

Clustering

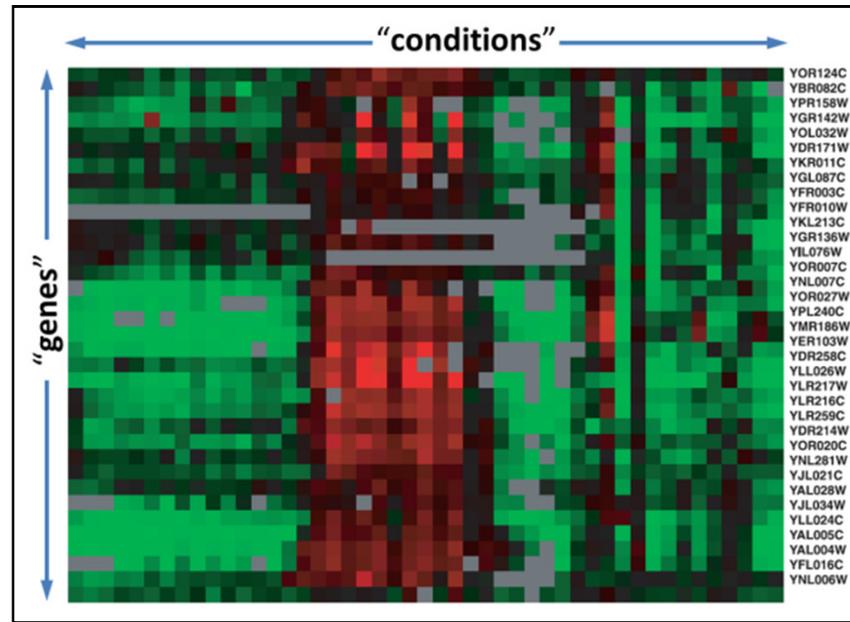
Genome 373

Genomic Informatics

Elhanan Borenstein

The clustering problem

- The goal of gene clustering process is to partition the genes into distinct sets such that genes that are assigned to the same cluster are “similar”, while genes assigned to different clusters are “non-similar”.

