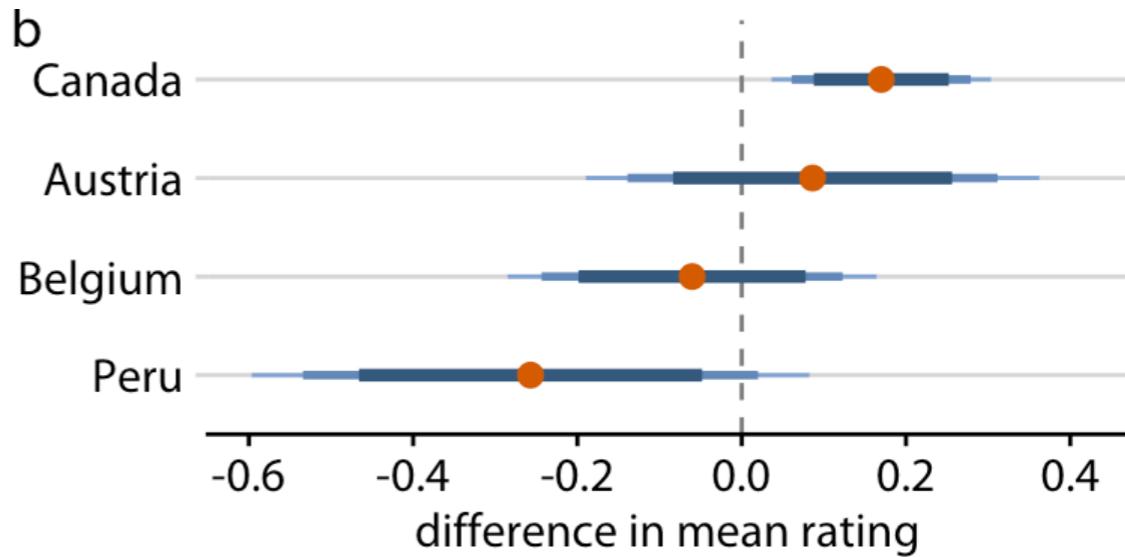
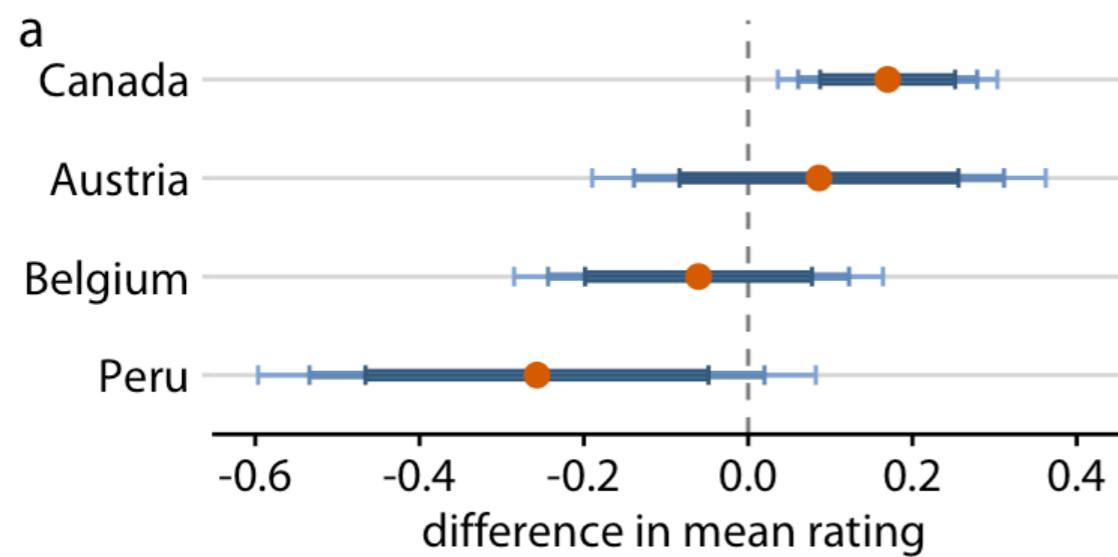
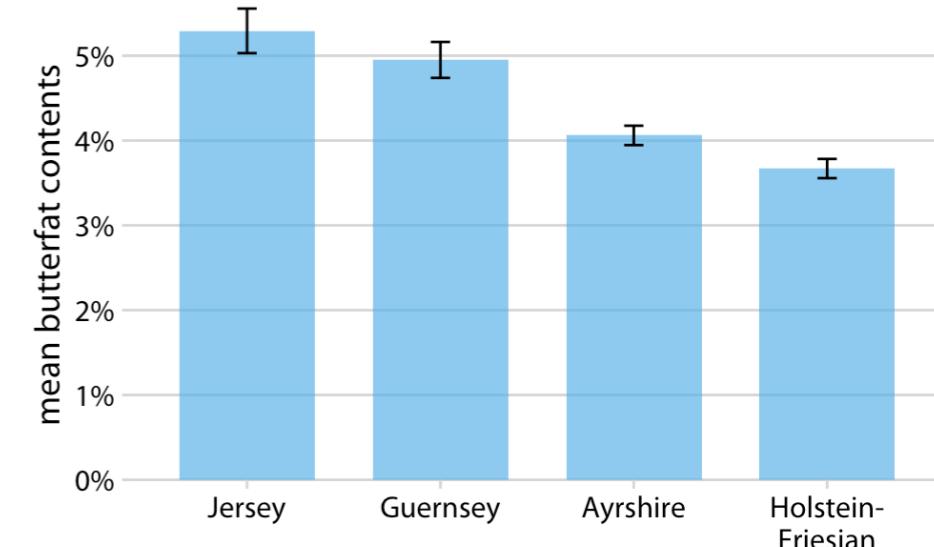


Issues with error bars: deterministic construal errors

With simple error bars, there is a risk that people perceive them deterministically:

- as representing min/max values;
- as precisely delineating the range of parameter estimates.

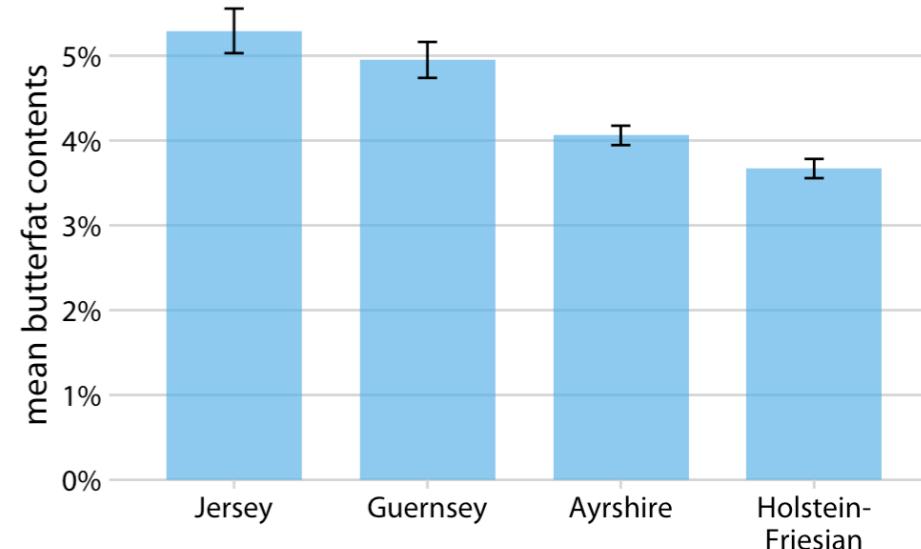
Visually convey the notion of uncertainty with, e.g., graded error bars or confidence strips.



Issues with error bars: can be misleading

Symmetric error bars are misleading if the data are skewed.

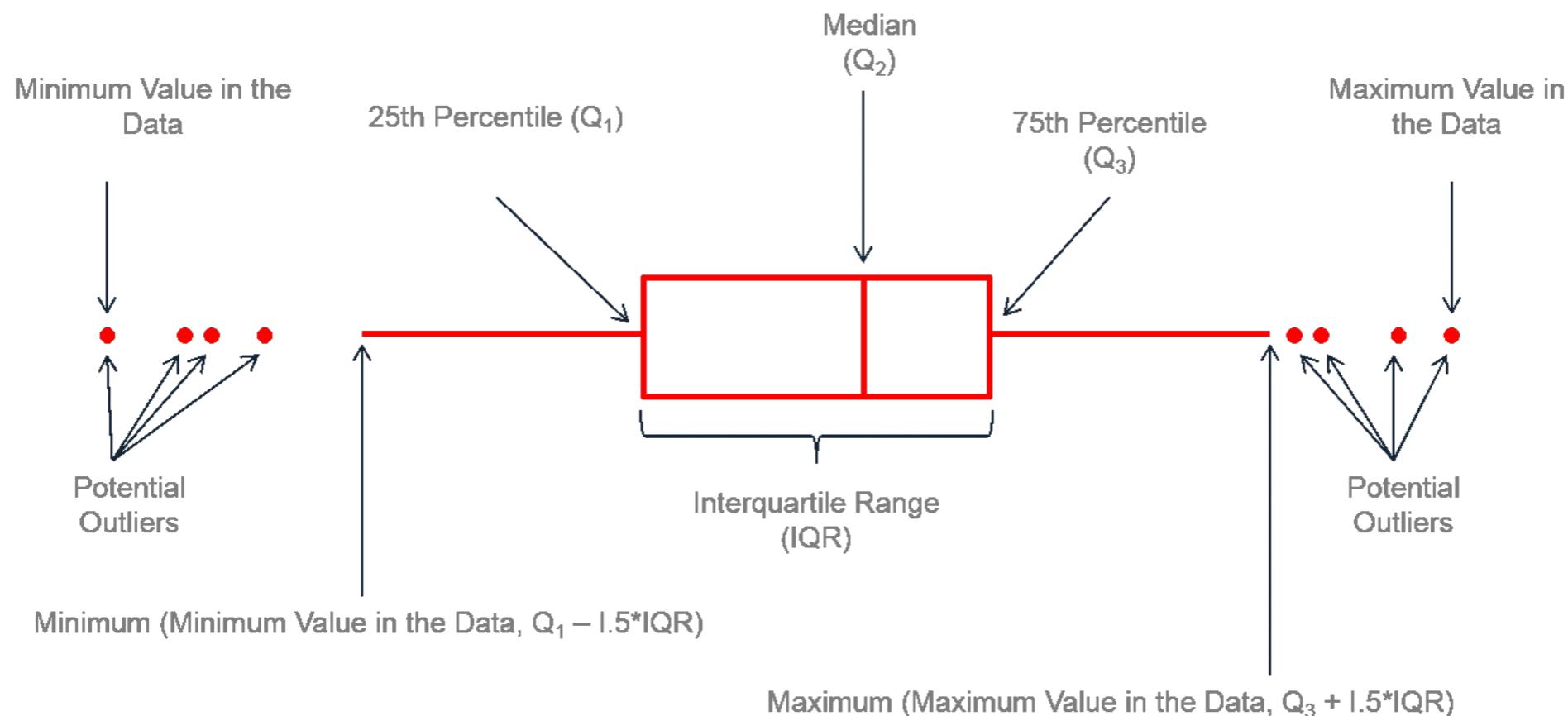
Absence of overlap of two error bars does not necessarily entail that the difference is significant at a given confidence level: *plot confidence intervals of the differences themselves.*



Interpretation of error bars representing confidence intervals (taking 95% as an example confidence level):

- ~~the bar delineates the interval that has 95% chance of containing the true value~~
- percentage of confidence intervals that would include the true value if an infinite number of random samples (same size) were pulled from the data and each time a 95% confidence interval was constructed.

