


www.cs.washington.edu/527



**University of Washington**  
Computer Science & Engineering

**CSE 527, Au '03: Computational Biology**

[CSE Home](#)
[About Us](#)
[Search](#)
[Contact Info](#)

**Administrative**  
[Syllabus](#)  
**Lecture Slides**  
[Overview](#)  
[Microarrays I](#)  
[Microarrays II](#)  
**Lecture Notes**  
[2. Microarrays I](#)  
[4. Microarrays III](#)  
**Assignments**  
[HW #1](#)  
[HW #2](#)  
**Notes on Readings**  
[HW #1: Primers](#)  
[HW #2: Microarrays](#)  
**Project Information**

**Time:** MW 12:00-1:20  
**Place:** MGH 284  
**Office Hours** **Phone**  
**Instructor:** Larry Ruzzo, [luzzo@cs,](mailto:luzzo@cs.washington.edu) M 1:30 554 Allen Center, 543-6298  
**TA:** Zizhen Yao, [yizizhen@cs,](mailto:yizizhen@cs.washington.edu)

An introduction to the use of computational biology to gain a better understanding of biological systems. Intended for graduate students in bioinformatics and computer science interested in applications of those fields to biology.

**Subscribe, if you haven't gotten any msgs**

**Mail archive** of all mail sent to [cse527@cs.](mailto:cse527@cs.washington.edu) Read it regularly or subscribe.

References:

The following slides partly from <http://staff.washington.edu/kayee/research.html>  
Errors are mine.

## Clustering 101

Ka Yee Yeung  
Center for Expression Arrays  
University of Washington

## Overview

- What is clustering?
- Similarity/distance metrics
- Hierarchical clustering algorithms
  - E.g. [Eisen et al. 1998]
- K-means
  - E.g. [MacQueen, 1965] [Tavazoie et al. 1999]
- Self-organizing map (SOM)
  - E.g. [Tamayo et al. 1999]

## What is clustering?

- Group *similar* objects together
- Objects in the same cluster (group) are more similar to each other than objects in different clusters
- Data exploratory tool

## Clustering Expression Data

- Why cluster gene expression data?
    - Tissue classification
    - Find biologically related genes
    - First step in inferring regulatory networks
    - Look for common promoter elements
- } columns  
} rows
- Hypothesis generation
  - One of the tools of choice for expression analysis

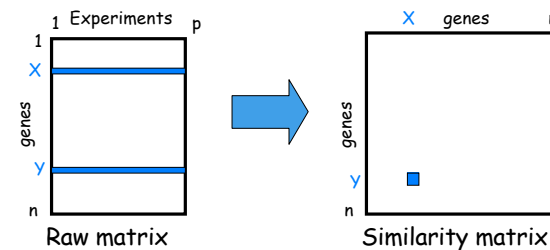
## Clustering Expression Data

- What has been done?
  - Partitional
    - CAST (Ben-Dor et al. 1999)
    - k-means, variously initialized (Hartigan 1975)
  - Hierarchical
    - single-, average-, complete-, centroid-link [Eisen et al. 98]
  - Self Organizing Maps (SOM) [Tamayo et al. 99]
  - Support Vector Machines (SVM) [Grundy et al. 00]
  - etc., etc., etc.

## Clustering Expression Data

- Why so many methods?
  - Clustering is NP-hard, even with simple objectives, data
  - Hard problem: high dimensionality, noise, ...
  - $\therefore$  many heuristic, local search, & approximation algorithms
  - No clear winner

## How to define similarity?



- Similarity metric:
  - A measure of pairwise similarity or dissimilarity
  - Examples:
    - Correlation coefficient
    - Euclidean distance

## Similarity metrics

- Euclidean distance

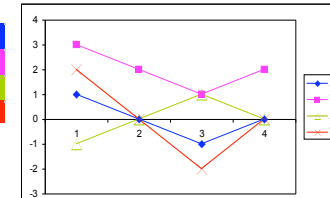
$$\sqrt{\sum_{j=1}^p (X[j] - Y[j])^2}$$

- Correlation coefficient

$$\frac{\sum_{j=1}^p (X[j] - \bar{X})(Y[j] - \bar{Y})}{\sqrt{\sum_{j=1}^p (X[j] - \bar{X})^2 \sum_{j=1}^p (Y[j] - \bar{Y})^2}}, \text{ where } \bar{X} = \frac{\sum_{j=1}^p X[j]}{p}$$