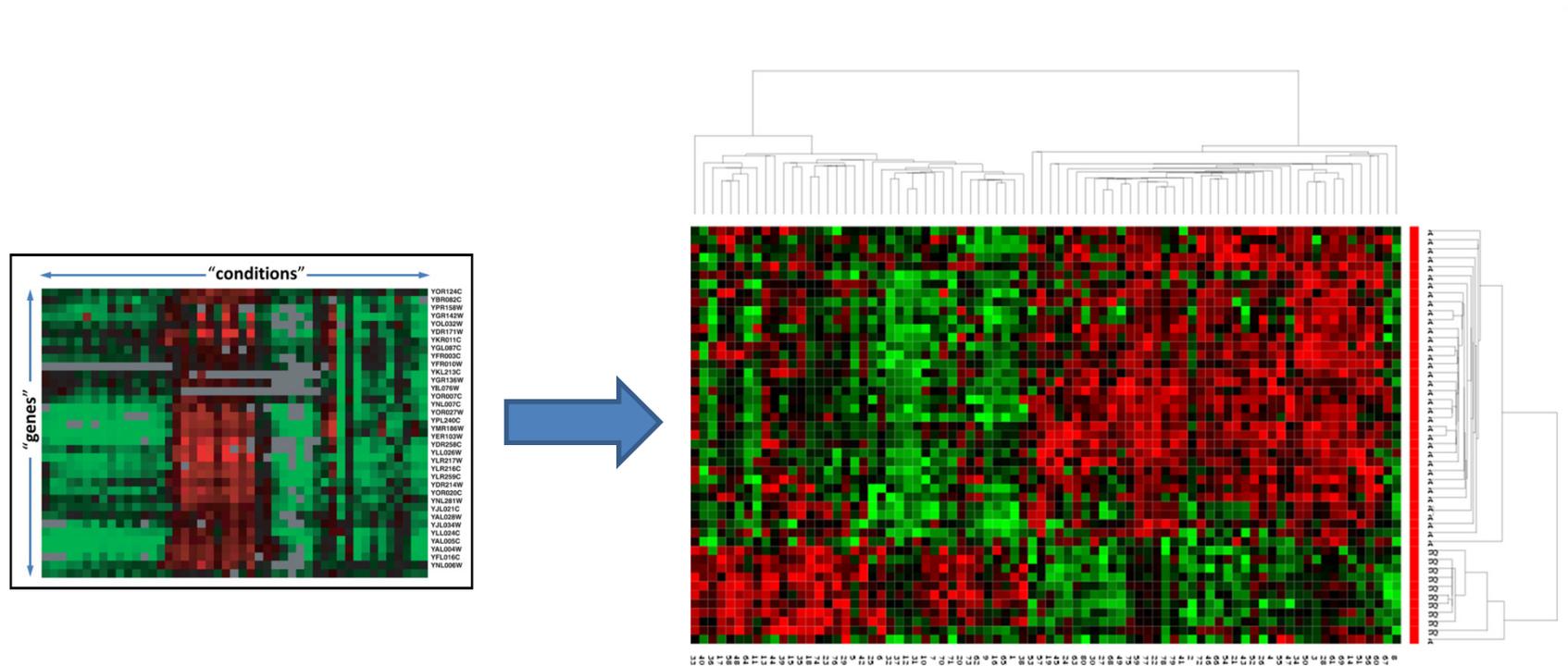


Clustering vs. Classification

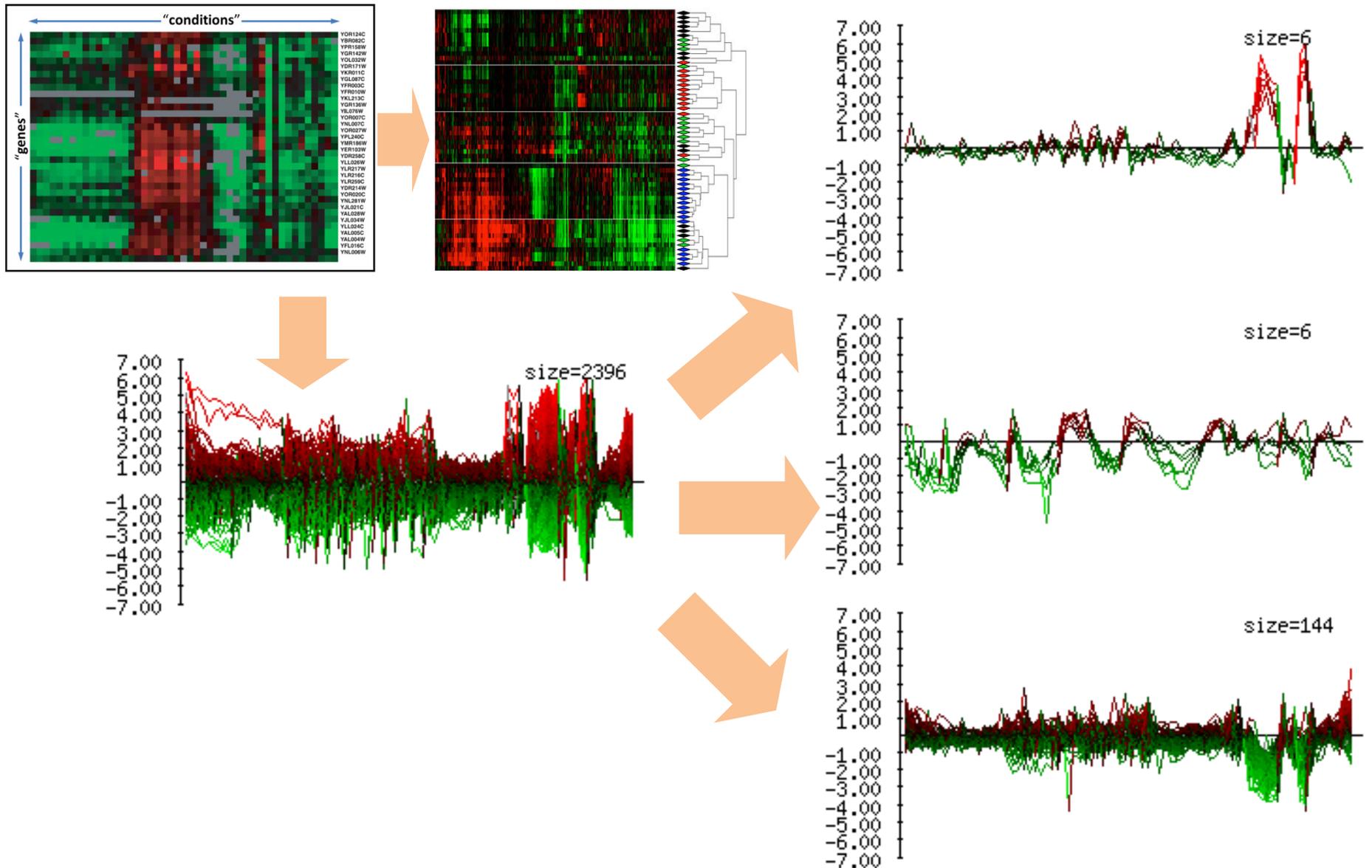
- Clustering is an **exploratory** tool: “who's running with who”.
- A very different problem from classification:
 - Clustering is about finding coherent groups
 - Classification is about relating such groups or individual objects to specific labels, mostly to support future prediction
- Most clustering algorithms are **unsupervised**
(in contrast, most classification algorithms are supervised)
- Clustering methods have been used in a vast number of disciplines and fields.