Providing further information about the distribution: Box Plots



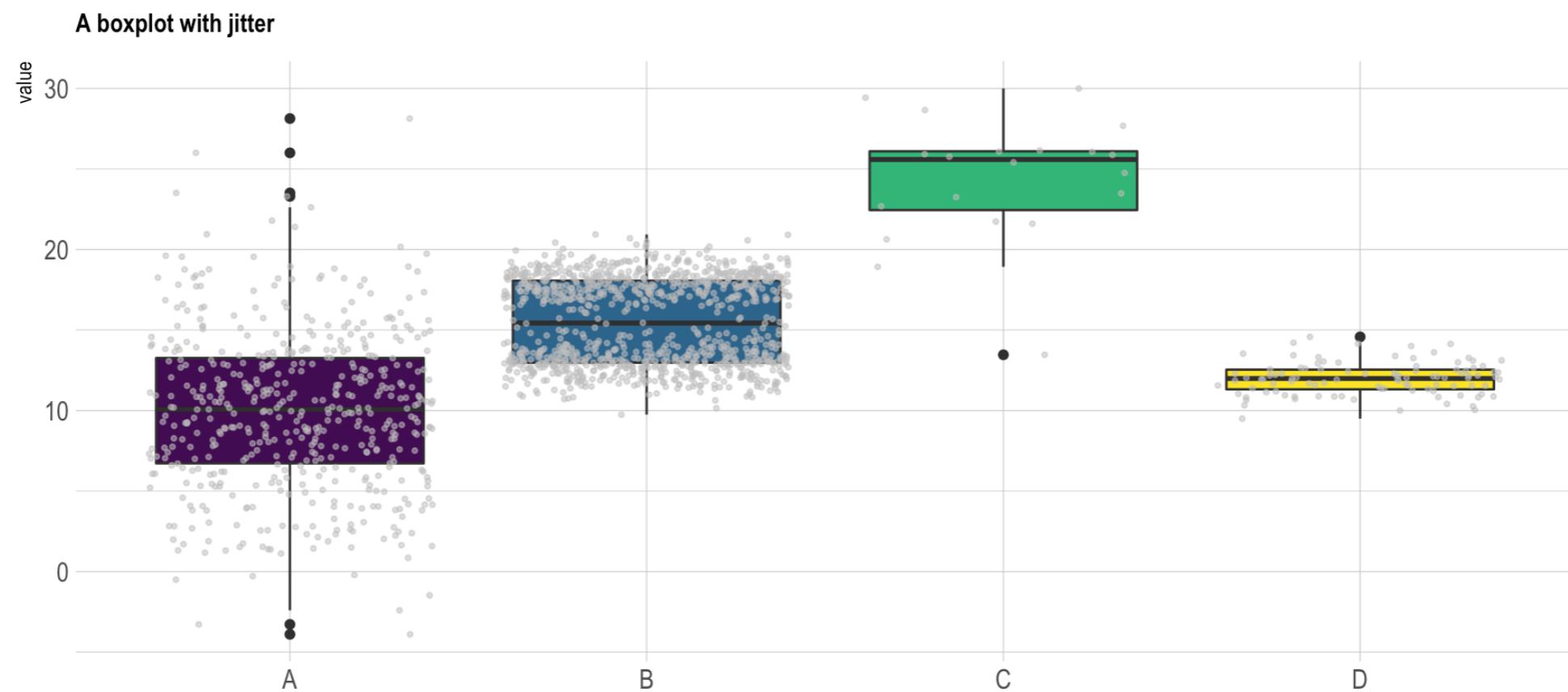
Provides more information, but still limited (missing sample size, distribution, ...)



Better: show the data *and* summary statistics

If not dealing with too many data points, showing them can be valuable.

An example using jitter:

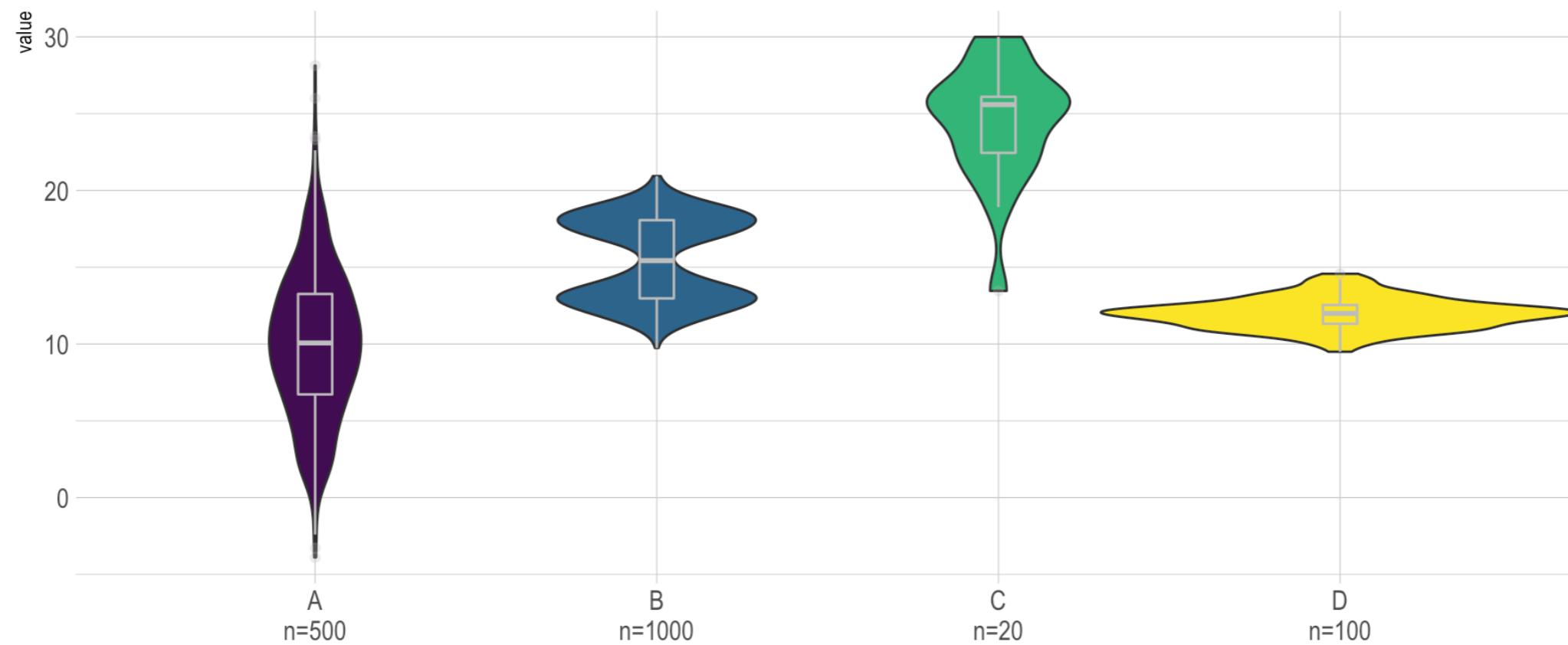


Shows that C's sample size is much smaller, that B seems to have a bimodal distribution.

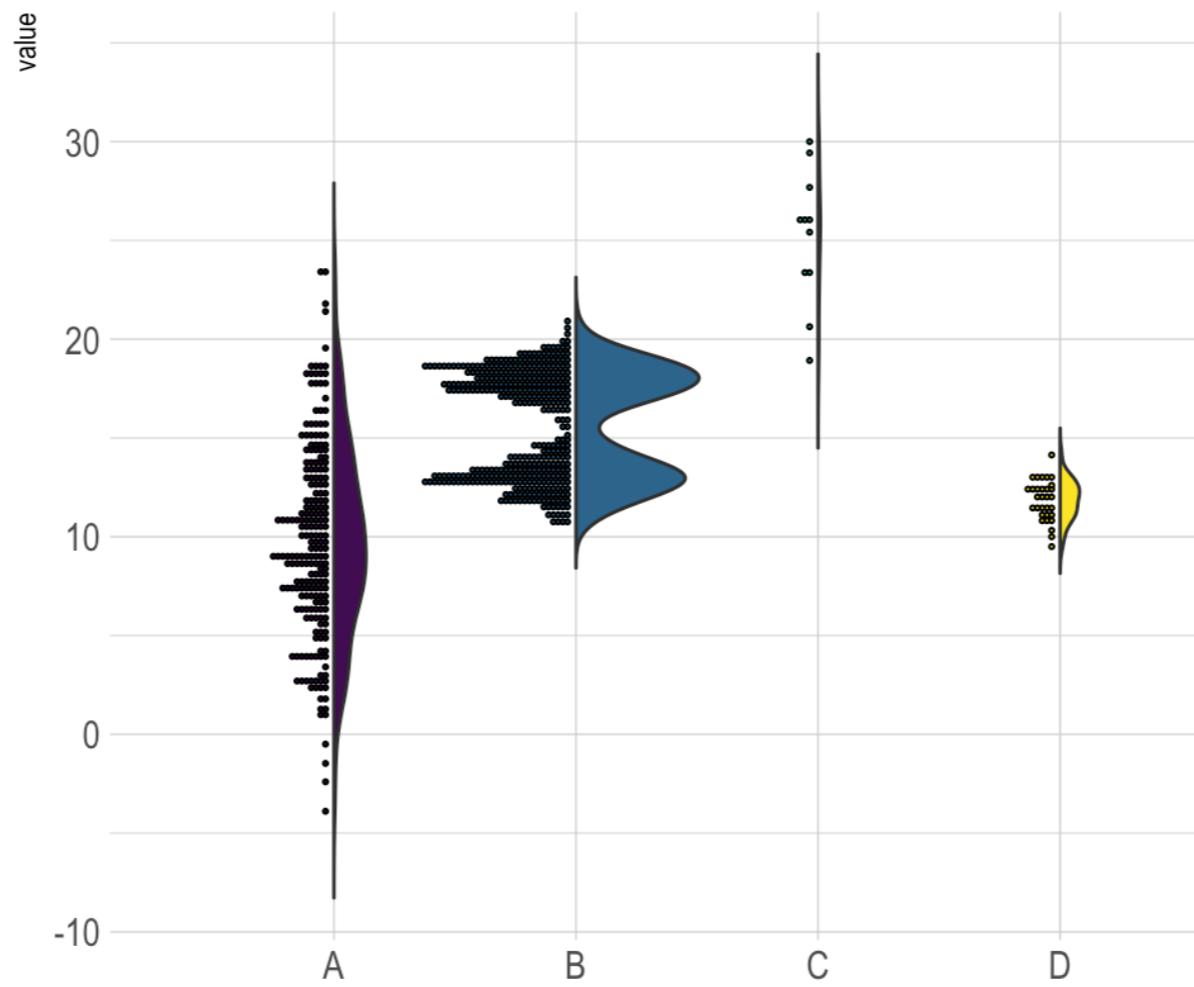
or show *distribution* and summary statistics

If dealing with many data points, showing them no longer works.

