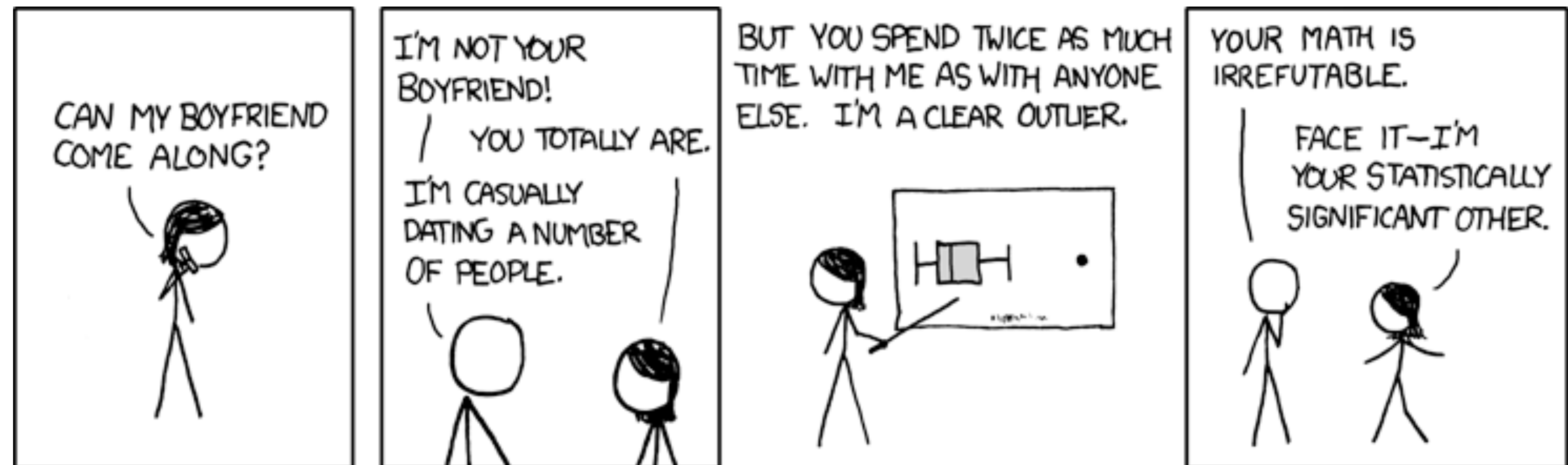


CS-5630 / CS-6630 Visualization

Filtering & Aggregation

Alexander Lex
alex@sci.utah.edu



Administrativa

This/Next Week

Tomorrow: Group announcement

Saturday: Homework 5

Monday: Project Proposal

Put proposal into github repo

Submit only link via canvas

Tuesday: Peer Feedback (mandatory, graded)

Thursday: Class cancelled

filter & Aggregate

Reducing Items and Attributes

➔ Filter

→ Items



→ Attributes



➔ Aggregate

→ Items



→ Attributes



Filter

elements are eliminated

What drives filters?

Any possible function that partitions a dataset into two sets

Bigger/smaller than x

Fold-change

Noisy/insignificant



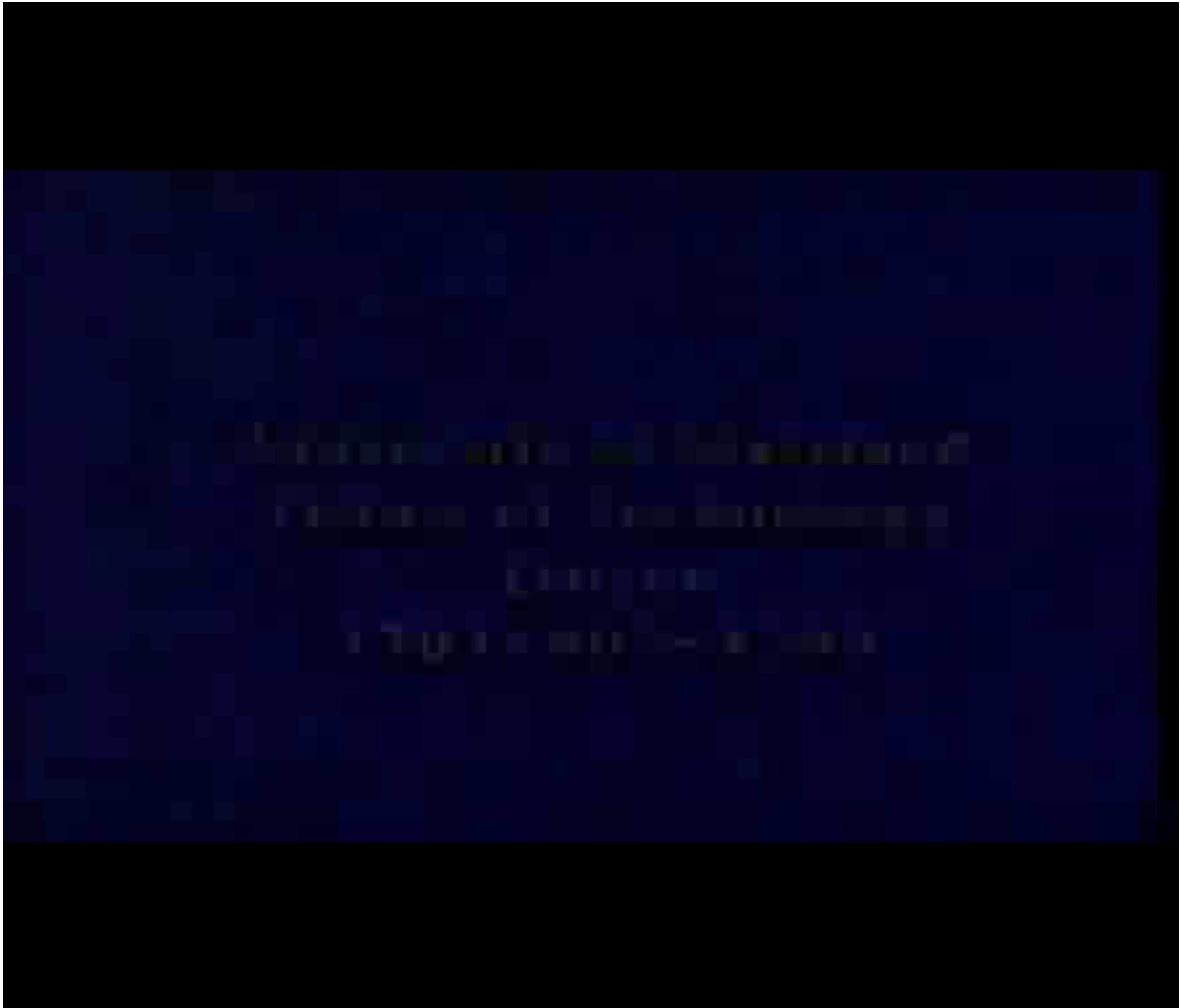
Dynamic Queries / Filters

coupling between encoding and interaction so that user can immediately see the results of an action

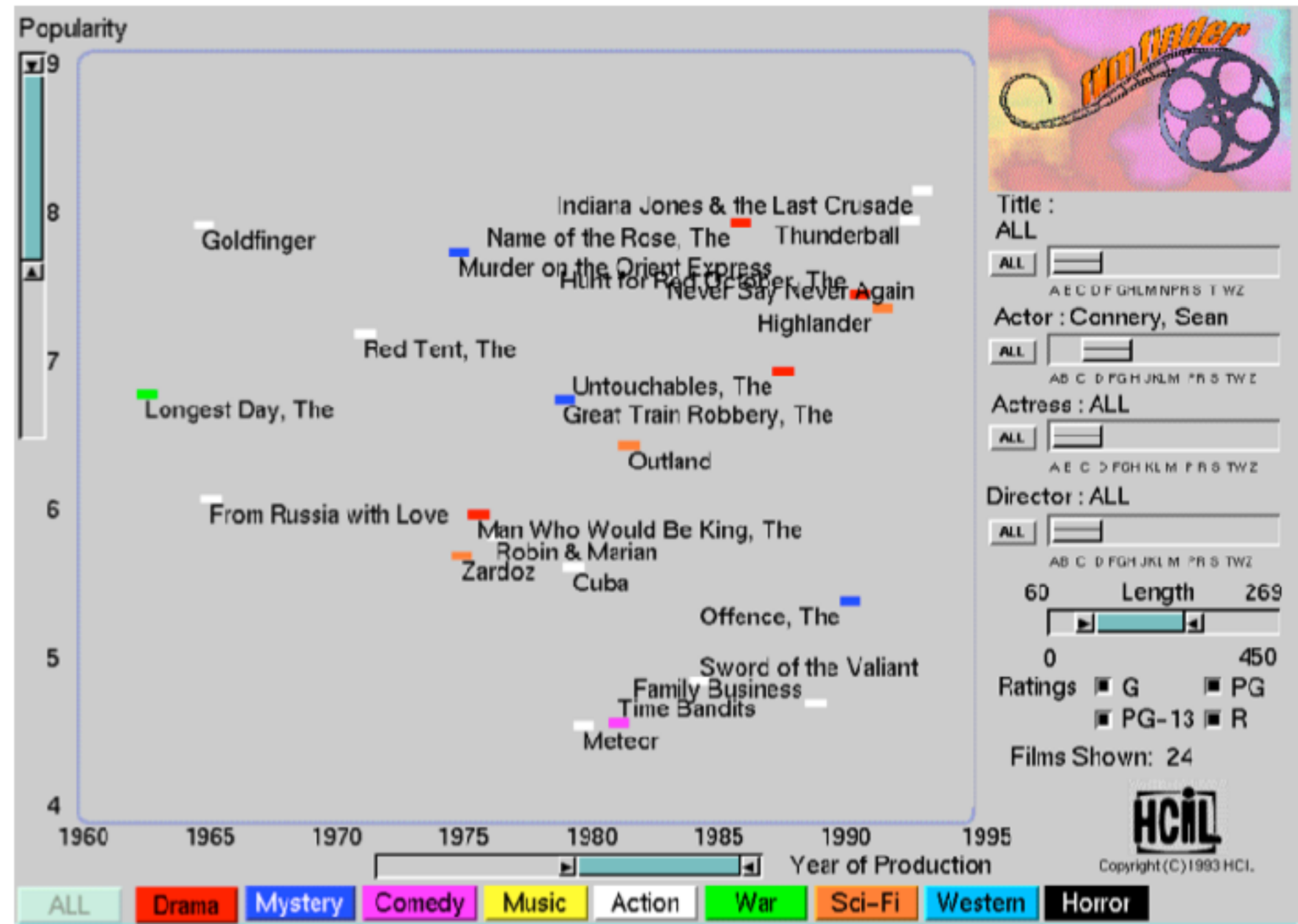
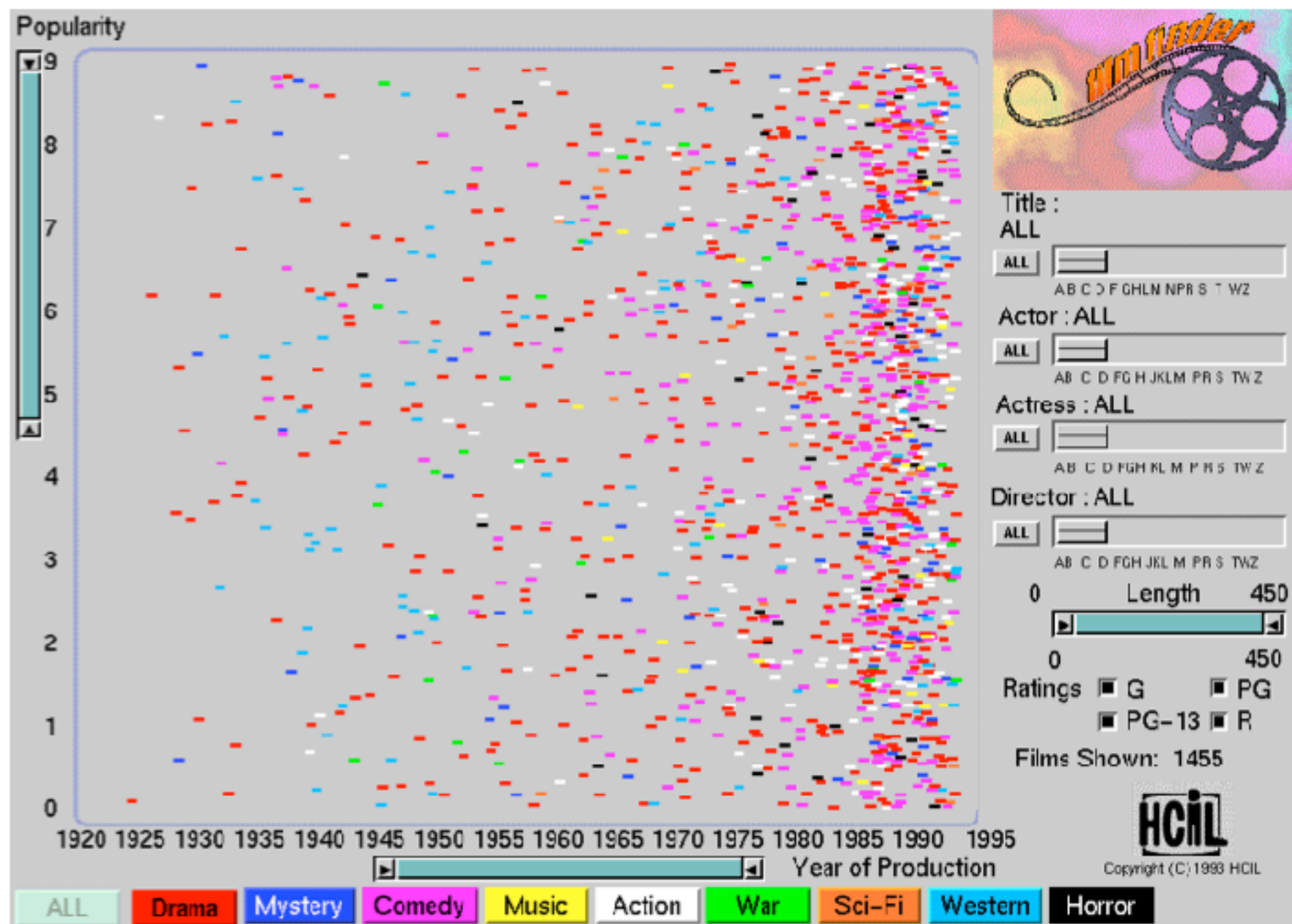
Queries: start with 0, add in elements

Filters: start with all, remove elements

Approach depends on dataset size




ITEM FILTERING





FIND A RESTAURANT

FIND A LOCATION

FILTER

 All grades

 All violations

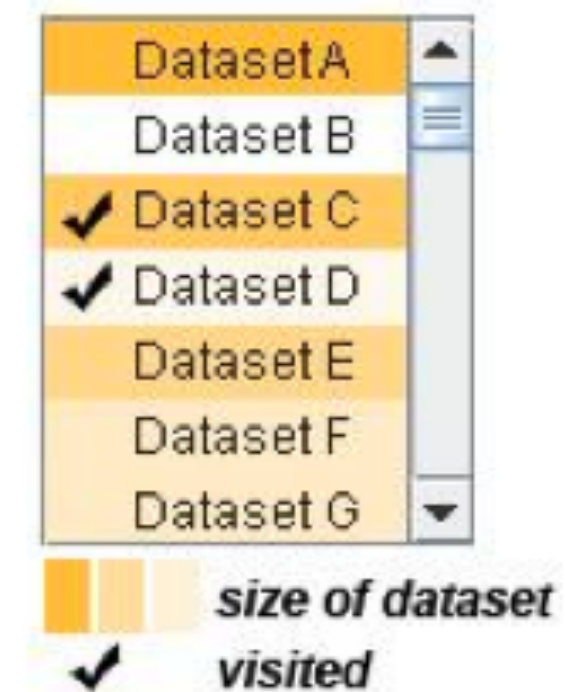
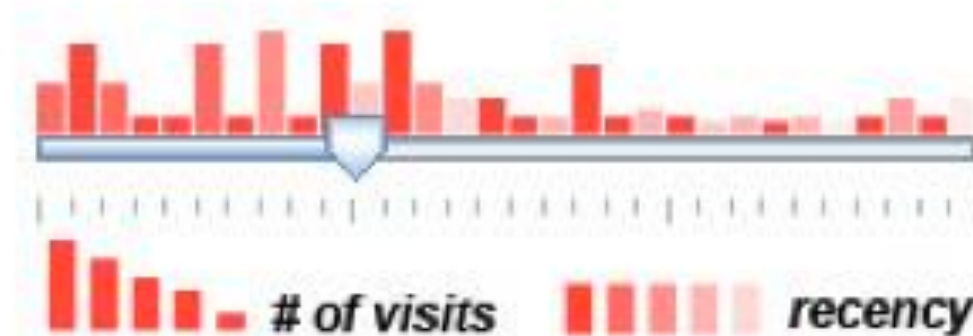
 All cuisines



Scented Widgets

information scent: user's (imperfect) perception of data

GOAL: lower the cost of information foraging through better cues



Aggregation

Aggregate

a group of elements is represented by a (typically smaller) number of derived elements

→ Items

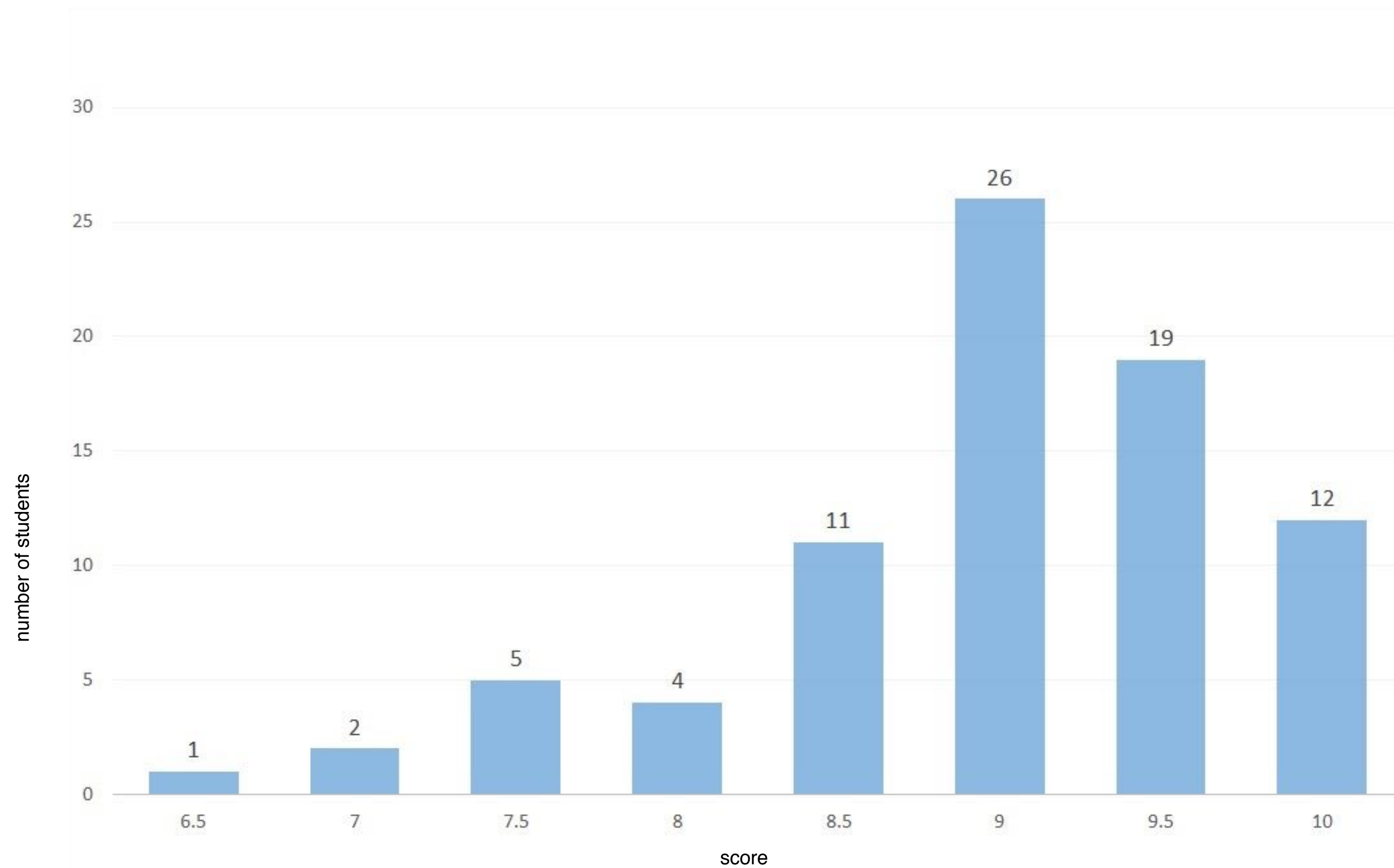


→ Attributes



Item Aggregation

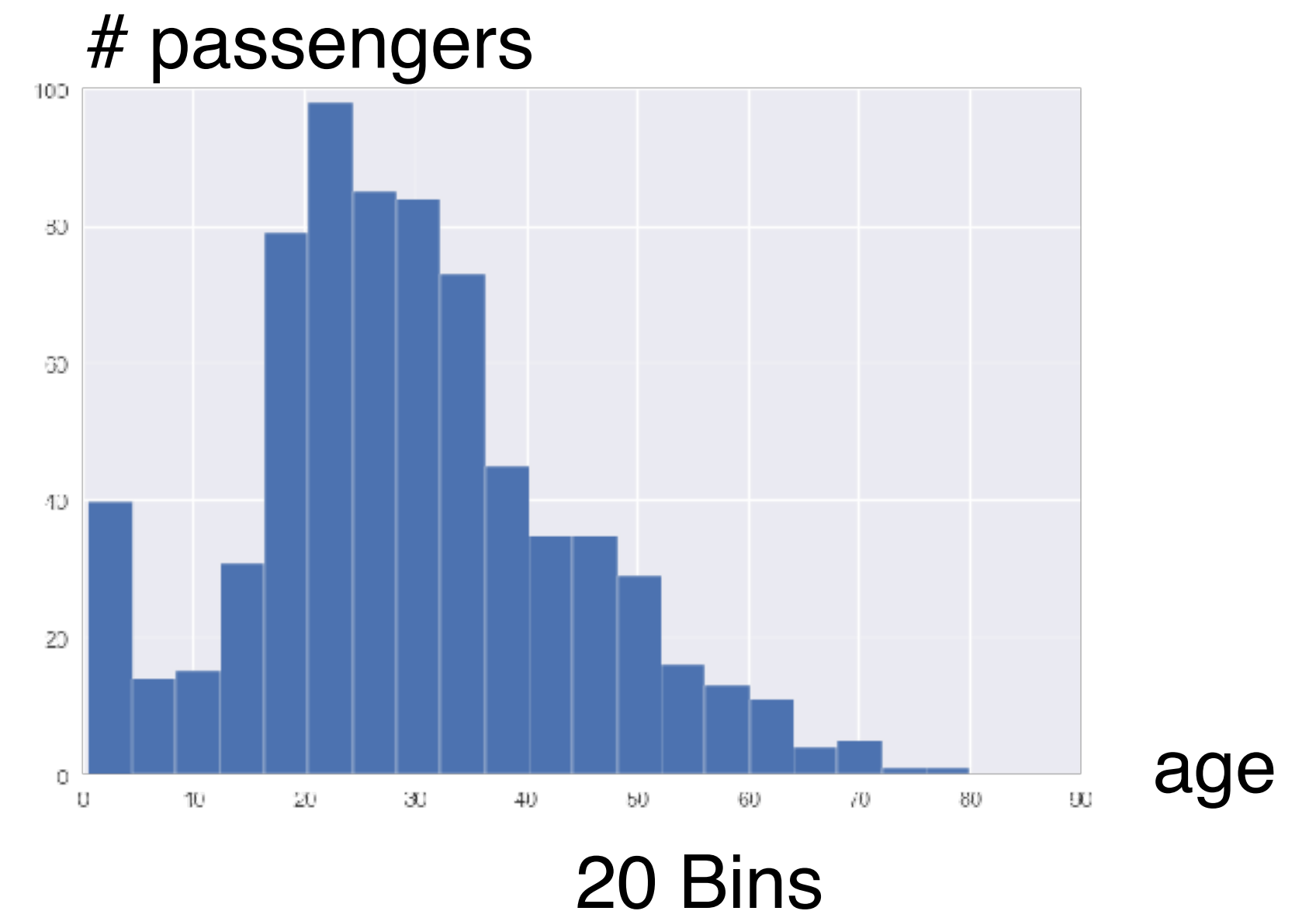
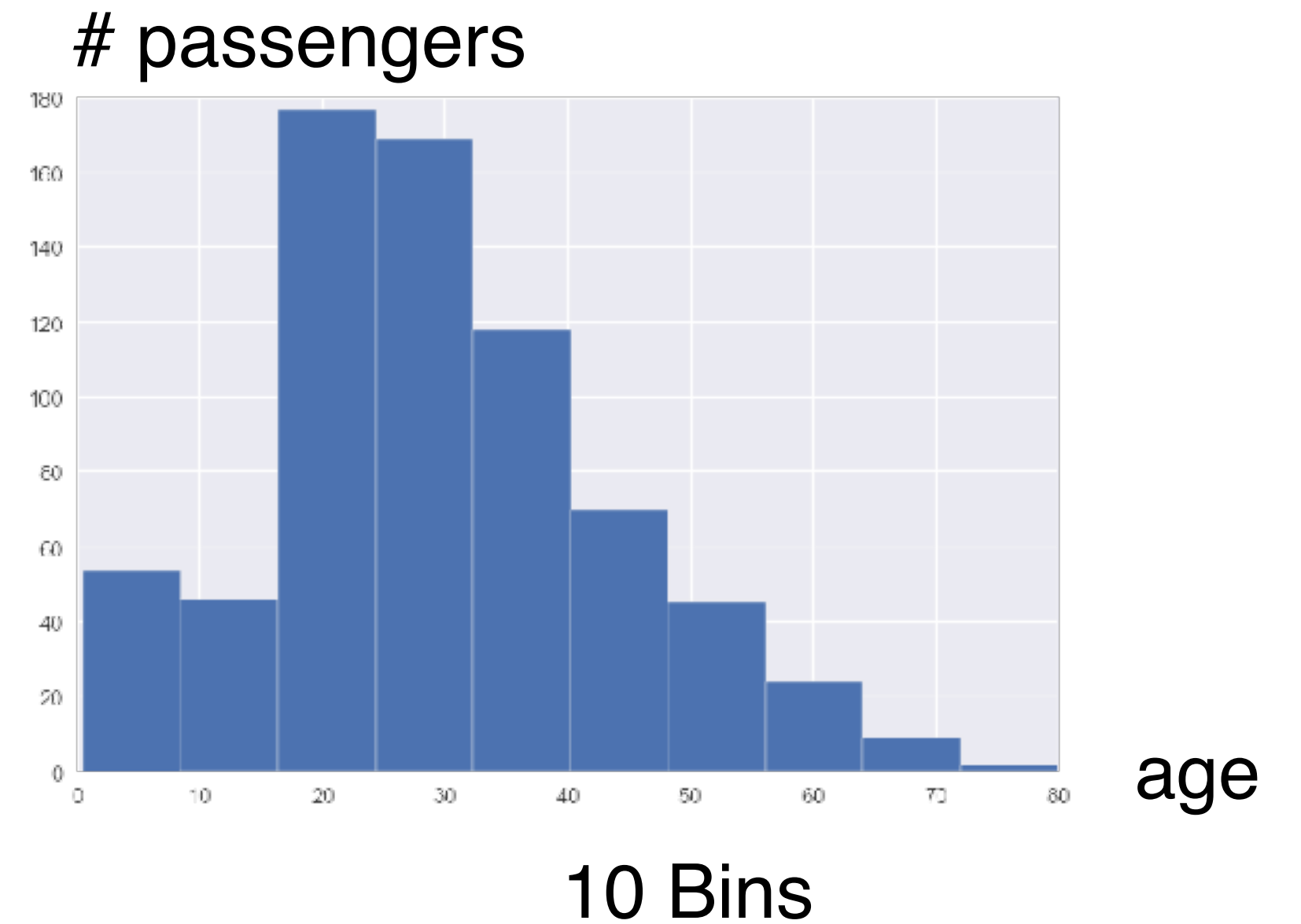
Histogram



