# Clustering in Bioinformatics

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# OVERVIEW

☐ Define the clustering problem

☐ Motivation: gene expression and microarrays

☐ Types of clustering

☐ Clustering algorithms

☐ Other applications of clustering

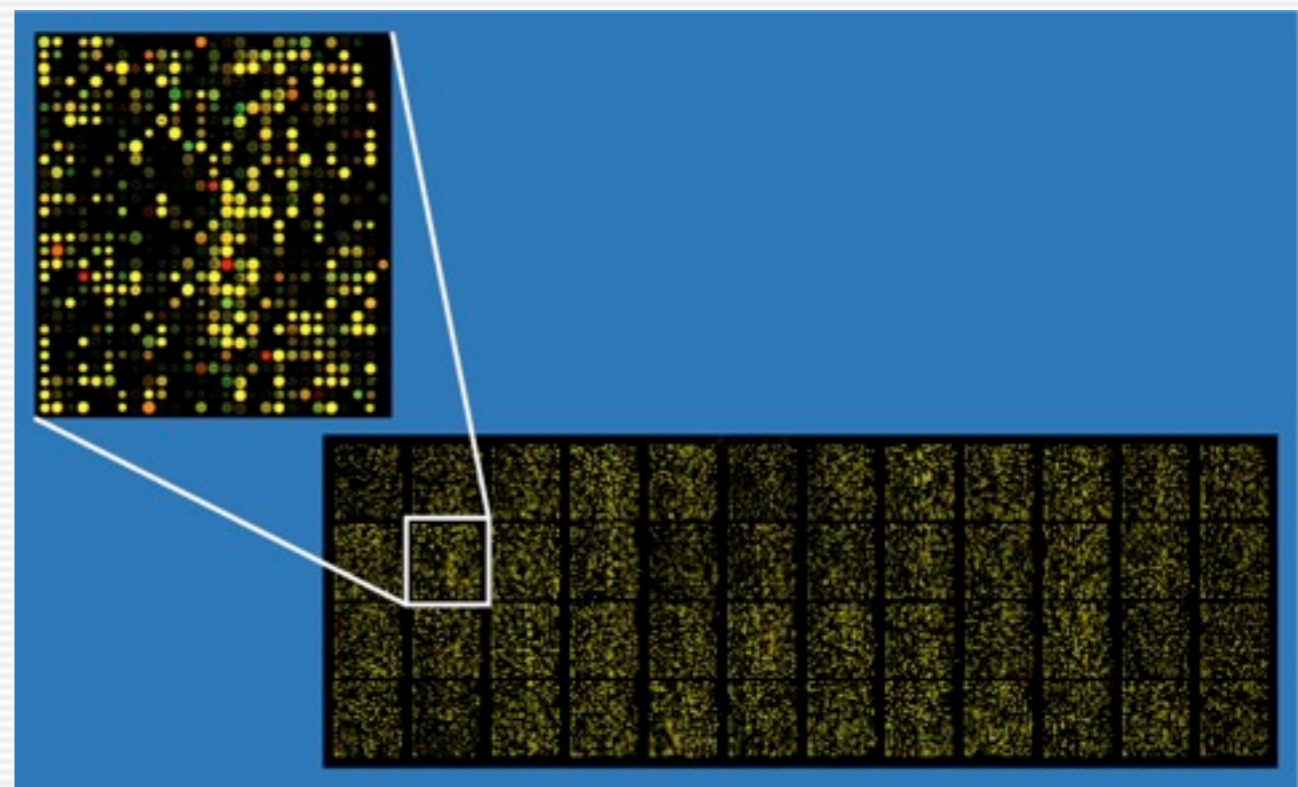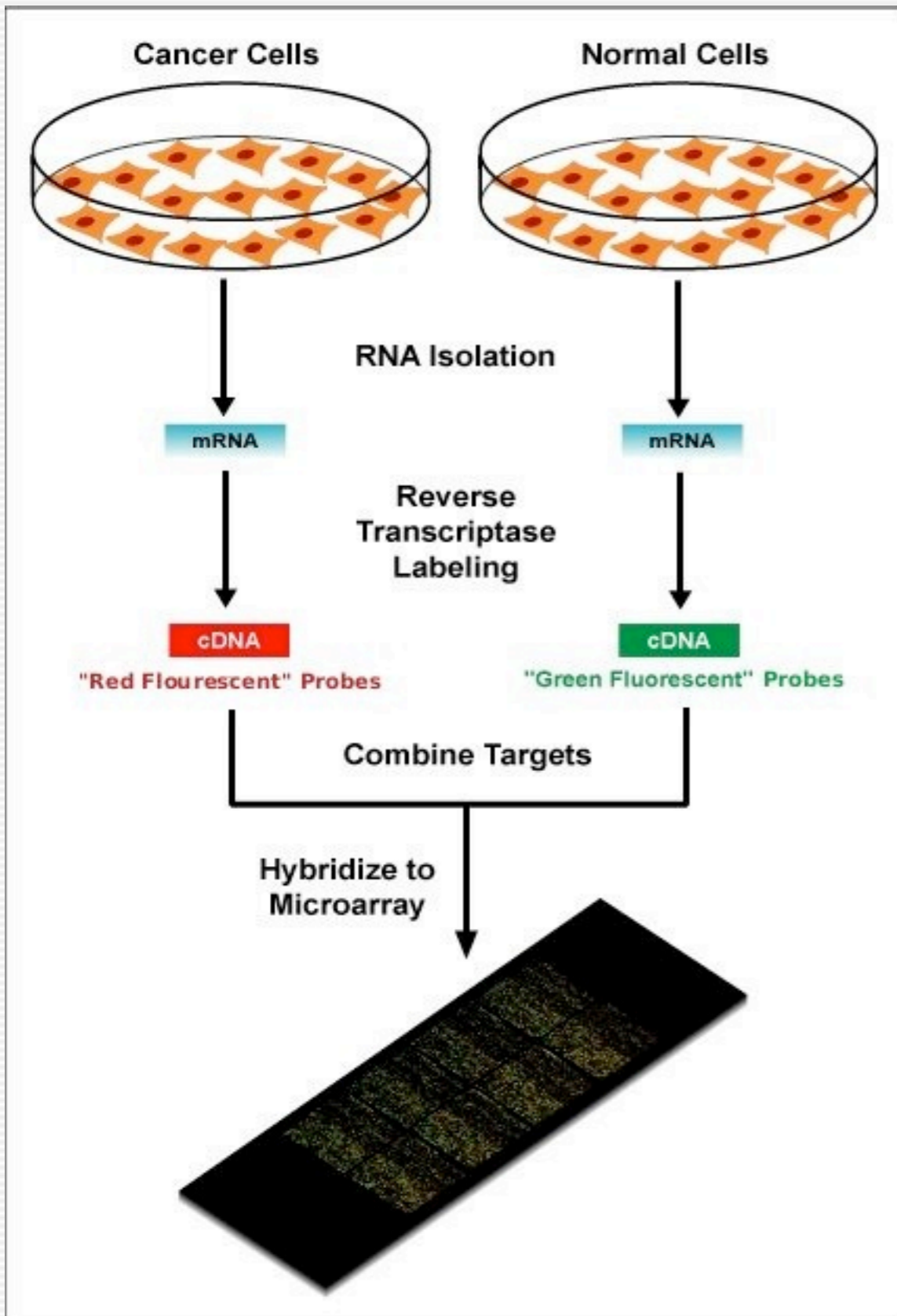SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# THE CLUSTERING PROBLEM