Show their distribution instead, using violin plots:

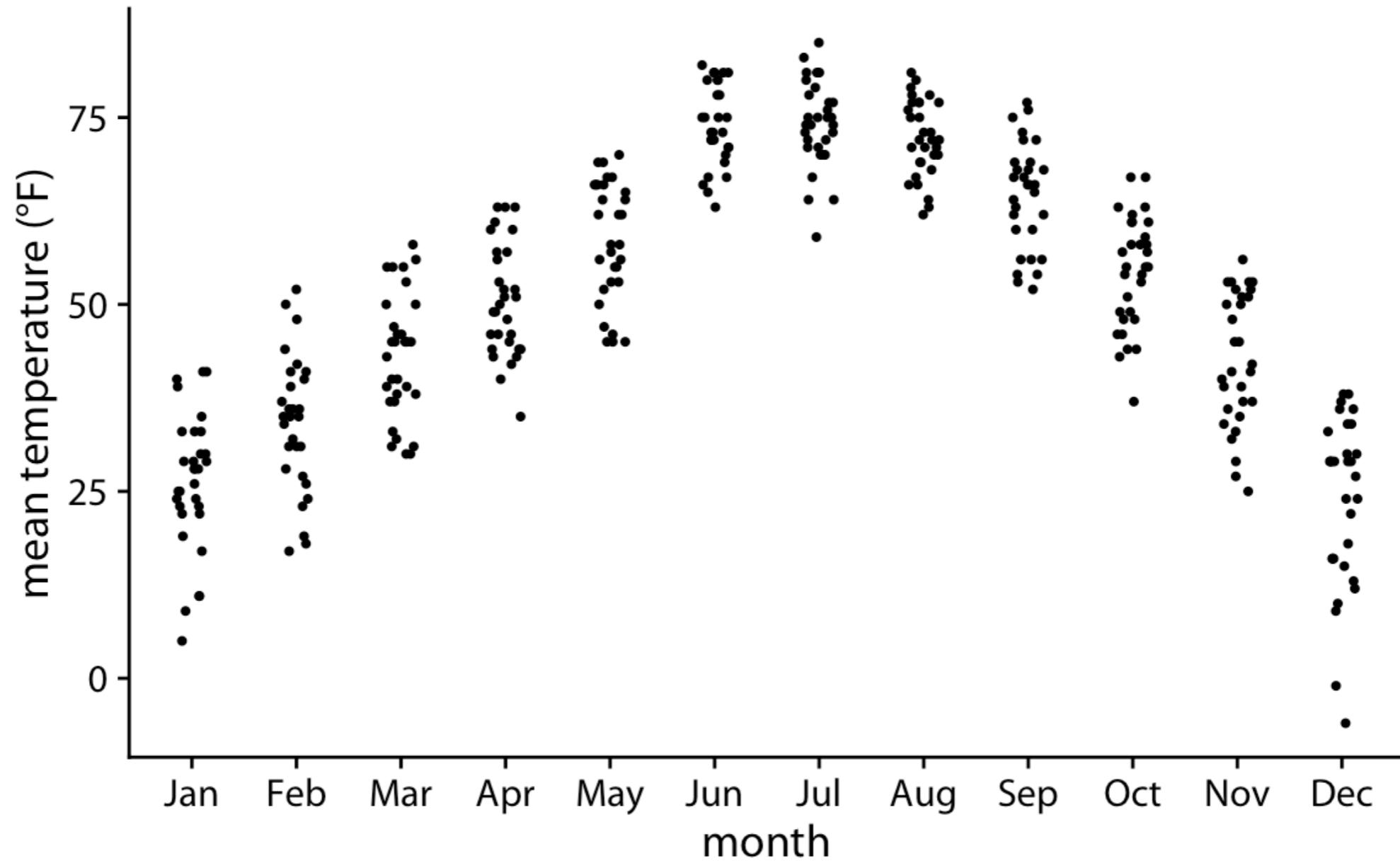


or give a sense of the raw data.



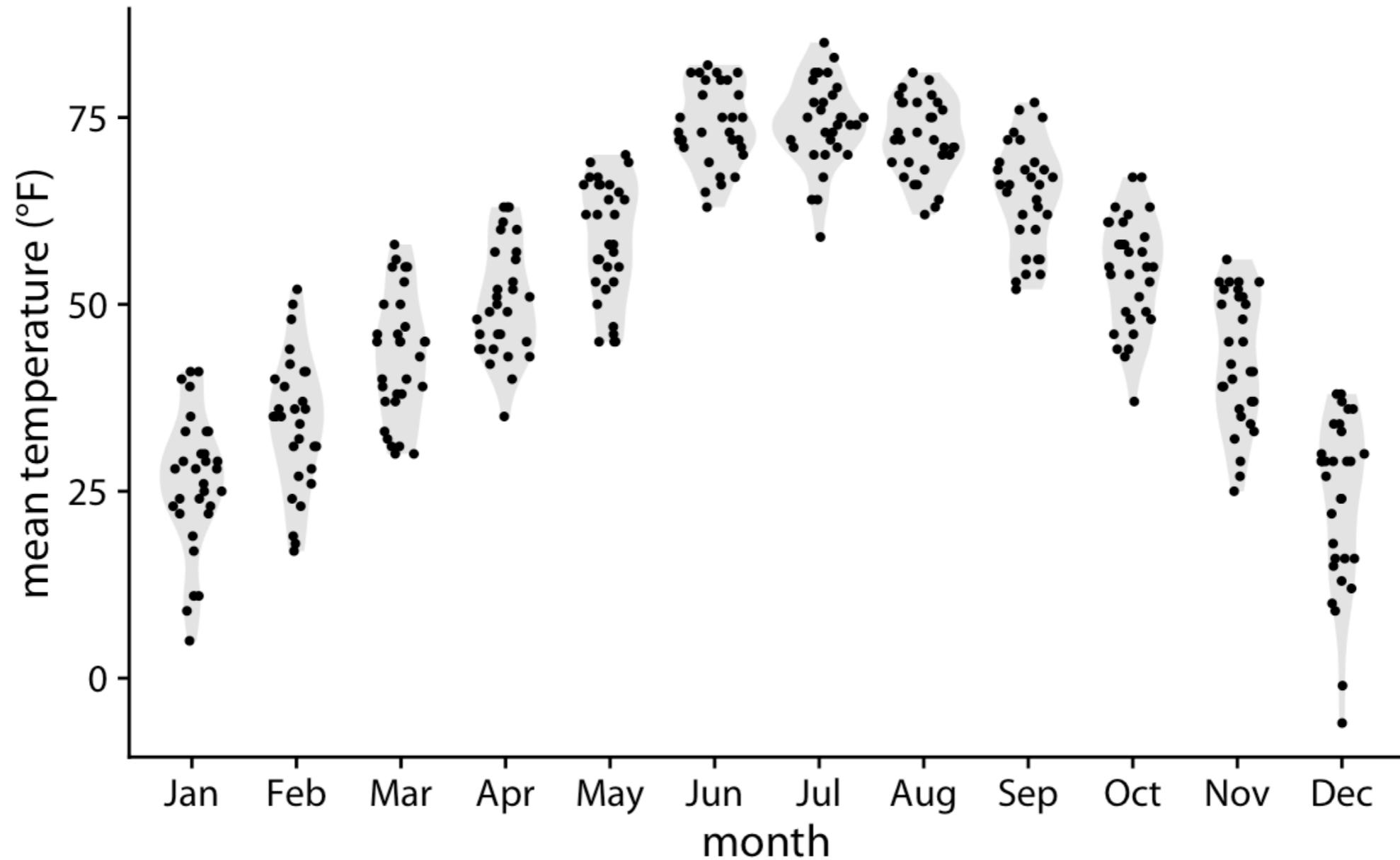
Sparse datasets

Plot the raw data points...

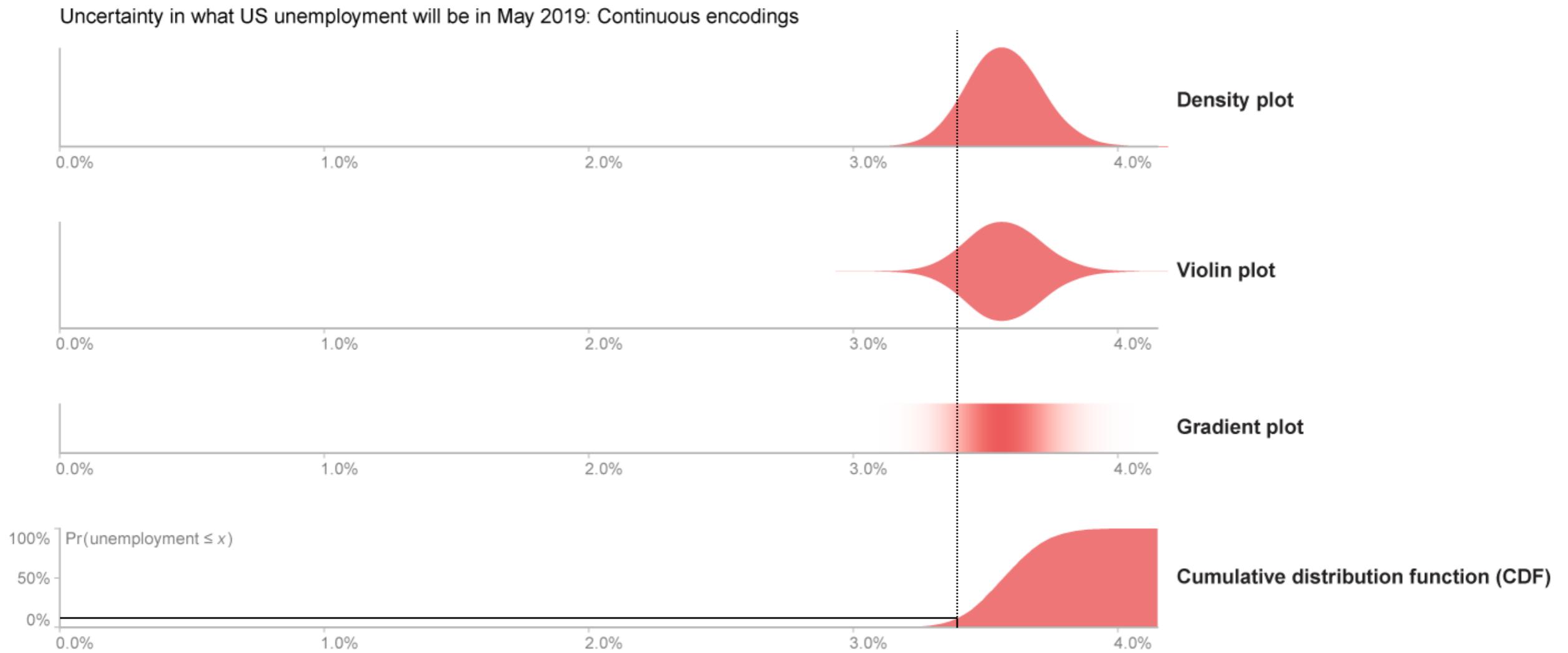


Sparse datasets

Plot the raw data points... and violins.