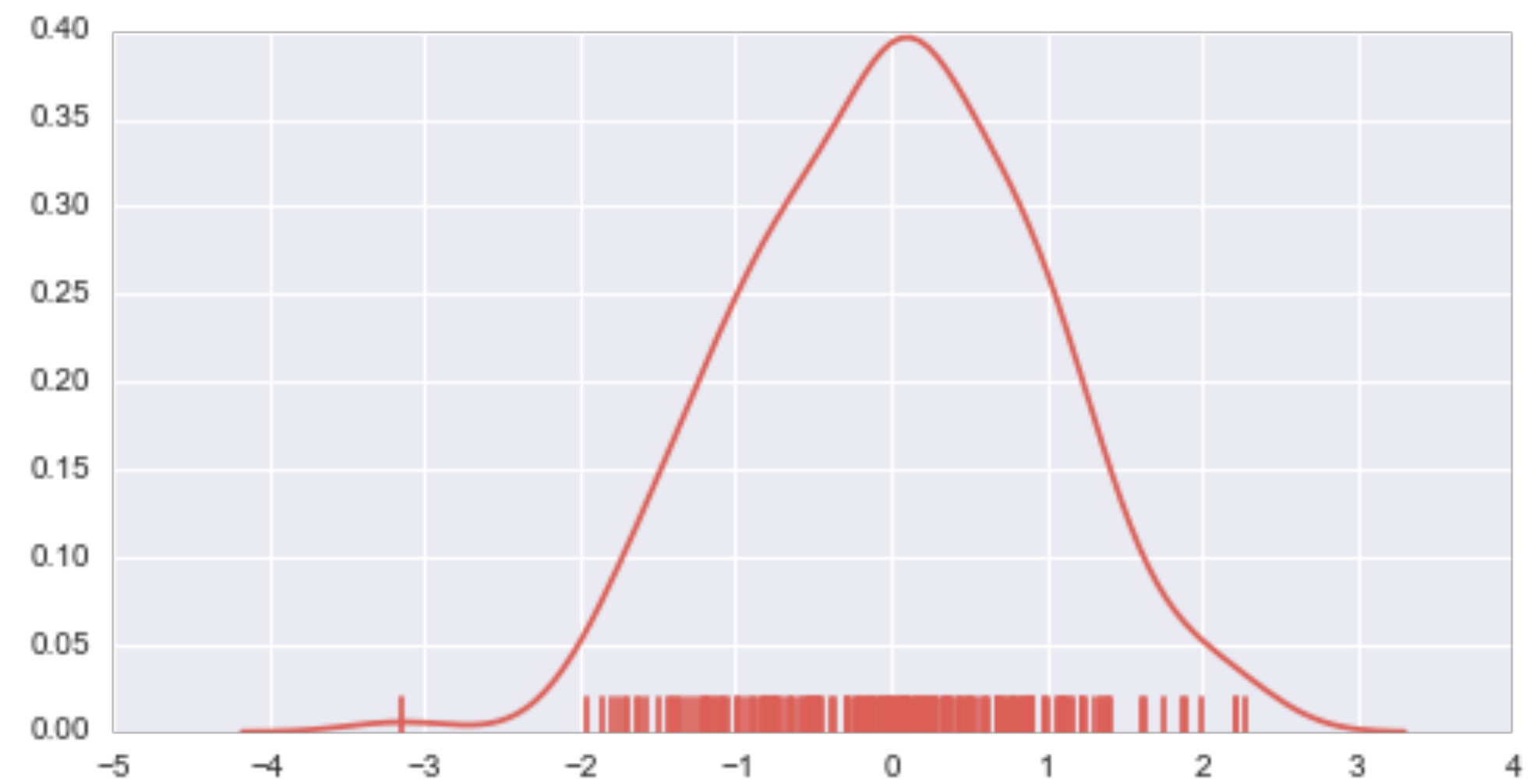
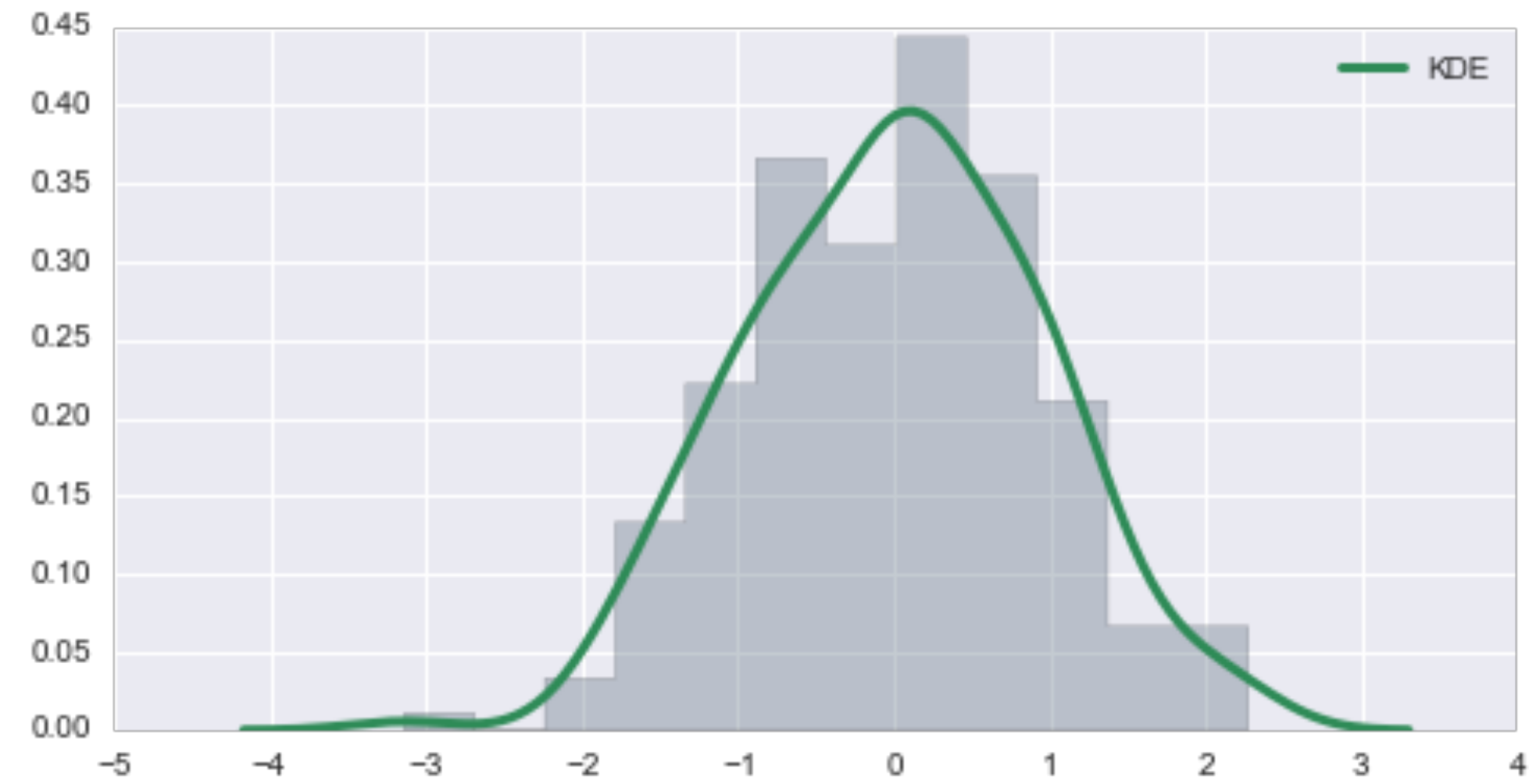
Histogram

Good #bins hard to predict
make interactive!

rule of thumb: #bins = \sqrt{n}



Density Plots



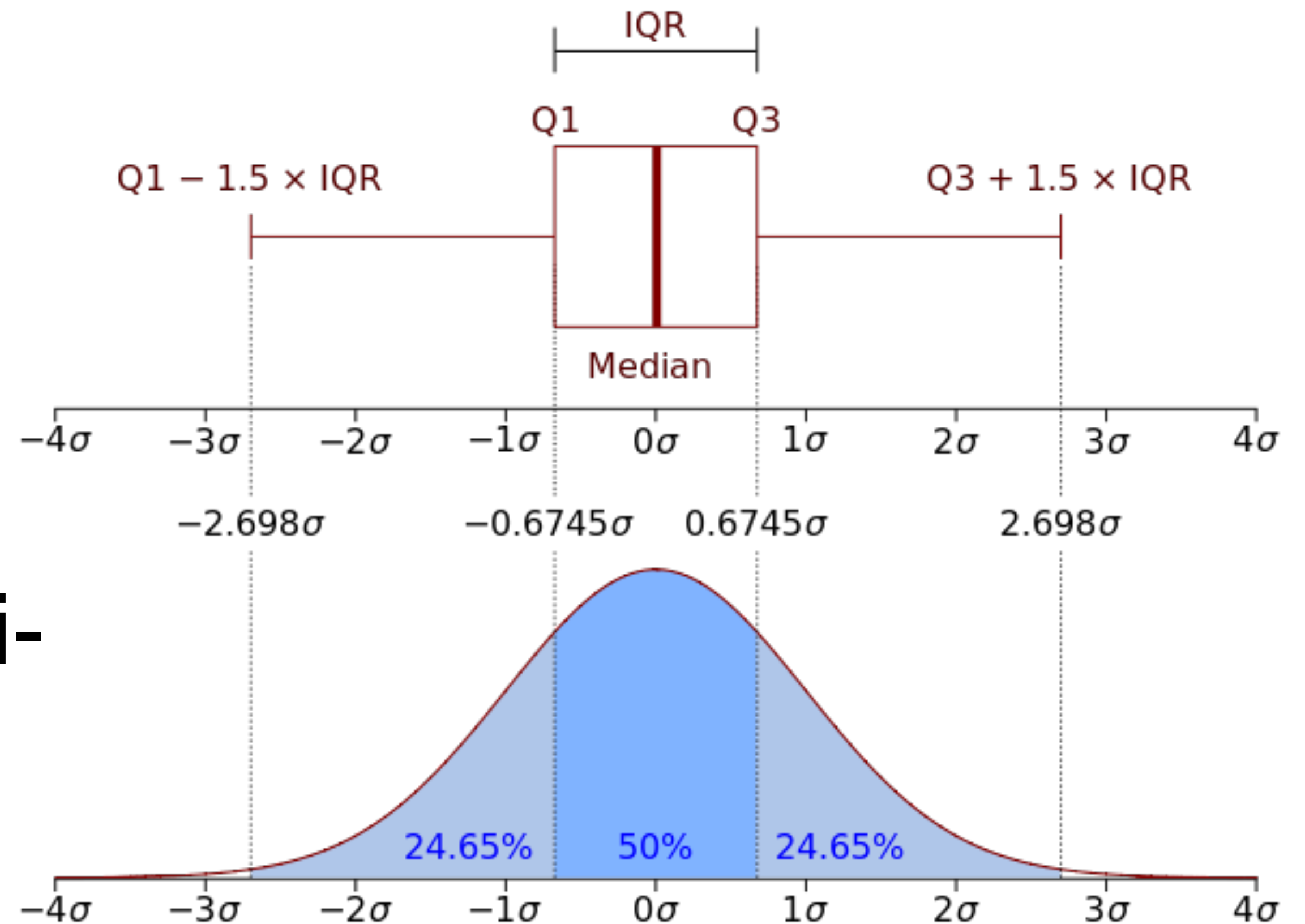
Box Plots

aka Box-and-Whisker Plot

Show outliers as points!

Not so great for non-normal distributed data

Especially bad for bi- or multi-modal distributions



One Boxplot, Four Distributions

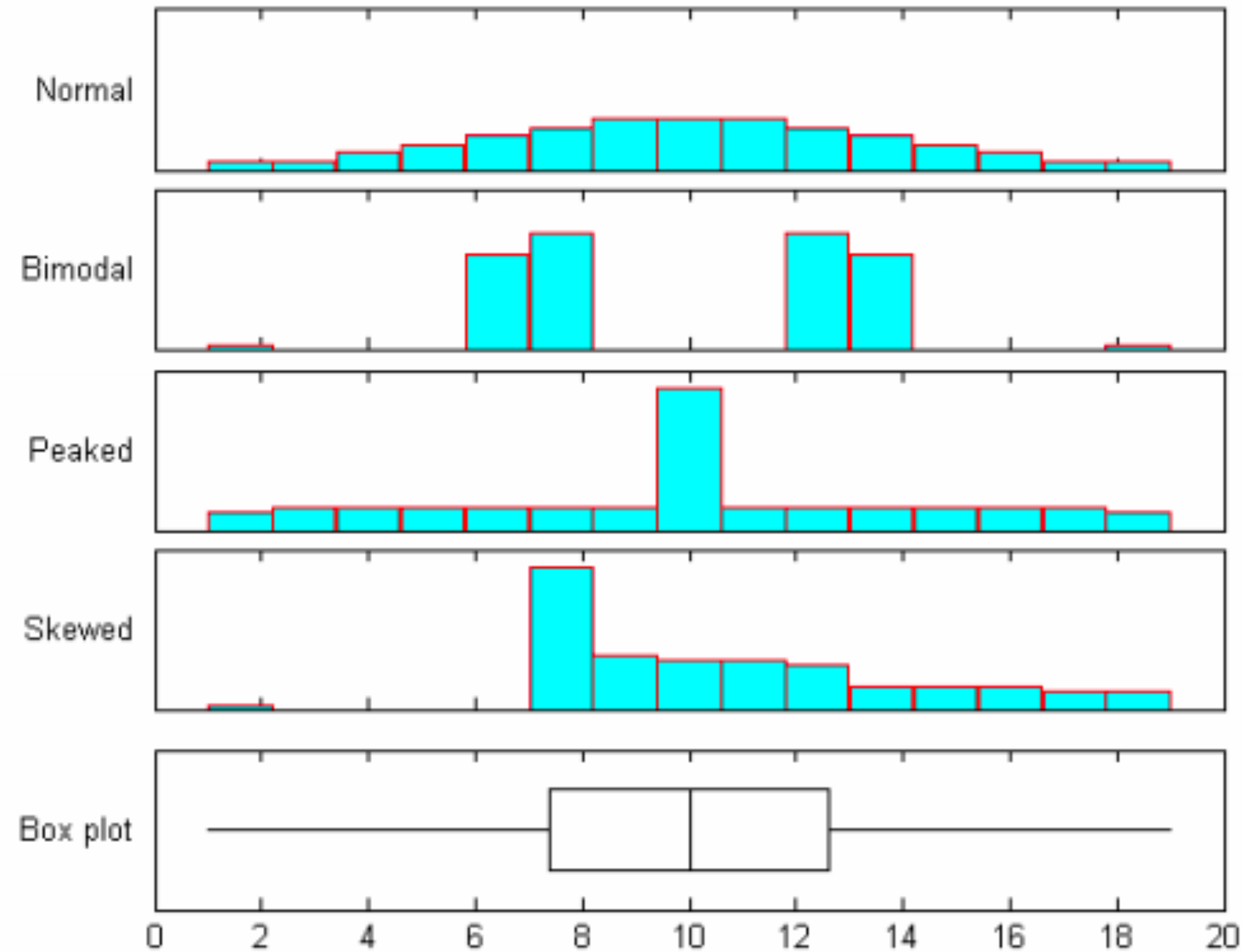
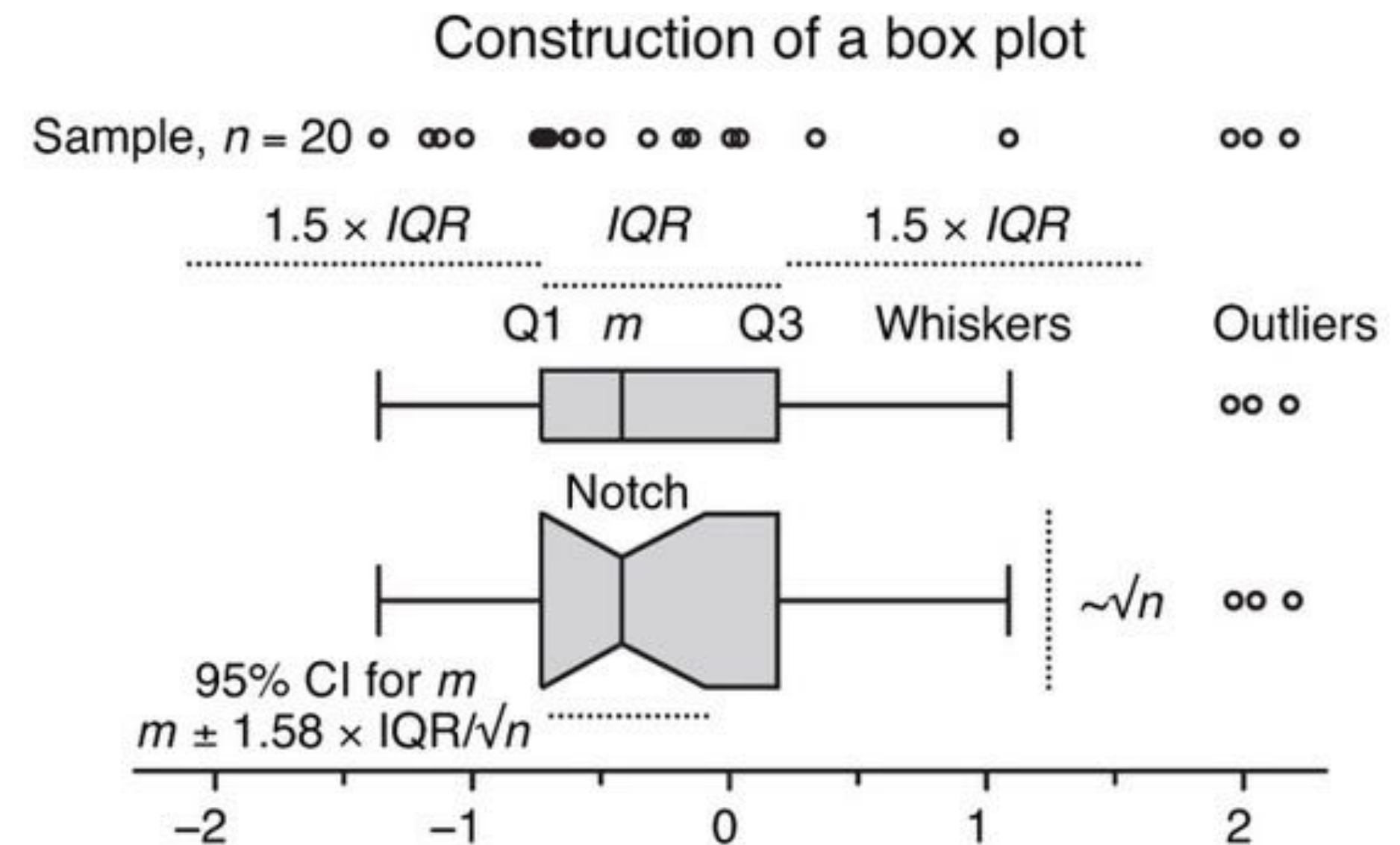


Figure 1: Histograms and box plot: four samples each of size 100

Notched Box Plots

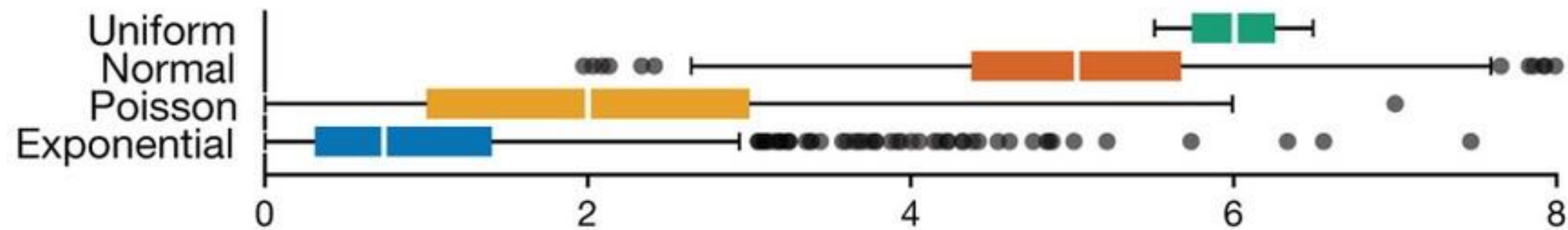
Notch shows
 $m \pm 1.58 \times IQR/\sqrt{n}$

A guide to statistical
significance.

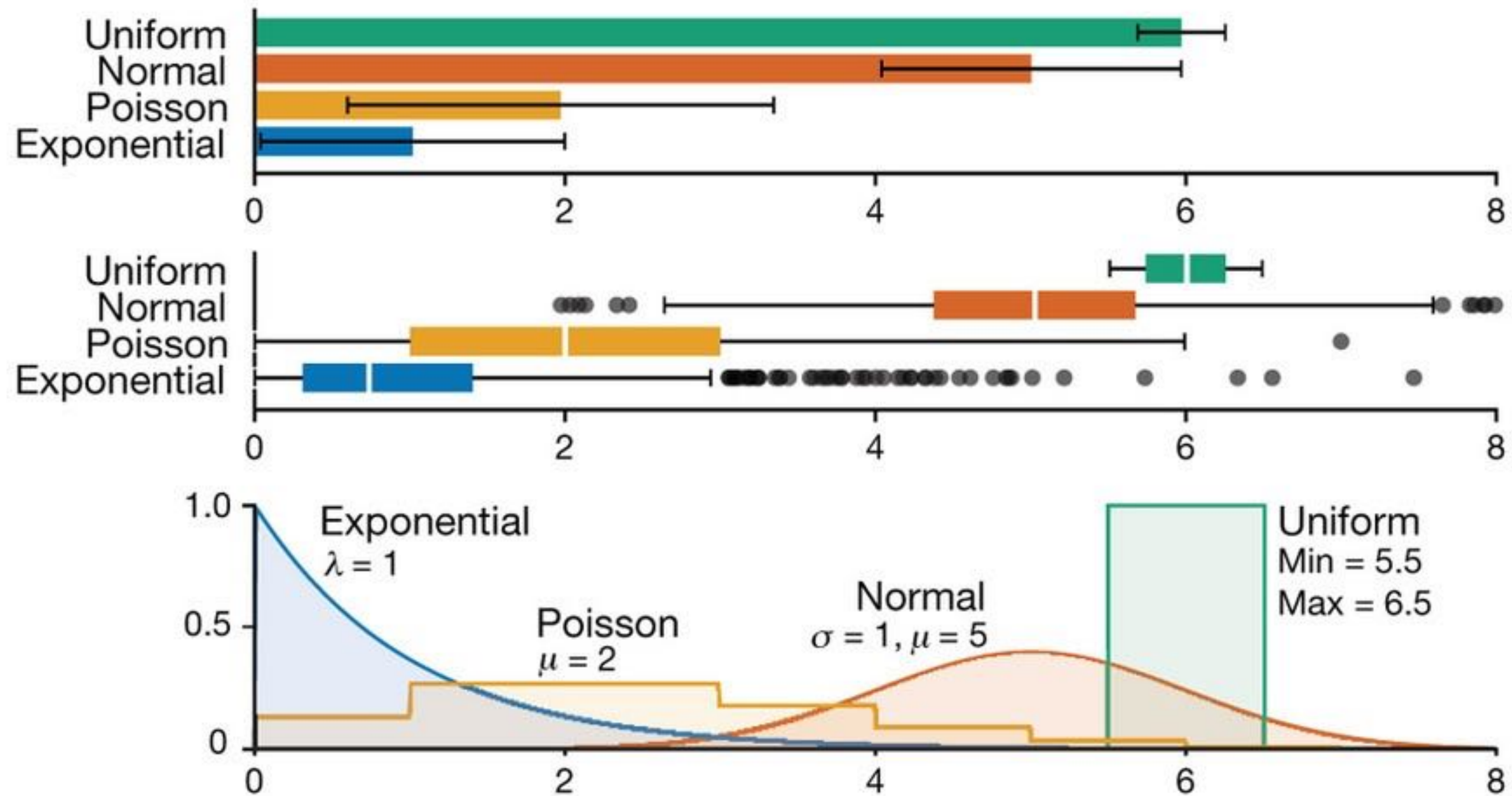


Box(and Whisker) Plots

<http://xkcd.com/539/>

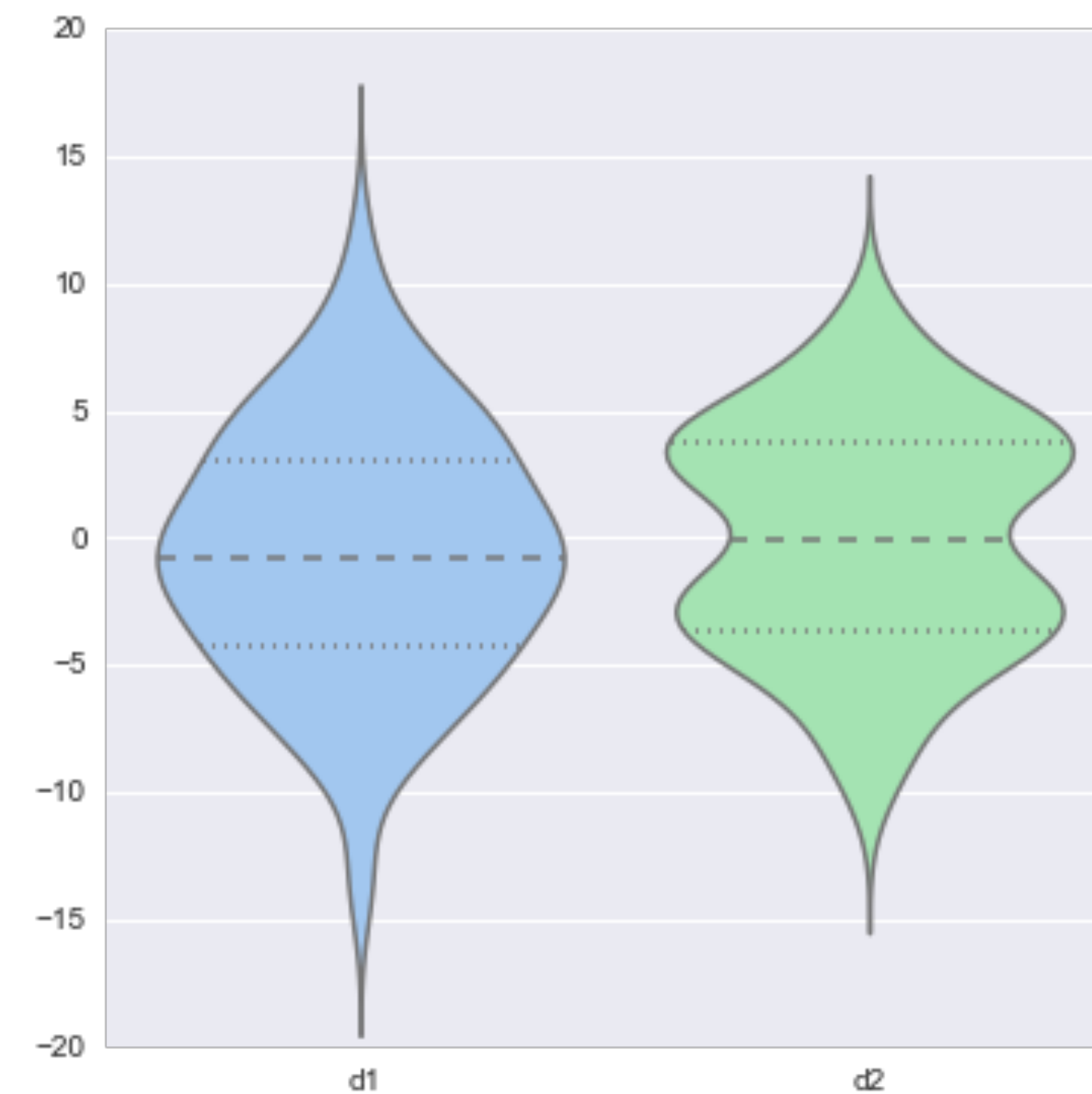
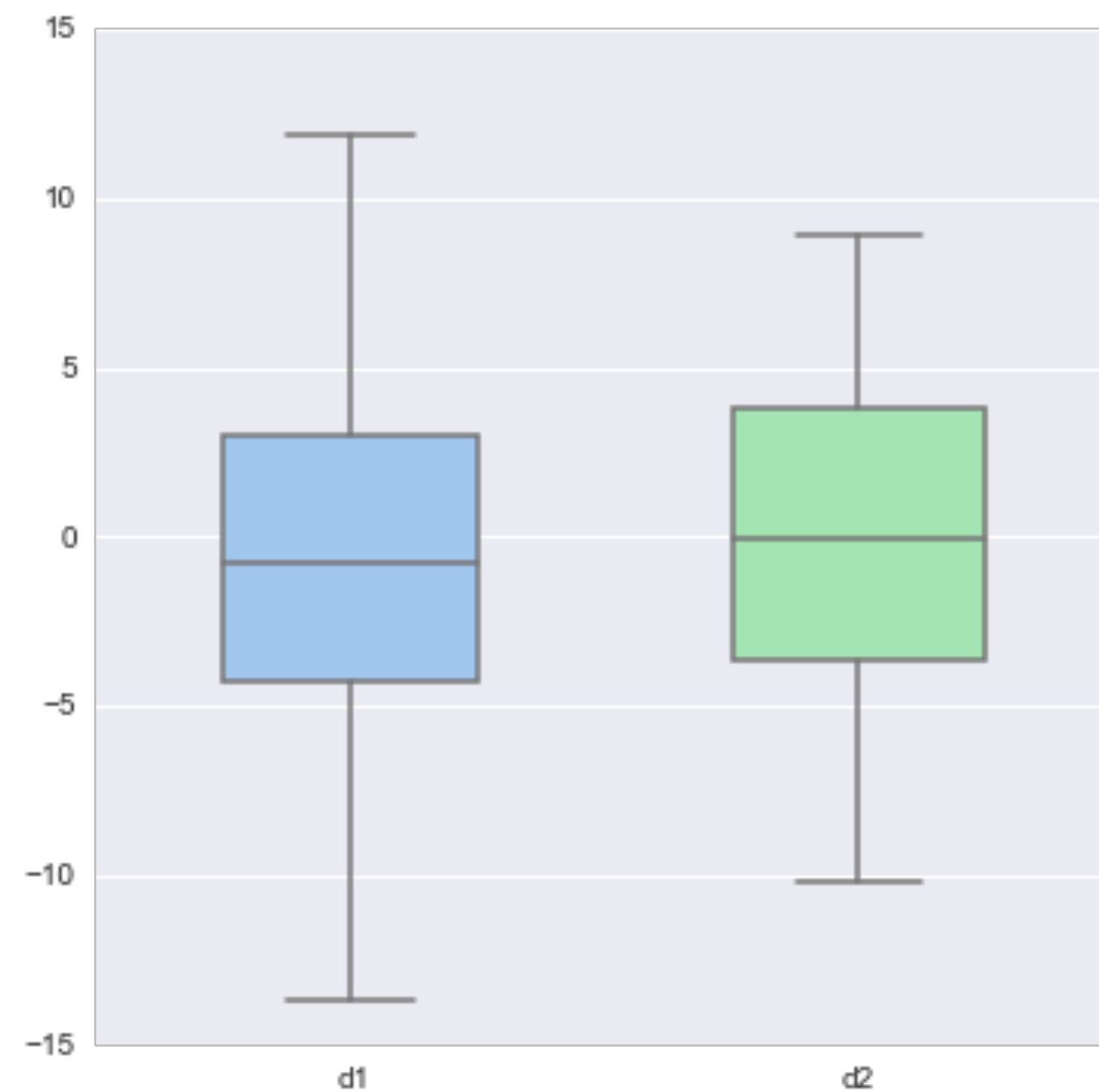