- Motivation: Find patterns in a sea of data

- Input:

  - A (large) number of datapoints: **N**

  - A measure of distance between any two data points $\mathbf{d_{ij}}$

- Output:

  - Groupings (**clustering**) of the elements into **K** (the number can be user-specified or automatically determined) 'similarity' classes

  - Sometimes there is also an objective measure that the obtained clustering seeks to minimize.

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# A MARQUEE APPLICATION: MICROARRAY ANALYSIS

- ☐ What do newly sequenced genes do?

- ☐ Simply comparing the new gene sequences to known DNA sequences often does not necessarily reveal the function of a gene: for 40% of sequenced genes, functionality cannot be ascertained by only comparing to sequences of other known genes

- ☐ Genes that perform similar or complementary function to known genes (reference) will be expressed (transcribed) at high levels together with known genes

- ☐ Genes that perform antagonistic functions (e.g. down-regulation) may be expressed at high levels at an earlier or later time point when compared to known genes

- ☐ E.g. what happens to gene expression in cancer cells?

- ☐ Expression level is estimated by measuring the amount of mRNA for that particular gene

  - ☐ A gene is active if it is being transcribed

  - ☐ More mRNA usually indicates more gene activity

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

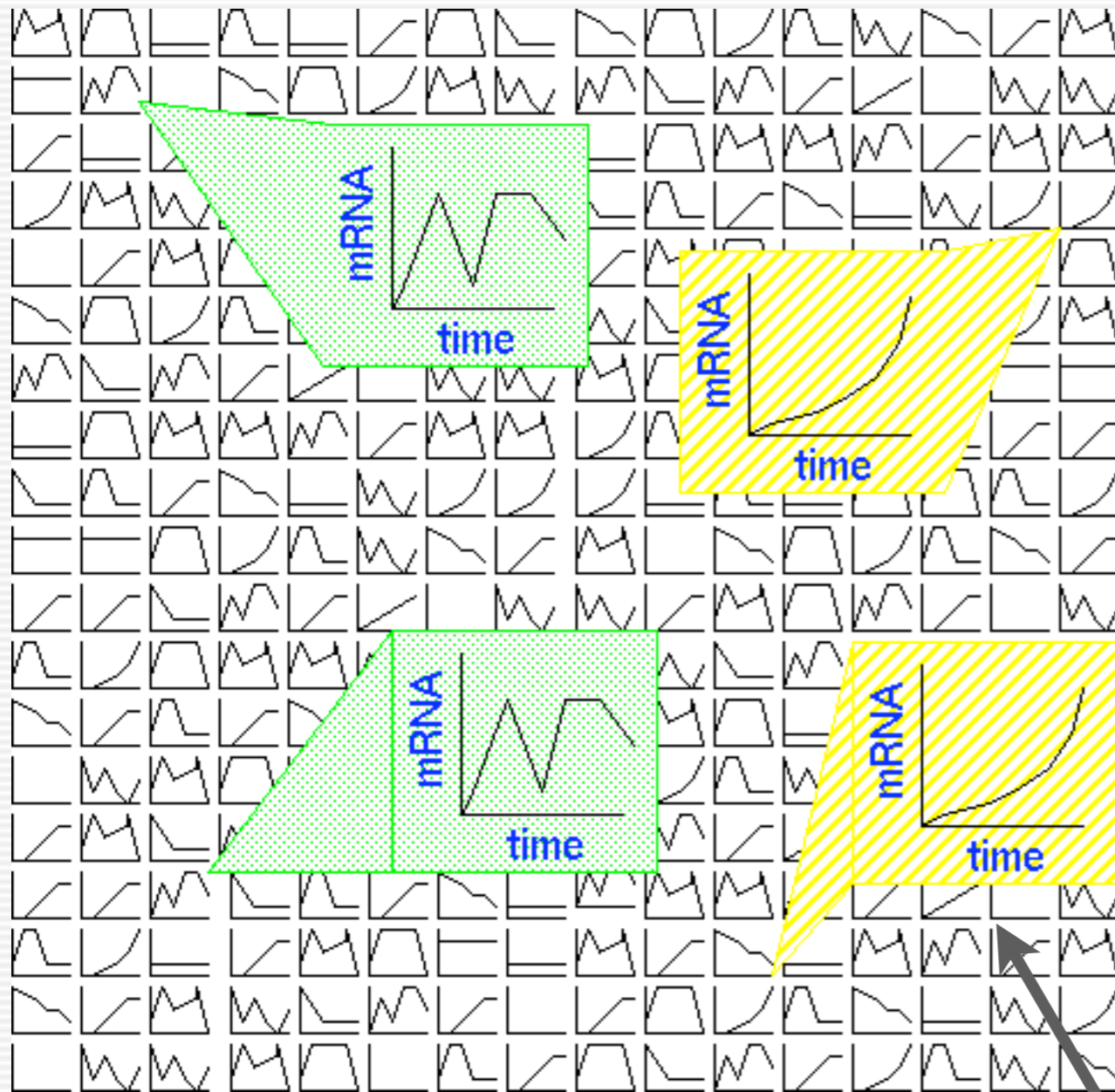# A MICROARRAY EXPERIMENT

- ☐ Produce cDNA from mRNA (cDNA is more stable)

- ☐ Label cDNA with a fluorescent dye or biotin for detection

- ☐ Different color labels are available to compare many samples at once

- ☐ Wash cDNA over the microarray containing thousands of high density **probes** that hybridize to complementary strands in the sample and immobilize them on the surface.