Visualizing uncertainty, one-sided



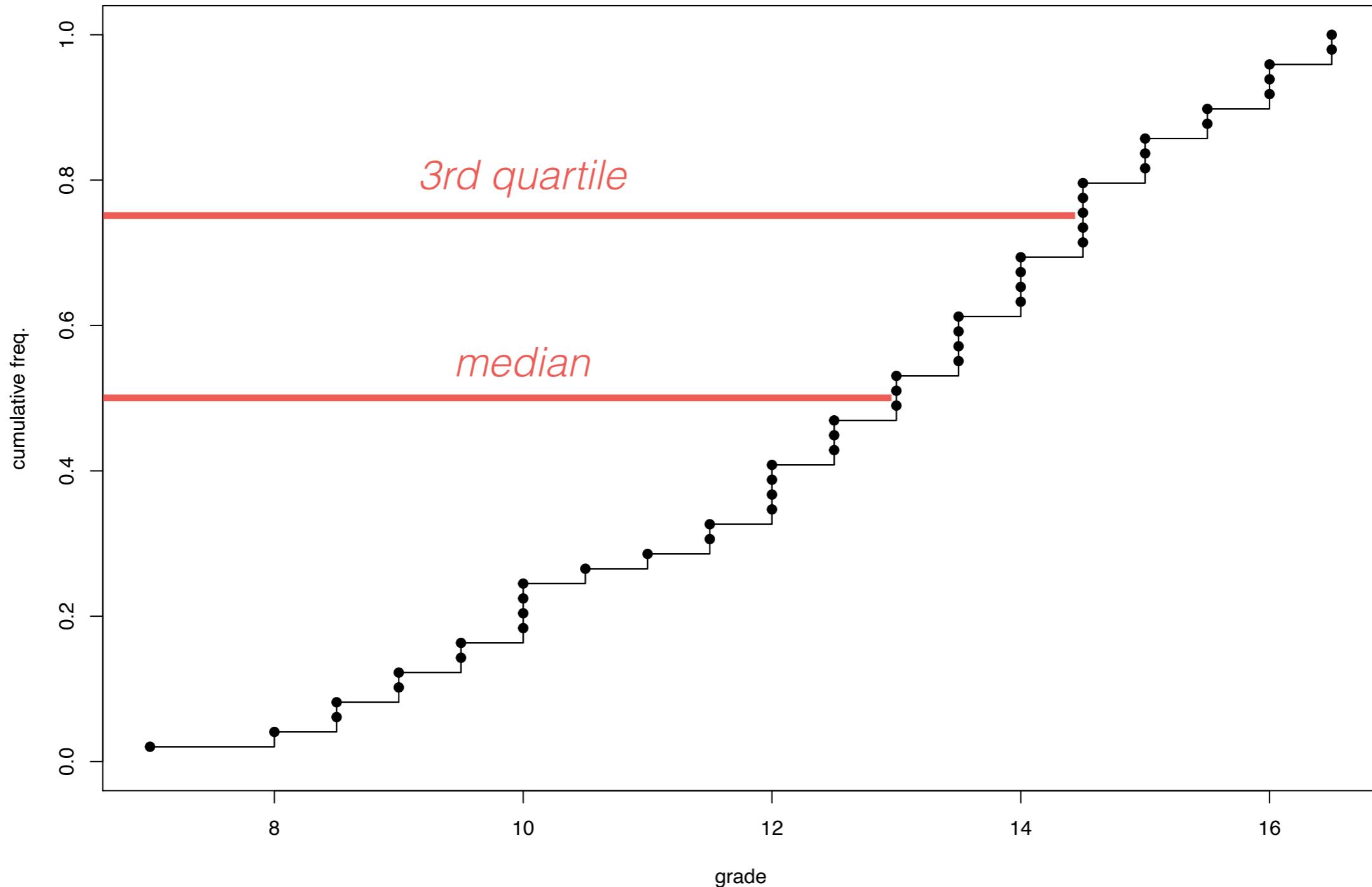
CDF effectively supports answering questions such as:
what is the probability that `var` will be less than `x`

Position judgment vs. area/color judgment with the other three encodings.

Mean, mode, etc. are more difficult to find though.

Empirical Cumulative Distribution Function (ECDF)

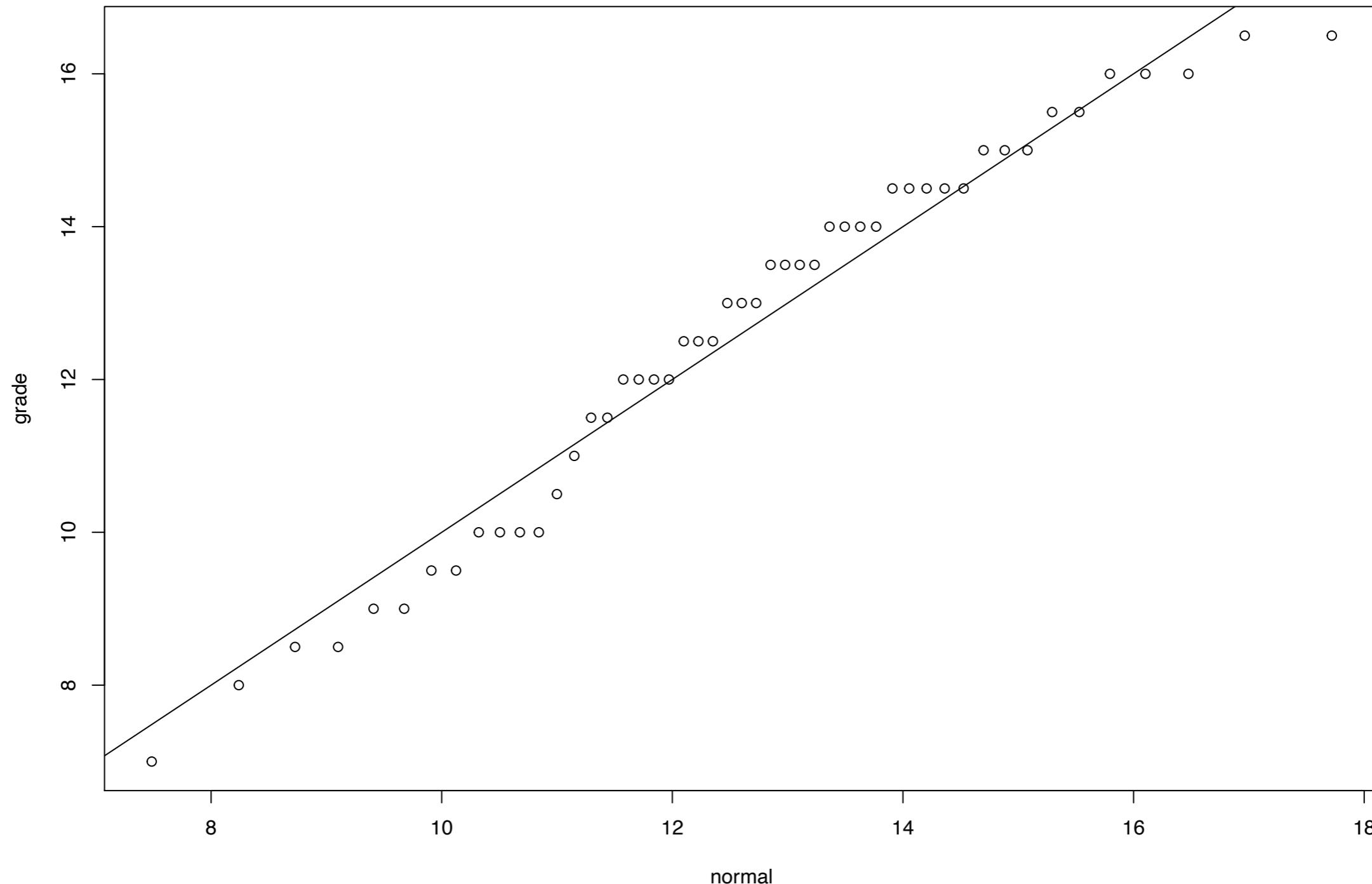
INF552–2018 Student Grades



Easy to find the median, quantiles, etc.

Quantile-quantile plot (Q-Q plot)

INF552–2018 Student Grade Distribution



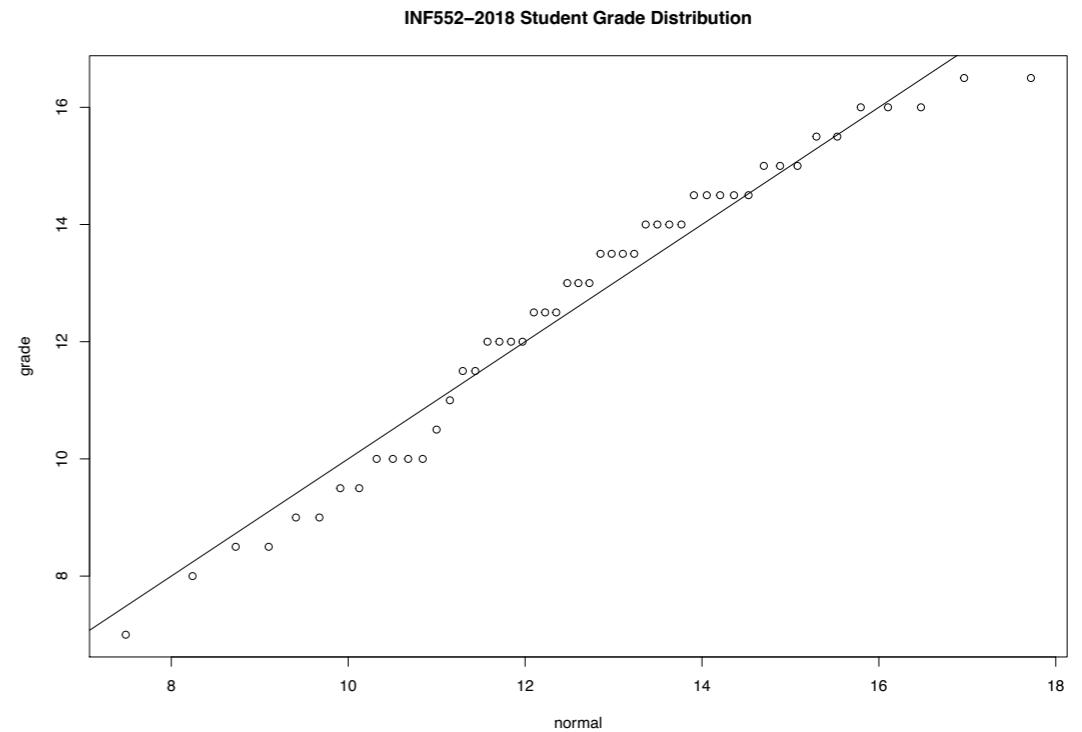
Visual representation of the similarity between two distributions:

- two distinct datasets;
- theoretical vs. observed.

Quantile-quantile plot (Q-Q plot)

= sample sizes:

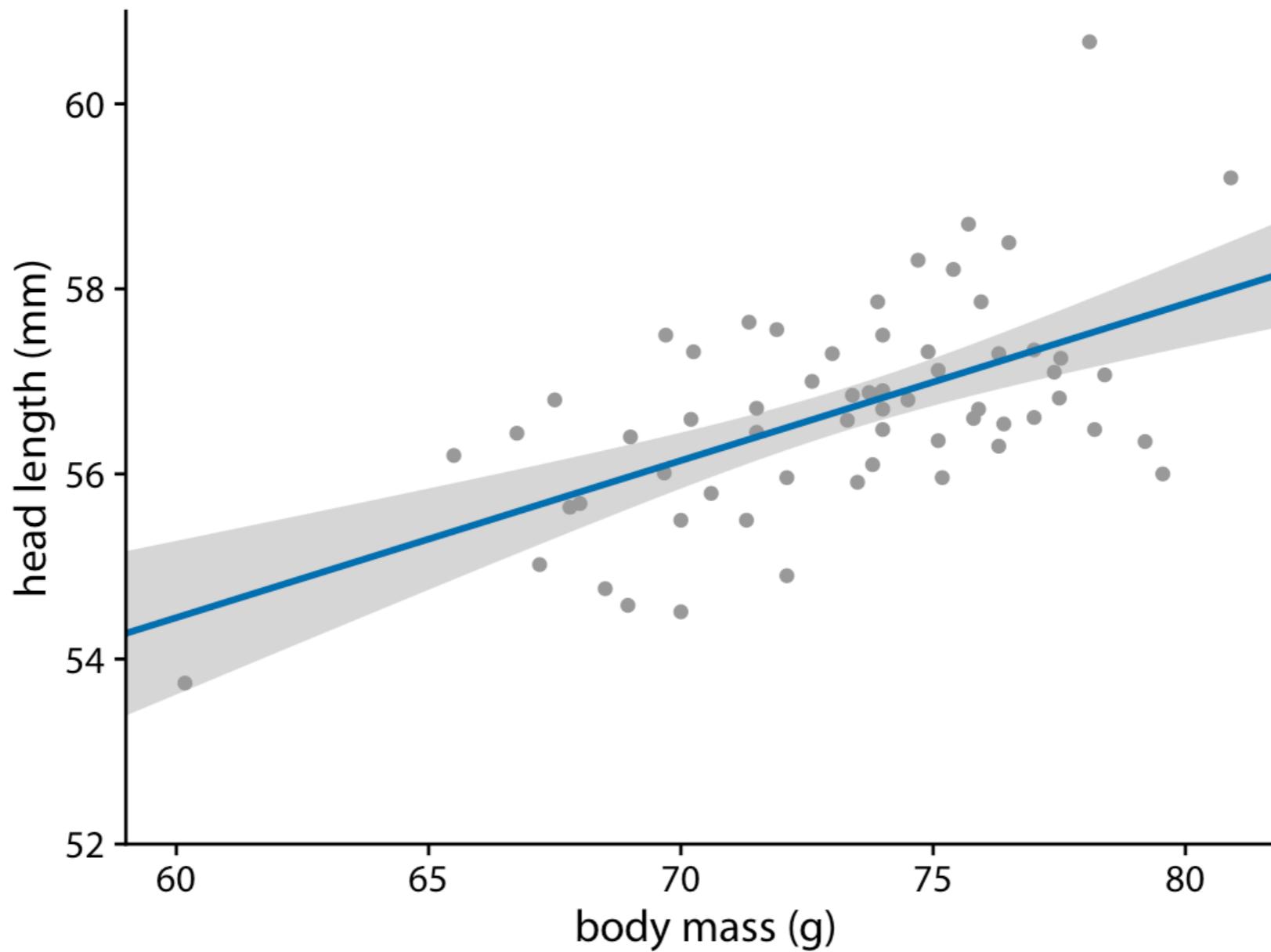
- plot the sorted data;