Comparison



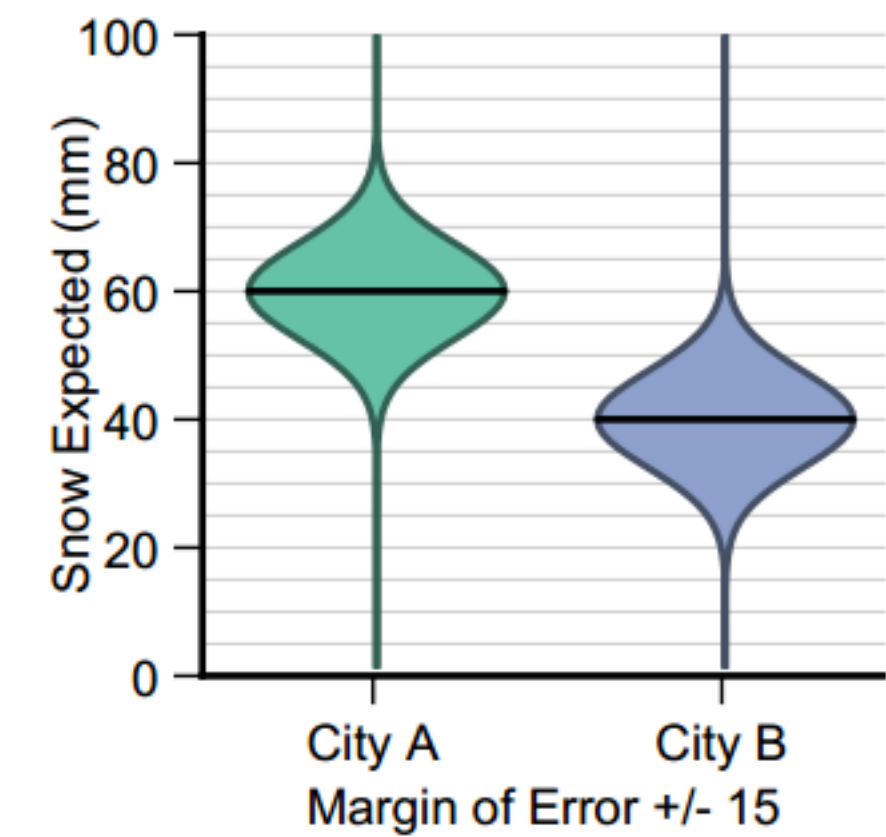
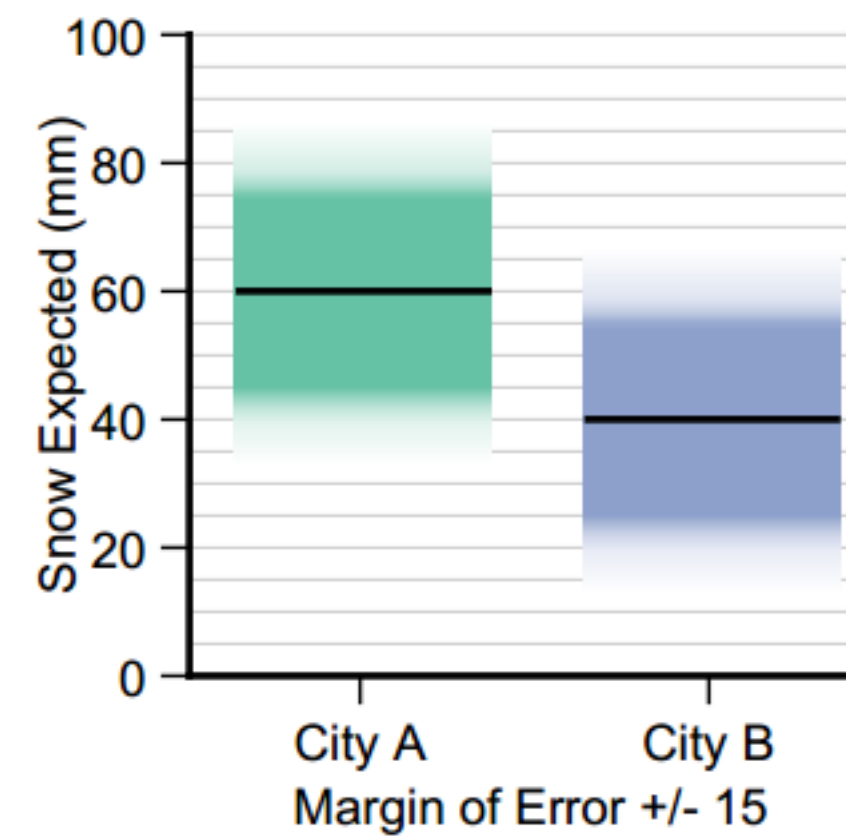
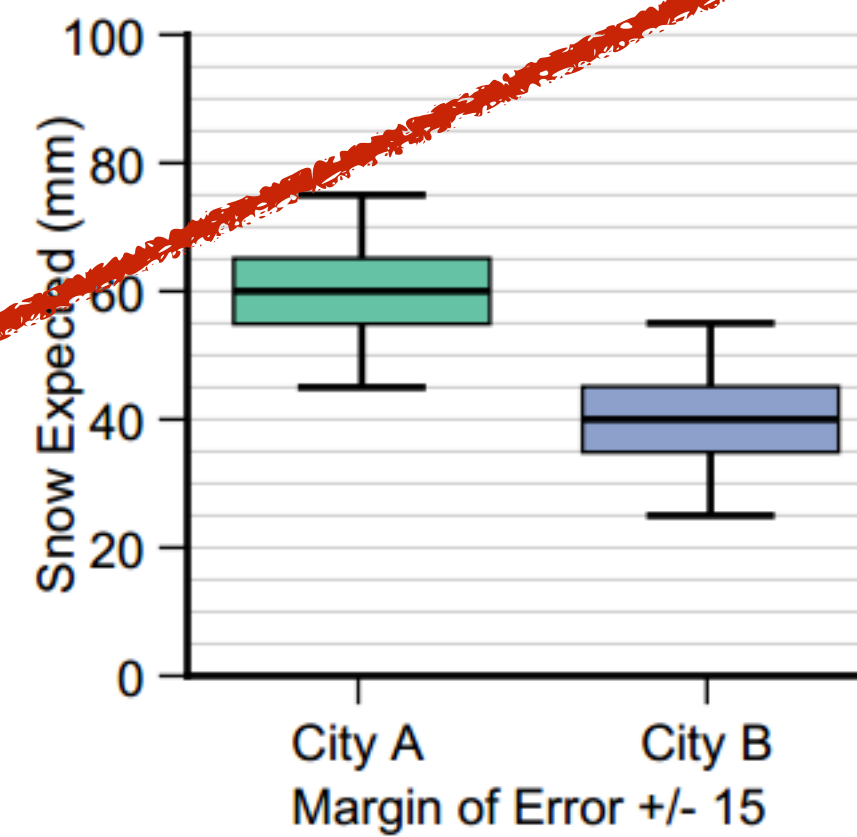
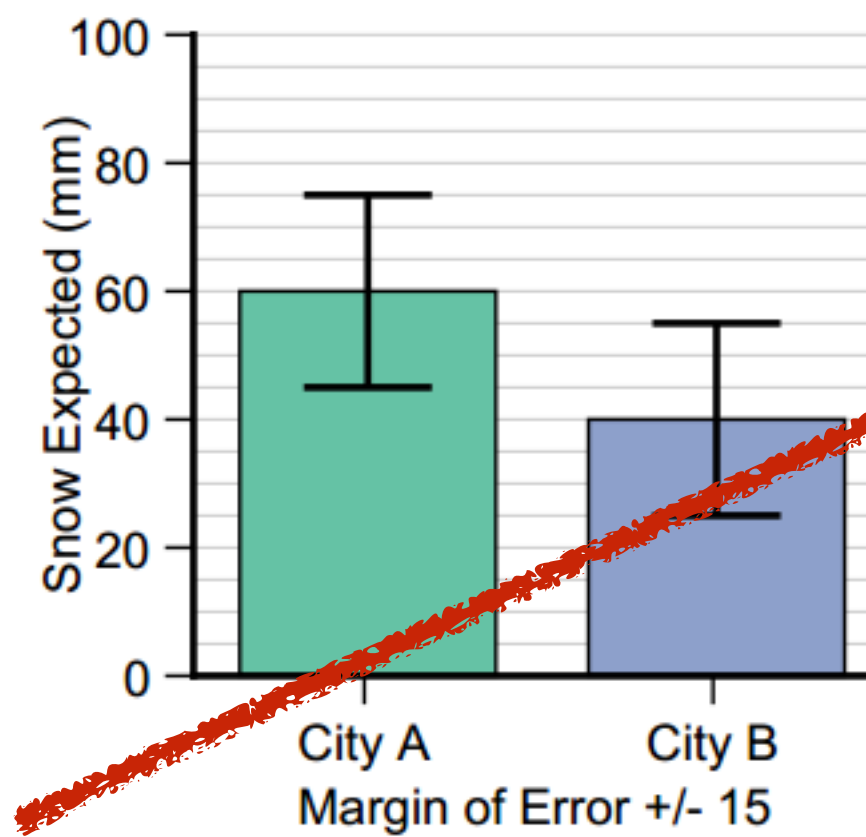
Violin Plot

= Box Plot + Probability Density Function



Showing Expected Values & Uncertainty

NOT a distribution!

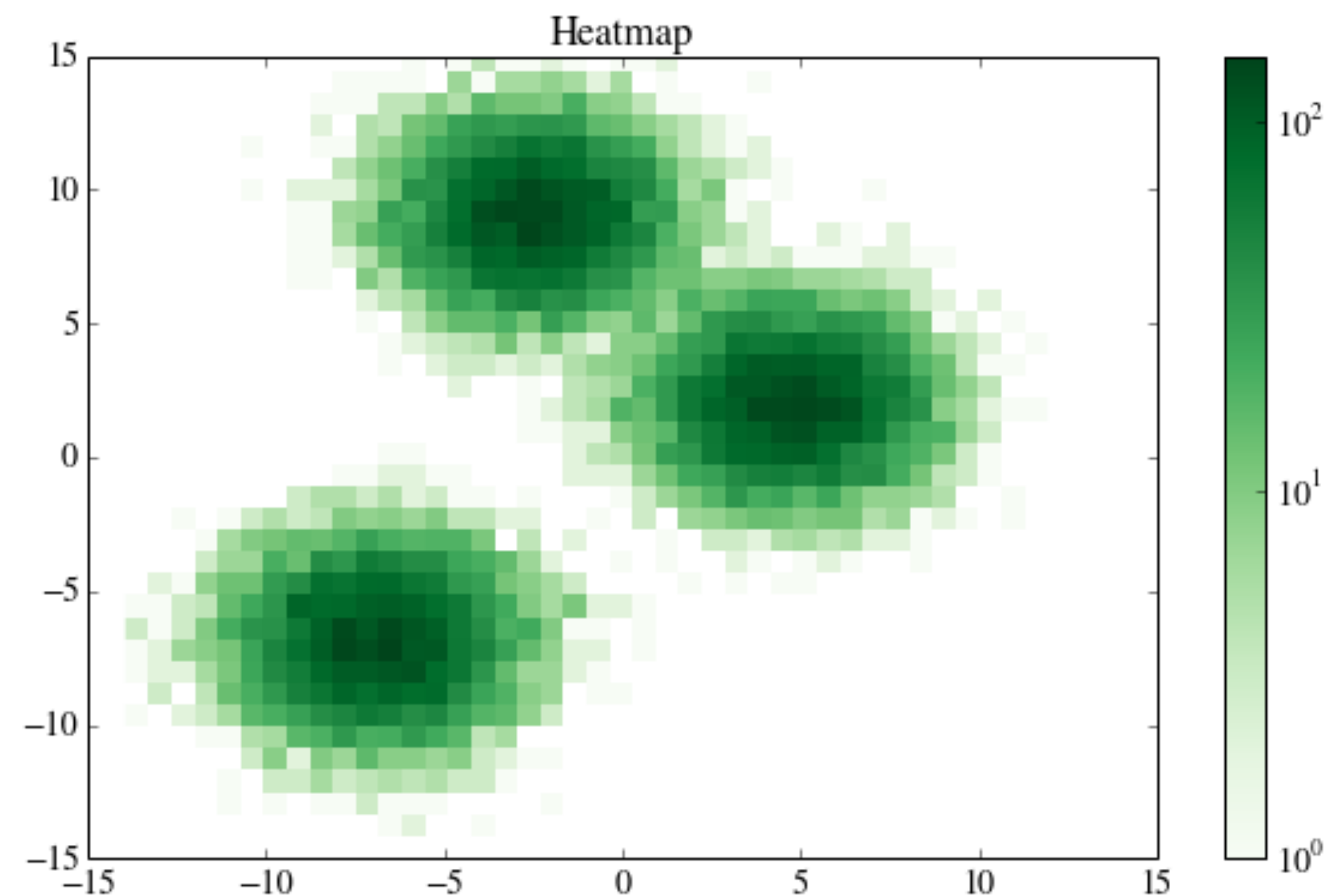
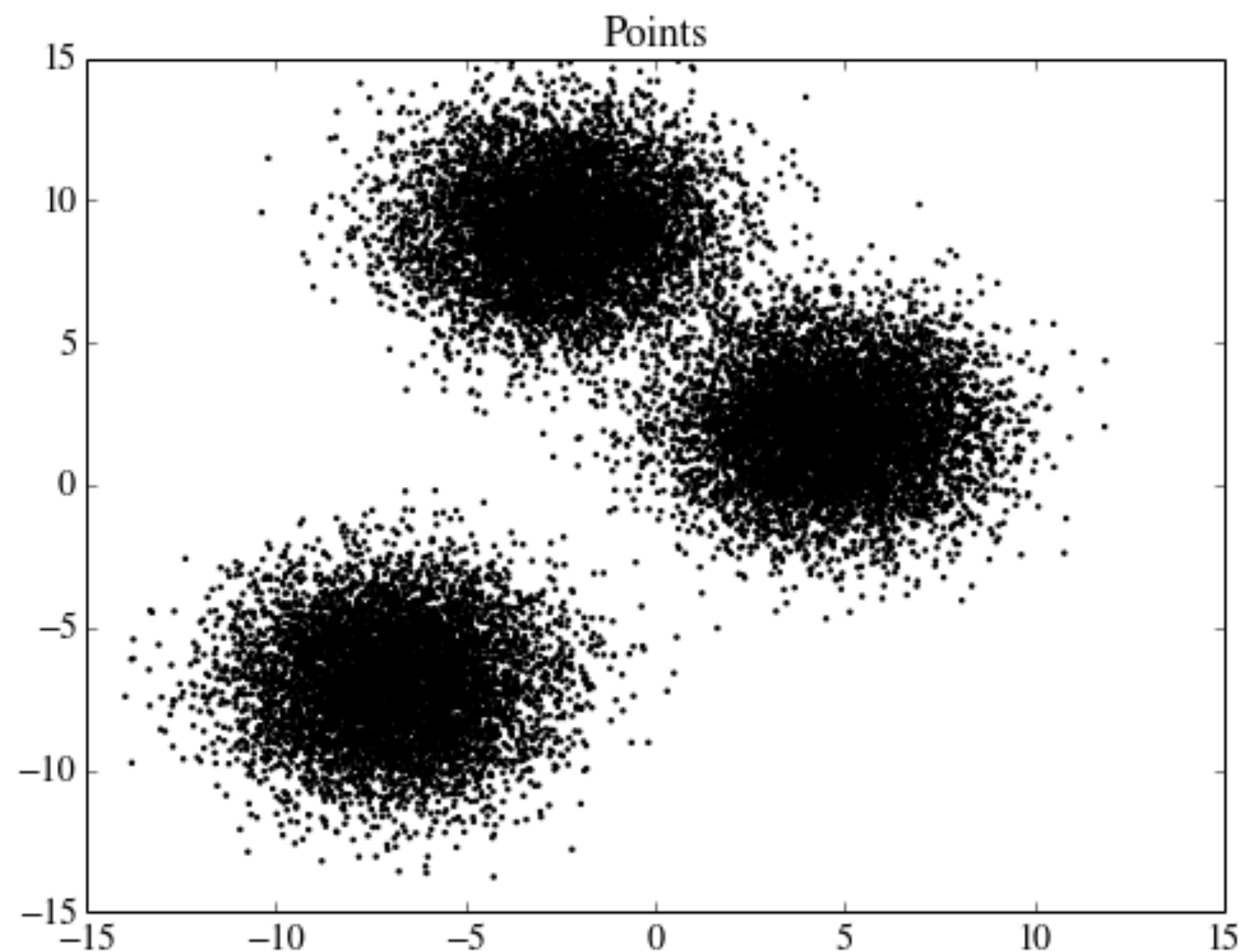


Error Bars Considered Harmful:
Exploring Alternate Encodings for Mean and Error
Michael Correll, and Michael Gleicher

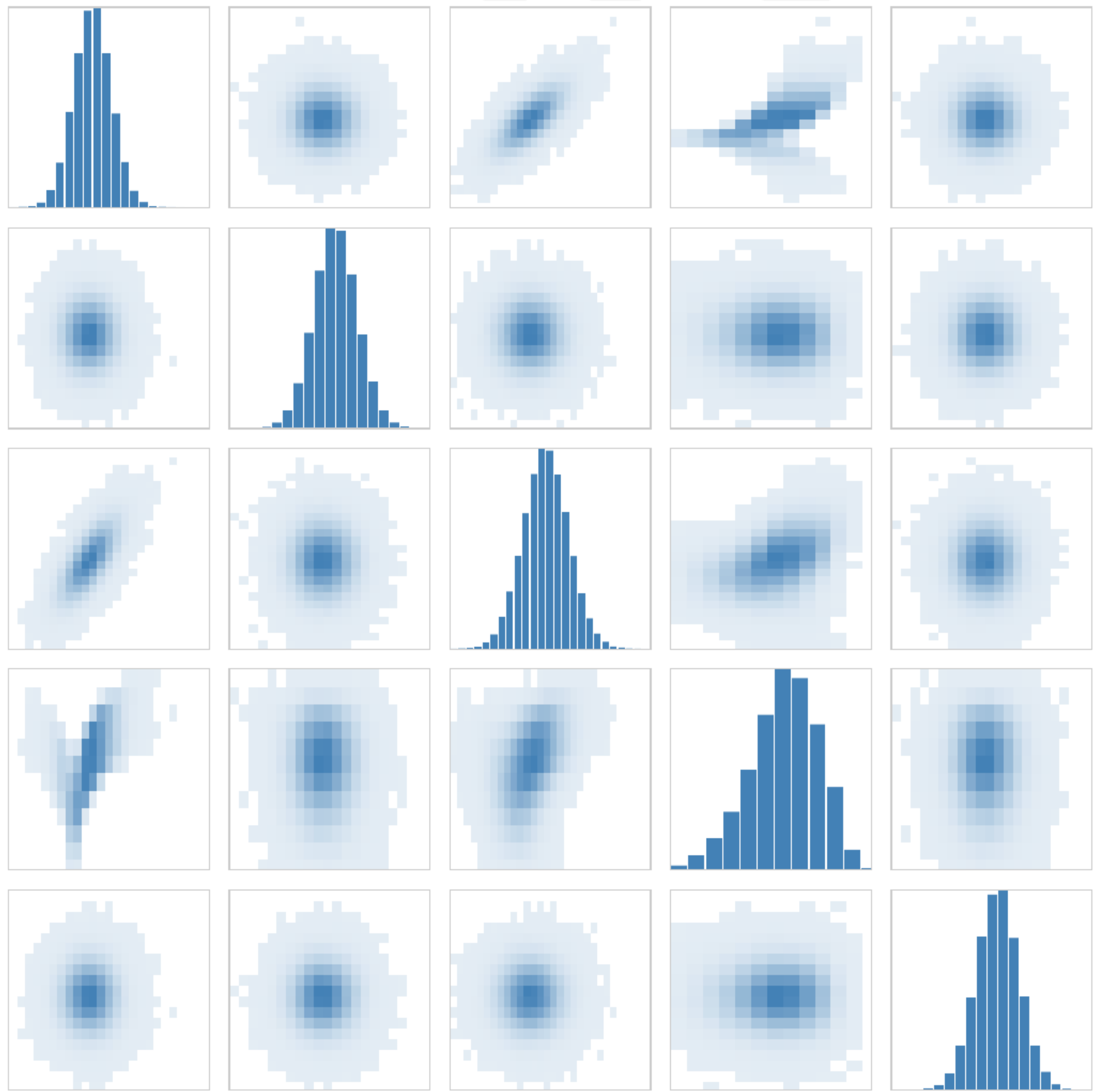
Heat Maps

binning of scatterplots

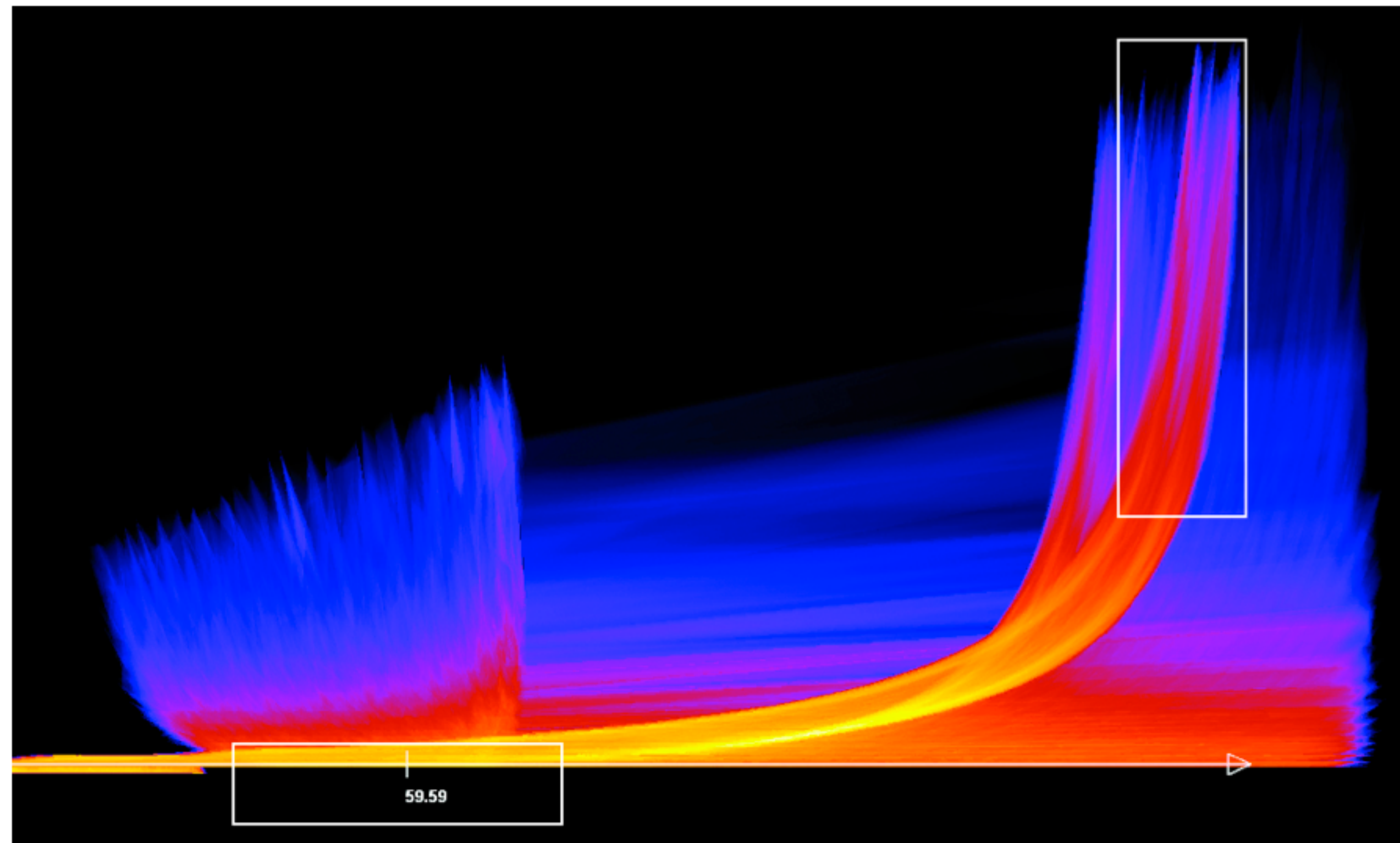
instead of drawing every point, calculate grid and intensities



Interactive Binned Scatterplot Matrix Dimensions: 5 Bins: 20 Data Points: 100k



Continuous Scatterplot



Spatial Aggregation

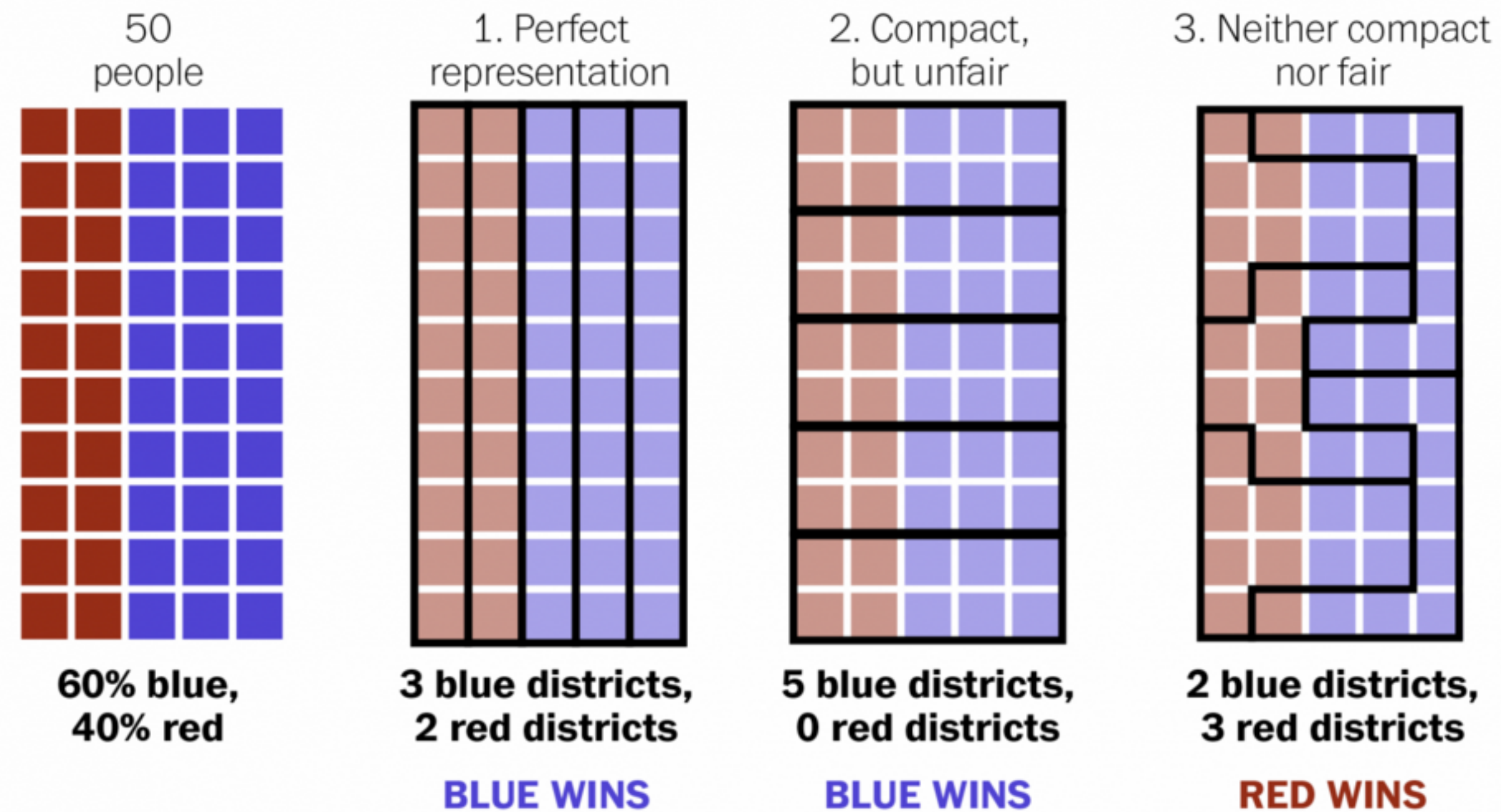
modifiable areal unit problem

in cartography, changing the boundaries of the regions used to analyze data can yield dramatically different results