- ☐ For biotin-labeled samples, stain with the biotin-specific fluorescently labeled antibody

- ☐ Read the microarray, using a laser or a high-resolution CCD

- ☐ Illumination reveals transcribed/co-expressed genes

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

- ☐ Green: expressed only in control

- ☐ Red: expressed only in an experimental cell

- ☐ Yellow: equally expressed in both samples

- ☐ Black: NOT expressed in either control or sample

☐ Track the sample over a period of time to observe changes in gene expression over time

☐ Track two samples under the same conditions to look for **differential expression**

Each box represents one gene's expression over time

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# MICROARRAY DATA

- Microarray data are usually transformed into a (relative, normalized) intensity matrix

- Can also be represented as a bit matrix ($\log_2$ of relative intensity)

- The intensity matrix allows biologists to infer correlations between different genes (even if they are dissimilar) and to understand how genes functions might be related

- Care must be taken to normalize the data appropriately, e.g. different time points can come from different arrays.

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

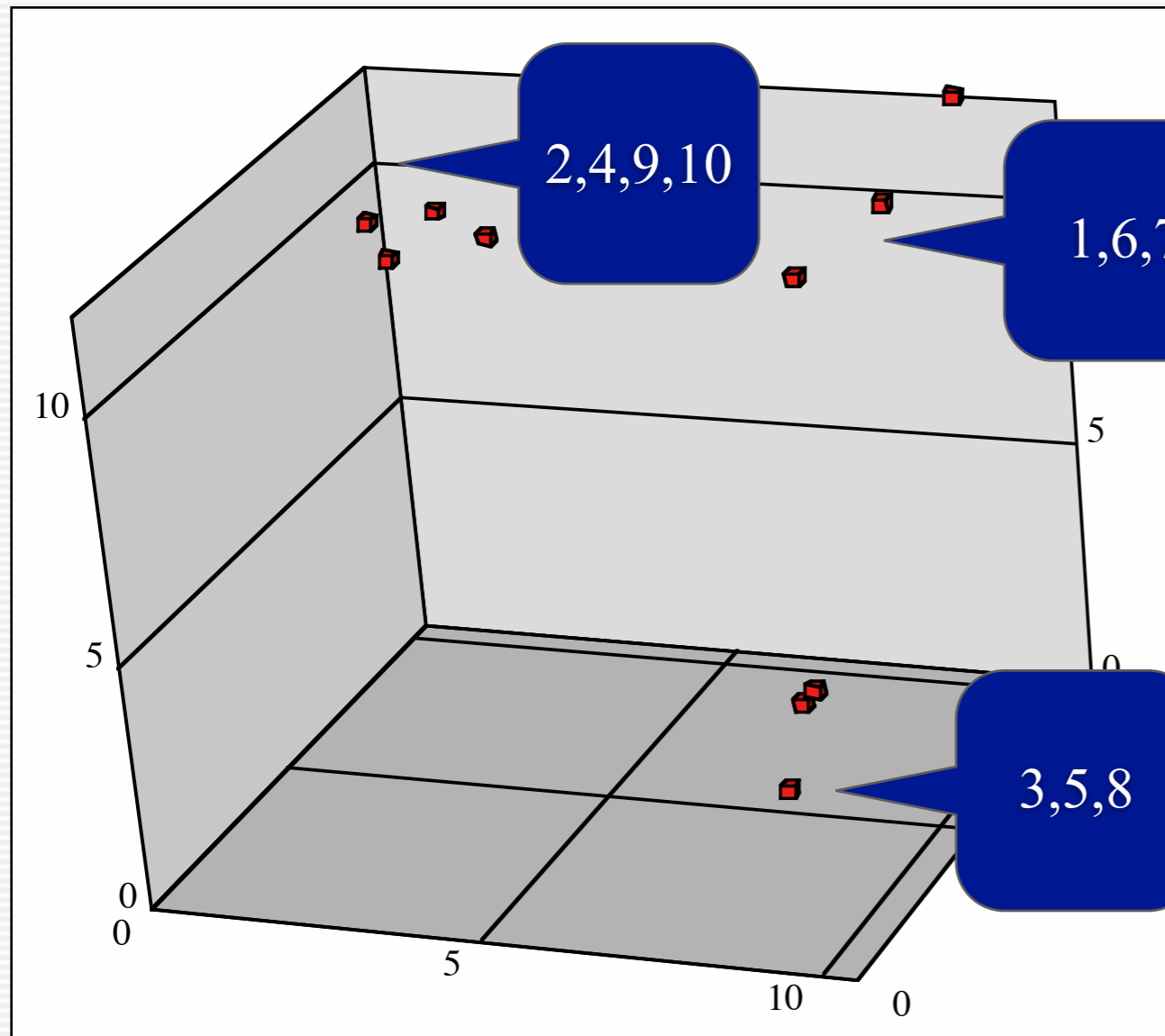| Gene | Time 1 | Time 2 | Time 3 |
|------|--------|--------|--------|
| 1 | 10 | 8 | 10 |
| 2 | 10 | 0 | 9 |
| 3 | 4 | 8.5 | 3 |
| 4 | 9.5 | 0.5 | 8.5 |
| 5 | 4.5 | 8.5 | 3 |
| 6 | 10.5 | 9 | 12 |
| 7 | 5 | 8.5 | 11 |
| 8 | 2.7 | 8.7 | 2 |
| 9 | 9.7 | 2 | 9 |
| 10 | 10.2 | 1 | 9.2 |

## INTENSITY TABLE

- Which genes are similar?

- What defines co-expression?

- How to measure the distance/similarity?

## EUCLIDEAN DISTANCE IN D-DIMENSIONS

$$D(x,y) = \sqrt{\sum_{i=1}^{d}(x_i - y_i)^2}$$

# Finding similar genes



## Pairwise distances

|    | 1    | 2    | 3    | 4    | 5    | 6    | 7    | 8    | 9    | 10   |
|----|------|------|------|------|------|------|------|------|------|------|
| 1  |      | 8.1  | 9.2  | 7.7  | 8.9  | 2.3  | 5.1  | 10.9 | 6.1  | 7.0  |
| 2  | 8.1  |      | 12.0 | 0.9  | 11.8 | 9.5  | 10.1 | 13.3 | 2.0  | 1.0  |
| 3  | 9.2  | 12.0 |      | 11.2 | 0.5  | 11.1 | 8.1  | 1.7  | 10.5 | 11.5 |
| 4  | 7.7  | 0.9  | 11.2 |      | 10.9 | 9.2  | 9.5  | 12.5 | 1.6  | 1.1  |
| 5  | 8.9  | 11.8 | 0.5  | 10.9 |      | 10.8 | 8.0  | 2.1  | 10.3 | 11.3 |
| 6  | 2.3  | 9.5  | 11.1 | 9.2  | 10.8 |      | 5.6  | 12.7 | 7.7  | 8.5  |
| 7  | 5.1  | 10.1 | 8.1  | 9.5  | 8.0  | 5.6  |      | 9.3  | 8.3  | 9.3  |
| 8  | 10.9 | 13.3 | 1.7  | 12.5 | 2.1  | 12.7 | 9.3  |      | 12.0 | 12.9 |
| 9  | 6.1  | 2.0  | 10.5 | 1.6  | 10.3 | 7.7  | 8.3  | 12.0 |      | 1.1  |
| 10 | 7.0  | 1.0  | 11.5 | 1.1  | 11.3 | 8.5  | 9.3  | 12.9 | 1.1  |      |