\neq sample sizes, *with* $\text{card}(d1) < \text{card}(d2)$:

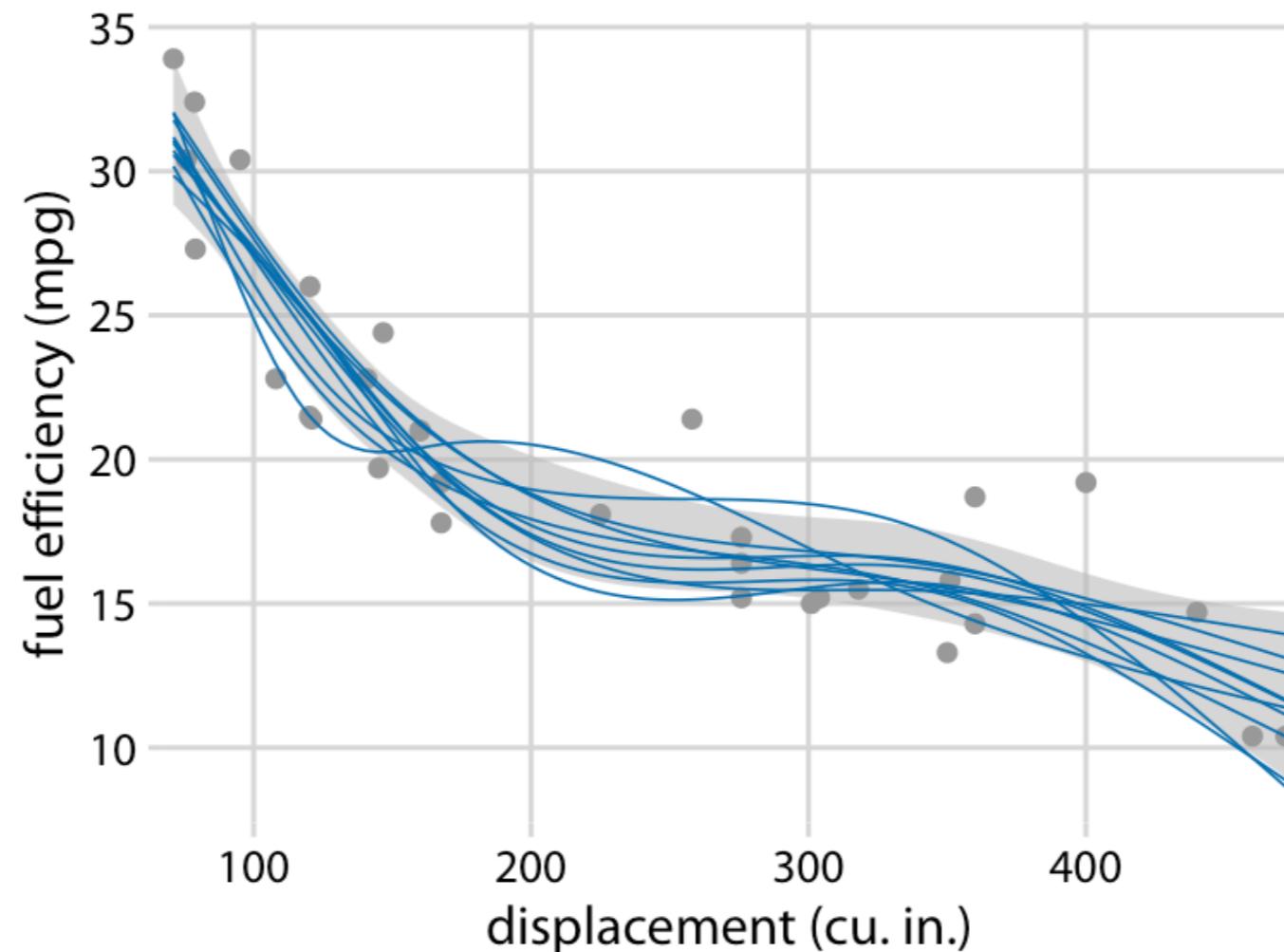
- quantiles correspond to the sorted data values from $d1$;
- quantiles for $d2$ are interpolated.

Uncertainty of curve fits: confidence bands (linear fits)



Uncertainty about both slope and intercept of the regression line.

Uncertainty of curve fits: confidence bands (non-linear fits)

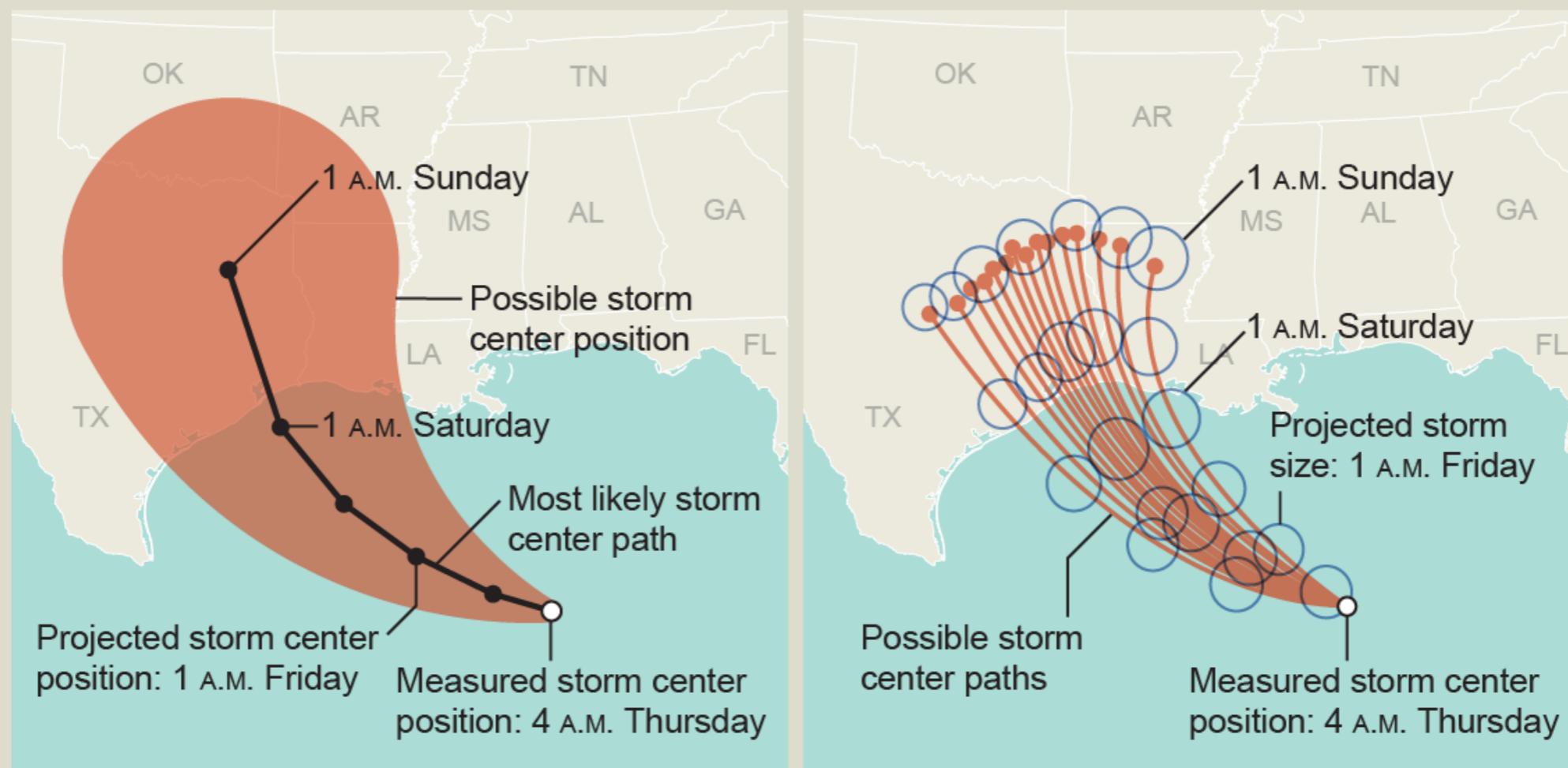


Uncertainty not only about position of the curve, but its actual shape (wiggliness).

