# Clustering

Katie Pollard

In this unit we will learn …

- Differences between supervised and unsupervised learning

- How commonly used distance measurements encode different notions of "close"

- Hierarchical versus partitioning algorithms for clustering bioinformatics data

- Strategies for statistical inference and assessing variability with clustering results

# Clustering vs. Classification

Clustering = unsupervised learning
- Classes: unknown a priori
- Goal: discover groups from the data

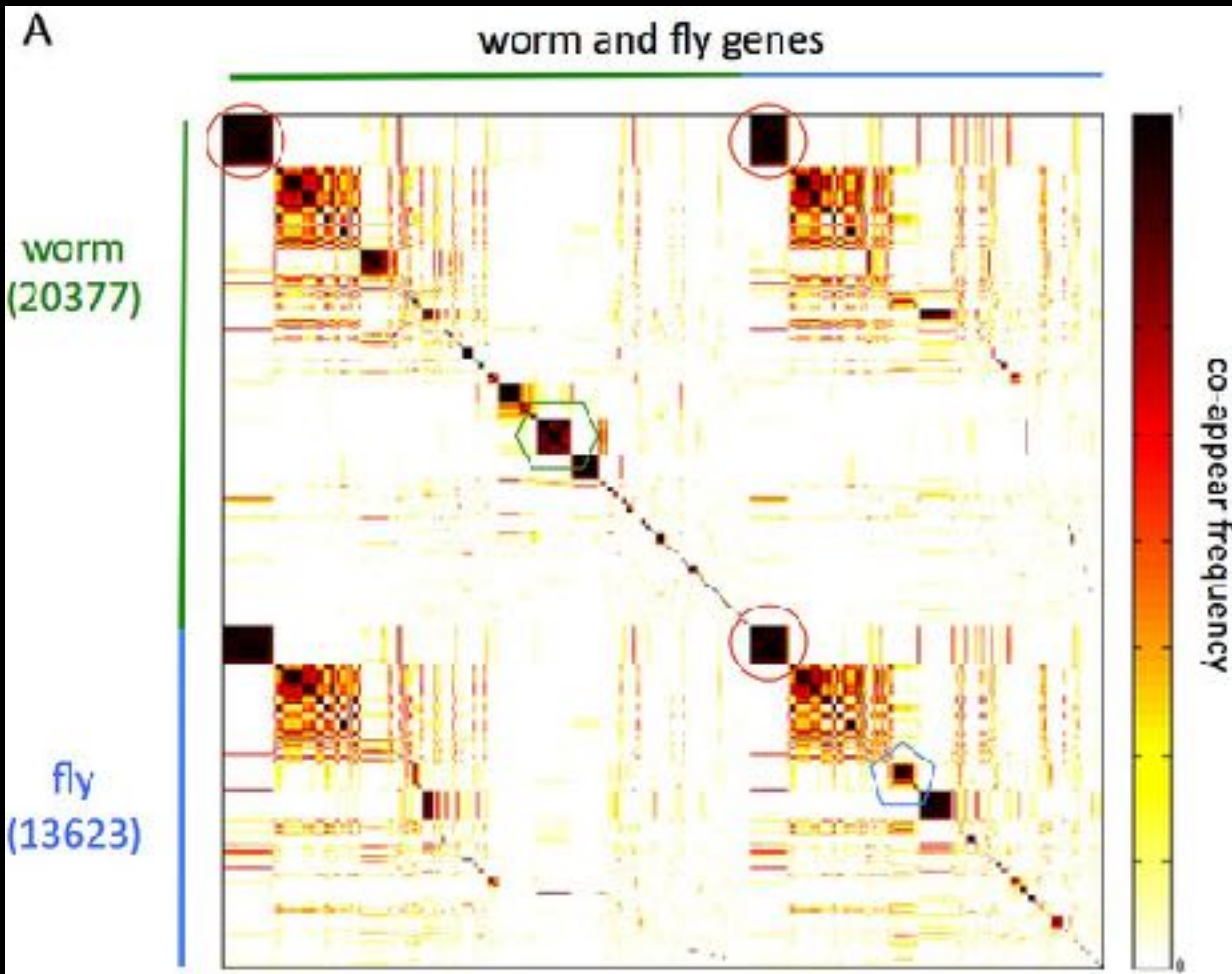Classification = supervised learning
- Classes: known/predefined
- Goals: understand the basis for the classes and build a predictor (classify new data)
- Prediction with a categorical outcome variable

# Cluster Analysis

Exploratory data analysis methods for:

- Discovering patterns

- Grouping
  - Variables
  - Samples
  - Both simultaneously

- Ordering and organizing

- Dimension reduction
  - How many distinct patterns?