## Rearranged distances

|    | 1    | 6    | 7    | 2    | 4    | 9    | 10   | 3    | 5    | 8    |
|----|------|------|------|------|------|------|------|------|------|------|
| 1  | 0.0  | 2.3  | 5.1  | 8.1  | 7.7  | 6.1  | 7.0  | 9.2  | 8.9  | 10.9 |
| 6  | 2.3  | 0.0  | 5.6  | 9.5  | 9.2  | 7.7  | 8.5  | 11.1 | 10.8 | 12.7 |
| 7  | 5.1  | 5.6  | 0.0  | 10.1 | 9.5  | 8.3  | 9.3  | 8.1  | 8.0  | 9.3  |
| 2  | 8.1  | 9.5  | 10.1 | 0.0  | 0.9  | 2.0  | 1.0  | 12.0 | 11.8 | 13.3 |
| 4  | 7.7  | 9.2  | 9.5  | 0.9  | 0.0  | 1.6  | 1.1  | 11.2 | 10.9 | 12.5 |
| 9  | 6.1  | 7.7  | 8.3  | 2.0  | 1.6  | 0.0  | 1.1  | 10.5 | 10.3 | 12.0 |
| 10 | 7.0  | 8.5  | 9.3  | 1.0  | 1.1  | 1.1  | 0.0  | 11.5 | 11.3 | 12.9 |
| 3  | 9.2  | 11.1 | 8.1  | 12.0 | 11.2 | 10.5 | 11.5 | 0.0  | 0.5  | 1.7  |
| 6  | 8.9  | 10.8 | 8.0  | 11.8 | 10.9 | 10.3 | 11.3 | 0.5  | 0.0  | 2.1  |
| 8  | 10.9 | 12.7 | 9.3  | 13.3 | 12.5 | 12.0 | 12.9 | 1.7  | 2.1  | 0.0  |

# CLUSTERING PRINCIPLES

- ☐ **Homogeneity**: elements of the same cluster are maximally close to each other

- ☐ **Separation**: elements in separate clusters are maximally far apart from each other

- ☐ One is actually implied by the other (in many cases)
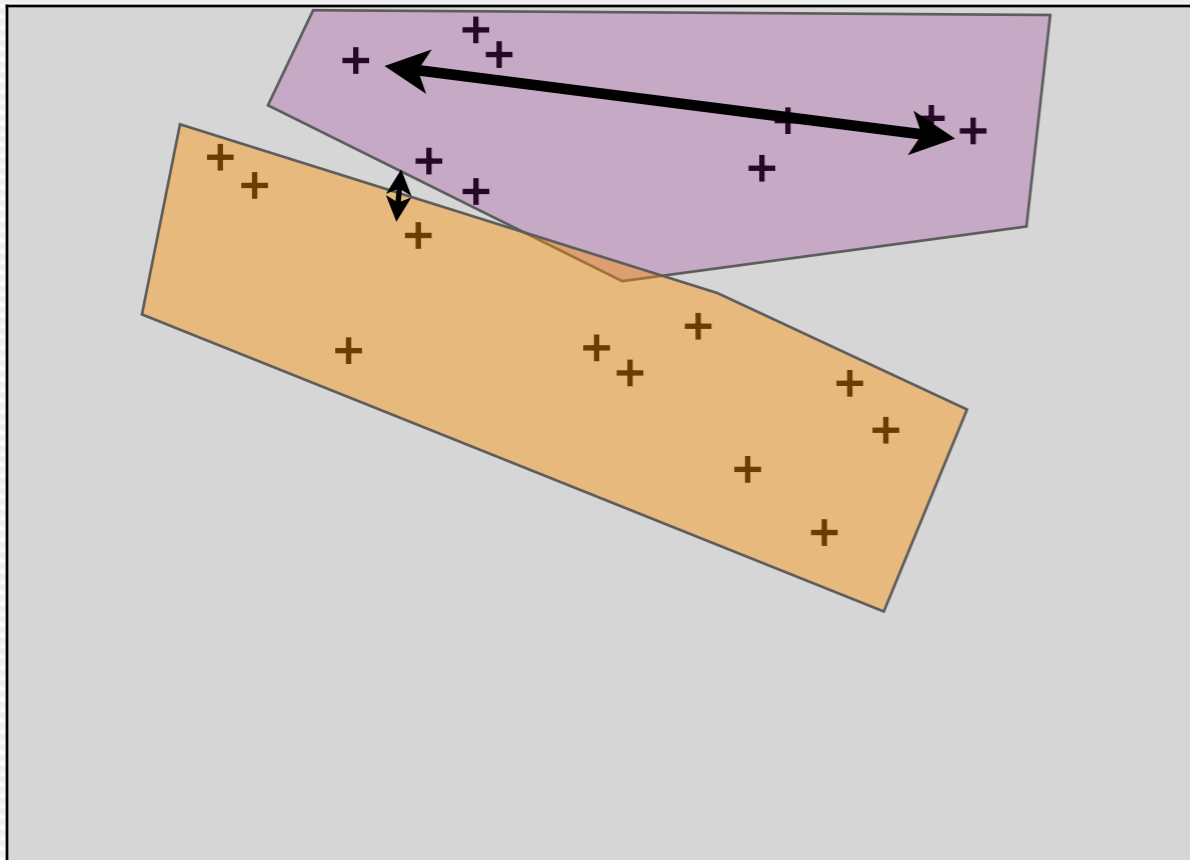
- ☐ Generally speaking, this is a hard problem.

$$\min_{\text{clustering}} \left[ \alpha \sum_{x,y \in \text{the same cluster}} d(x,y) - \beta \sum_{x,y \in \text{different clusters}} d(x,y) \right]$$

RELATIVE IMPORTANCE

# Because

$$\sum_{x,y \in \text{the same cluster}} d(x,y) + \sum_{x,y \in \text{different clusters}} d(x,y)$$

$$= \sum_{x,y} d(x,y) = D = const$$

# We can simplify

$$\min_{\text{clustering}} \left[ \alpha \sum_{x,y \in \text{the same cluster}} d(x,y) - \beta \sum_{x,y \in \text{different clusters}} d(x,y) \right]$$

# To an equivalent expression that only depends on intRA-cluster distances

$$(\alpha + \beta) \min_{\text{clustering}} \left[ \sum_{x,y \in \text{ the same cluster}} d(x,y) \right] - \beta D$$

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# Poor clustering example



- This clustering violates both principles:

  - Points in the same cluster are far apart

  - Points in different cluster are close together

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# BETTER CLUSTERING EXAMPLE



☐ This clustering appears sensible.