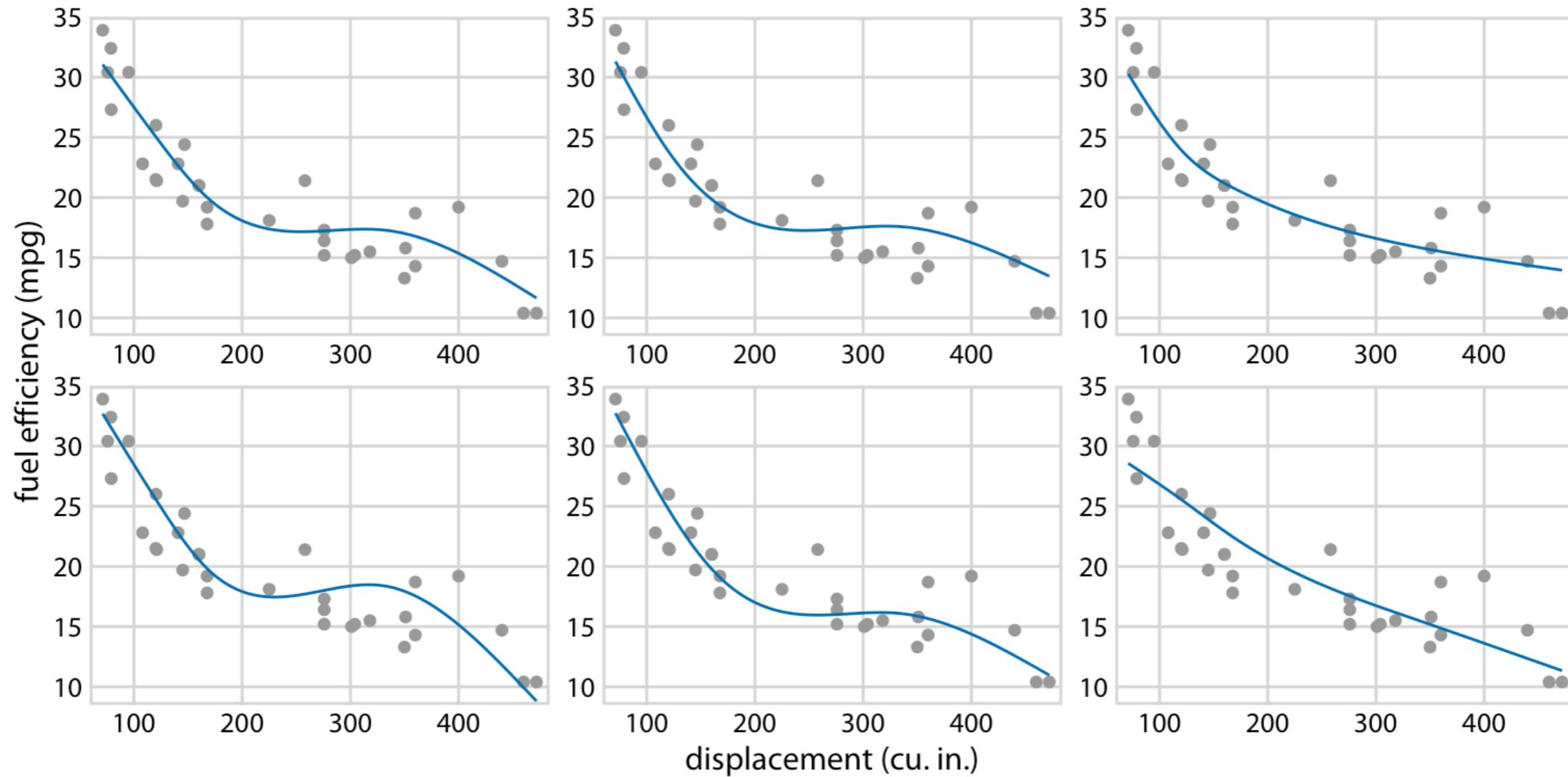
Issue of deterministic construal error again.

Deterministic construal errors

“Cone of uncertainty” (*left*) shows where a hurricane may head, according to a group of forecasts. An alternative is to show the specific path predicted by each forecast (*right*). Both approaches have pros and cons in helping people judge the risk they may face, but the one on the right makes it clearer that the path is difficult to predict.

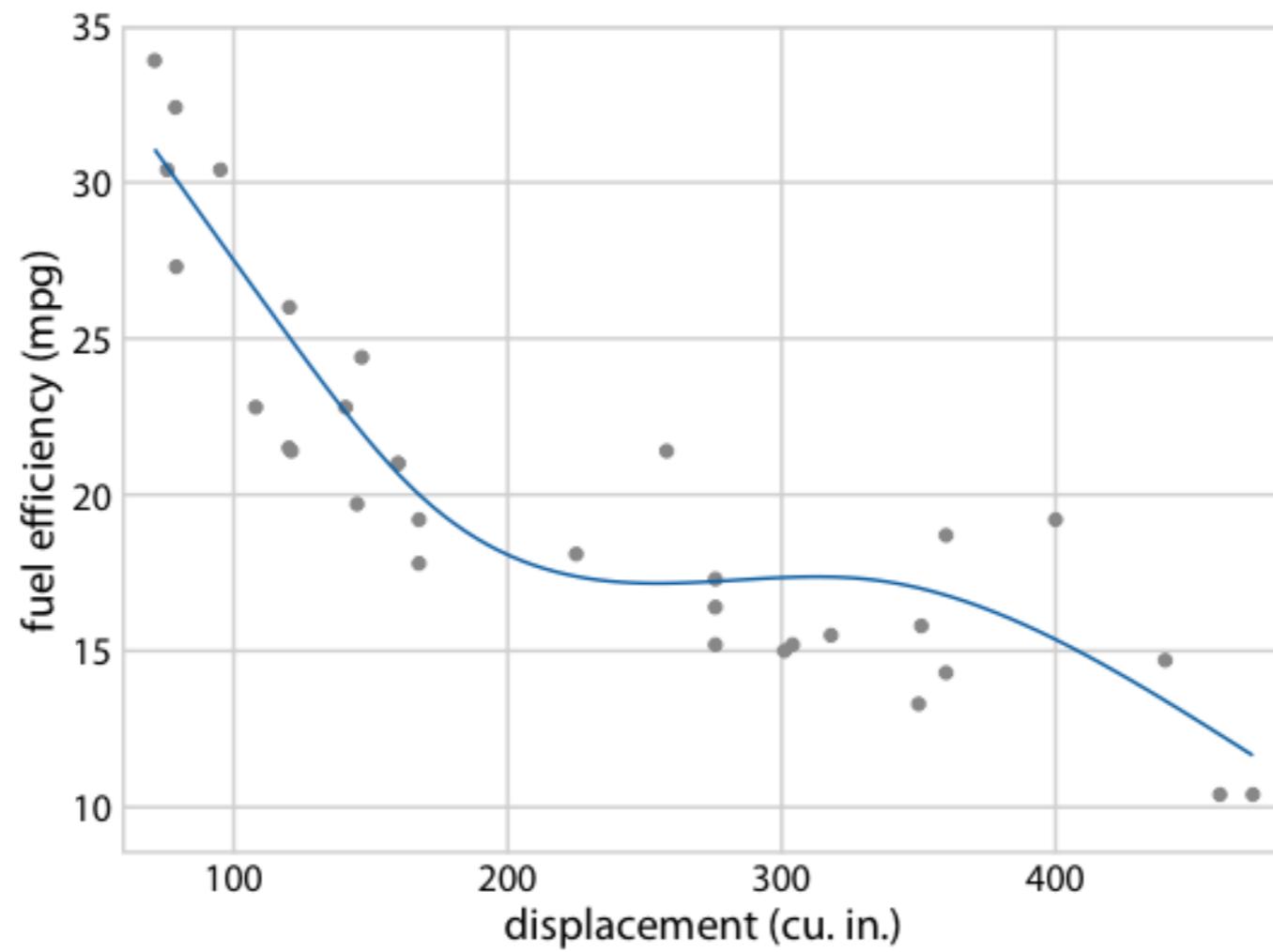


Show different possible outcomes



Plot multiple (hypothetical) outcomes, either as small multiples...

Show different possible outcomes: Hypothetical Outcome Plots

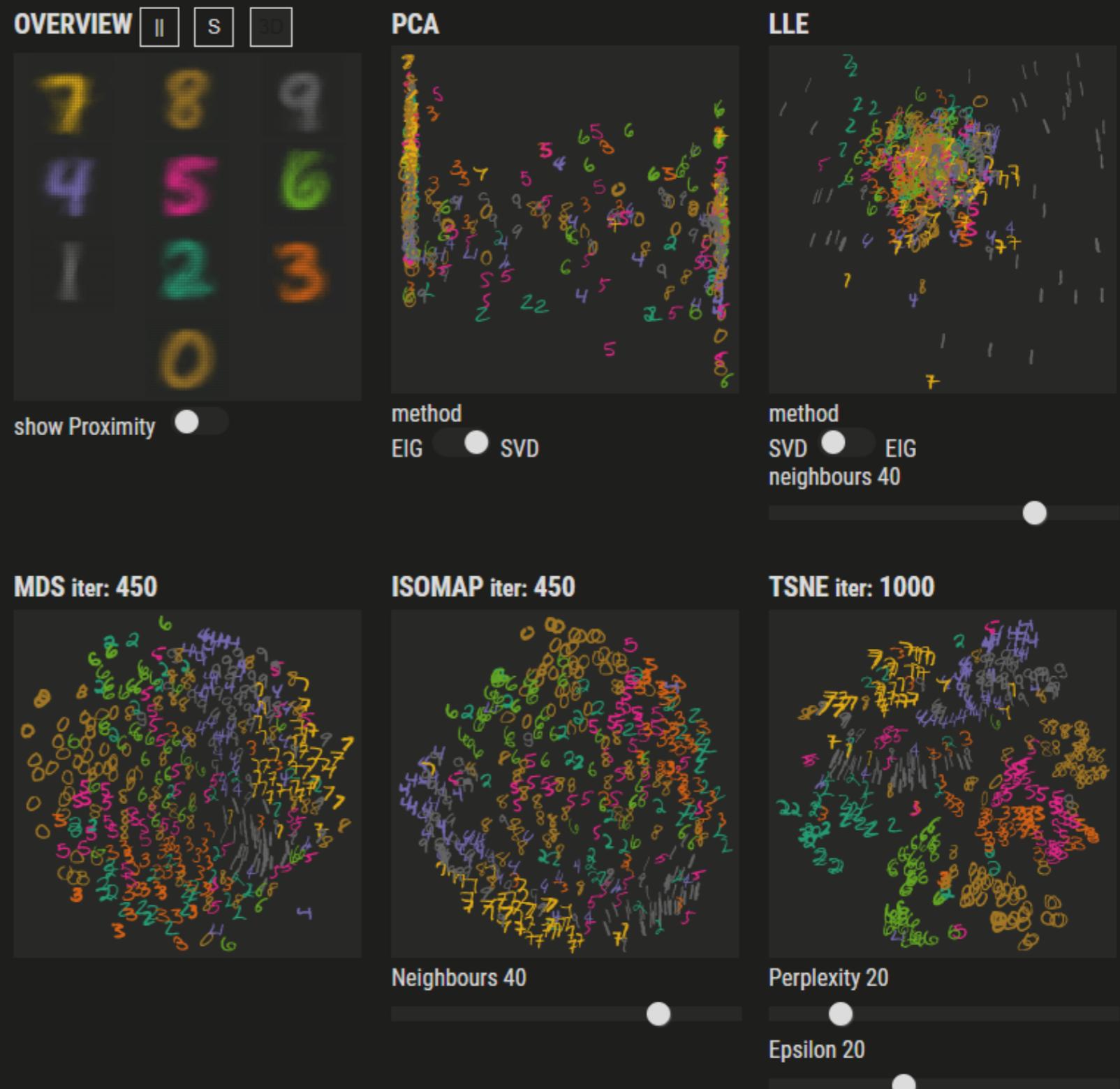


Plot multiple (hypothetical) outcomes, either as small multiples... or animating between them.