What are we clustering?

- We can cluster genes, conditions (samples), or both.



Clustering gene expression profiles



Why clustering

- Clustering genes or conditions is a basic tool for the analysis of expression profiles, and can be useful for many purposes, including:
 - Inferring functions of unknown genes
(assuming a similar expression pattern implies a similar function).
 - Identifying disease profiles
(tissues with similar pathology should yield similar expression profiles).
 - Deciphering regulatory mechanisms: co-expression of genes may imply co-regulation.
 - **Reducing dimensionality.**

The clustering problem

- A good clustering solution should have two features:
 1. **High homogeneity:** homogeneity measures the similarity between genes assigned to the same cluster.
 2. **High separation:** separation measures the distance/dissimilarity between clusters.
(If two clusters have similar expression patterns, then they should probably be merged into one cluster).
- Note that there is usually a tradeoff between these two features:
 - More clusters → increased homogeneity but decreased separation
 - Less clusters → Increased separation but reduced homogeneity

One problem, numerous solutions

- No single solution is necessarily the true/correct mathematical solution!
- There are many formulations of the clustering problem; most of them are **NP-hard**.
- Therefore, in most cases, heuristic methods or approximations are used.
 - Hierarchical clustering
 - k-means
 - self-organizing maps (SOM)
 - Knn
 - PCC
 - CAST
 - CLICK

One problem, numerous solutions

- The results (i.e., obtained clusters) can vary drastically depending on:
 - Clustering method
 - Similarity or dissimilarity metric
 - Parameters specific to each clustering method (e.g. number of centers for the k-mean method, agglomeration rule for hierarchical clustering, etc.)

Clustering methods

- We can distinguish between two types of clustering methods:
 1. **Agglomerative:** These methods build the clusters by examining small groups of elements and merging them in order to construct larger groups.
 2. **Divisive:** A different approach which analyzes large groups of elements in order to divide the data into smaller groups and eventually reach the desired clusters.
- There is another way to distinguish between clustering methods:
 1. **Hierarchical:** Here we construct a hierarchy of clusters in order to examine the relationship between entities.
 2. **Non-Hierarchical:** In non-hierarchical methods, the elements are partitioned into non-overlapping groups.