## Example

X	1	0	-1	0
Y	3	2	1	2
Z	-1	0	1	0
W	2	0	-2	0



Correlation (X,Y) = 1

Distance (X,Y) = 4

Correlation (X,Z) = -1

Distance (X,Z) = 2.83

Correlation (X,W) = 1

Distance (X,W) = 1.41

## Lessons from the example

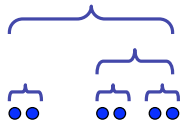
- Correlation – direction only
- Euclidean distance – magnitude & direction
- Min # attributes (experiments) to compute pairwise similarity
  - $\geq 2$  attributes for Euclidean distance
  - $\geq 3$  attributes for correlation
- Array data is noisy  $\rightarrow$  need many experiments to robustly estimate pairwise similarity

## Clustering algorithms

- Inputs:
  - Raw data matrix or similarity matrix
  - Number of clusters or some other parameters
- Many different classifications of clustering algorithms:
  - Hierarchical vs partitional
  - Heuristic-based vs model-based
  - Soft vs hard

# Hierarchical Clustering

[Hartigan 1975]

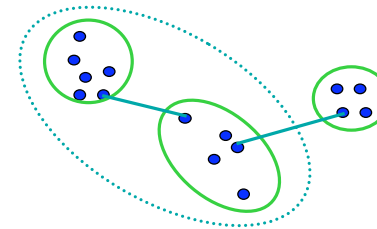


dendrogram

- Agglomerative (bottom-up)
- Algorithm:
  - Initialize: each item a cluster
  - Iterate:
    - select two most similar clusters
    - merge them
  - Halt: when required number of clusters is reached

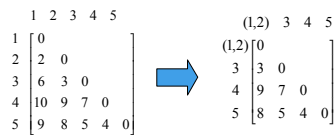
# Hierarchical: Single Link

- cluster similarity = similarity of two most similar members

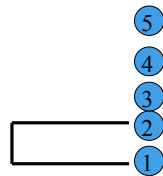


- Potentially long and skinny clusters
- + Fast

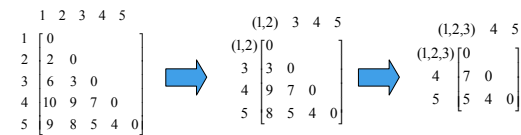
# Example: single link



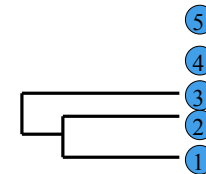
$$\begin{aligned}
 d_{(1,2),3} &= \min\{d_{1,3}, d_{2,3}\} = \min\{6,3\} = 3 \\
 d_{(1,2),4} &= \min\{d_{1,4}, d_{2,4}\} = \min\{10,9\} = 9 \\
 d_{(1,2),5} &= \min\{d_{1,5}, d_{2,5}\} = \min\{9,8\} = 8
 \end{aligned}$$



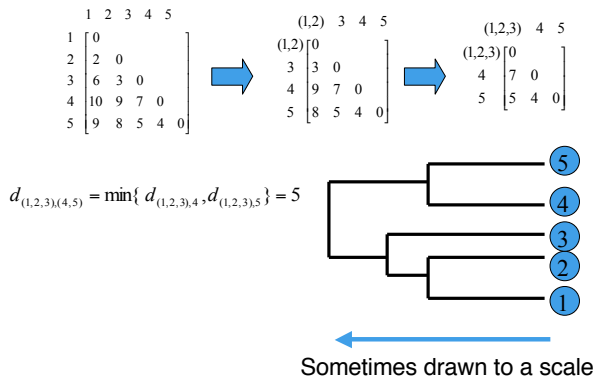
# Example: single link



$$\begin{aligned}
 d_{(1,2,3),4} &= \min\{d_{(1,2),4}, d_{3,4}\} = \min\{9,7\} = 7 \\
 d_{(1,2,3),5} &= \min\{d_{(1,2),5}, d_{3,5}\} = \min\{8,5\} = 5
 \end{aligned}$$