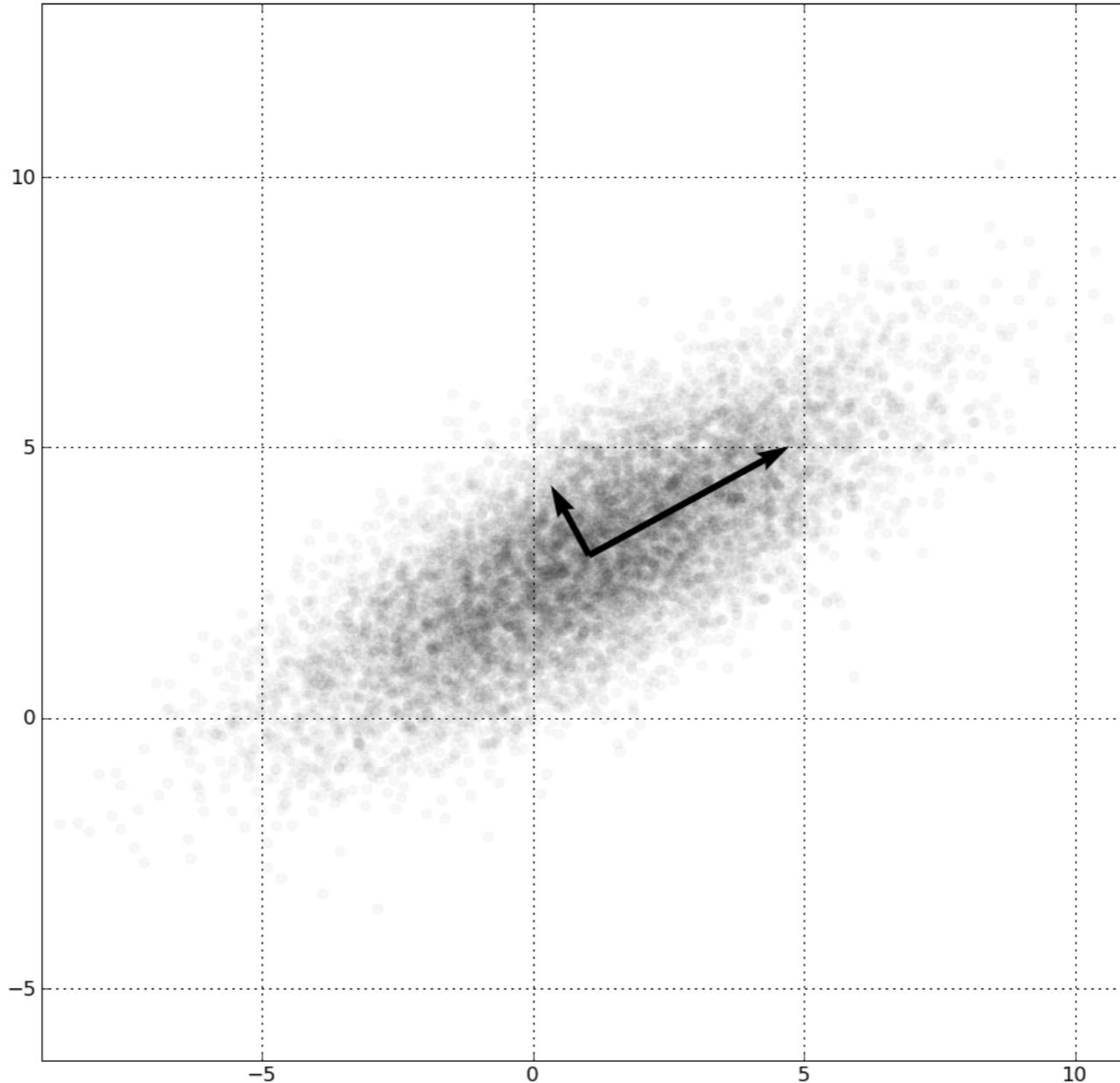
- Very effective at conveying uncertainty.
- Hard switch between them, or smooth (quick) fade in/out, but no smooth shape morphing.
- Show outcomes that are representative of the true distribution of possible outcomes.

Dimensionality Reduction

Given a set of data points that lie in a high-dimensional space, find a representation of those points in a lower-dimensional space (2D space, 3D volume).



Dimensionality Reduction: PCA



Principal Component Analysis

PCs are **linear** combinations of original variables

PCs are uncorrelated

PCs are ordered by max. variation they yield (decreasingly)

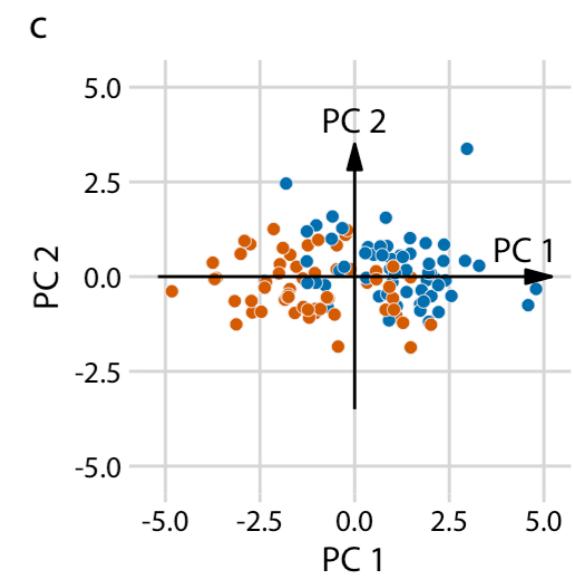
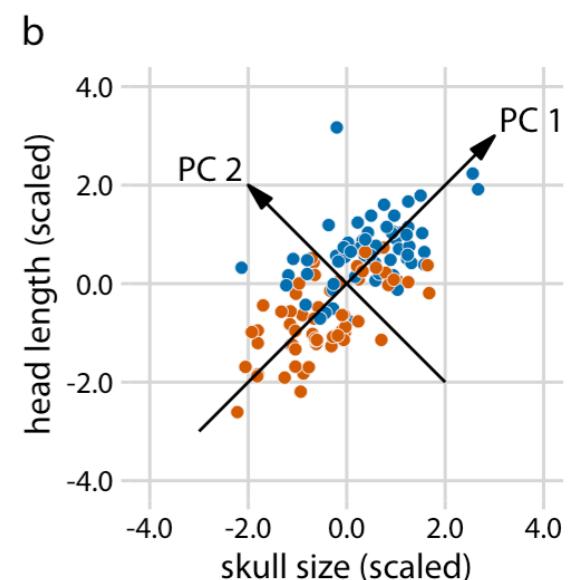
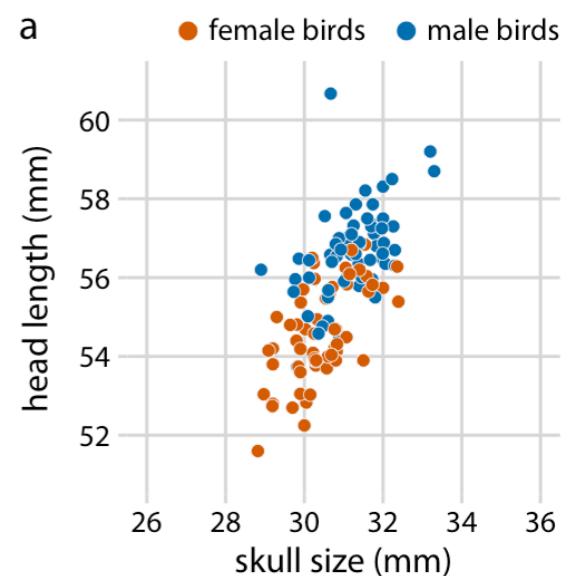


Dimensionality Reduction: PCA

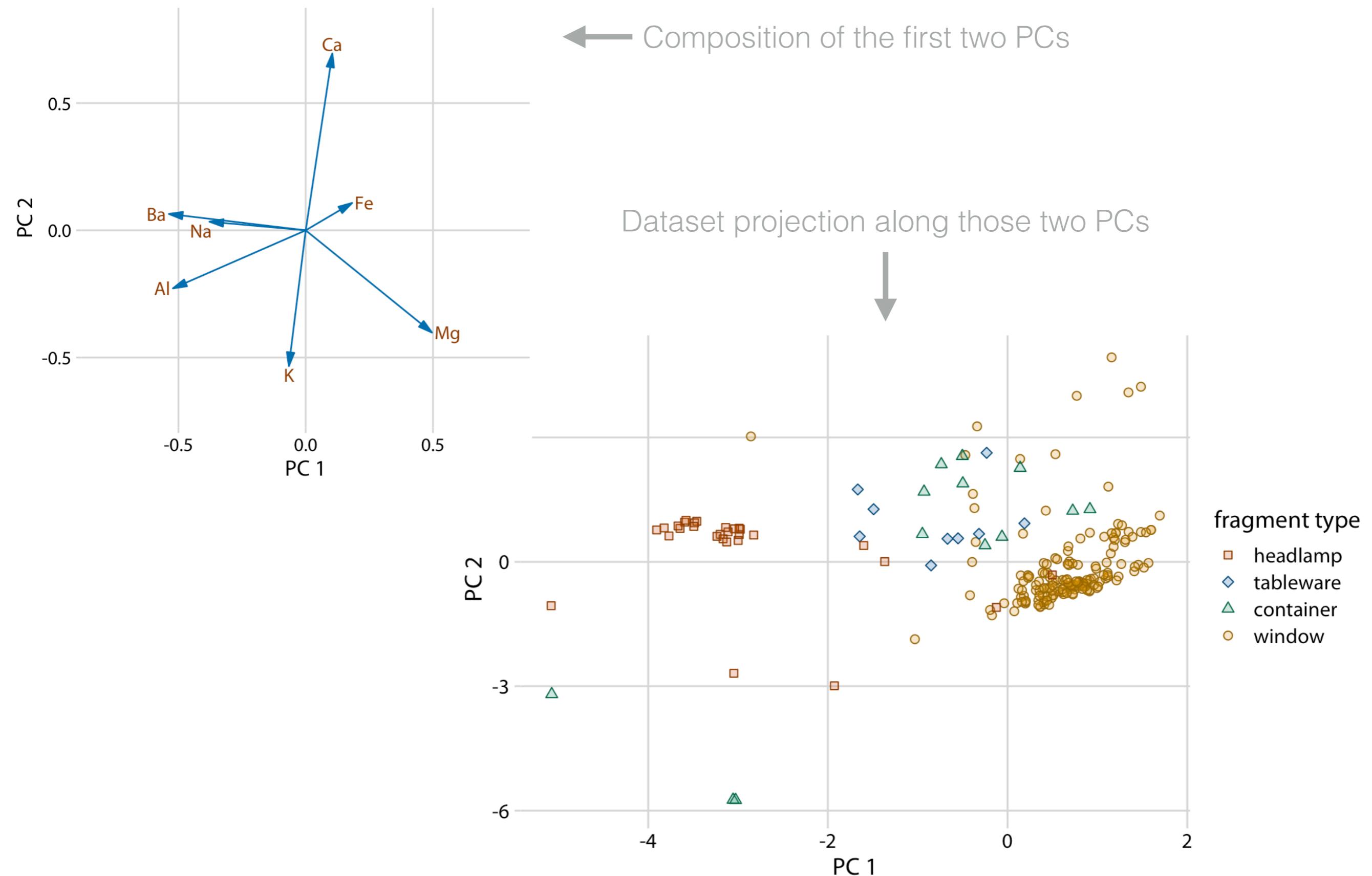
Standardize variables to zero mean and unit variance.

Define PCs along directions of maximum variation.