Hierarchical
clustering

K-mean
clustering

Hierarchical clustering

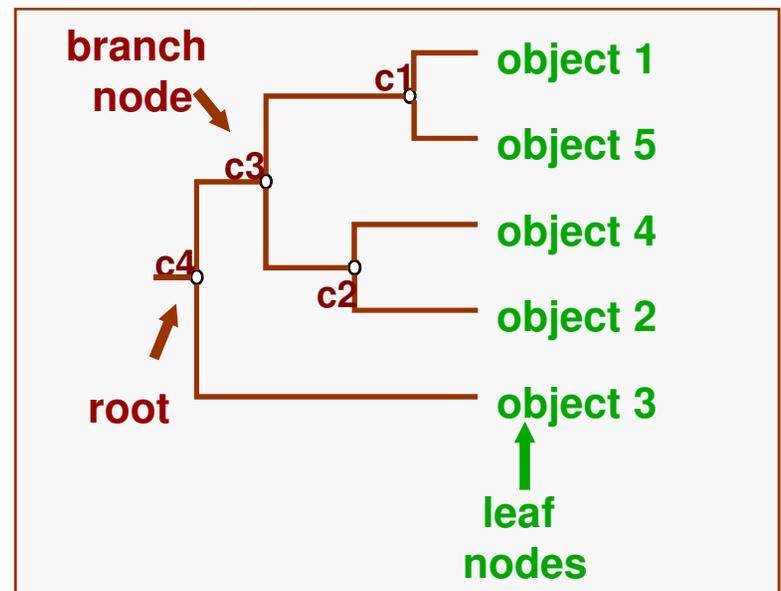
Hierarchical clustering

- **Hierarchical** clustering is an **agglomerative** clustering method
 - Takes as input a distance matrix
 - Progressively regroups the closest objects/groups

Distance matrix

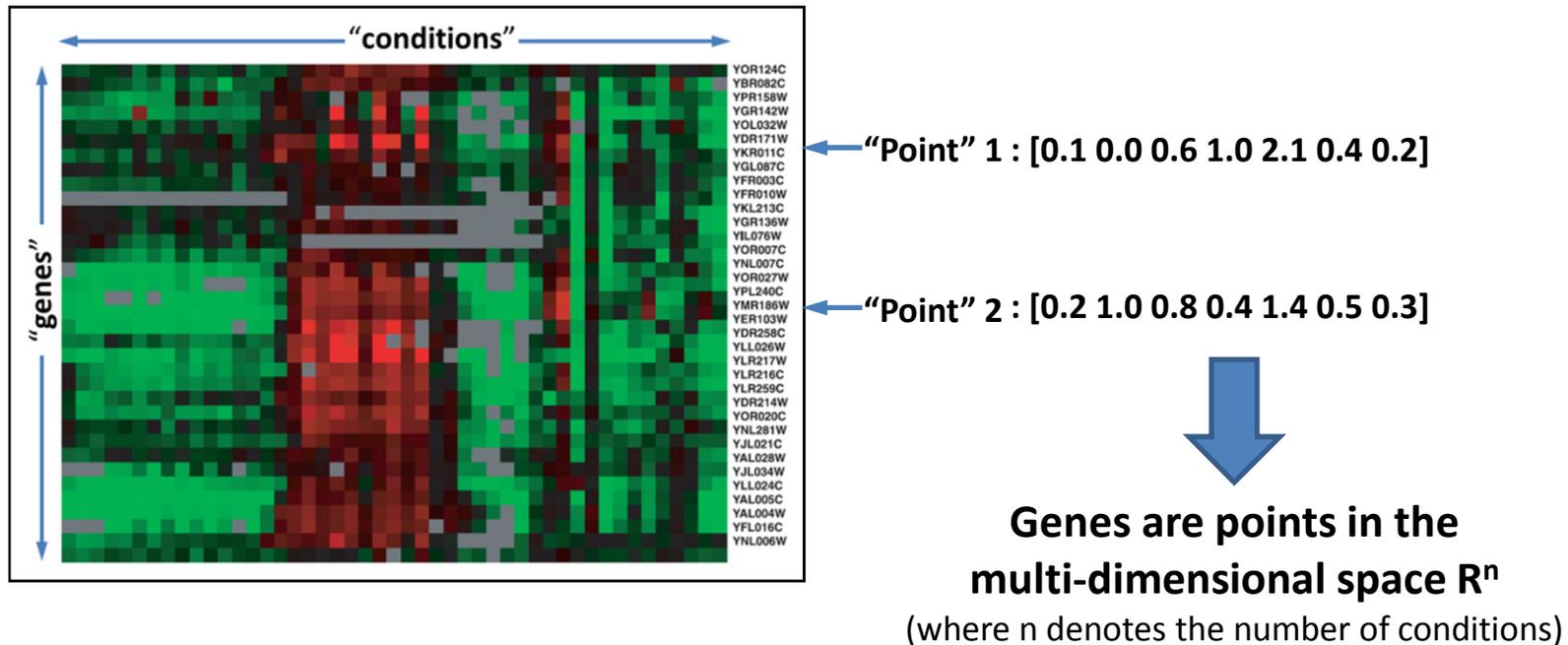
	object 1	object 2	object 3	object 4	object 5
object 1	0.00	4.00	6.00	3.50	1.00
object 2	4.00	0.00	6.00	2.00	4.50
object 3	6.00	6.00	0.00	5.50	6.50
object 4	3.50	2.00	5.50	0.00	4.00
object 5	1.00	4.50	6.50	4.00	0.00