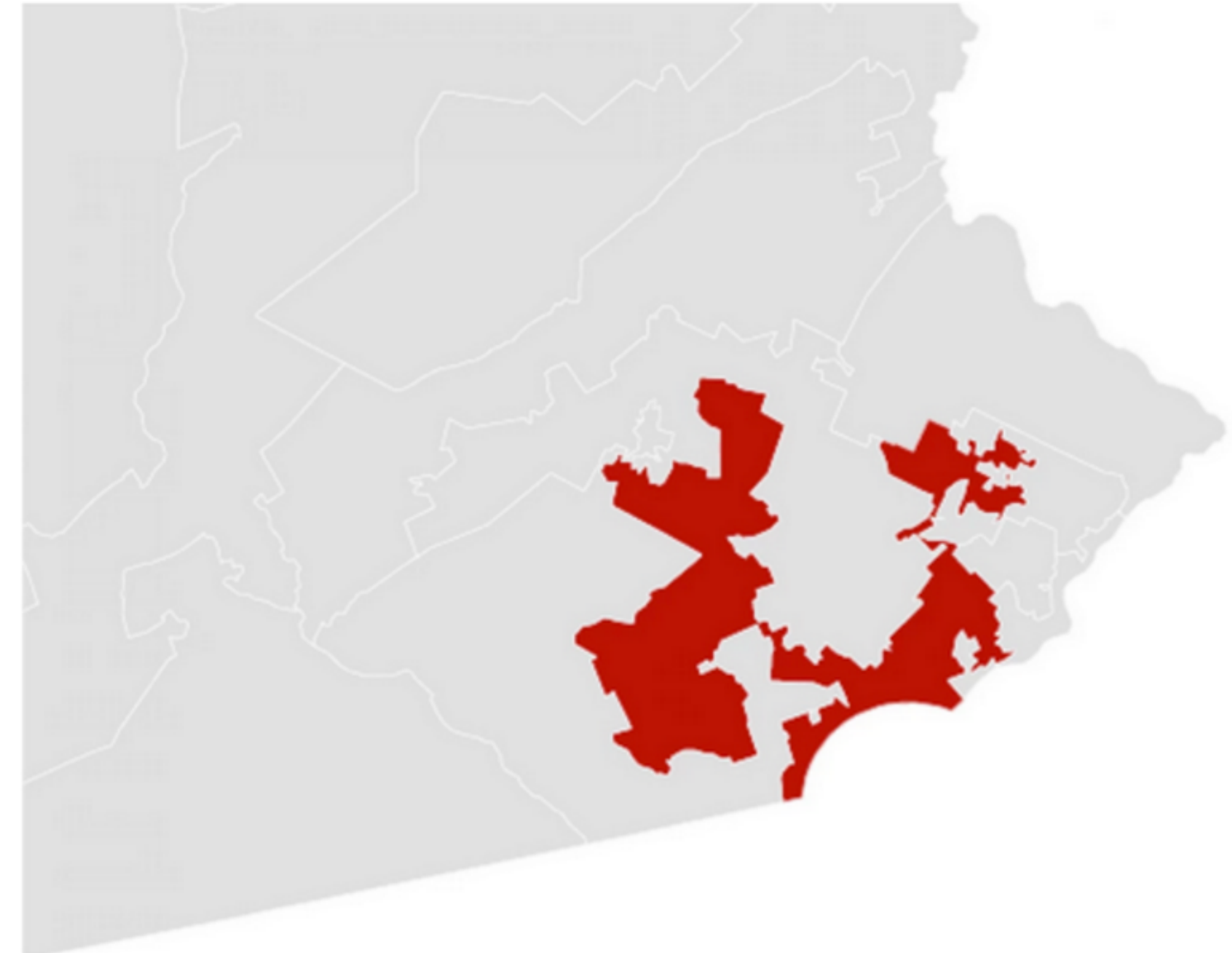
Gerrymandering, explained

Three different ways to divide 50 people into five districts



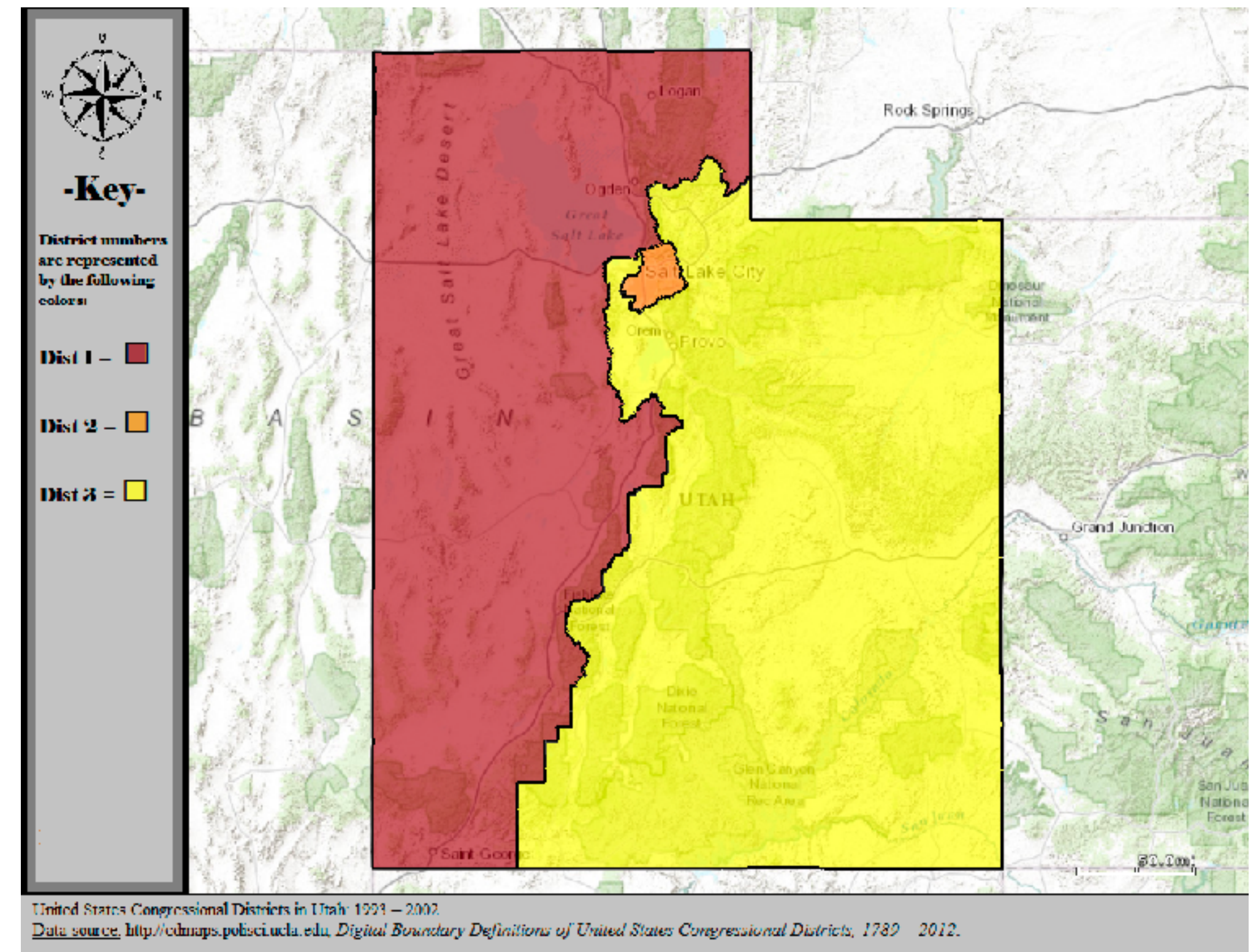
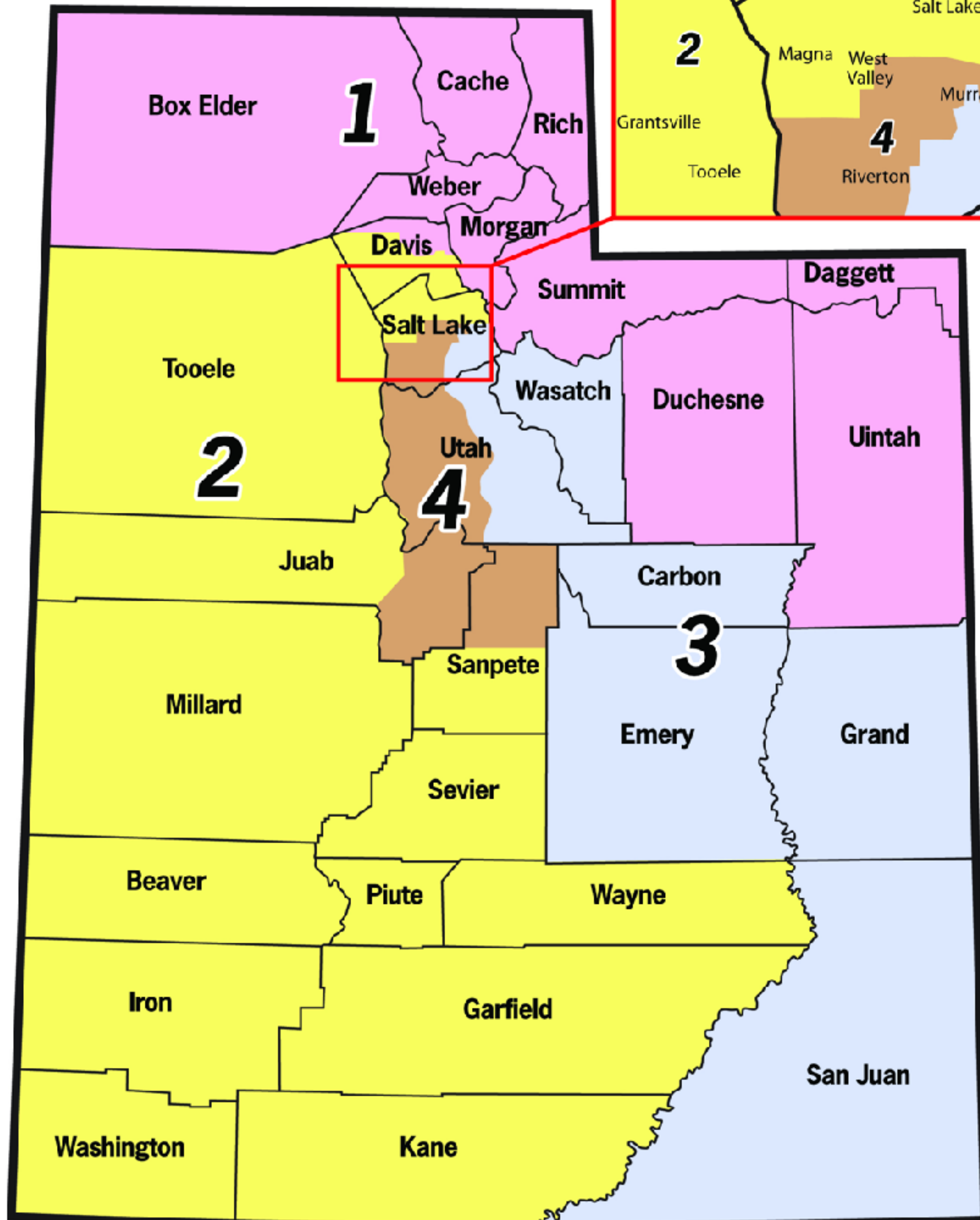
WASHINGTONPOST.COM/**WONKBLOG**

Adapted from Stephen Nass



A real district in Pennsylvania
Democrats won 51% of the vote
but only 5 out of 18 house seats

Congressional Districts



Valid till 2002

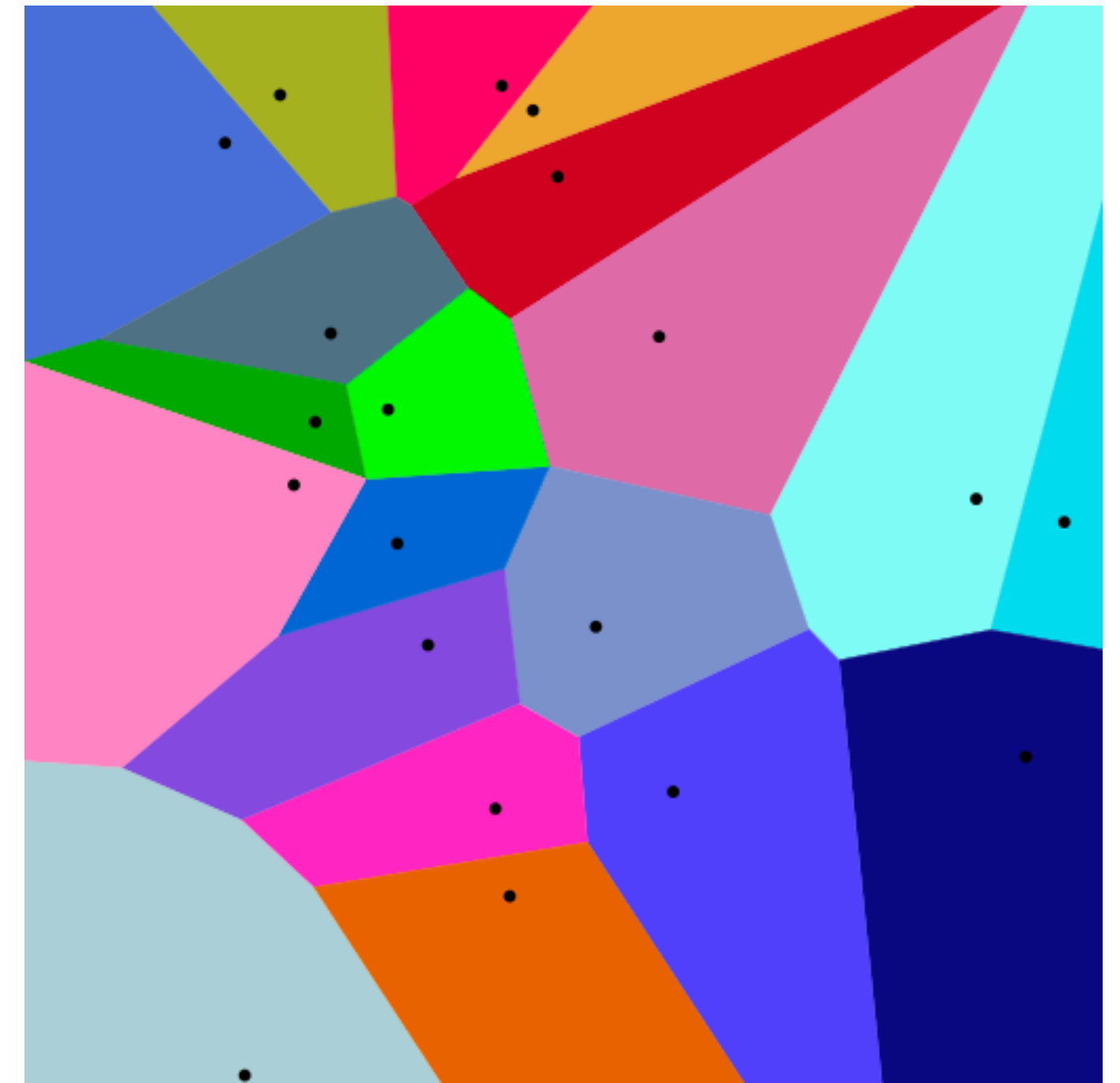
<http://www.sltrib.com/opinion/1794525-155/lake-salt-republican-county-http-utah>

Voronoi Diagrams

Given a set of locations, for which area is a location n closest?

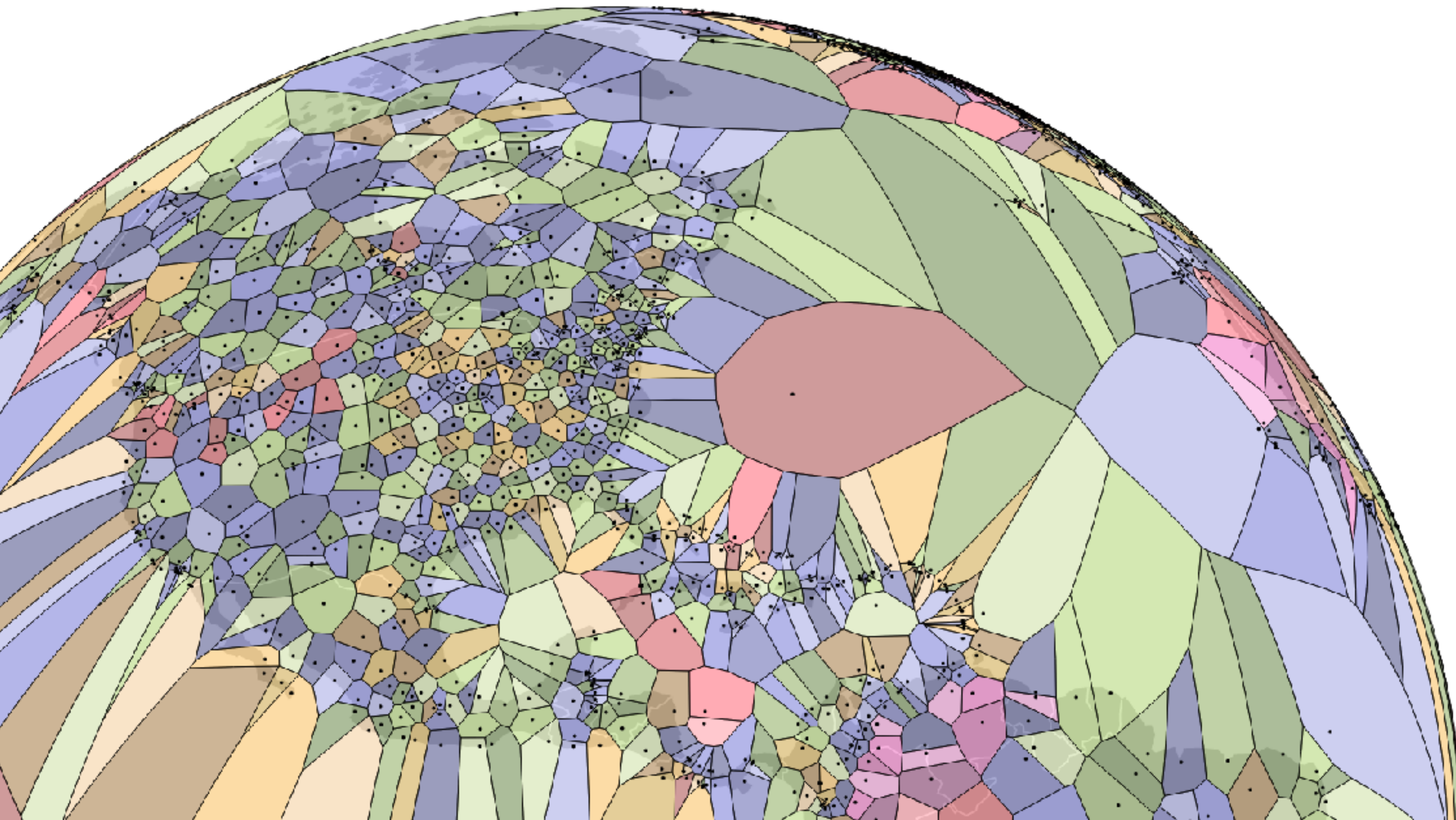
D3 Voronoi Layout:

<https://github.com/mbostock/d3/wiki/Voronoi-Geom>



Voronoi Examples

World Airports Voronoi

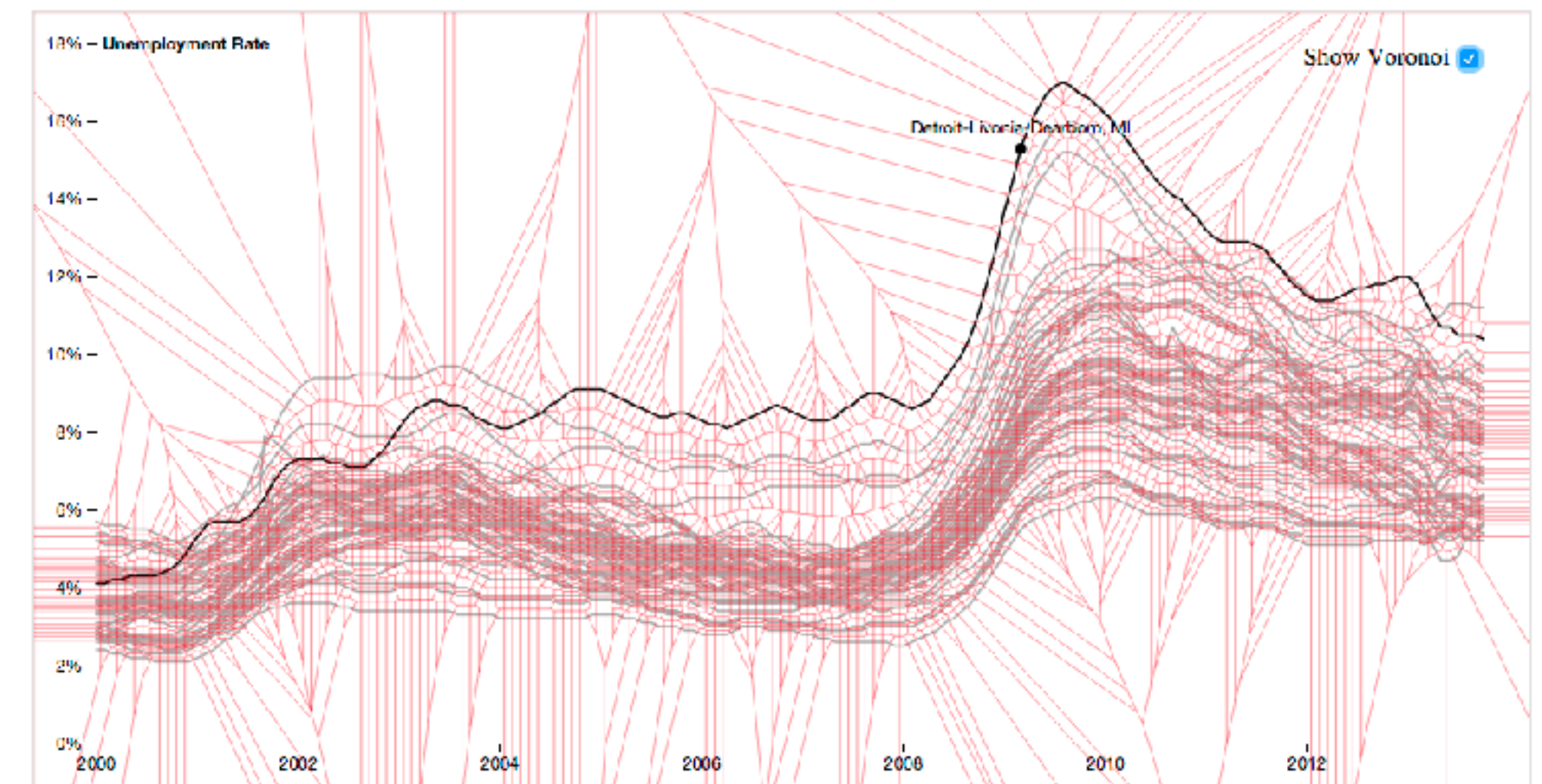
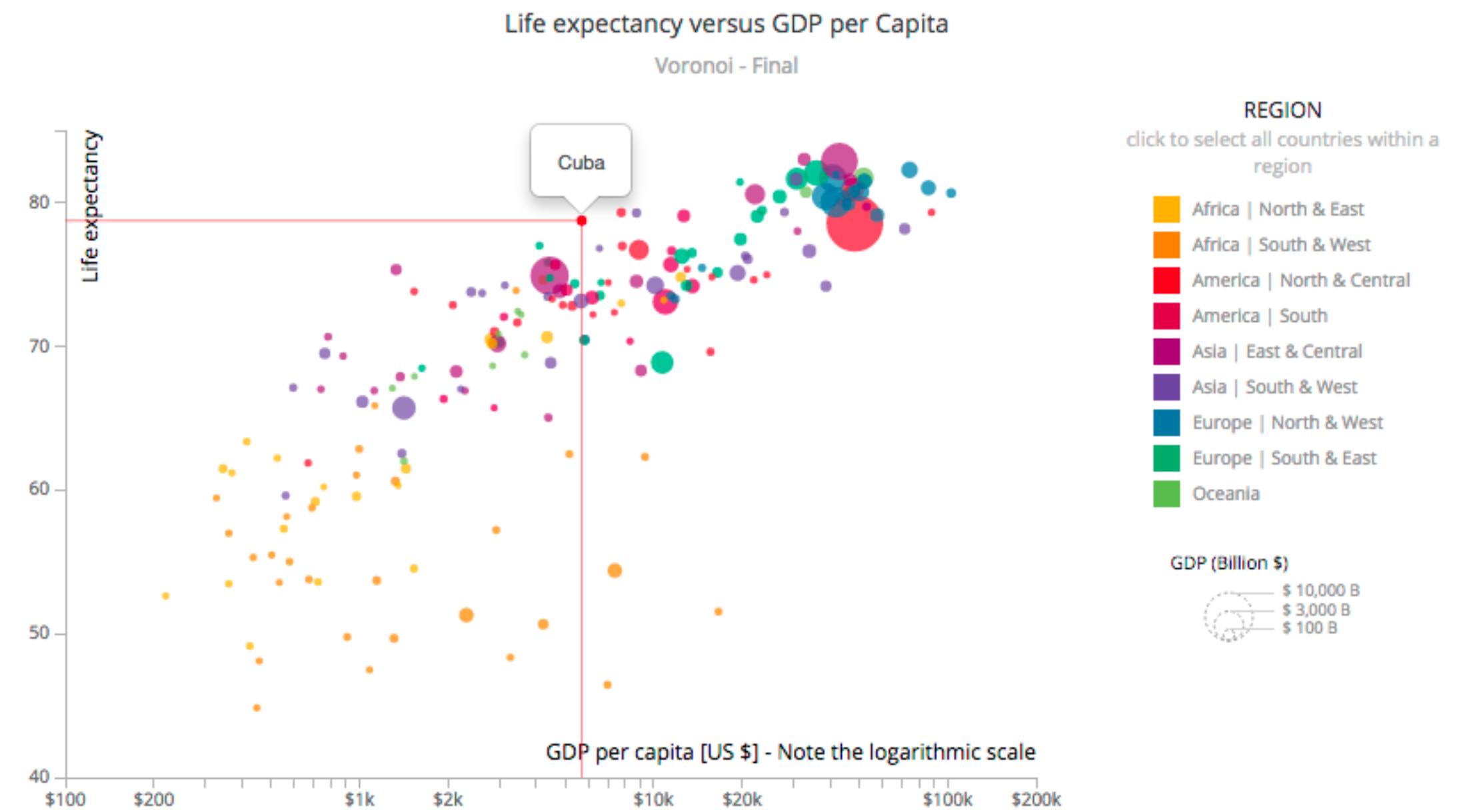


Voronoi for Interaction

Useful for interaction:
Increase size of target area to
click/hover

Instead of clicking on point,
hover in its region

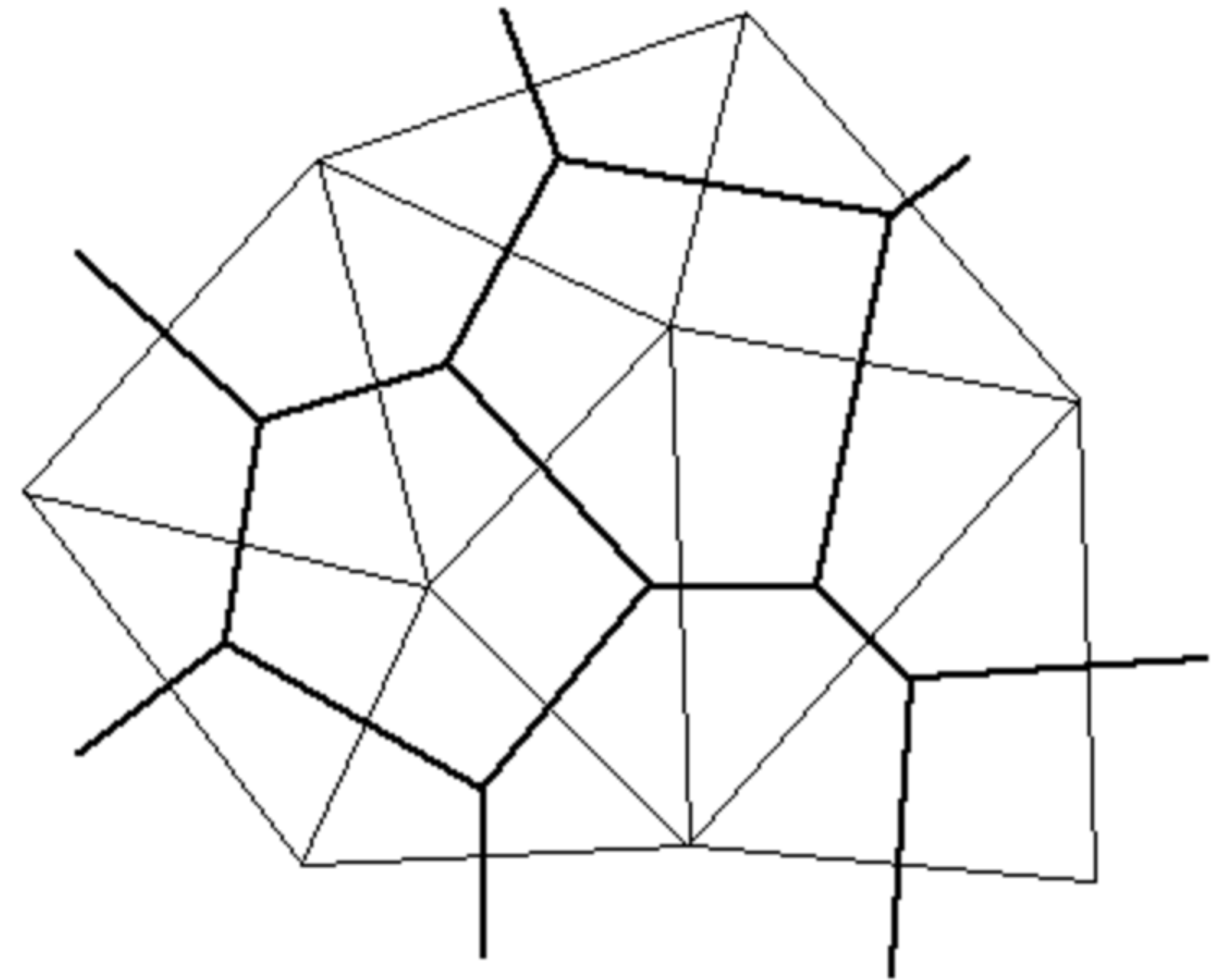
<https://github.com/d3/d3-voronoi/>



Constructing a Voronoi Diagram

Calculate a Delauney triangulation

Voronoi edges are perpendicular to triangle edges.