☐ But we need to use an **objective** metric to optimize cluster assignment.

# CLUSTERING TECHNIQUES

☐ **Agglomerative**: Start with every element in its own cluster, and iteratively join clusters together

☐ **Divisive**: Start with one cluster and iteratively divide it into smaller clusters

☐ **Hierarchical**: Organize elements into a tree, leaves represent genes and the length of the paths between leaves represents the distances between genes. Similar genes lie within the same subtrees

☐ Generally, finding the exact solution to a clustering problem is NP hard.

# K-MEANS CLUSTERING

- A technique to partition a set of **N** points into **K** clusters

- Each cluster is represented with a **mean** (a centroid) – hence 'K-means'

- **Input**: A set **V** with **N** points ($v_1$, $v_2$ ... $v_n$), the desired number of clusters **K** and a distance measure between any two points **d(v,w)**

- **Output**: A set **X** of **K** cluster centers that minimize the <u>squared error distortion</u> D(V,X) over all possible choices of X.

$$D(V, X) = \frac{1}{N} \sum_{i=1}^{N} \min_{k} d^2(v_i, x_k)$$

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# K-MEANS CLUSTERING

☐ For K=1, the problem becomes trivial: the centroid of all points is the solution for Euclidean distances.

$$x = \frac{1}{N} \sum_i v_i$$

☐ For K≥2 the problem becomes NP-complete

☐ An efficient heuristic exists

☐ **Lloyd's algorithm**.

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# LLOYD'S ALGORITHM

1. Arbitrarily assign the **K** cluster centers (this can significantly influence the outcome)

2. **while** cluster centers keep changing

   A. Compute the distance from each data point to the current cluster center $C_i$ ($1 \leq i \leq K$) and assign the point to the nearest cluster

   B. After the assignment of all data points, compute new centers for each cluster by taking the centroid of all the points in that cluster