Tree representation



Measuring similarity/distance

- An important step in many clustering methods is the selection of a distance measure (**metric**), defining the distance between 2 data points (e.g., 2 genes)



Measuring similarity/distance

- So ... how do we measure the distance between two point in a multi-dimensional space?

p-norm

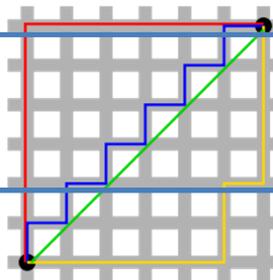
$$\|x\|_p := \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$$

- Common distance functions:

- The **Euclidean** distance $\|x\| := \sqrt{x_1^2 + \dots + x_n^2}$ ← *2*-norm
(a.k.a “distance as the crow flies” or distance).

- The **Manhattan** distance ← *1*-norm
(a.k.a **taxicab** distance)

- The **maximum** norm ← *infinity*-norm
(a.k.a **infinity** distance)

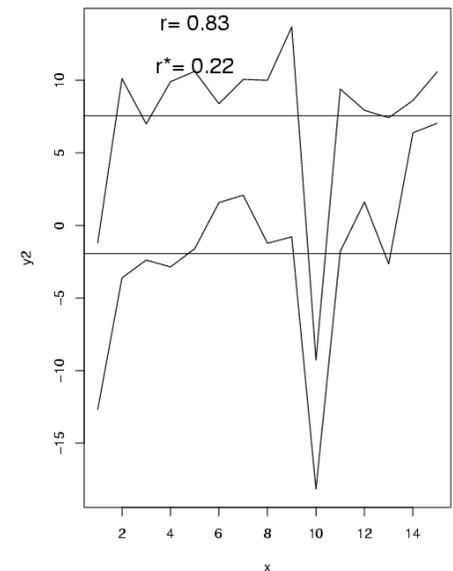
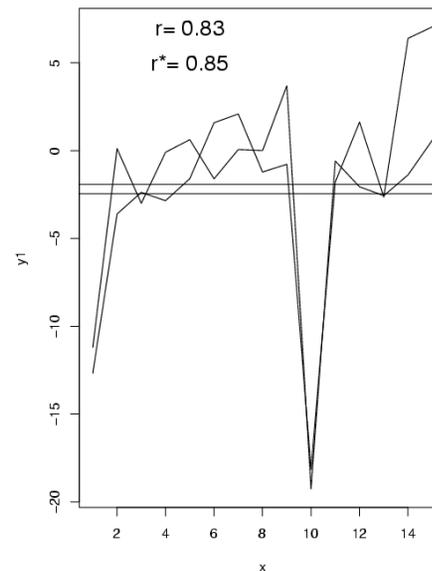


- The **Hamming** distance
(number of substitutions required to change one point into another).

- **Symmetric vs. asymmetric distances.**

Correlation as distance

- Another approach is to use the correlation between two data points as a distance metric.
 - Pearson Correlation
 - Spearman Correlation
 - Absolute Value of Correlation



Metric matters!

- The metric of choice has a marked impact on the shape of the resulting clusters:
 - Some elements may be close to one another in one metric and far from one another in a different metric.
- Consider, for example, the point $(x=1, y=1)$ and the origin.
 - What's their distance using the 2-norm (Euclidean distance)?
 - What's their distance using the 1-norm (a.k.a. taxicab/ Manhattan norm)?
 - What's their distance using the infinity-norm?

Hierarchical clustering algorithm

1. Assign each object to a separate cluster.
2. Find the pair of clusters with the shortest distance, and regroup them into a single cluster.
3. Repeat 2 until there is a single cluster.