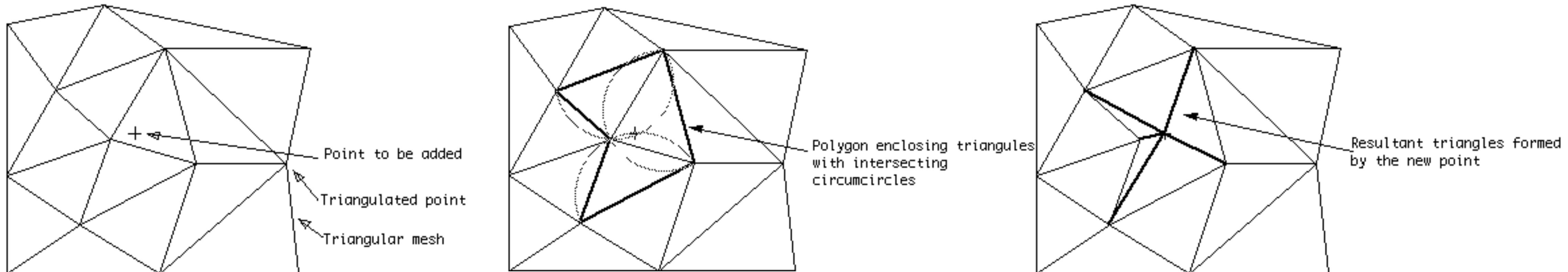
Delauney Triangulation

Start with all-encompassing fake triangle

For existing triangles: check if circumcircle contains new point

Outer edges of triangles form polygon, delete all inner edges

Create triangle connecting all outer edges to new point.



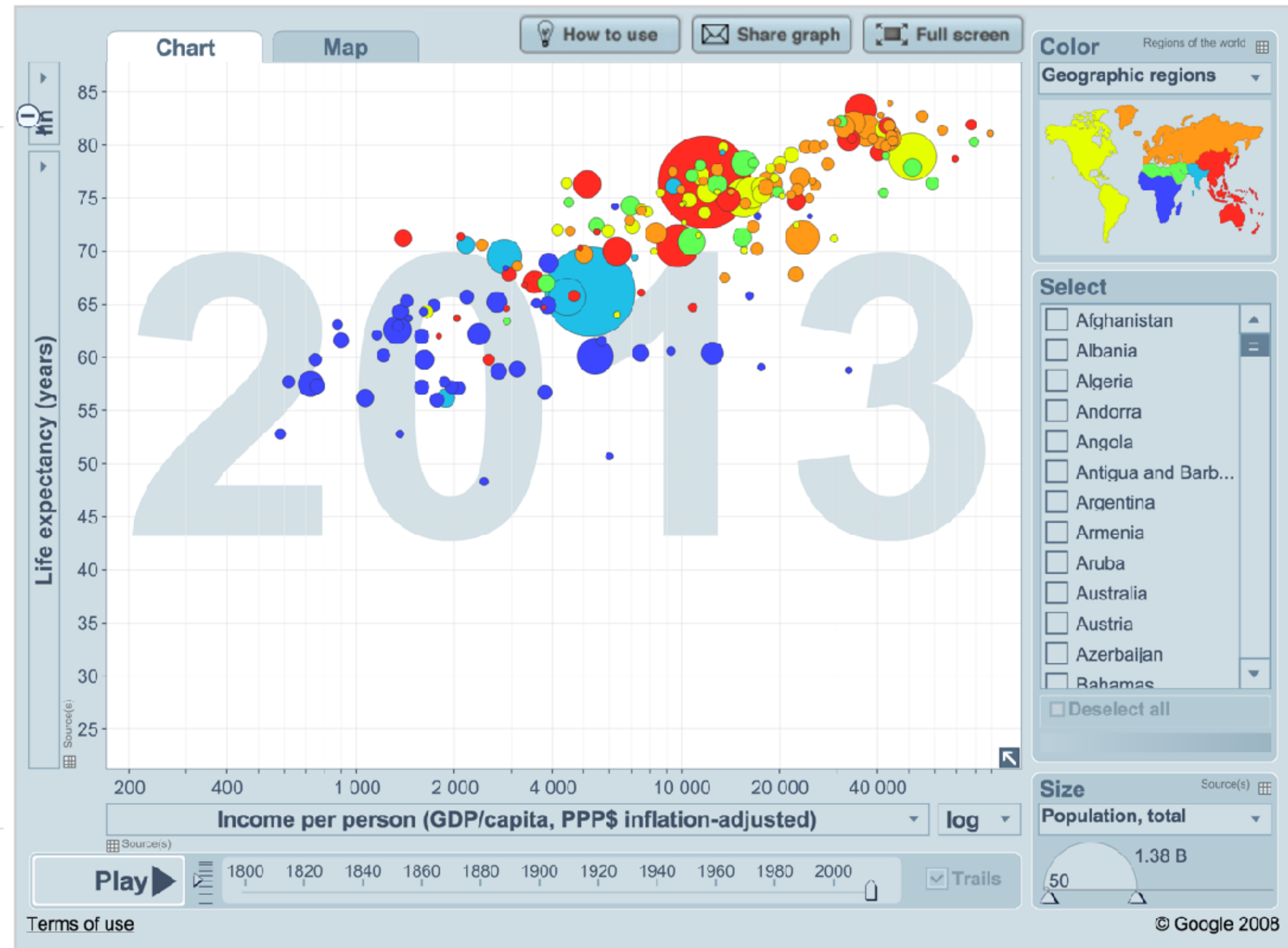
Design Critique

GapMinder

<https://goo.gl/Fcx28n>

Tool:

<http://goo.gl/jWNOUb>



Attribute aggregation

- 1) group attributes and compute a similarity score across the set
- 2) dimensionality reduction, to preserve meaningful structure

Attribute aggregation

1) group attributes and compute a similarity score across the set

2) dimensionality reduction,
to preserve meaningful structure

Clustering

Classification of items into “similar” bins

Based on similarity measures

Euclidean distance, Pearson correlation, ...

Partitional Algorithms

divide data into set of bins

bins either manually set (e.g., k-means) or automatically determined (e.g., affinity propagation)

Hierarchical Algorithms

Produce “similarity tree” – dendrogram

Bi-Clustering

Clusters dimensions & records

Fuzzy clustering

allows occurrence of elements in multiples clusters

Clustering Applications

Clusters can be used to

- order (pixel based techniques)

- brush (geometric techniques)

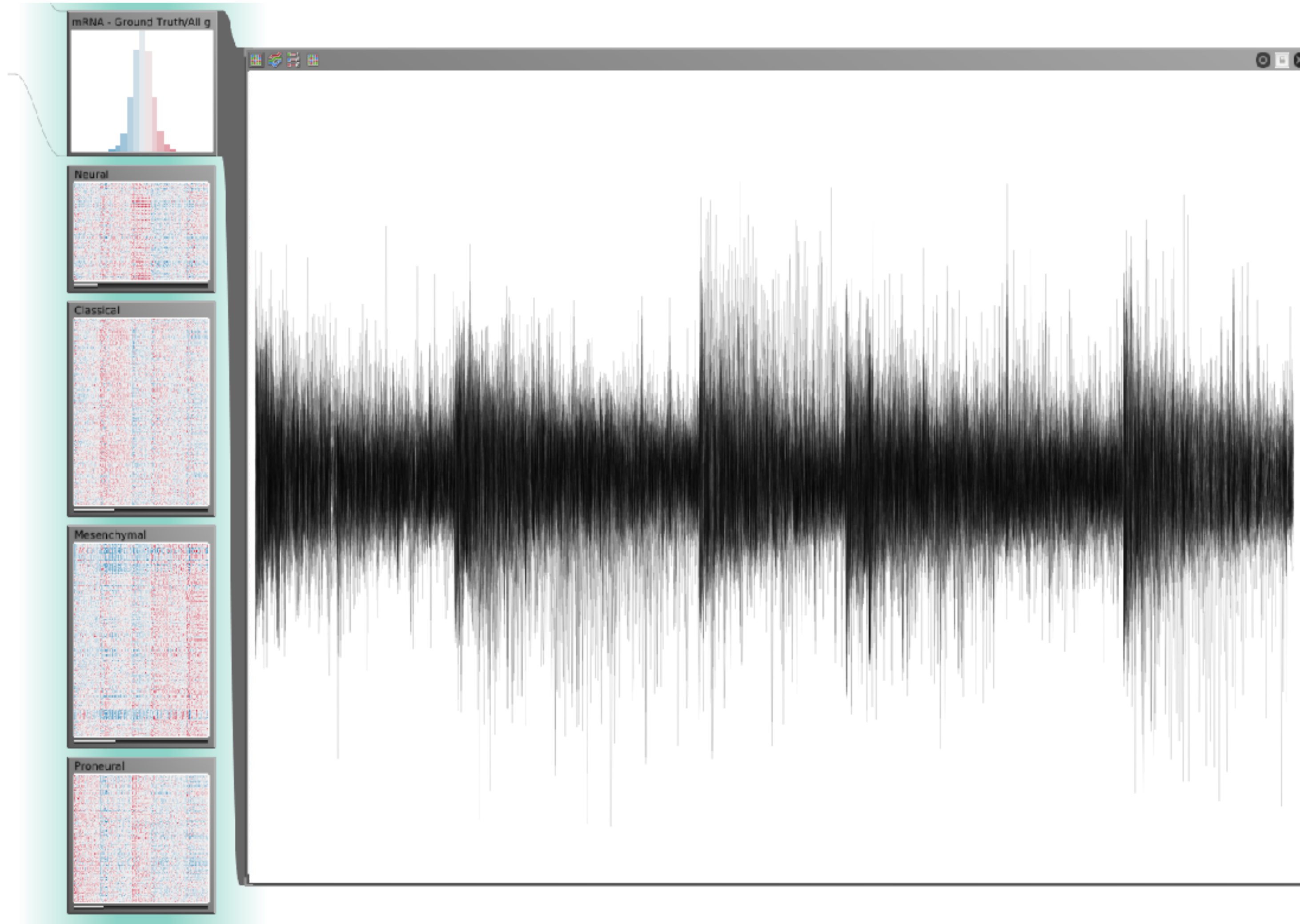
- aggregate

Aggregation

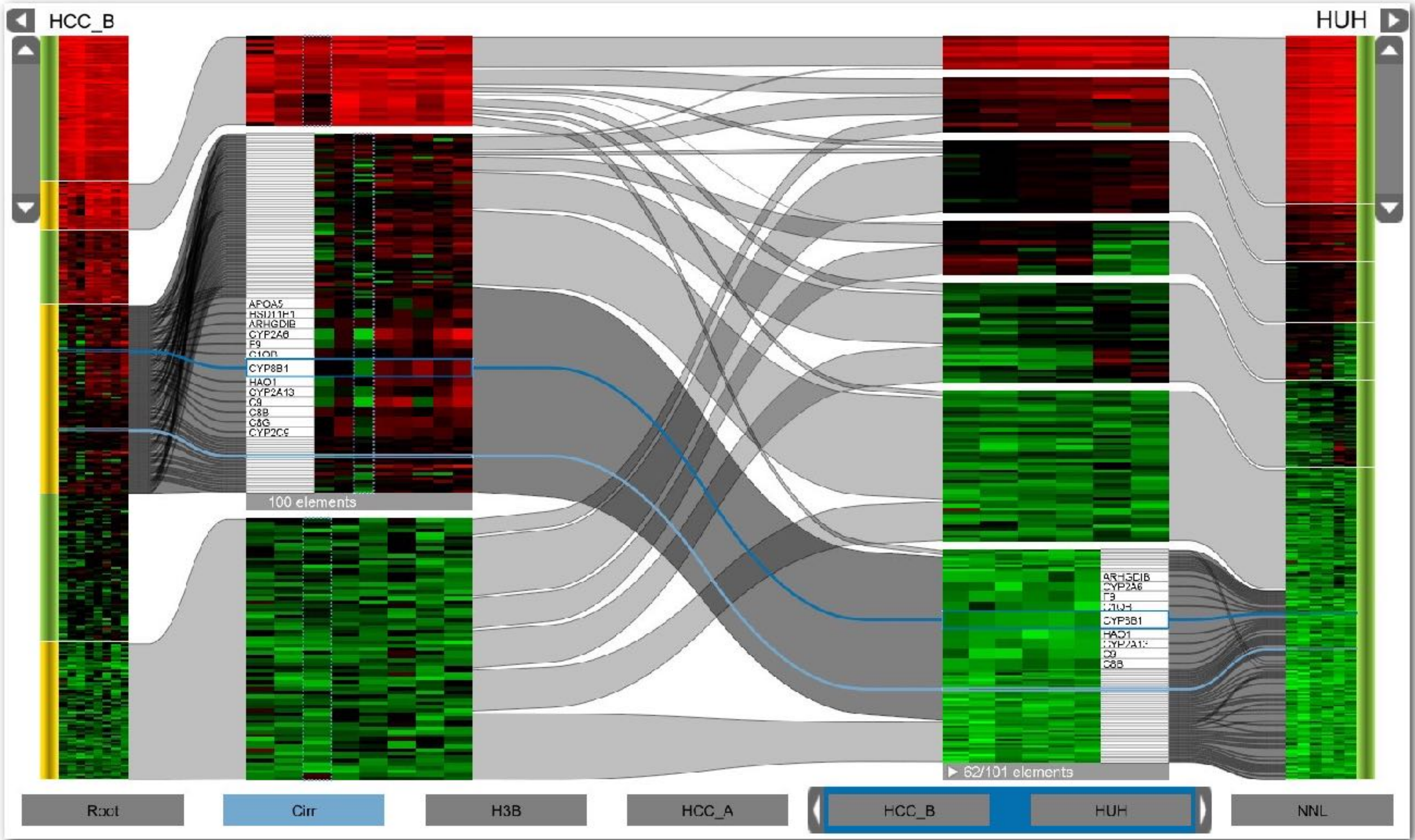
- cluster more homogeneous than whole dataset

- statistical measures, distributions, etc. more meaningful

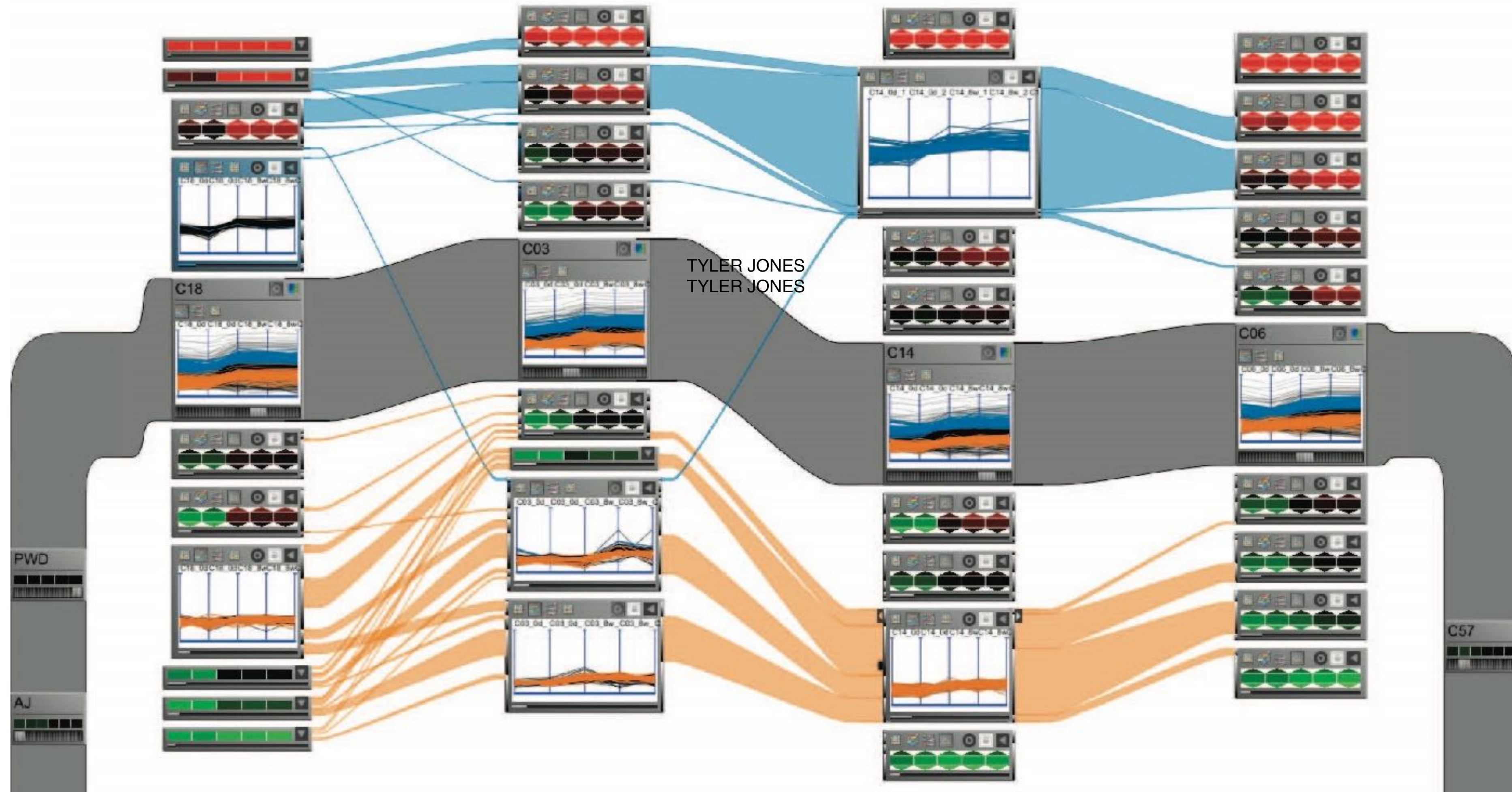
Clustered Heat Map



Cluster Comparison



Aggregation



Example: K-Means

Goal: Minimize aggregate intra-cluster distance (*inertia*)

$$\underset{C}{\operatorname{argmin}} \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

total squared distance from point to center of its cluster

for euclidian distance: this is the variance

measure of how internally coherent clusters are

Lloyd's Algorithm

Input: set of records $x_1 \dots x_n$, and k (nr clusters)

Pick k starting points as centroids $c_1 \dots c_k$

While not converged:

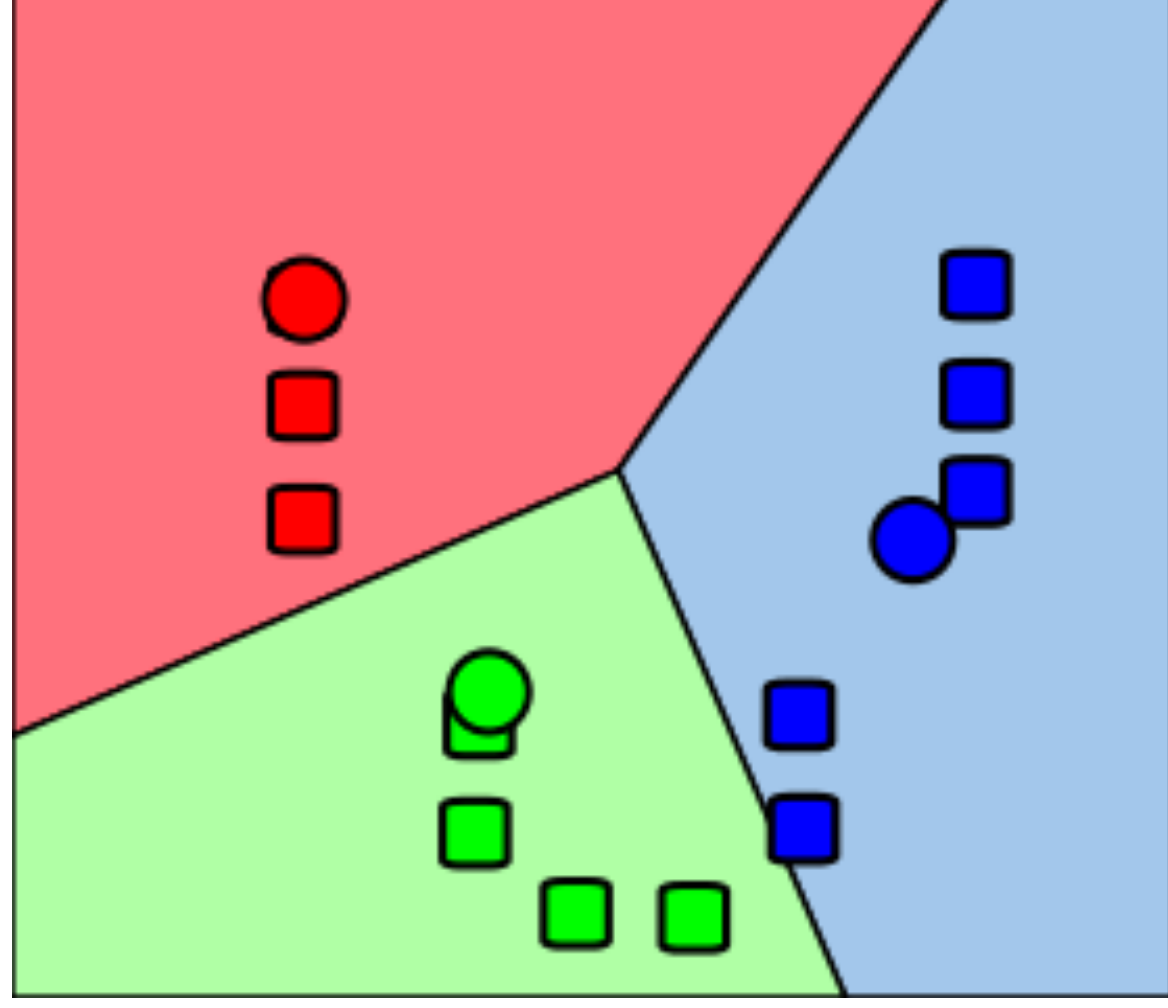
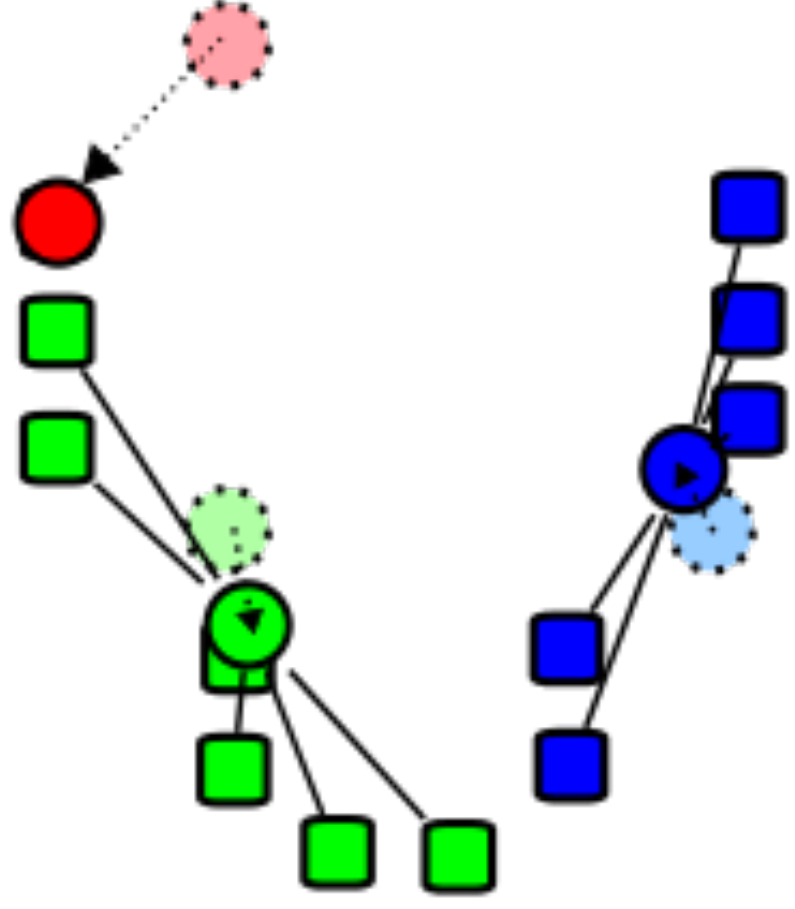
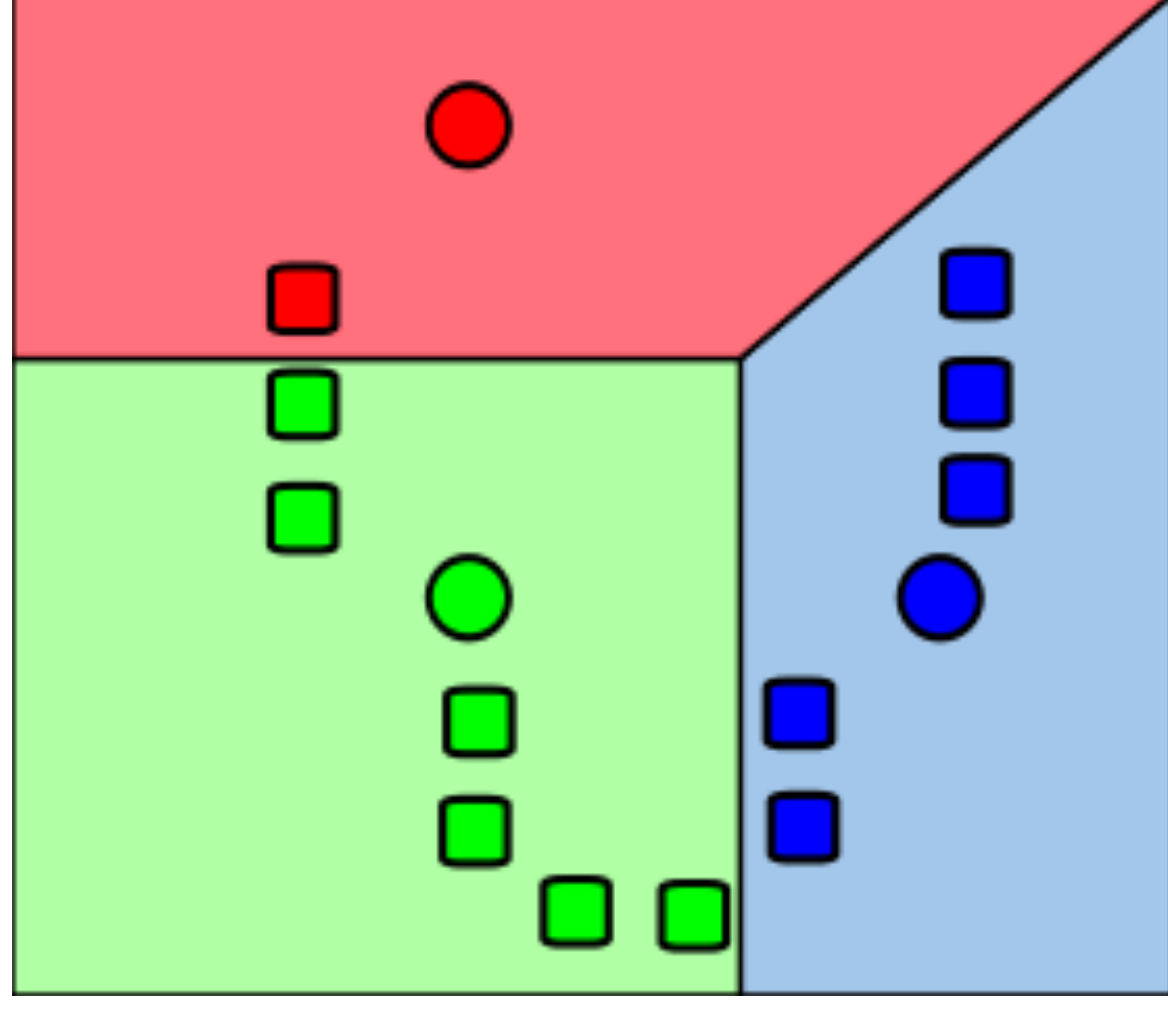
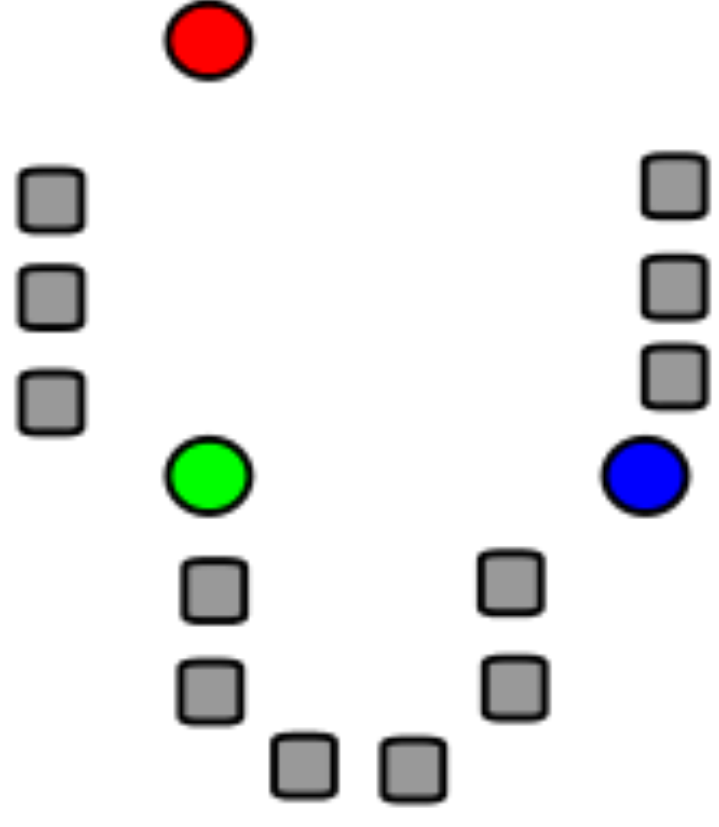
1. for each point x_i find closest centroid c_j
 - for every c_j calculate distance $D(x_i, c_j)$
 - assign x_i to cluster j defined by smallest distance
2. for each cluster j , compute a new centroid c_j by calculating the average of all x_i assigned to cluster j