3. Output cluster centers and assignments

# K-means execution example

## Step 1
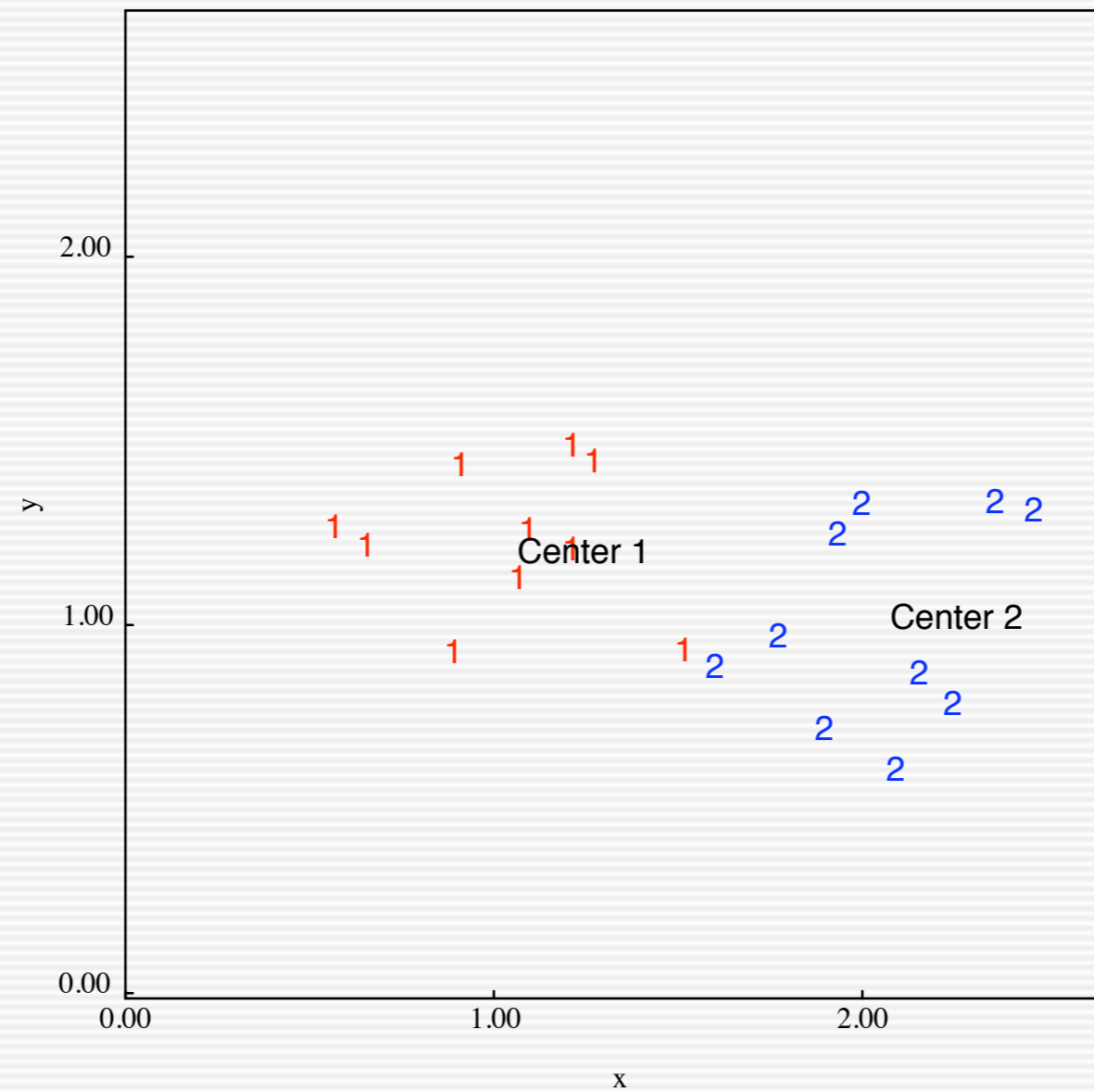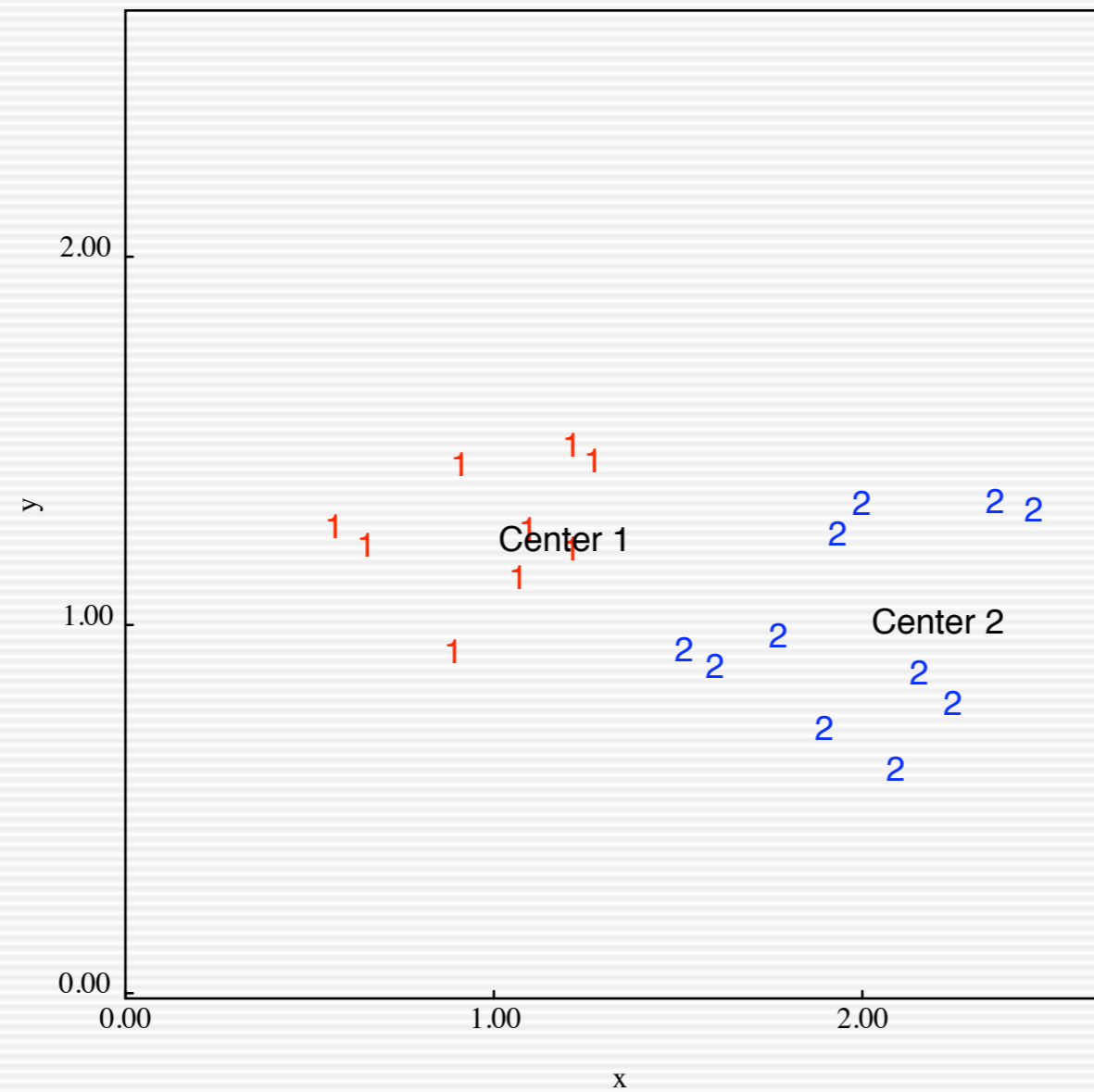
K-MEANS EXECUTION EXAMPLE

STEP 2

# K-MEANS EXECUTION EXAMPLE

## STEP 3

# K-MEANS EXECUTION EXAMPLE

## STEP 4

# K-means execution example

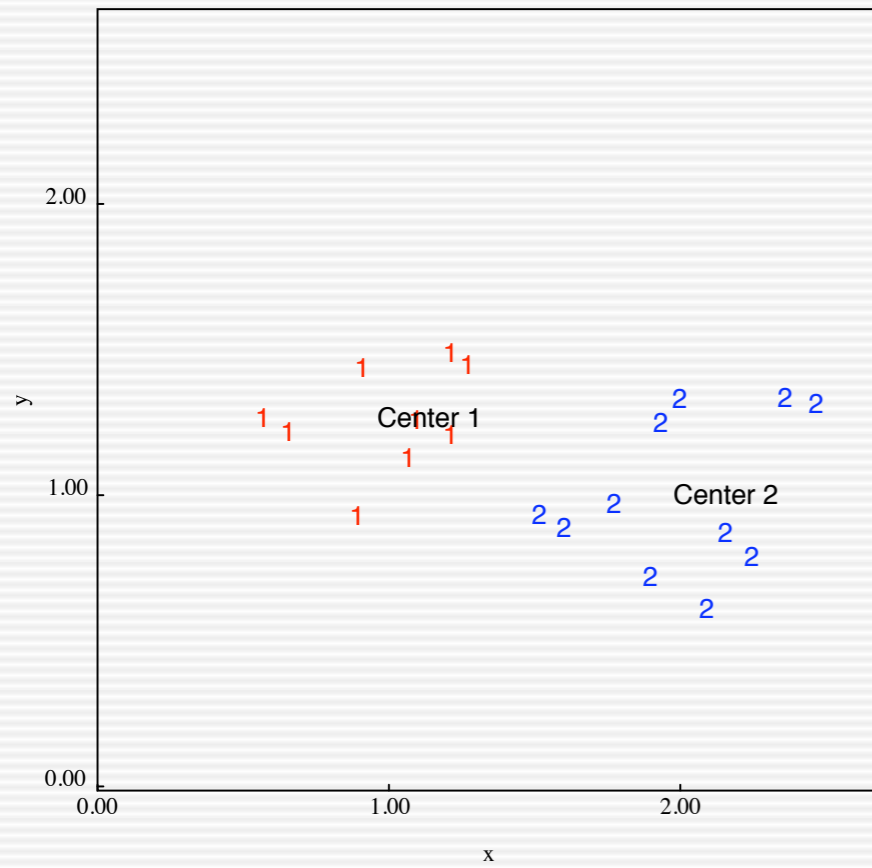## Step 5

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# K-means execution example

## Step 6

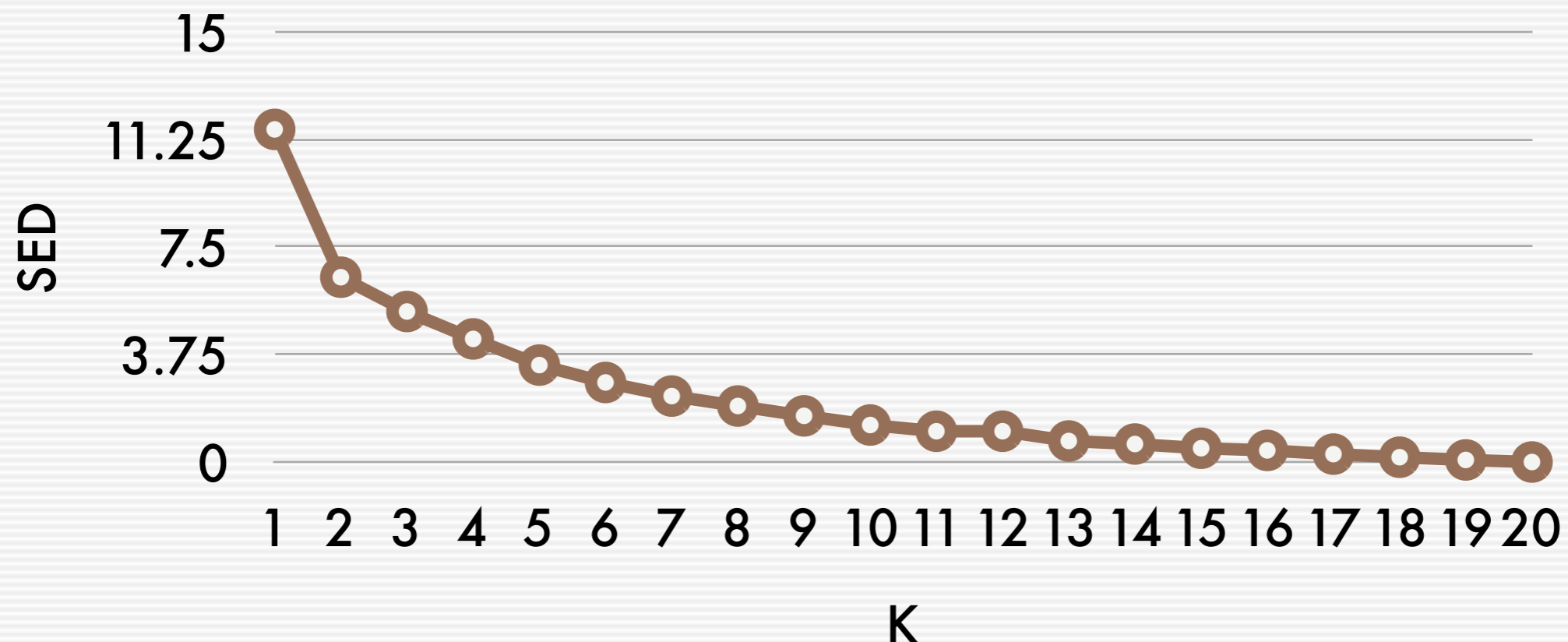Sergei L Kosakovsky Pond [spond@ucsd.edu]

# K-MEANS EXECUTION EXAMPLE

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# HOW TO CHOOSE K?

☐ The simplest approach is to start with K=1 and increase K until the squared error distortion (SED) stops decreasing

 ☐ The problems is that K=N always achieves the value of **0** (each point is a cluster), so we always keep increasing K.

 ☐ Generally, need to add further constraints (e.g. model complexity) to obtain non-trivial results

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# Conservative K-Means Algorithm