- The result is a tree, whose intermediate nodes represent clusters
- Branch lengths represent distances between clusters

Hierarchical clustering

1. Assign each object to a separate cluster.
 2. **Find the pair of clusters with the shortest distance, and regroup them into a single cluster.**
 3. Repeat 2 until there is a single cluster.
- One needs to define a (dis)similarity metric between two **groups**. There are several possibilities
 - **Average linkage:** the average distance between objects from groups A and B
 - **Single linkage:** the distance between the closest objects from groups A and B
 - **Complete linkage:** the distance between the most distant objects from groups A and B

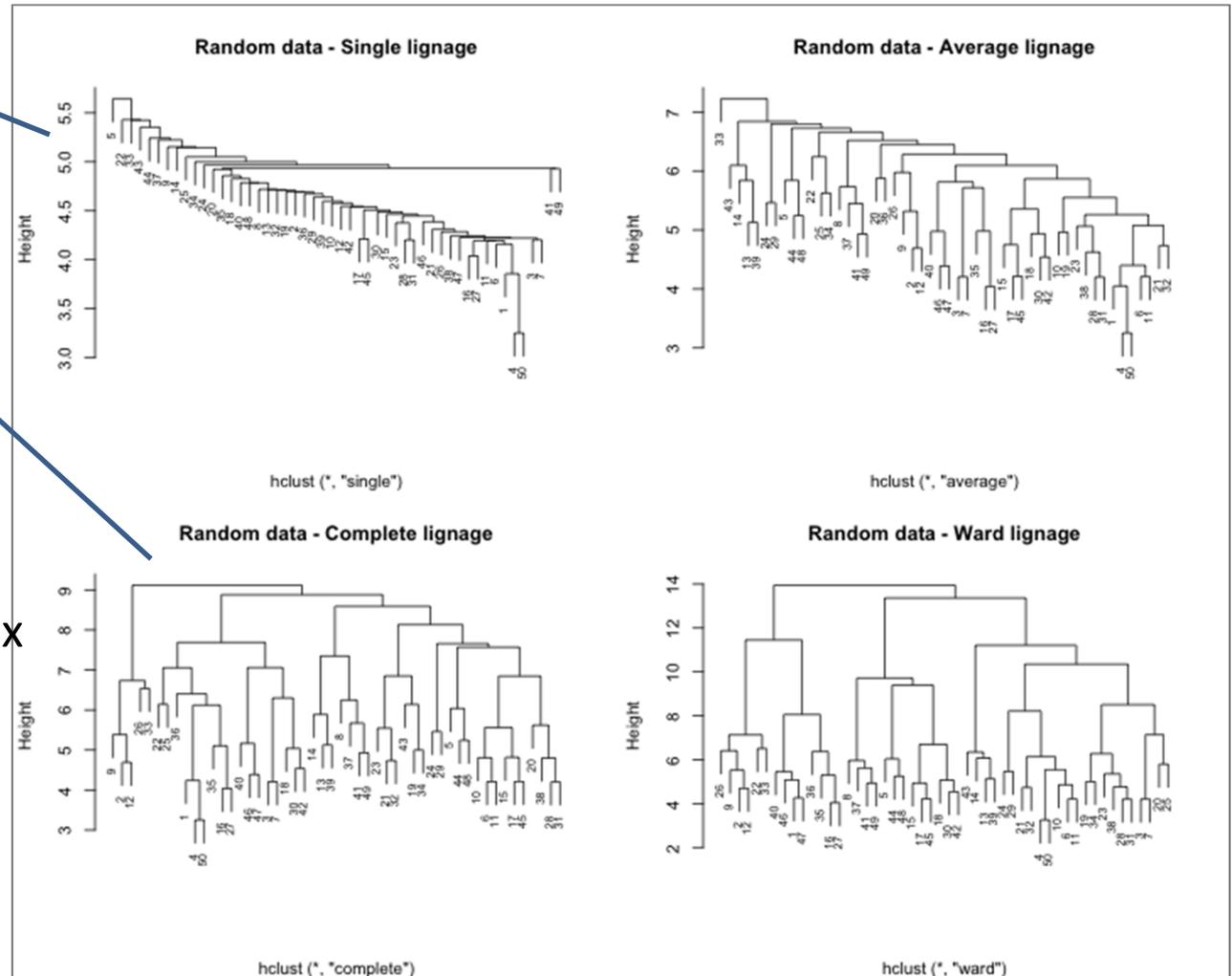
Impact of the agglomeration rule

- These four trees were built from the same distance matrix, using 4 different agglomeration rules.