Repeat until convergence, e.g.,

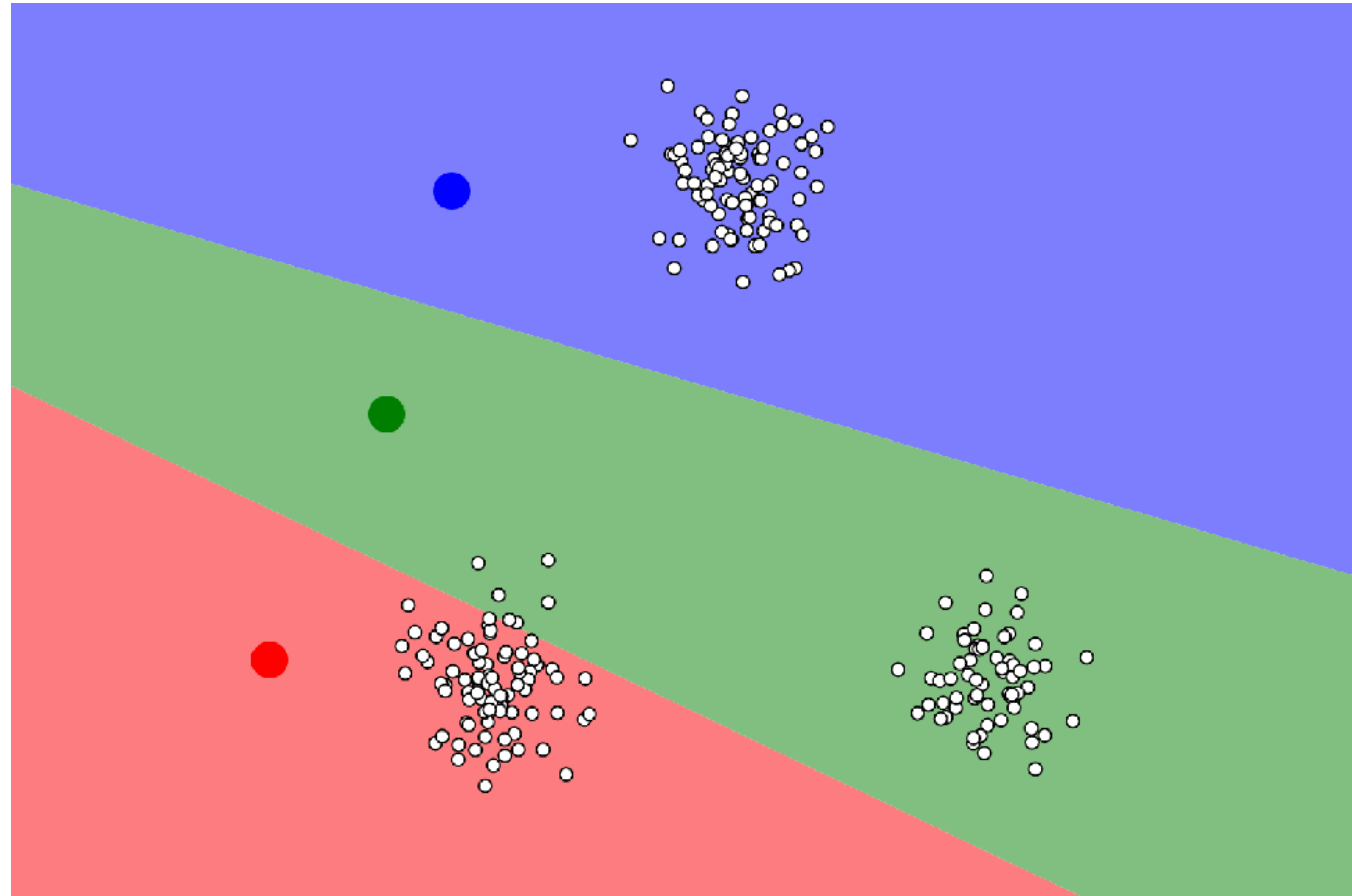
no point has changed cluster

distance between old and new centroid below threshold

number of max iterations reached

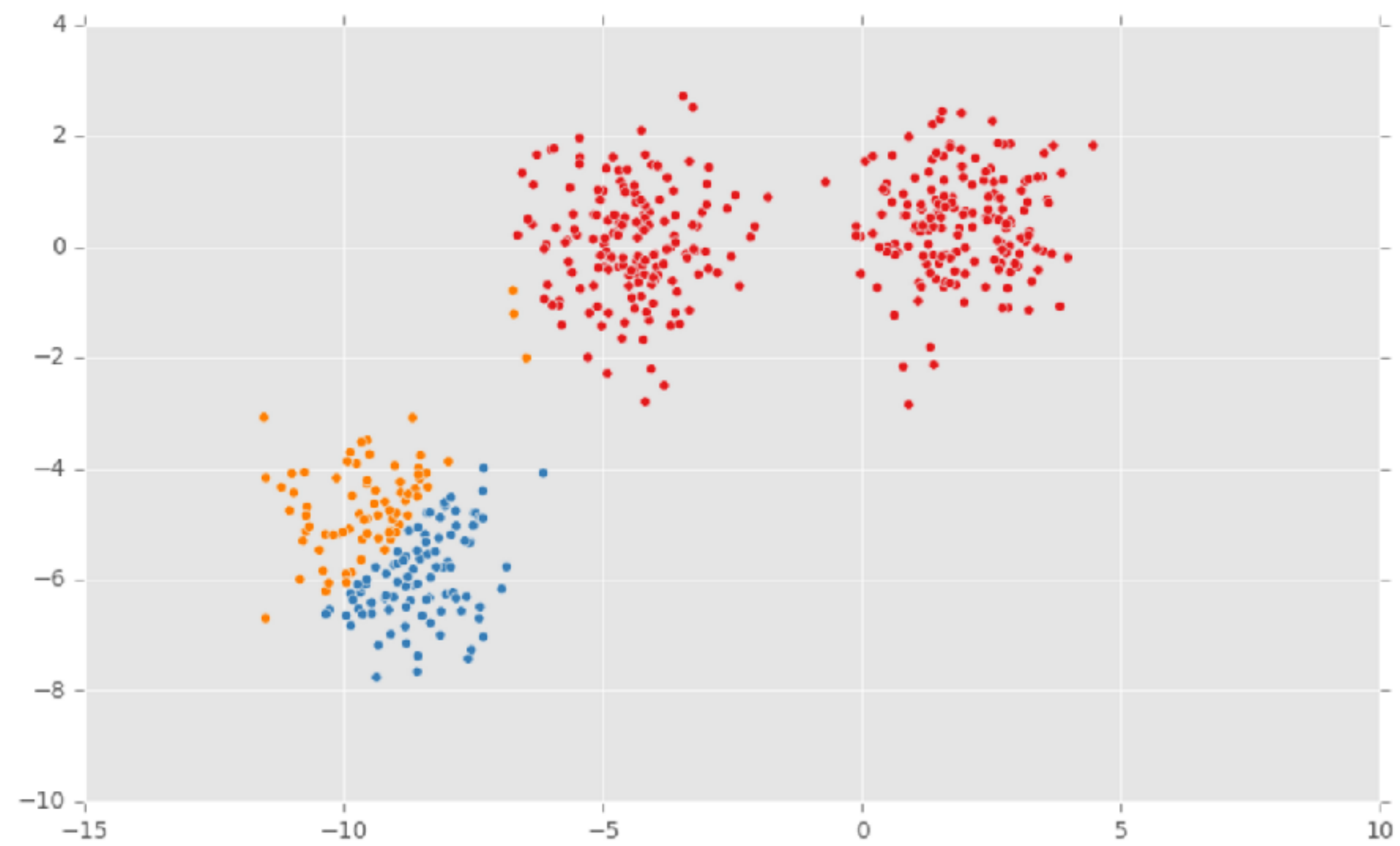
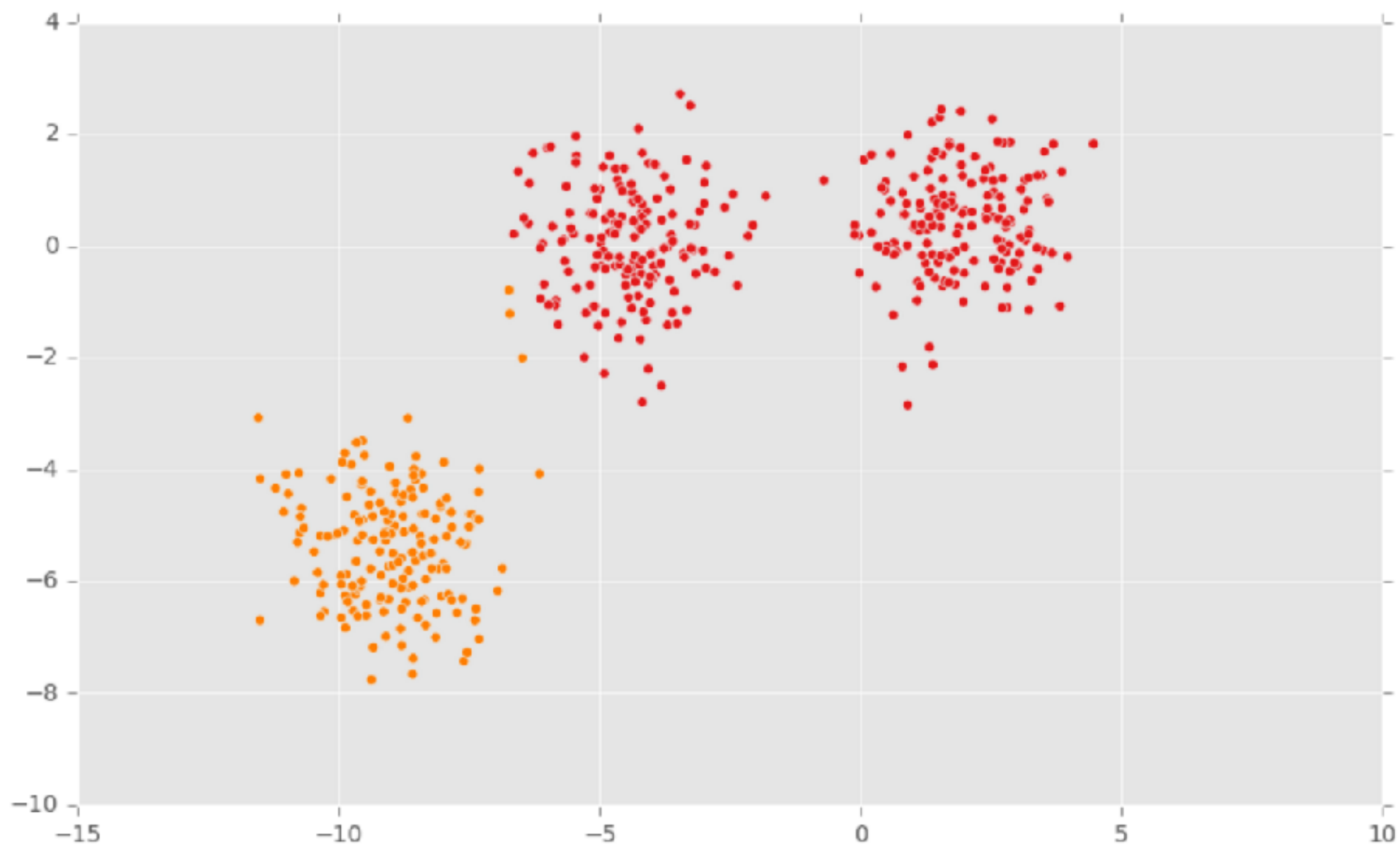


Illustrated



<https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

Choosing K



Properties

Lloyds algorithm doesn't find a global optimum

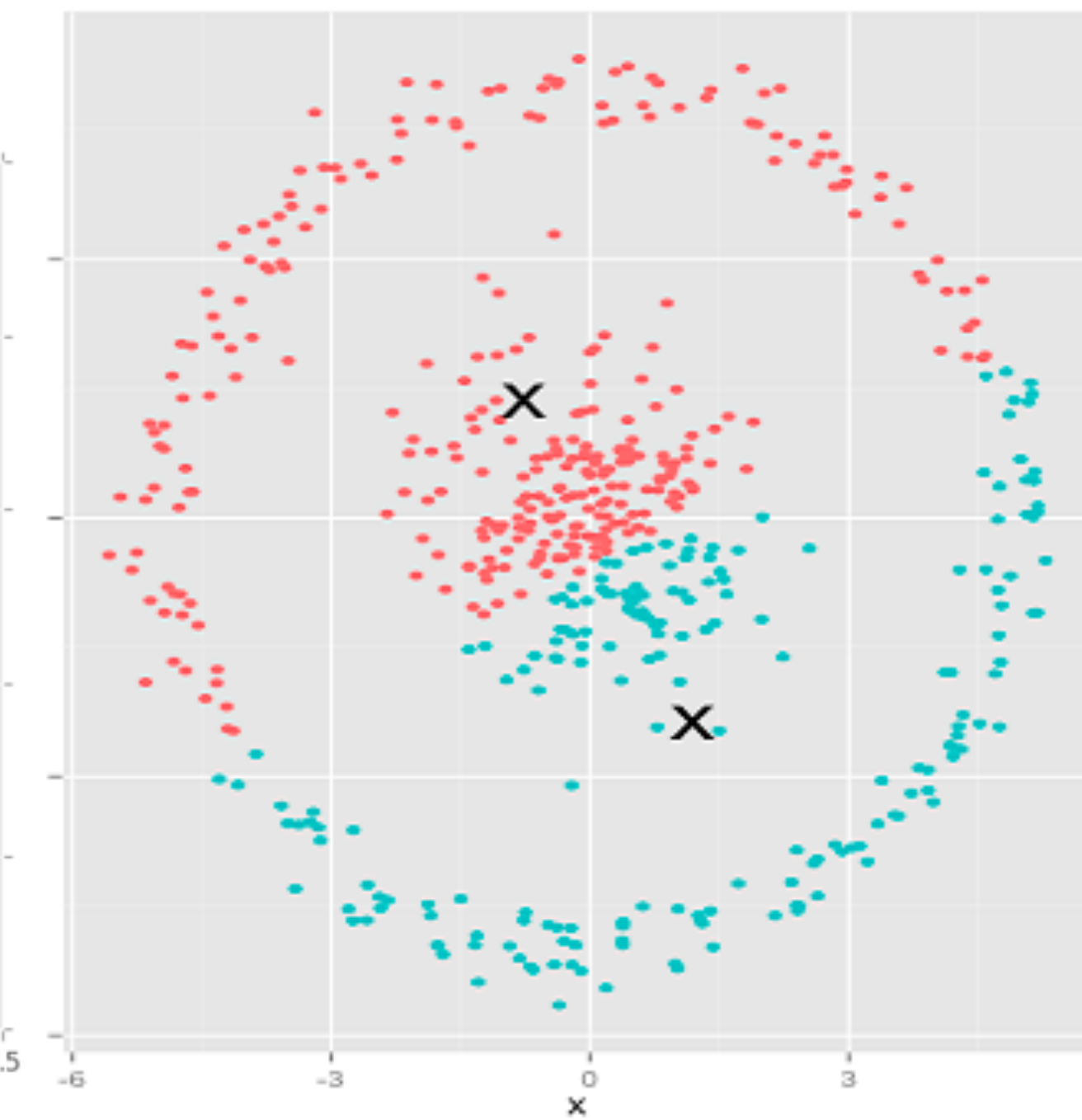
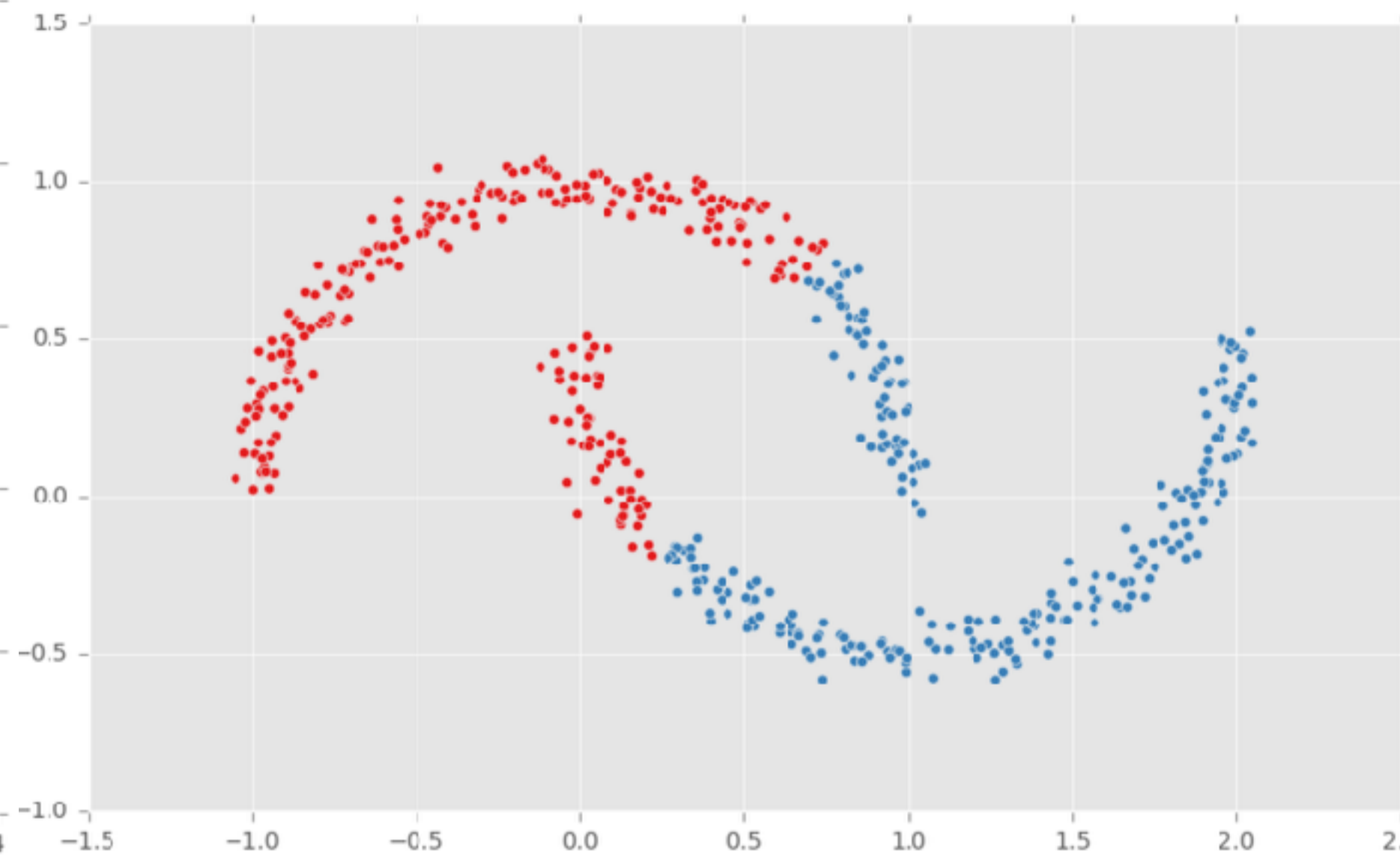
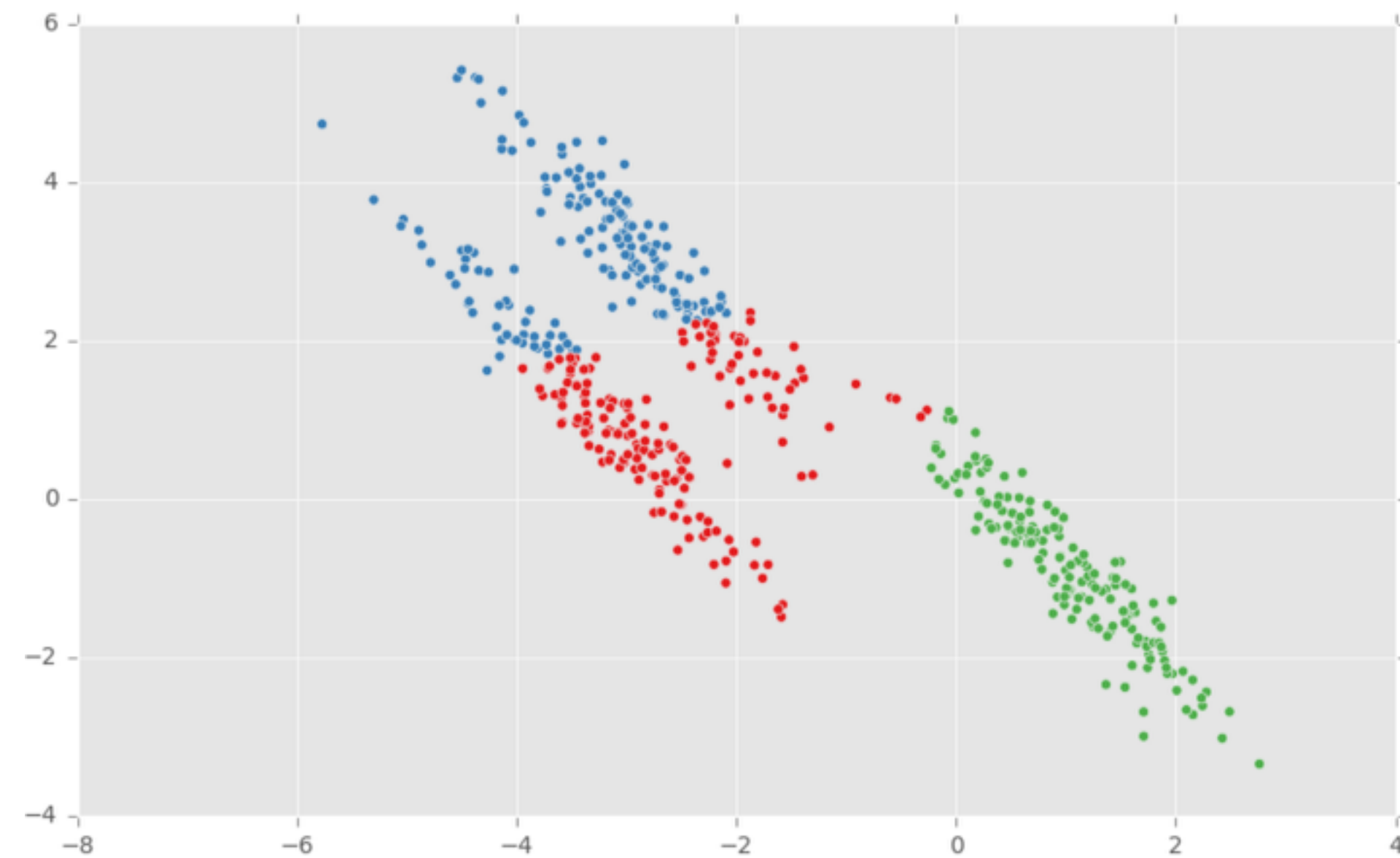
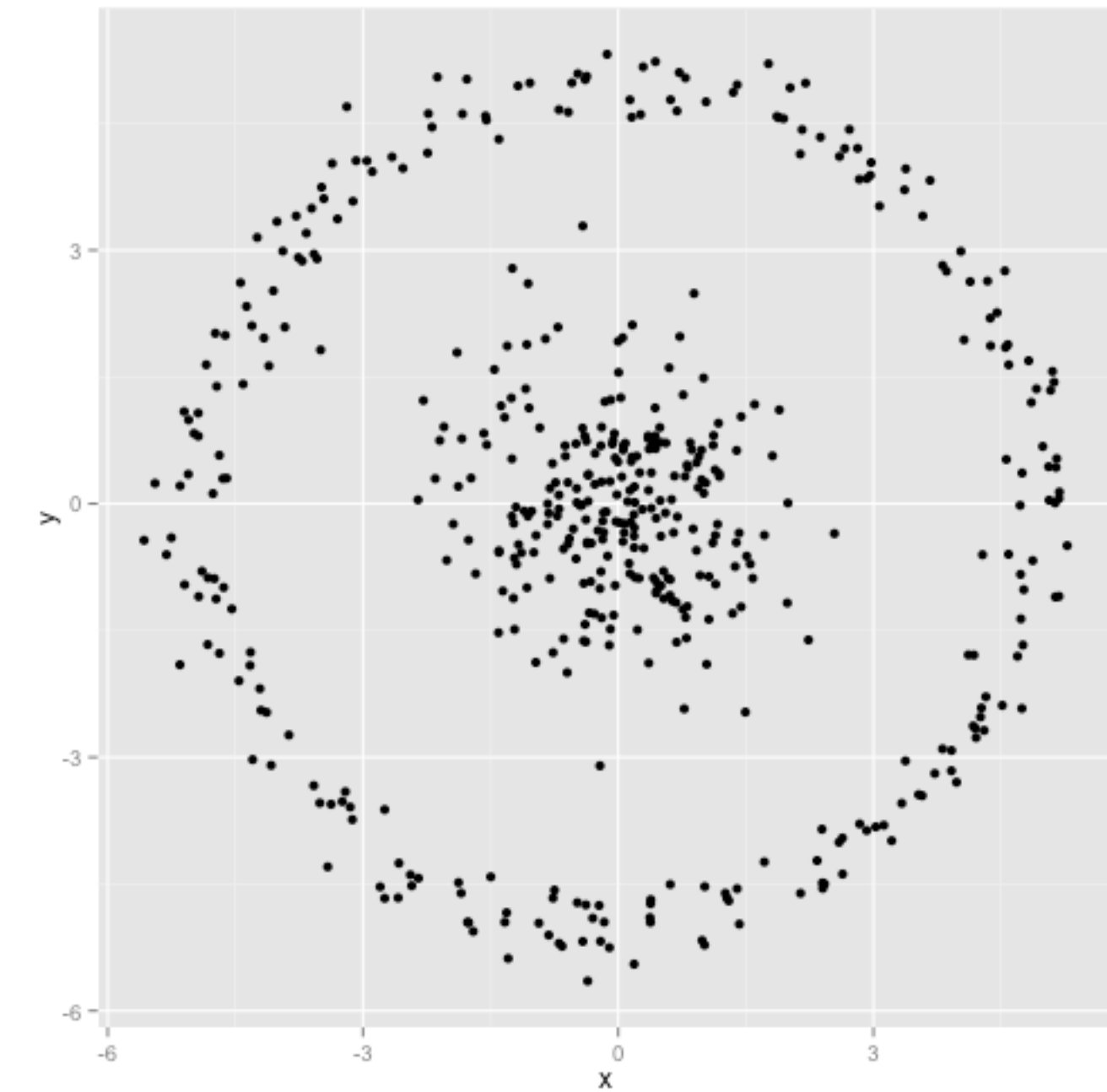
Instead it finds a local optimum

It is very fast:

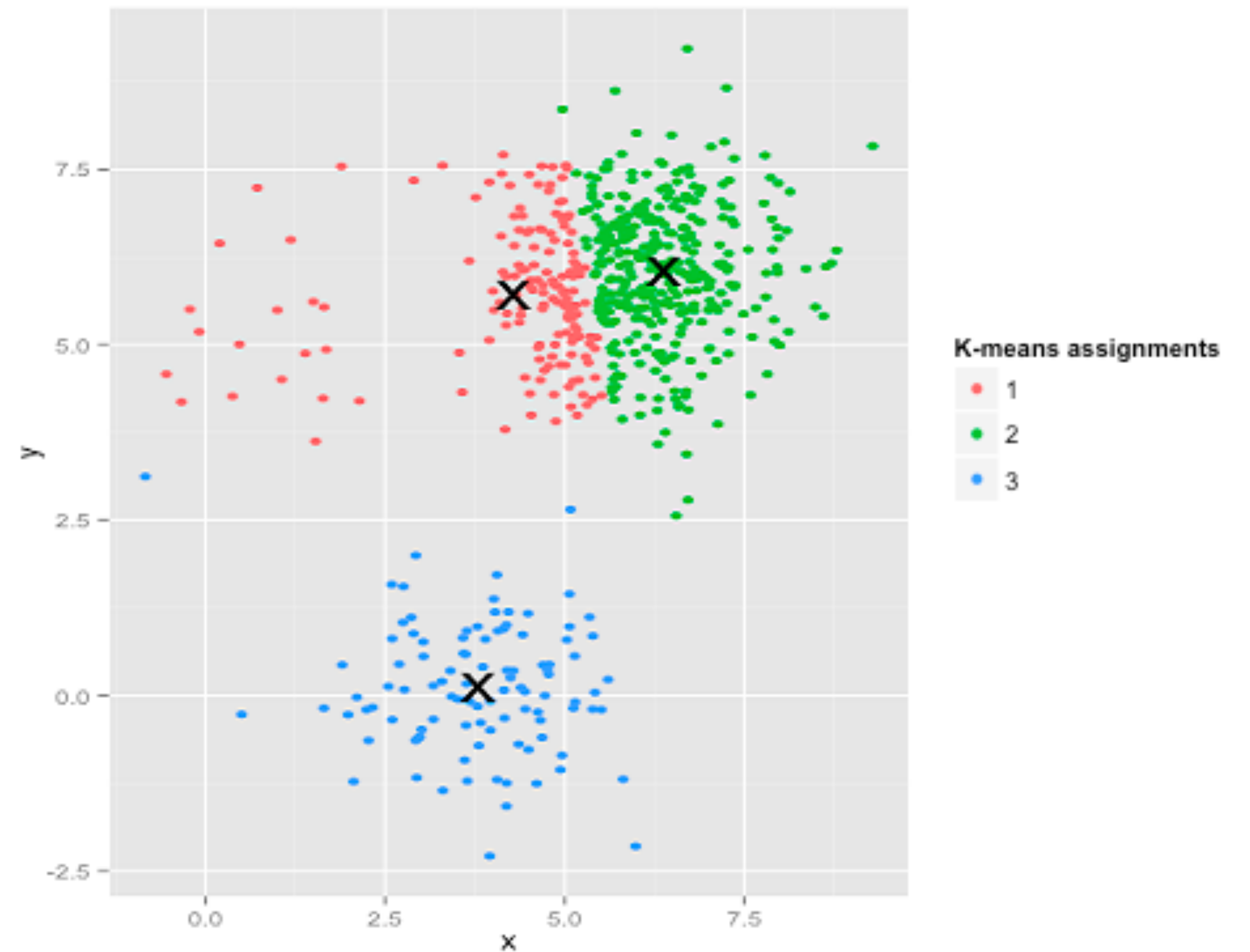
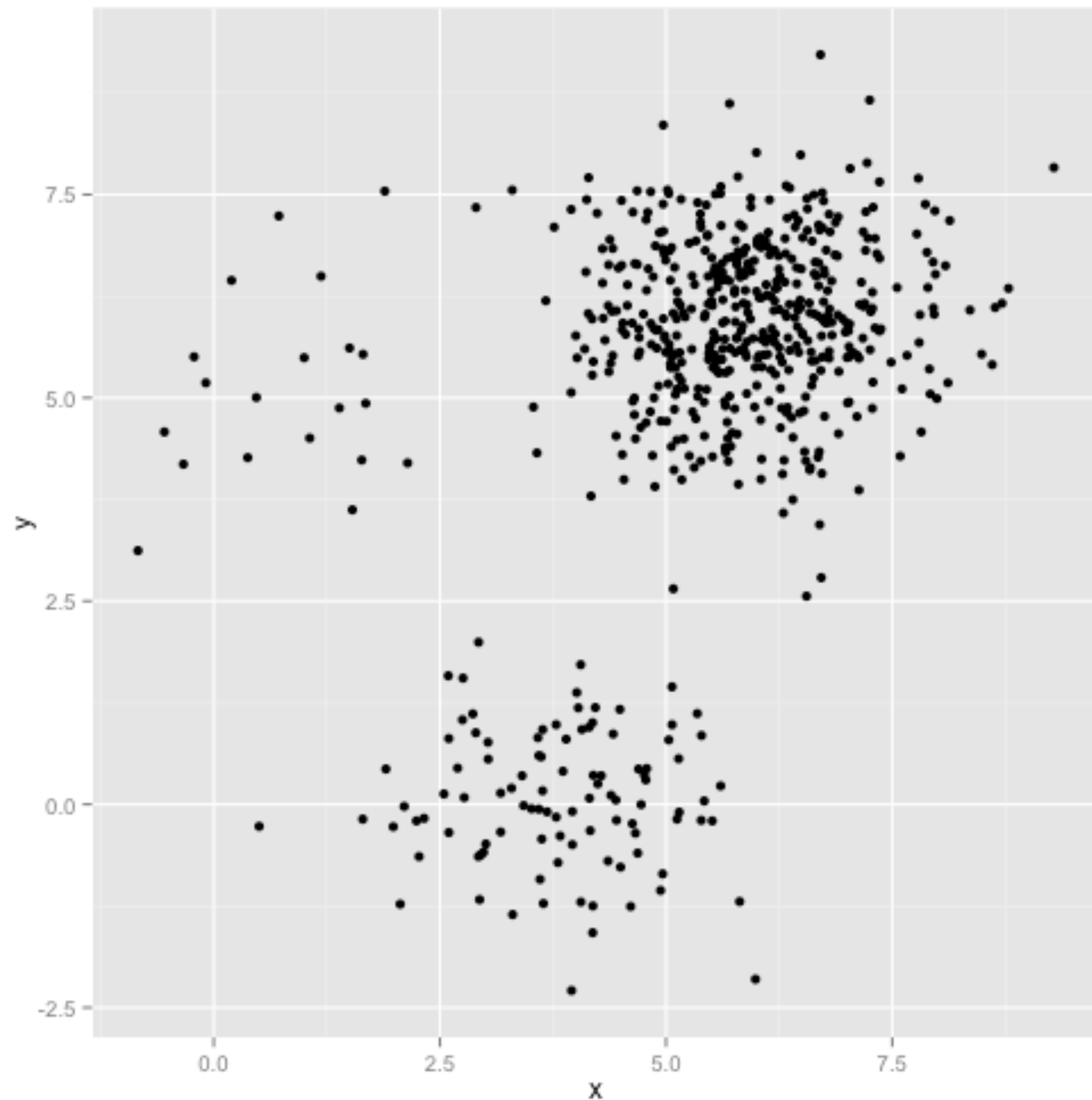
common to run multiple times and pick the solution with the minimum inertia

K-Means Properties

Assumptions about data:
roughly “circular” clusters of
equal size



K-Means Unequal Cluster Size



Hierarchical Clustering

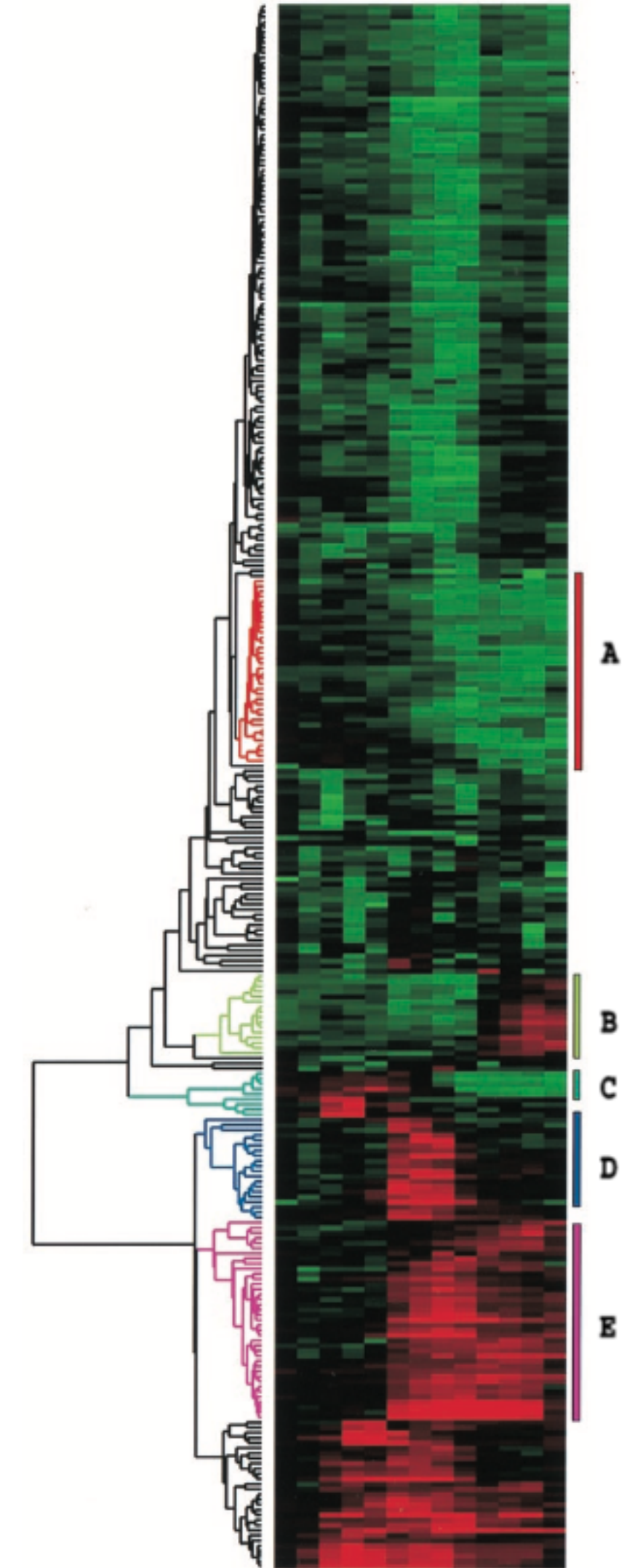
Two types:

agglomerative clustering

start with each node as a cluster and merge

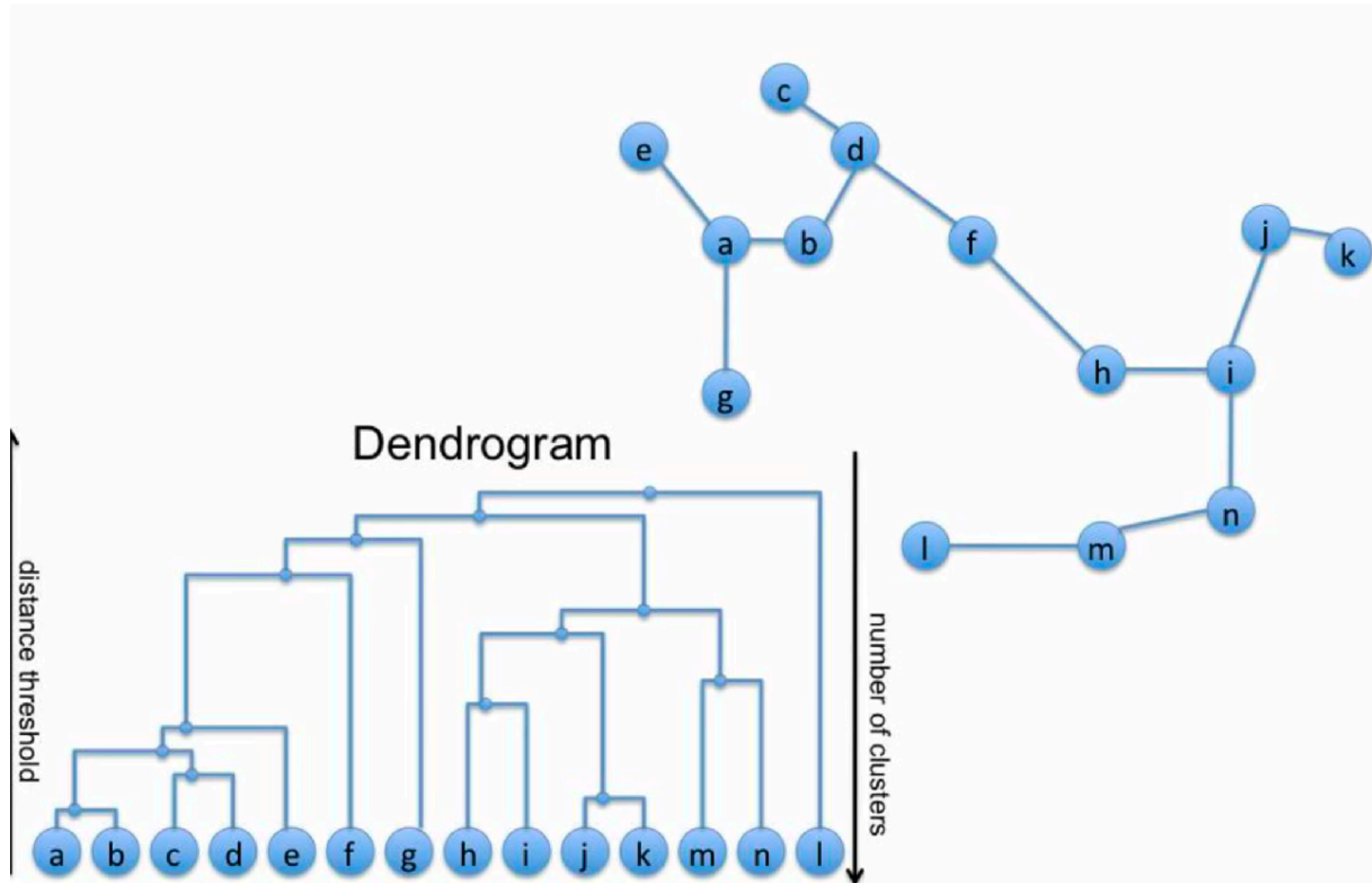
divisive clustering

start with one cluster, and split



Agglomerative Clustering Idea

Agglomerative Clustering Idea



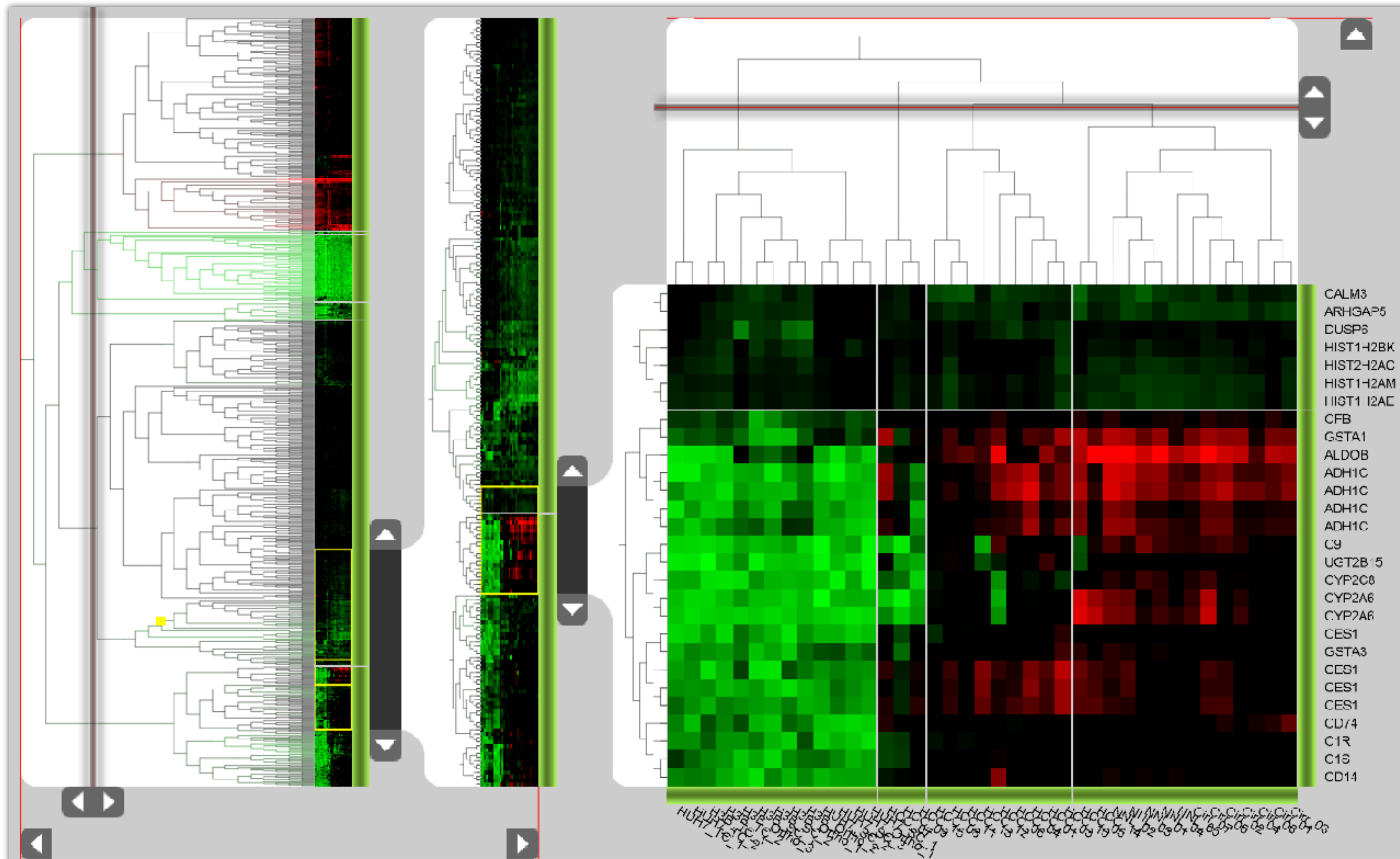
Linkage Criteria

How do you define similarity between two clusters to be merged (A and B)?

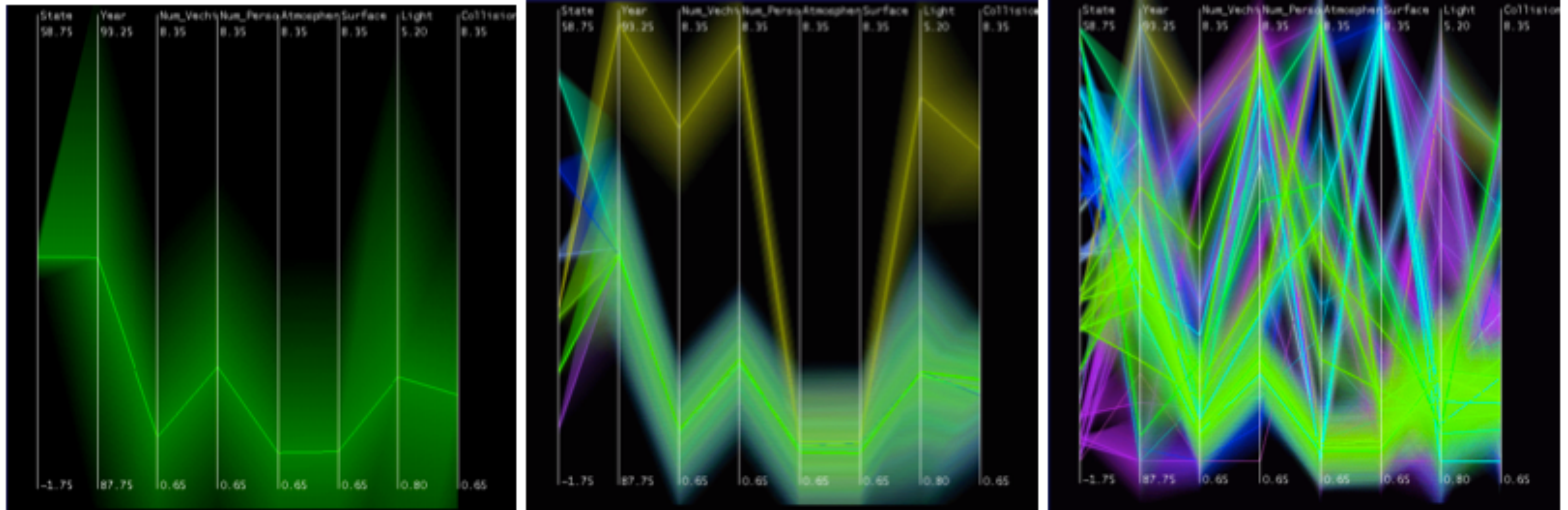
- use maximum linkage distance
- use minimum linkage distance
- use average linkage distance
- use centroid distance

Names	Formula
Maximum or complete-linkage clustering	$\max \{ d(a, b) : a \in A, b \in B \}.$
Minimum or single-linkage clustering	$\min \{ d(a, b) : a \in A, b \in B \}.$
Mean or average linkage clustering, or UPGMA	$\frac{1}{ A B } \sum_{a \in A} \sum_{b \in B} d(a, b).$
Centroid linkage clustering, or UPGMC	$\ c_s - c_t\ $ where c_s and c_t are the centroids of clusters s and t , respectively.

F+C Approach, with Dendrograms



Hierarchical Parallel Coordinates



Attribute aggregation

- 1) group attributes and compute a similarity score across the set
- 2) dimensionality reduction,
to preserve meaningful structure**

Dimensionality Reduction

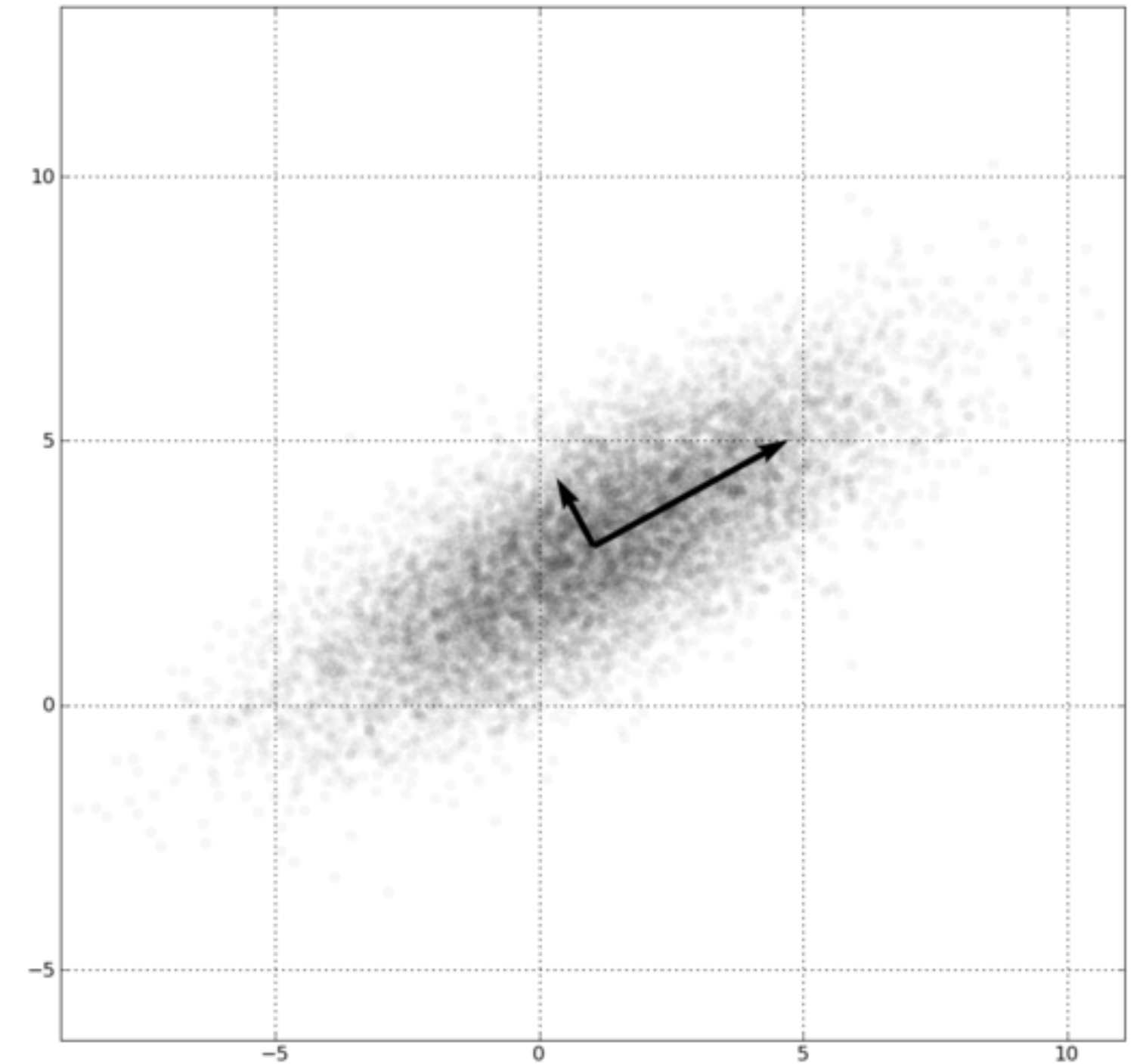
Reduce high dimensional to lower dimensional space

Preserve as much of variation as possible

Plot lower dimensional space

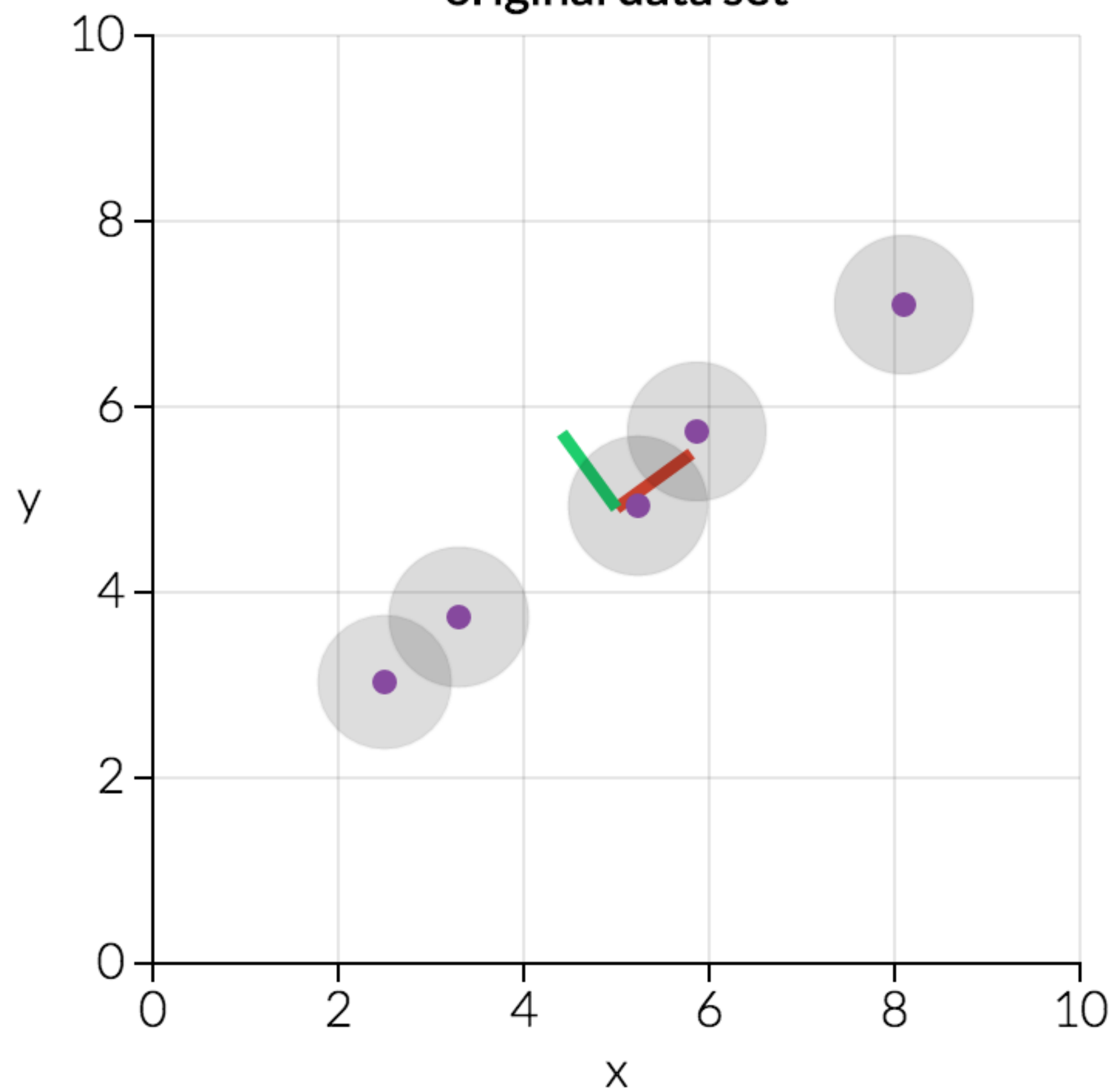
Principal Component Analysis (PCA)