# Clustering in Bioinformatics



A — worm and fly genes

worm (20377)

fly (13623)

co-appear frequency

Worm orthologs of co-expressed genes in fly are some times also co-expressed.

Yan et al. (2014) Genome Biology

# Clustering Methods

Two main components:

1. Distance measure
2. Algorithm

These produce a mapping from data to parameters of interest:

– Cluster labels, sizes, profiles
– Hierarchical tree structure, ordering
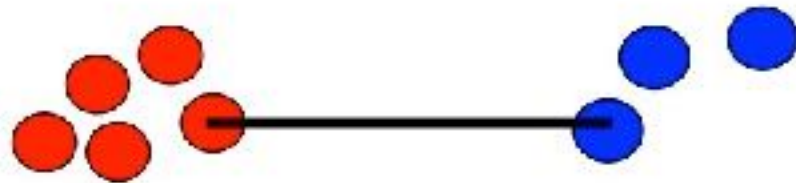– Number of groups

# Distances

Clustering algorithms require a notion of pairwise distance between objects.

- Minkowski metrics: magnitude
- Correlation distances: pattern (or both)
- The absolute value of any distance can also be used, e.g.

$$d(x,y) = 1 - \left|r(x,y)\right| \in (0,1)$$

# Distance Between Clusters
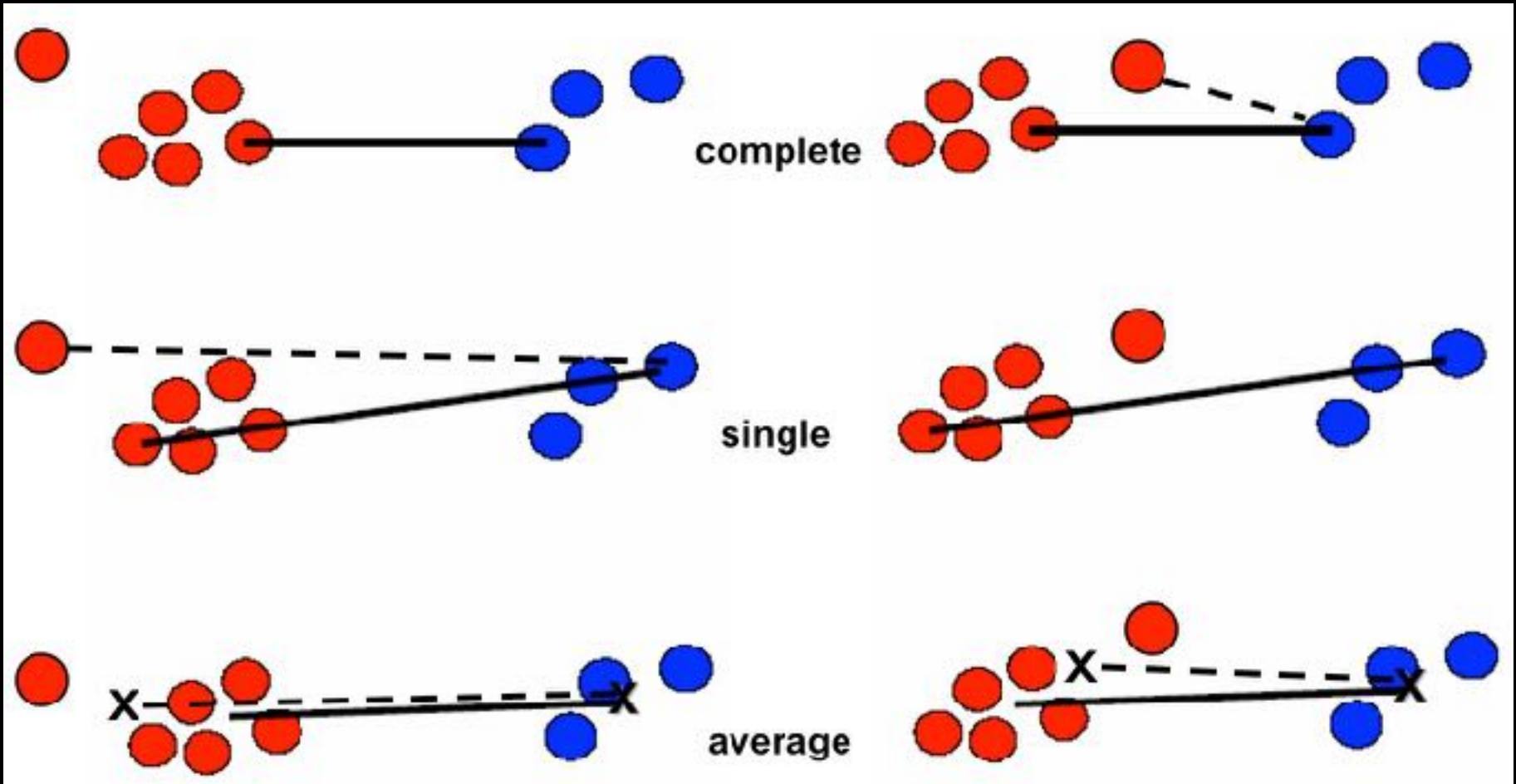


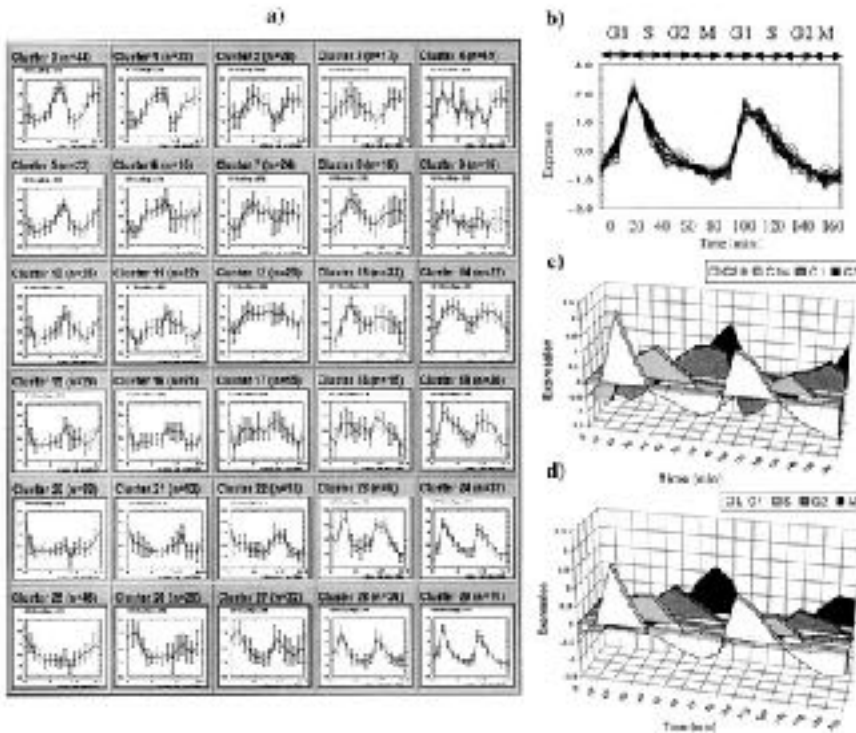Complete (minimum)

Single (maximum)