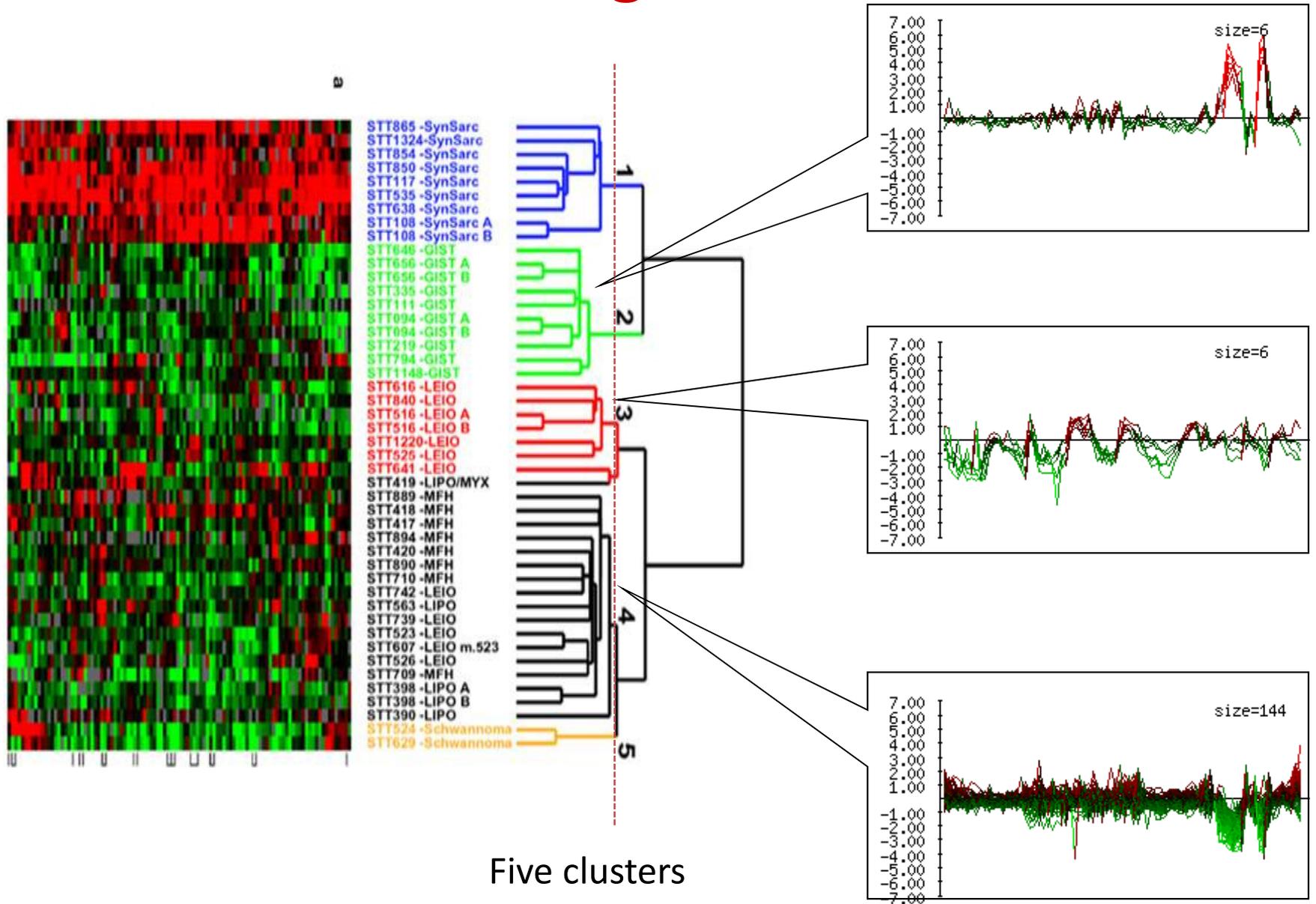
Single-linkage typically creates nesting clusters

Complete linkage create more balanced trees.

Note: these trees were computed from a matrix of random numbers. The impression of structure is thus a complete artifact.



Hierarchical clustering result



Clustering in both dimensions