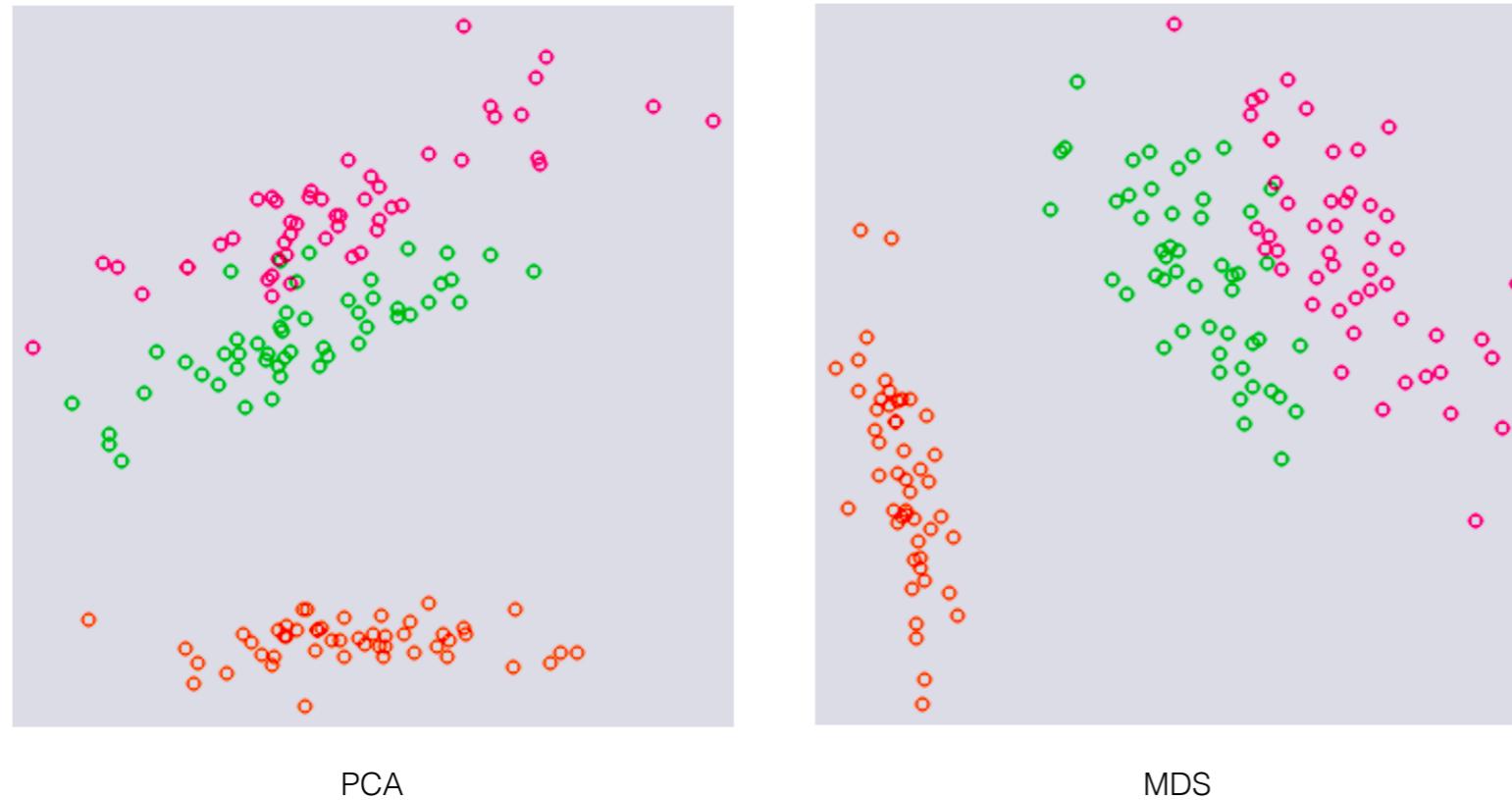
Project the data in the new coordinates system.



Dimensionality Reduction: PCA



Dimensionality Reduction: MDS



Multi-Dimensional Scaling (**non-linear**).

PCA considers the data points themselves, while MDS considers the pairwise distances between data points: it tries to preserve distances from the n D space in the 2D plane projection.

Preserves global structure (larger distances) at the expense of local structure (smaller distances).

Analogy with force-directed graph layout algorithms.

ISOMAP only considers the k-nearest-neighbors of each data point.

Dimensionality Reduction: t-SNE



t-distribution Stochastic Neighborhood Embedding (**non-linear**).

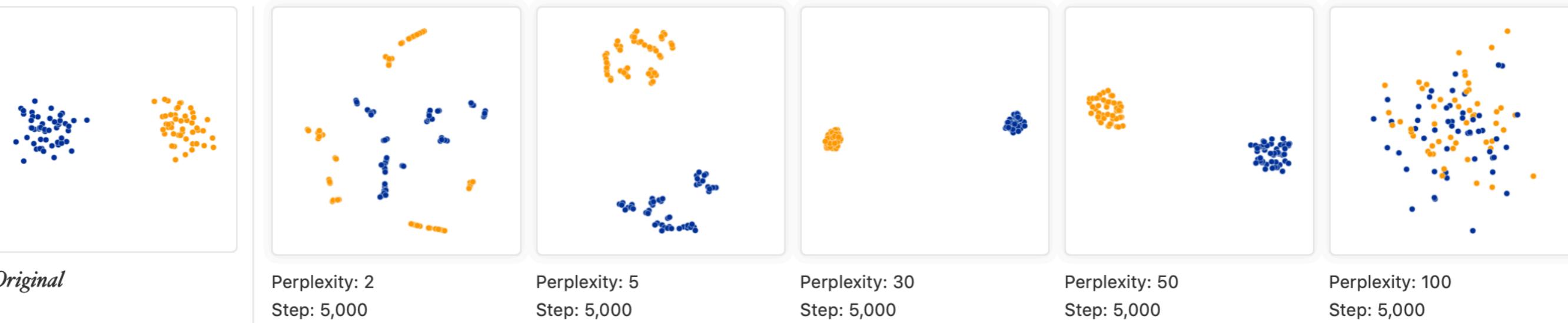
Analogy with force-directed graph layout algorithms (again).

Better preservation of both global *and* local structure.

Can work well with a large number of dimensions (1000+).

Different runs can yield different outputs.

Dimensionality Reduction: t-SNE



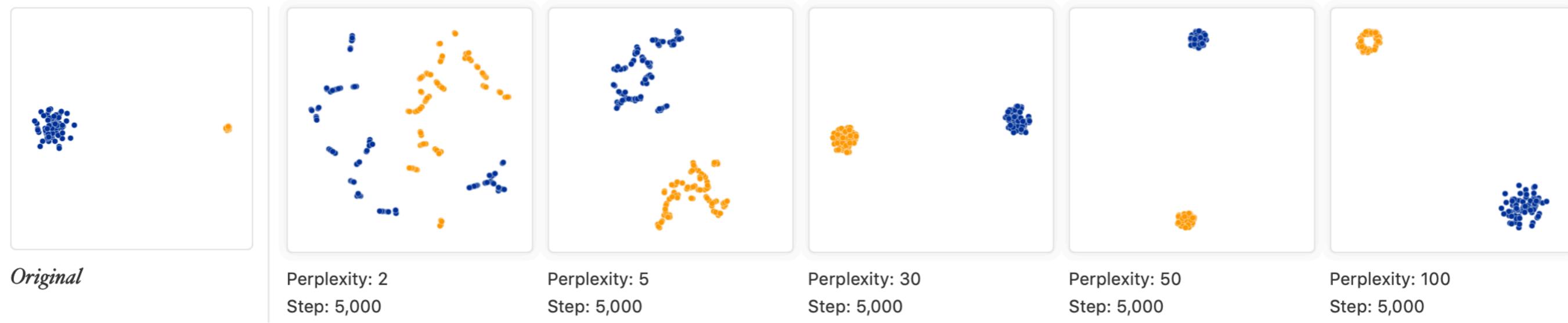
perplexity must be < number of points (here 50)



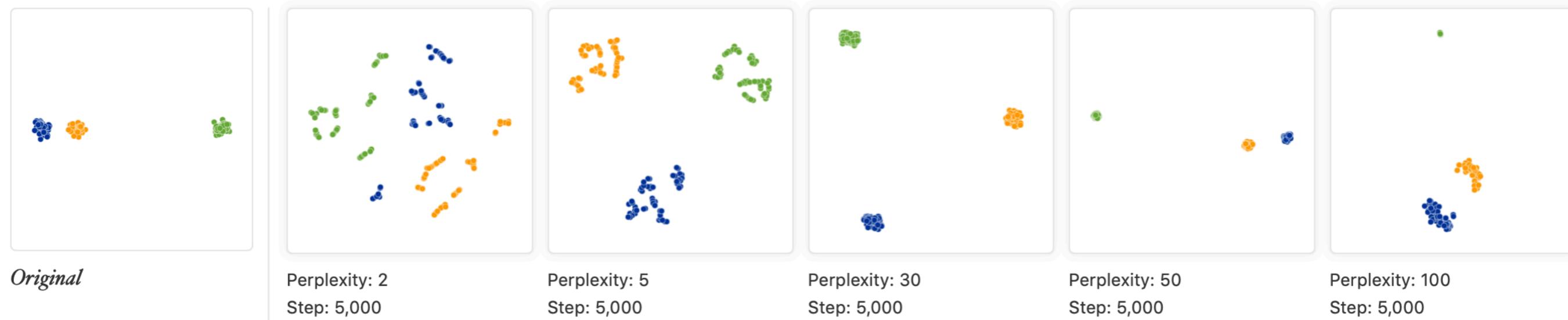
iterate until a stable configuration is reached

Dimensionality Reduction: t-SNE

t-SNE adapts the notion of distance to regional density



⇒ cannot compare relative size of clusters



distance is not reliable

HOW TO USE

A button is causing an action in the visualizations.

Underlined words highlight something in the visualizations.

EXAMPLE

IRIS is loading the dataset Iris.

t-SNE highlights the visualization of the t-SNE algorithm.

Activate the proximity view here, or temporary by pressing **P**.

Activate the brush here, or temporary by pressing **B**.

Activate the component planes to show the distribution of a dimension in the projection, by clicking on the respective density plot here.

You can hover over a projected point to see their positions in the other projections and in the overview, and the details of the data point in the DIMENSIONS section.

WHO AND WHY

Who is using Dimensional Reduction (DR) techniques?

Why they use it?

Data represent some characteristics of the real world as numbers and categories. *Analysts, data scientists or domain experts* seek for patterns in the data that might reveal interesting properties of the real world. A dataset

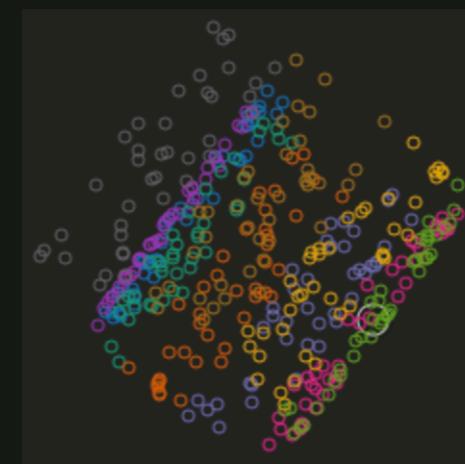
OVERVIEW

II S 3D



B brush
P proximity view,

PCA



method
EIG SVD

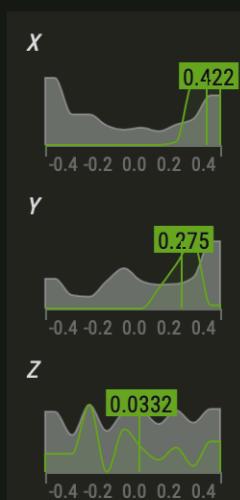
LLE



method
SVD EIG
neighbours 12

DIMENSIONS

[g] 263



MDS iter: 450

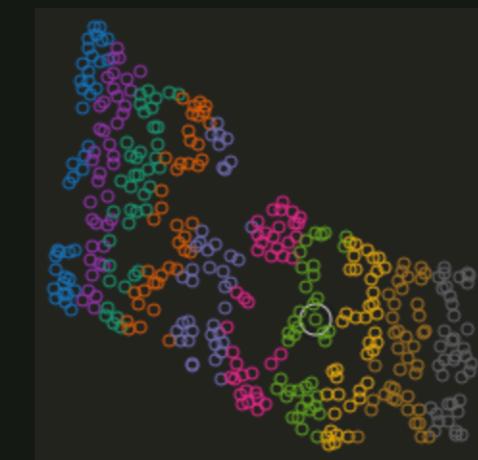


ISOMAP iter: 450



Neighbours 15

TSNE iter: 1000



Perplexity 40

Epsilon 5