Average
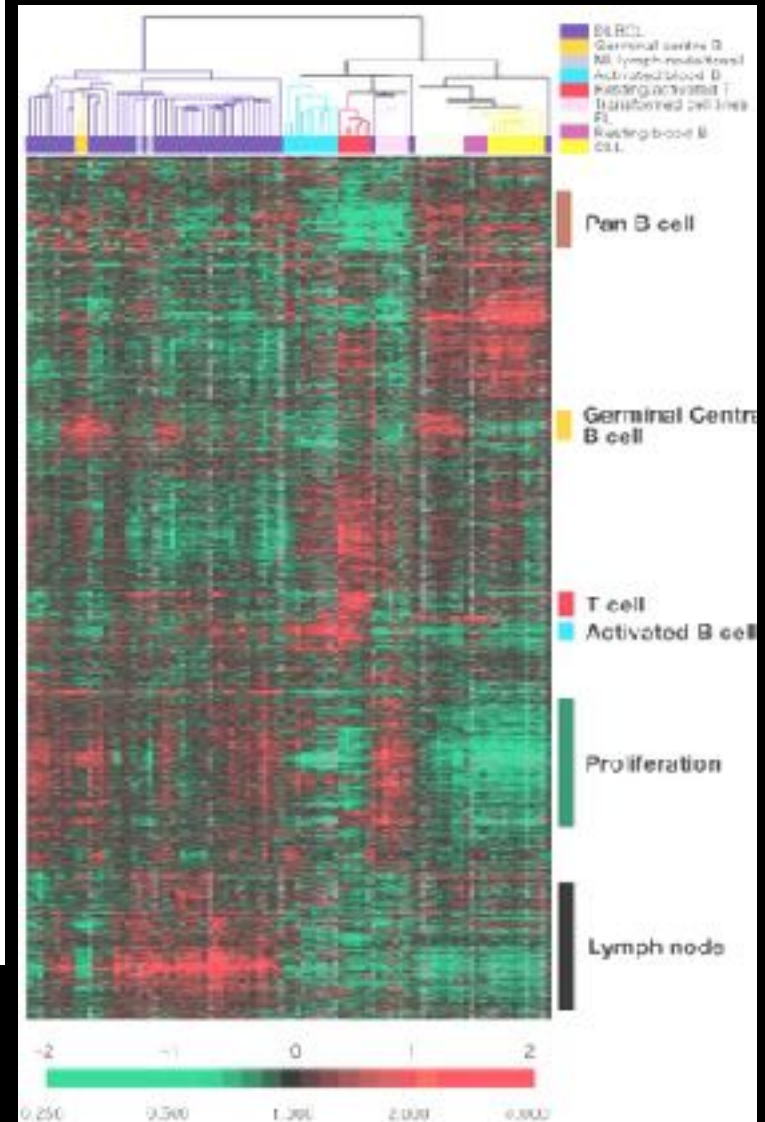
# Effects of Outliers

# Clustering Algorithms

- Model-based (AUTOCLASS,SNOB)
  vs. Non-parametric
- Partitioning (SOMs, PAM, KMEANS)
  vs. Hierarchical
  - Agglomerative (CLUSTER, AGNES)
    - Linkage: single, complete, average
  - Devisive (SOTA,DIANA,TSVQ)
  - Hybrid (HOPACH, MUTUAL CLUSTERS)

# Partitioning vs. Hierarchical



Cho *et al.* (1998) Molecular Cell, 2: 65-73

Tamayo *et al.* (1999) PNAS, 96: 2907-2912

Alizadeh *et al.* (2000) Nature, 403: 503-511

# How Many Clusters?

- Relevant for both partitioning and hierarchical algorithms (with pruning)

- Level of structure: global vs. detailed

- Two main approaches:

    - Direct Methods (criteria)

        e.g. sums of squares, silhouettes

    - Resampling Methods (testing)

        e.g. Clest, gap statistic, bootstrap

# Silhouette-based Criteria

- The silhouette for j'th object (e.g., gene):

avg distance own cluster

avg distance next closest cluster

$$S_j = \frac{a_j - b_j}{\max(a_j, b_j)}$$

- Cluster average silhouette: mean $S_j$ per cluster
- Average silhouette: overall mean $S_j$
- Median split silhouette (MSS): split each cluster and see if silhouettes get smaller

# Inference for Clustering

- How reliable and repeatable are cluster results from a single data set?

- Can view output (e.g. gene cluster labels) as a parameter estimate.

- Use resampling methods to estimate the variability of this estimator (since no closed form typically).

Example: Bootstrap cluster memberships

# Issues in Cluster Analysis

- Results can be very sensitive to input (i.e., pre-processing and filtering).
- What method fits your application?
  - Distance: capture what "close" means
  - Algorithm: the kind of clusters you seek
- Clustering methods will <u>always</u> return some output, but is it meaningful?
  - Evaluate variability
  - Assess biological relevance
  - Confirm hypotheses with experiments

# Clustering Algorithms in R

| Package | Functions | Type |
|---------|-----------|------|
| `stats` | `kmeans` `hclust` | partitioning divisive |
| `class` | `SOM` | partitioning |
| `cluster` | `pam` `agnes` `diana` | partitioning agglomerative divisive |
| `hopach` | `hopach` | hybrid |