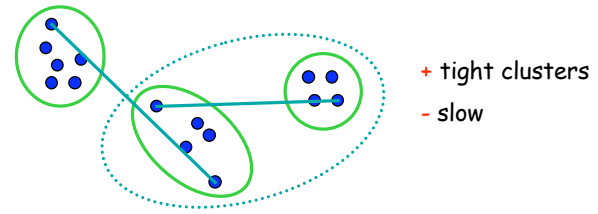


### Example: single link

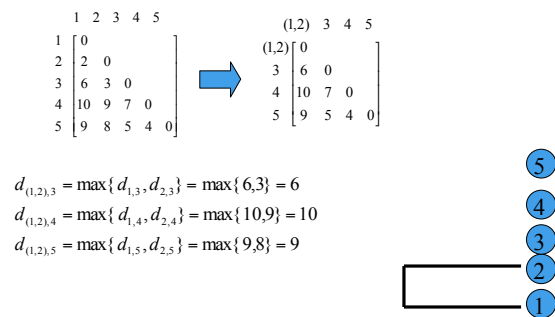


### Hierarchical: Complete Link

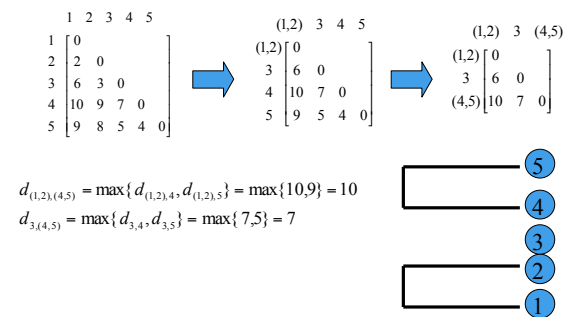
- cluster similarity = similarity of two **least** similar members



### Example: complete link

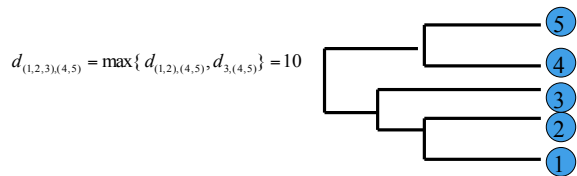


### Example: complete link



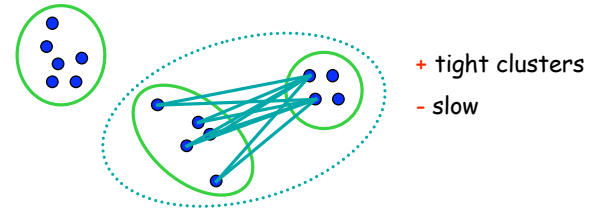
### Example: complete link

$$\begin{array}{c} 1 & 2 & 3 & 4 & 5 \\ \begin{bmatrix} 0 & & & & \\ 2 & 0 & & & \\ 6 & 3 & 0 & & \\ 10 & 9 & 7 & 0 & \\ 9 & 8 & 5 & 4 & 0 \end{bmatrix} & \rightarrow & \begin{array}{c} (1,2) & 3 & 4 & 5 \\ \begin{bmatrix} 0 & & & \\ 3 & 6 & 0 & \\ 4 & 10 & 7 & 0 \\ 5 & 9 & 5 & 4 & 0 \end{bmatrix} & \rightarrow & \begin{array}{c} (1,2) & 3 & (4,5) \\ \begin{bmatrix} 0 & & \\ 3 & 6 & 0 \\ (4,5) & 10 & 7 & 0 \end{bmatrix} \end{array}
 \end{array}$$