☐ Lloyd algorithm is fast but in each iteration it moves many data points, not necessarily causing better convergence.

☐ A more conservative method would be to move one point at a time only if it improves the overall **clustering cost**

  ☐ The smaller the clustering cost of a partition of data points is the better that clustering is

  ☐ Different methods (e.g. the squared error distortion) can be used to measure this clustering cost

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# K-Means "Greedy" Algorithm

ProgressiveGreedyK-Means(k)
Select an arbitrary partition P into k clusters
**while** forever
    bestChange ← 0
    **for** every cluster C
        **for** every element i not in C
            **if** $\text{cost}(P) - \text{cost}(P_{i \to C}) > \text{bestChange}$
                $\text{bestChange} \leftarrow \text{cost}(P) - \text{cost}(P_{i \to C})$
                $i^* \leftarrow I$
                $C^* \leftarrow C$
    if bestChange > 0
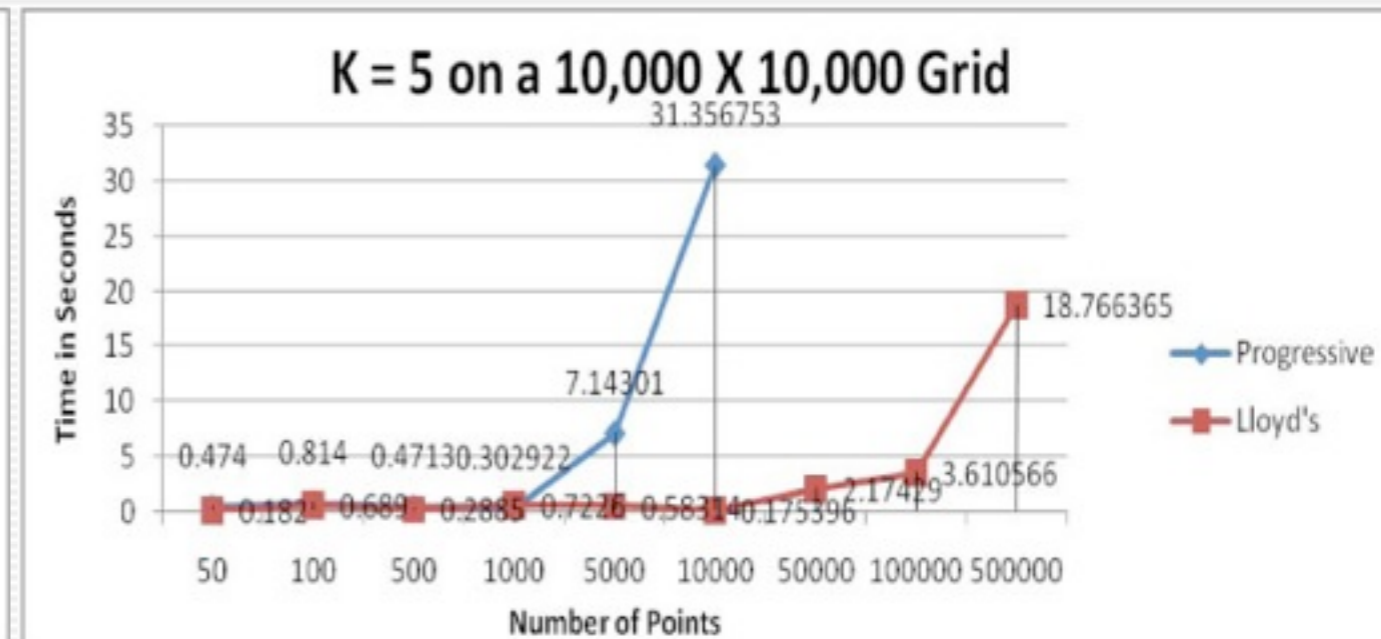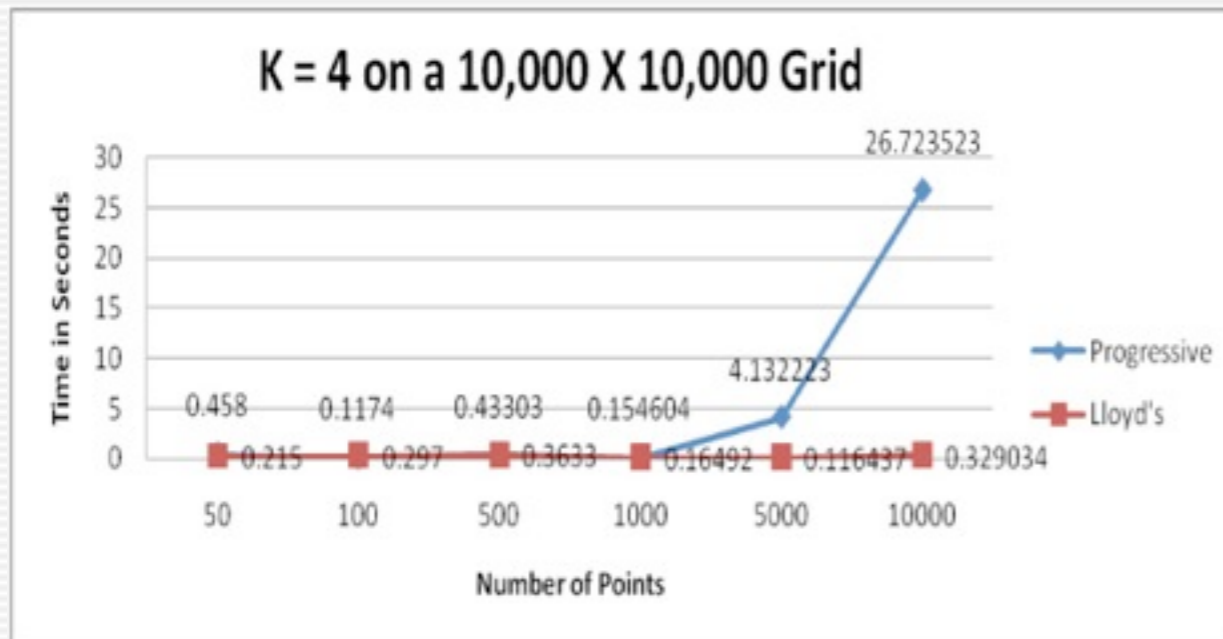        Change partition **P** by moving $i^*$ to $C^*$
    else
      return **P**