Other packages: `cclust`, `e1071`, `flexmix`, `fpc`, `mclust`

# Comparing Clustering

- Any distance measure?
- Reaction to noise: robust vs. efficient
- Results reproducible?
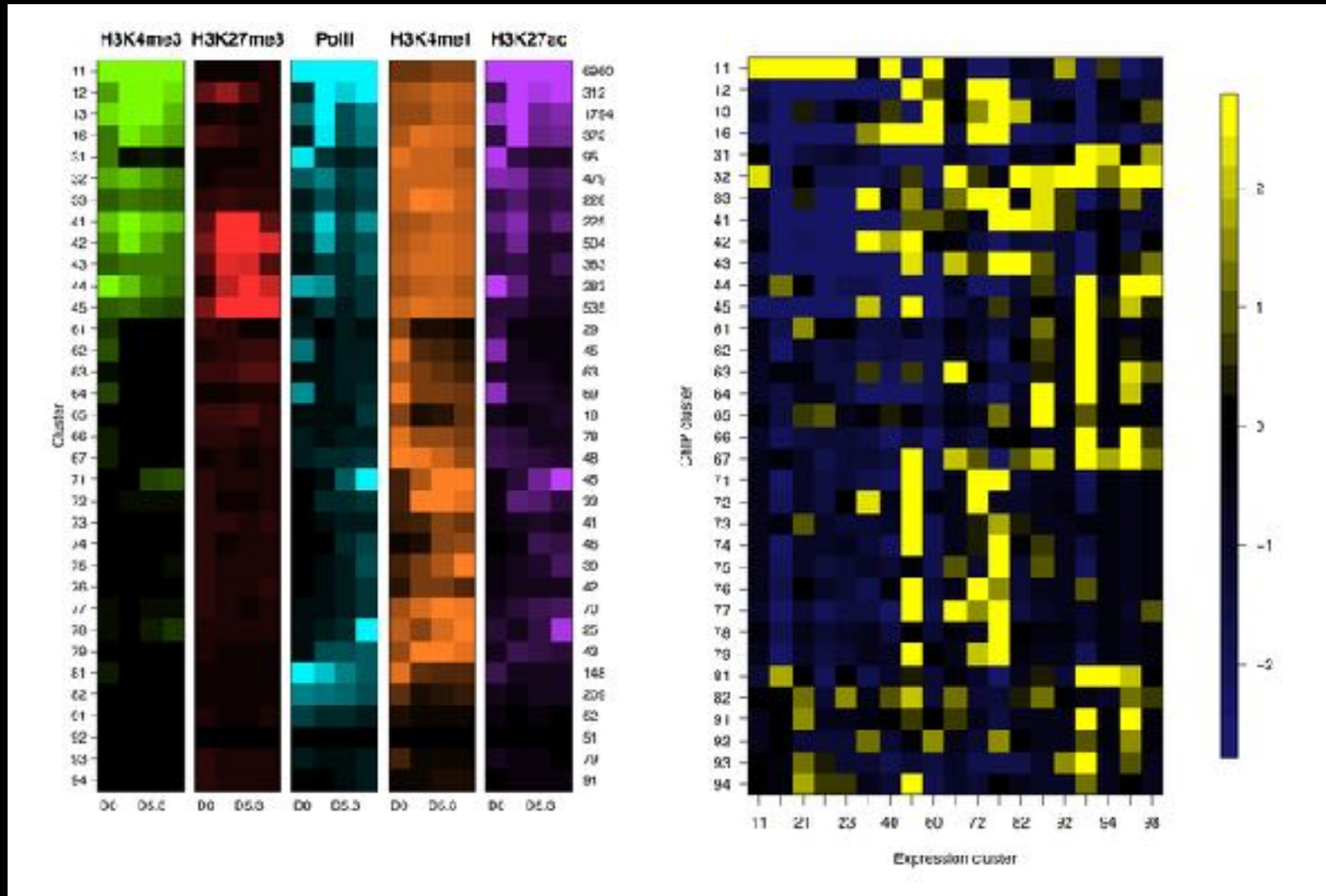- Biological relevance

Partitioning

Variety of cluster sizes

Overlapping clusters

Hierarchical

Sensible ordering

# Hybrid tree of partitions

HOPACH: Pruned hierarchical tree produces nested clusters



Wamstad *et al.* (2012) Cell, 151: 206-220

# Bootstrap Fuzzy Clustering



**boothopach** function

*Pyrobaculum aerophilum* array data from Lowe Lab (UCSC)