linear mapping, by order of variance

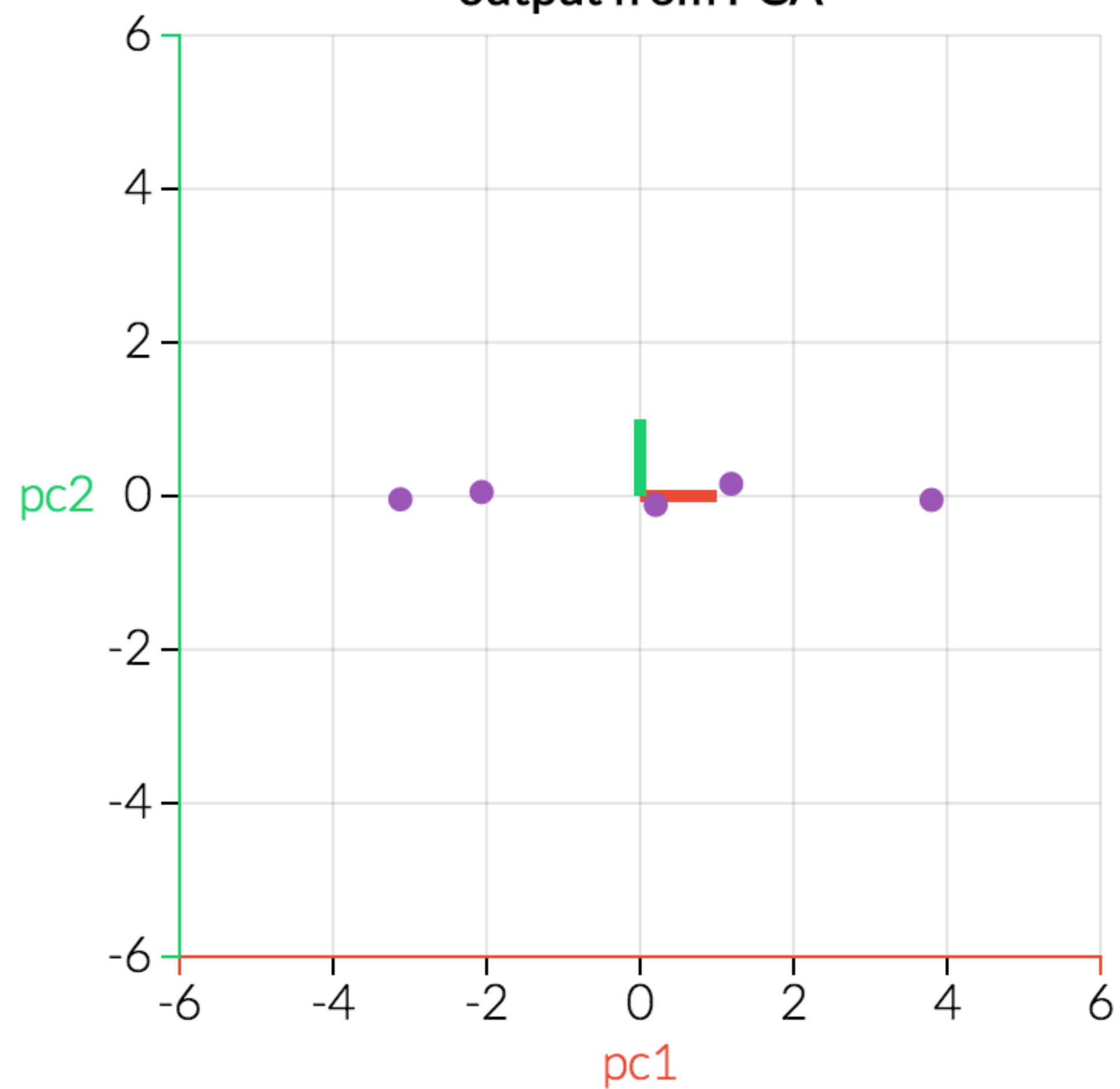


PCA

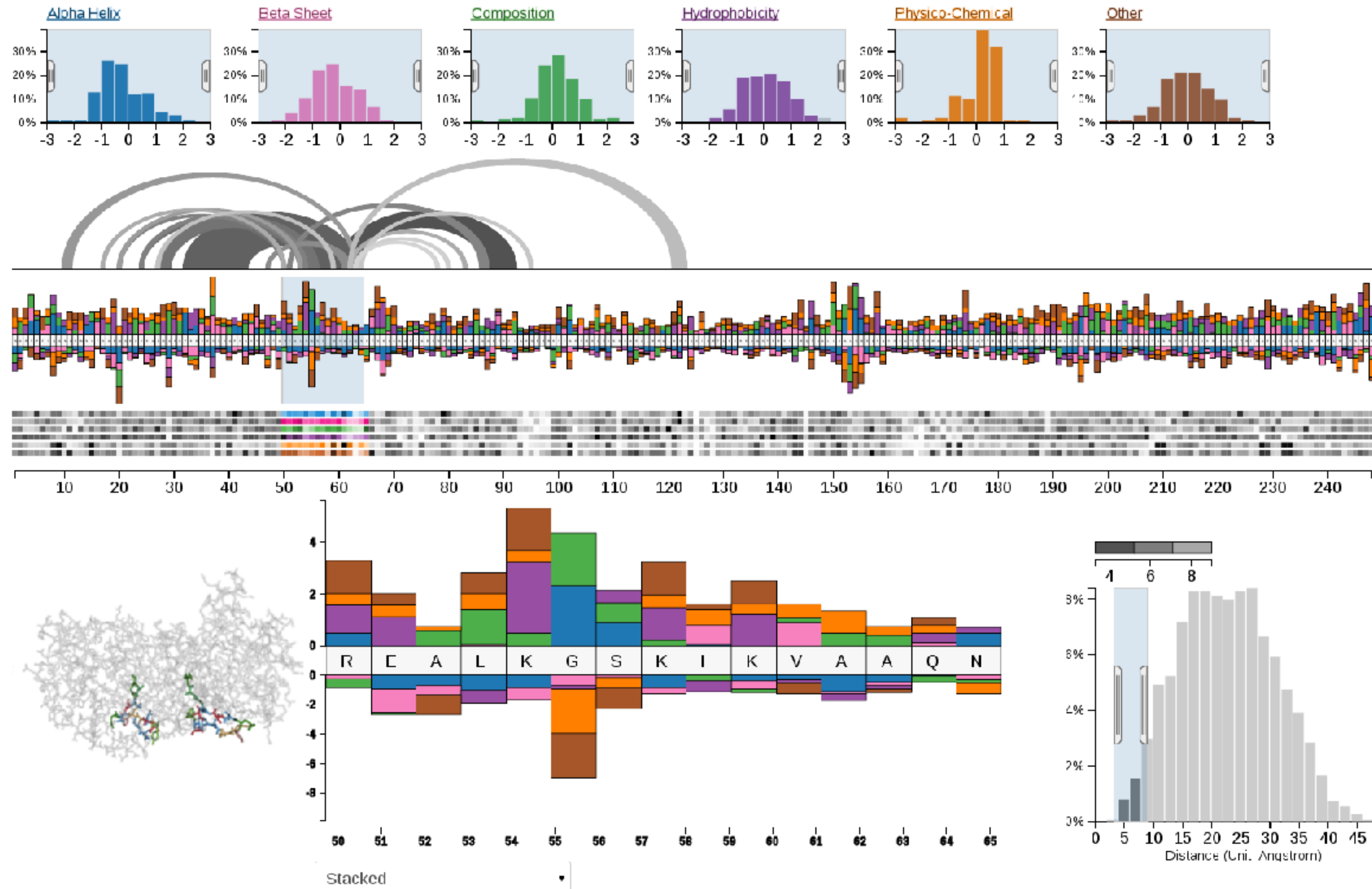
original data set



output from PCA



PCA Example – Class Project 2013



Multidimensional Scaling

Nonlinear, better suited for some DS

Multiple approaches

Works based on projecting a similarity matrix

How do you compute similarity?

How do you project the points?

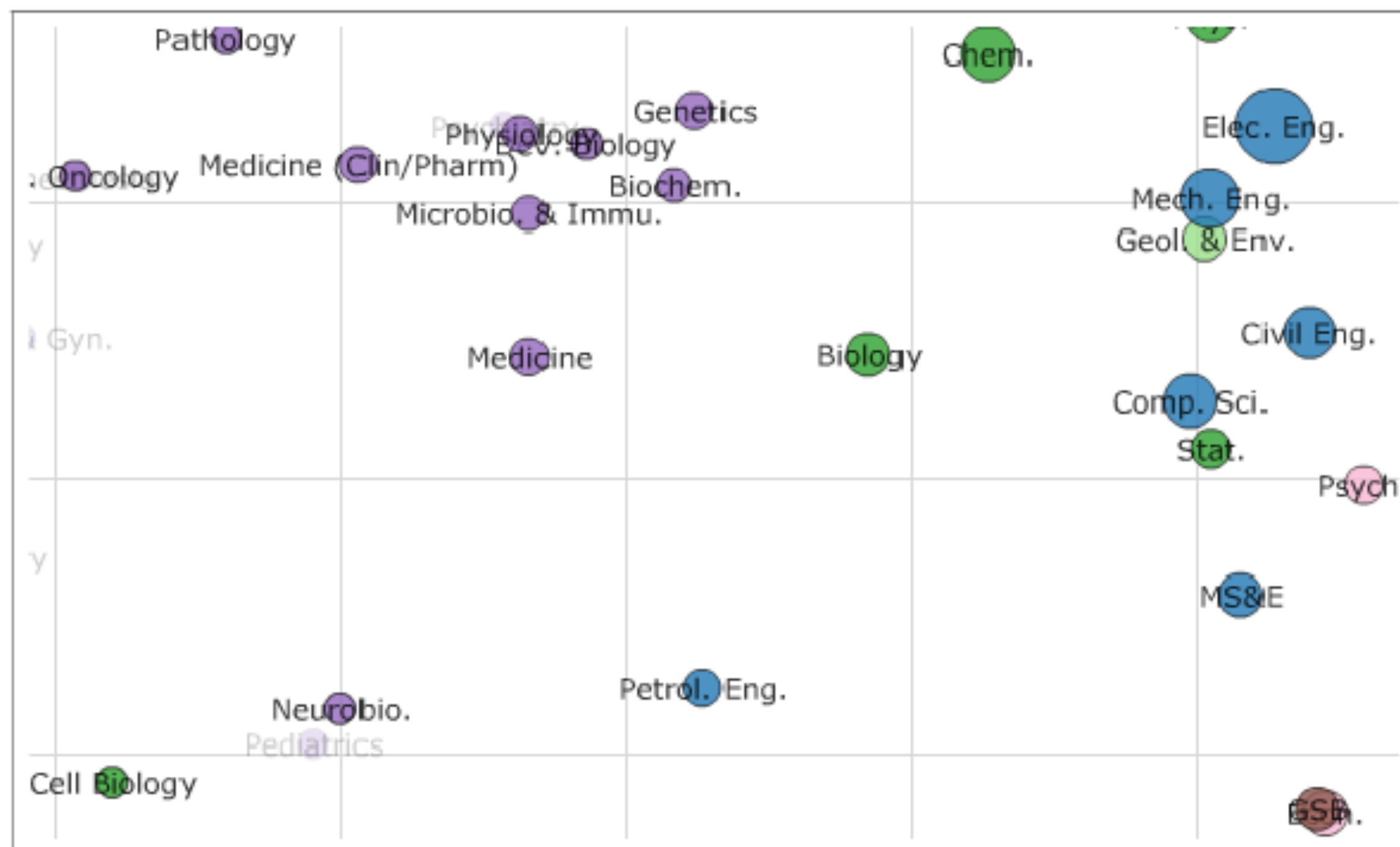
Popular for text analysis



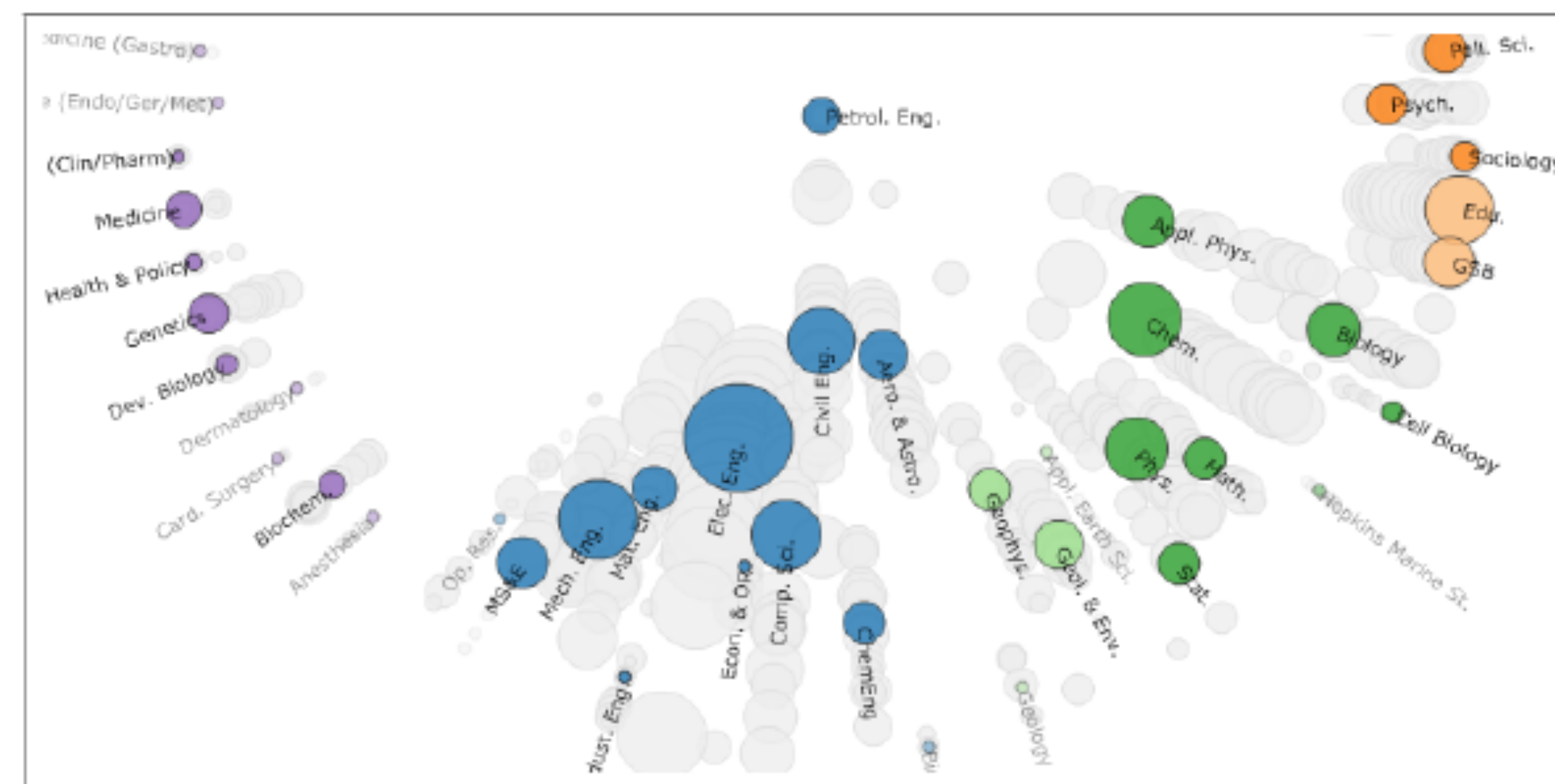
[Doerk 2011]

Can we Trust Dimensionality Reduction?

Topical distances between departments in a 2D projection

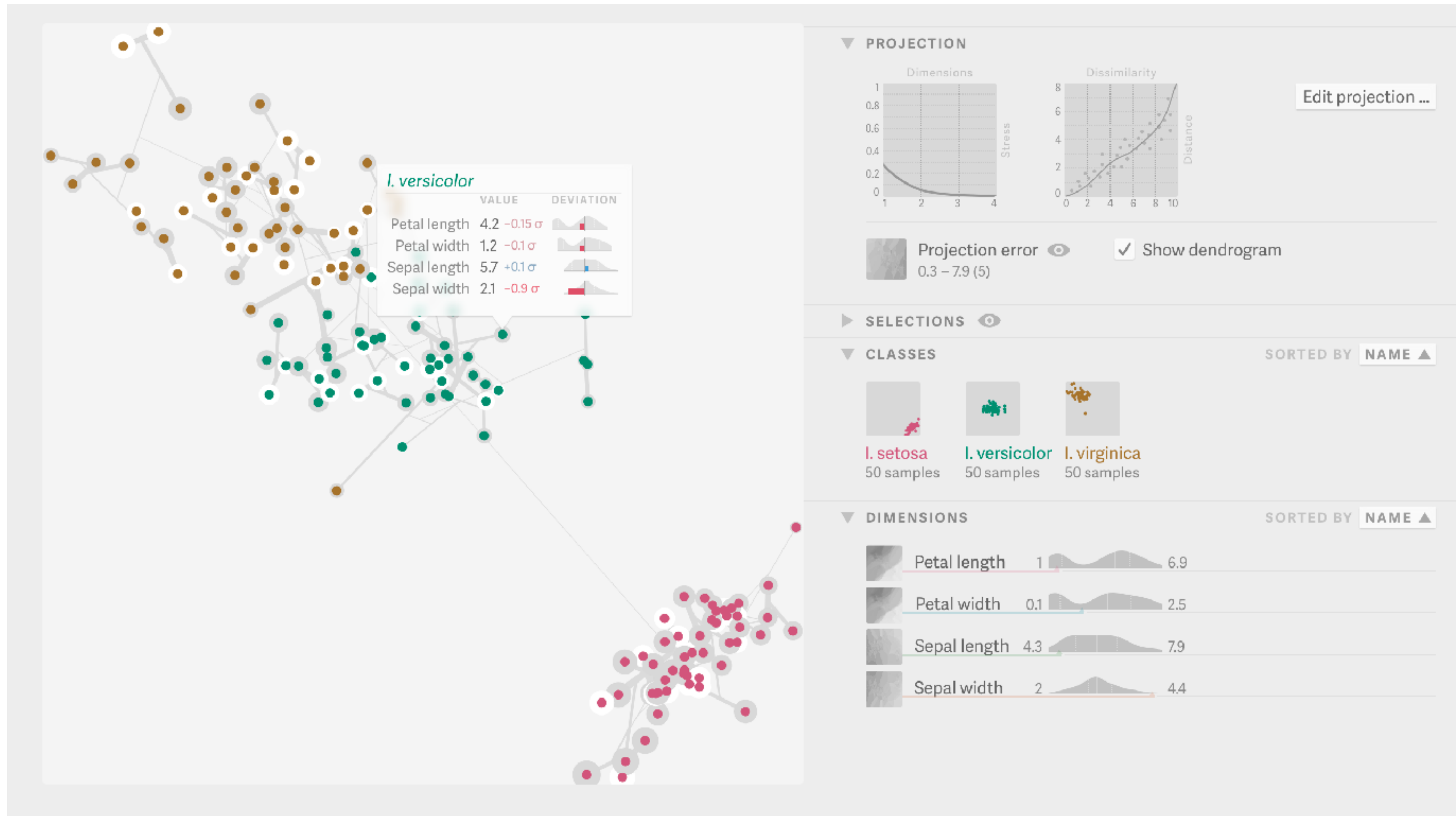


Topical distances between the selected Petroleum Engineering and the others.

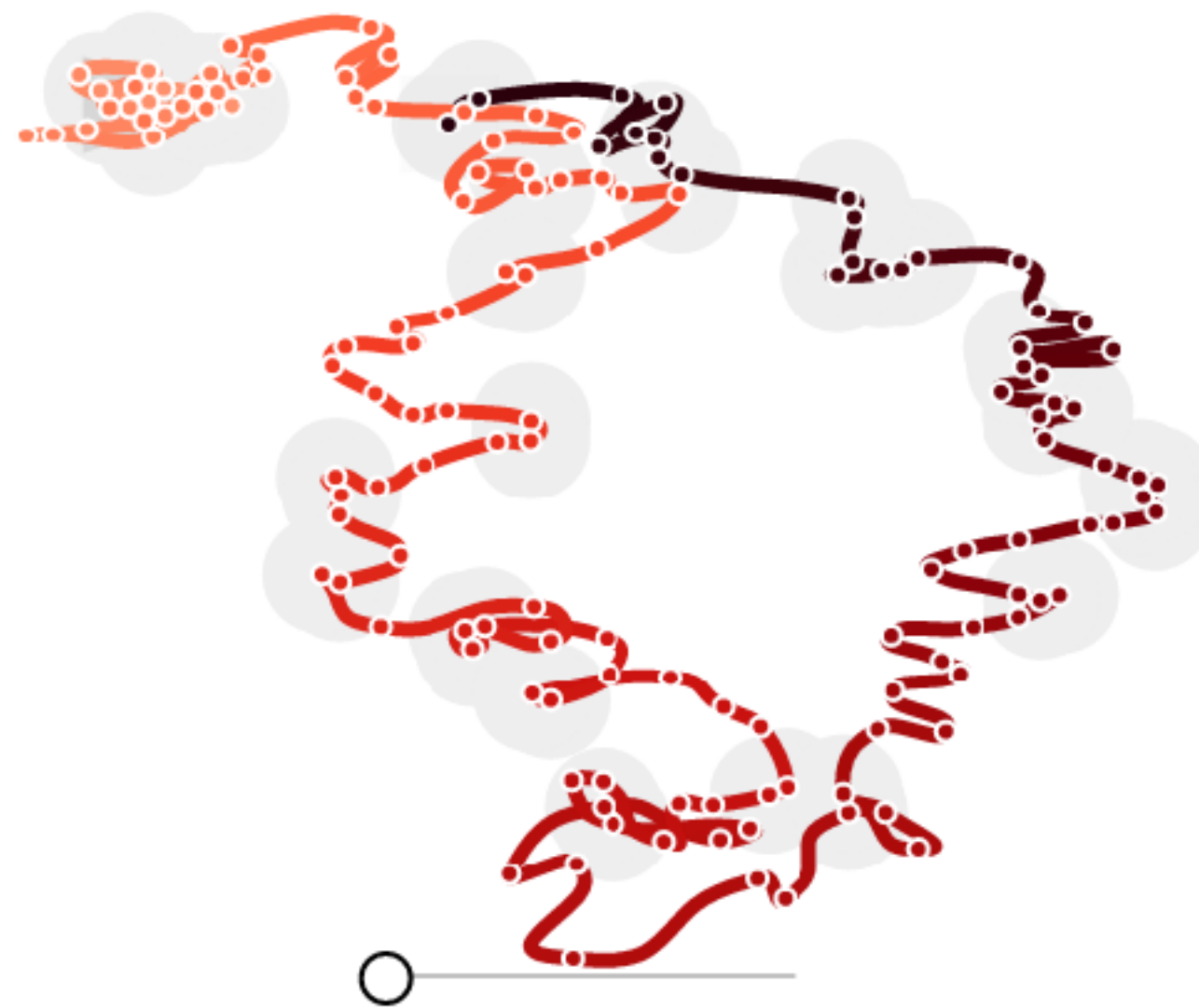


[Chuang et al., 2012]

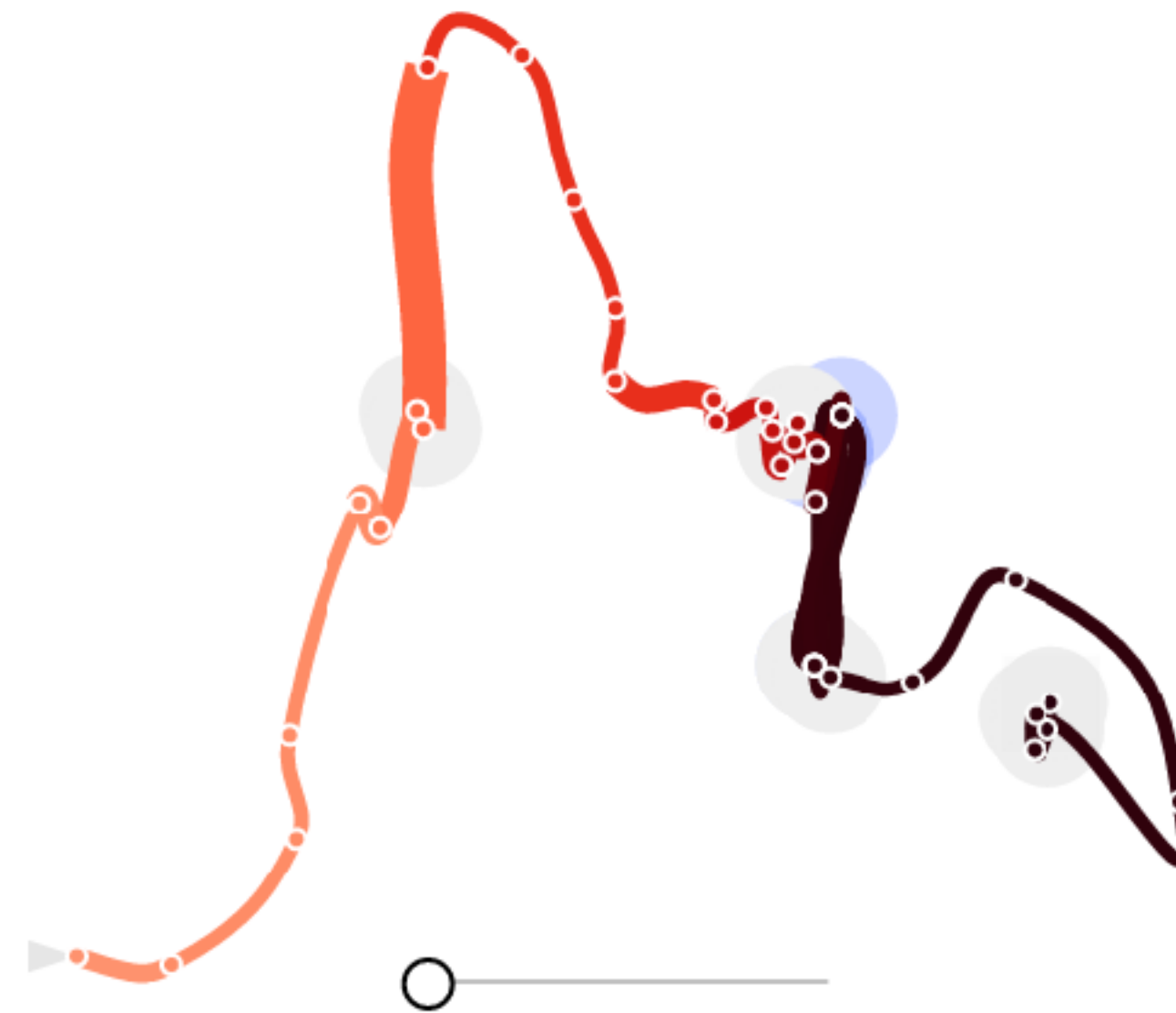
Probing Projections



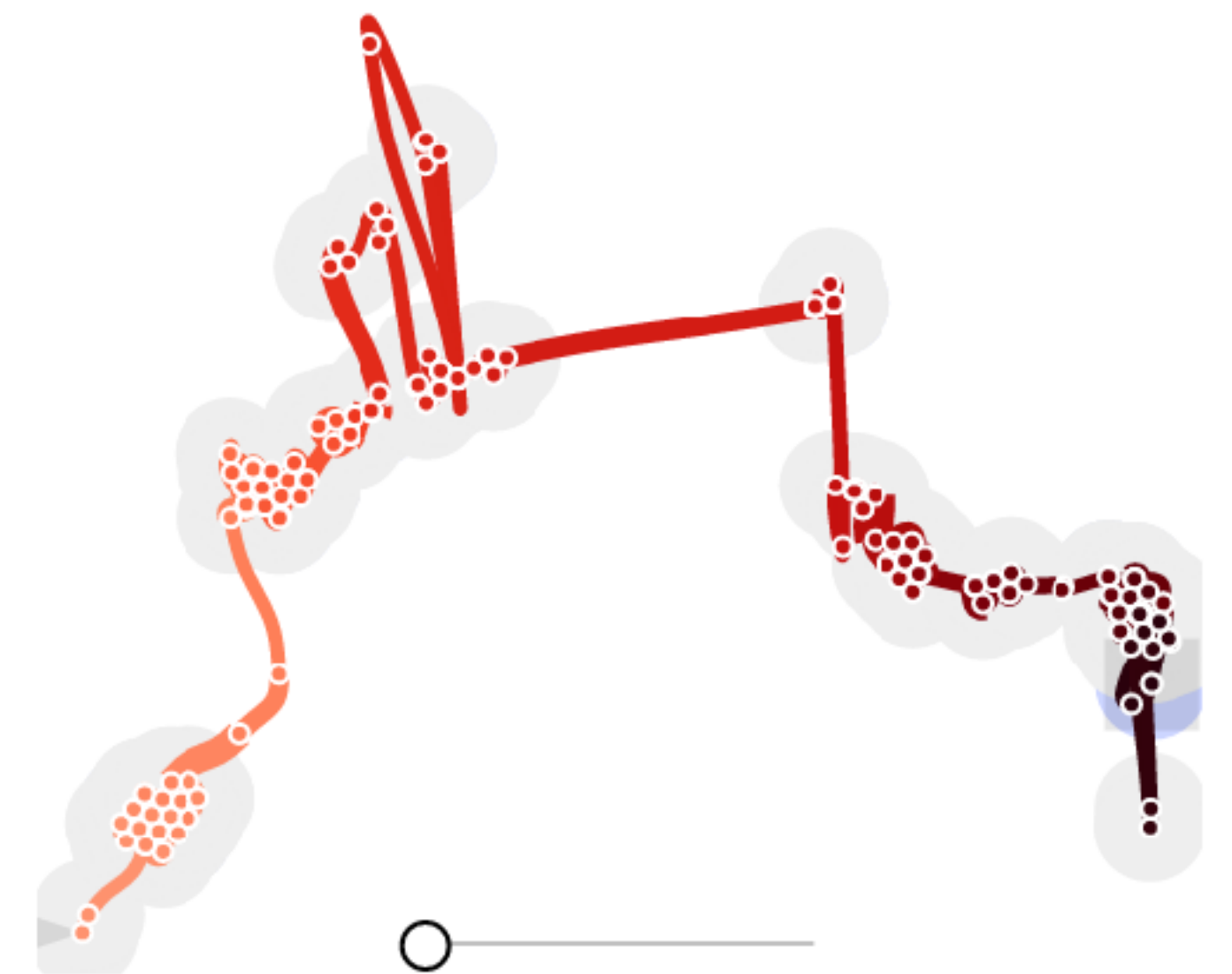
MDS for Temporal Data: TimeCurves



Video: Global Cloud Circulation (146)



Wikipedia: Chocolate (46) 



Wikipedia: Palestine 200 1 (200) 