### Hierarchical: Average Link

- cluster similarity = average similarity of all pairs



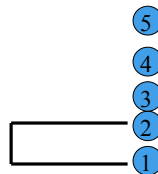
### Example: average link

$$\begin{array}{c} 1 & 2 & 3 & 4 & 5 \\ \begin{bmatrix} 0 & & & & \\ 2 & 0 & & & \\ 6 & 3 & 0 & & \\ 10 & 9 & 7 & 0 & \\ 9 & 8 & 5 & 4 & 0 \end{bmatrix} & \rightarrow & \begin{array}{c} (1,2) & 3 & 4 & 5 \\ \begin{bmatrix} 0 & & & \\ 3 & 4.5 & 0 & \\ 4 & 9.5 & 7 & 0 \\ 5 & 8.5 & 5 & 4 & 0 \end{bmatrix} \end{array}
 \end{array}$$

$$d_{(1,2),3} = \frac{1}{2}(d_{1,3} + d_{2,3}) = \frac{6+3}{2} = 4.5$$

$$d_{(1,2),4} = \frac{1}{2}(d_{1,4} + d_{2,4}) = \frac{10+9}{2} = 9.5$$

$$d_{(1,2),5} = \frac{1}{2}(d_{1,5} + d_{2,5}) = \frac{9+8}{2} = 8.5$$

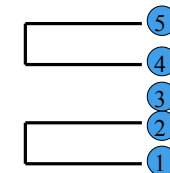


### Example: average link

$$\begin{array}{c} 1 & 2 & 3 & 4 & 5 \\ \begin{bmatrix} 0 & & & & \\ 2 & 0 & & & \\ 6 & 3 & 0 & & \\ 10 & 9 & 7 & 0 & \\ 9 & 8 & 5 & 4 & 0 \end{bmatrix} & \rightarrow & \begin{array}{c} (1,2) & 3 & 4 & 5 \\ \begin{bmatrix} 0 & & & \\ 3 & 4.5 & 0 & \\ 4 & 9.5 & 7 & 0 \\ 5 & 8.5 & 5 & 4 & 0 \end{bmatrix} & \rightarrow & \begin{array}{c} (1,2) & 3 & (4,5) \\ \begin{bmatrix} 0 & & \\ 3 & 4.5 & 0 \\ (4,5) & 9 & 6 & 0 \end{bmatrix} \end{array}
 \end{array}$$

$$d_{(1,2),(4,5)} = \frac{1}{4}(d_{1,4} + d_{1,5} + d_{2,4} + d_{2,5}) = 9$$

$$d_{3,(4,5)} = \frac{1}{2}(d_{3,4} + d_{3,5}) = 6$$



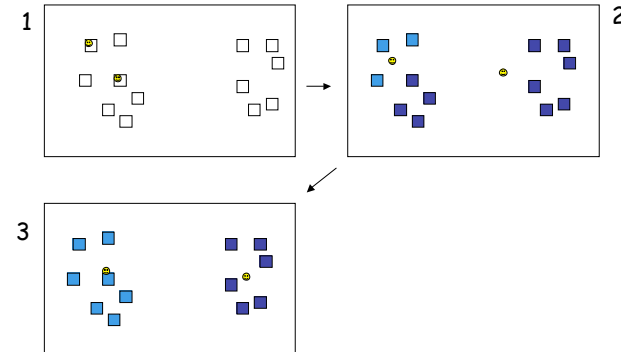


## Hierarchical divisive clustering algorithms

- Top down
  - Start with all the objects in one cluster
  - Successively split into smaller clusters
- Tend to be less efficient than agglomerative
- Resolver implemented a deterministic annealing approach from [Alon et al. 1999]

## Partitional: K-Means

[MacQueen 1965]



## Details of k-means

- Iterate until converge:
  - Assign each data point to the closest centroid
  - Compute new centroid

### Objective function:

Minimize

$$\sum_x (x - \text{Centroid}(\text{Cluster}(x)))^2$$

## Properties of k-means

- Fast
- Proved to converge to local optimum
- In practice, converge quickly
- Tend to produce spherical, equal-sized clusters
- Related to the model-based approach (next lecture)



## Summary

- Definition of clustering
- Pairwise similarity:
  - Correlation
  - Euclidean distance
- Clustering algorithms:
  - Hierarchical (single-, complete-, average-, centroid-link)
  - K-means
  - SOM
- Different clustering algorithms → different clusters

## Misc Notes

- Greedy algorithms. Can get trapped in local minima. Can be sensitive to addition of new points, order of points,...
- + simple, intuitive algorithms, reasonably fast, ok on simple data, no obvious preconception about structure
- no model of structure; biases unclear