# BMC Bioinformatics

**BioMed** Central

Research

**Open Access**

## A practical comparison of two K-Means clustering algorithms

Gregory A Wilkin[1] and Xiuzhen Huang*[2]

K = 4 on a 10,000 X 10,000 Grid

K = 5 on a 10,000 X 10,000 Grid

CONCLUSION: LLOYD'S IN MORE EFFICIENT, BOTH IN RUN-TIME AND IN BEST FOUND SED

# Distance Measures in DNA Microarray Data Analysis.

## R. Gentleman, B. Ding, S. Dudoit, and J. Ibrahim

☐ Euclidean distance is not necessarily the best measure for co-expression.

---

**MAN** Manhattan metric

$$d_{man}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{m} |x_i - y_i|.$$

---

**SPEAR** Spearman sample correlation distance

$$d_{spear}(\mathbf{x}, \mathbf{y}) = 1 - \frac{\sum_{i=1}^{m} (x_i' - \overline{x'})(y_i' - \overline{y'})}{\sqrt{\sum_{i=1}^{m} (x_i' - \overline{x'})^2 \sum_{i=1}^{m} (y_i' - \overline{y'})^2}}.$$

where $x_i' = rank(x_i)$ and $y_i' = rank(y_i)$.

**TAU** Kendall's $\tau$ sample correlation

$$d_{tau}(\mathbf{x}, \mathbf{y}) = 1 - \tau(\mathbf{x}, \mathbf{y})| = 1 - \frac{|\sum_{i=1}^{m} \sum_{j=1}^{m} C_{x_{ij}} C_{y_{ij}}|}{m(m-1)}$$

where $C_{x_{ij}} = sign(x_i - x_j)$ and $C_{y_{ij}} = sign(y_i - y_j)$.

---

**COR** Pearson sample correlation distance

$$d_{cor}(\mathbf{x}, \mathbf{y}) = 1 - r(\mathbf{x}, \mathbf{y}) = 1 - \frac{\sum_{i=1}^{m} (x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum_{i=1}^{m} (x_i - \overline{x})^2 \sum_{i=1}^{m} (y_i - \overline{y})^2}}.$$
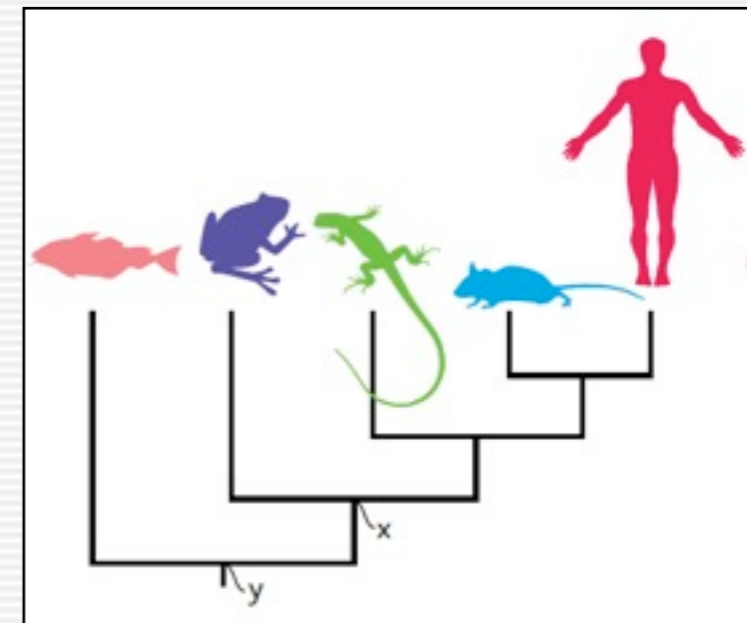
**EISEN** Cosine correlation distance

$$d_{eisen}(\mathbf{x}, \mathbf{y}) = 1 - \frac{\mathbf{x}'\mathbf{y}}{\|\mathbf{x}\|\|\mathbf{y}\|} = 1 - \frac{|\sum_{i=1}^{m} x_i y_i|}{\sqrt{\sum_{i=1}^{m} x_i^2 \sum_{i=1}^{m} y_i^2}}$$

which is a special case of Pearson's correlation with $\overline{x}$ and $\overline{y}$ both replaced by zero.

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# Hierarchical clustering

☐ Instead of grouping into discrete clusters, produces a 'classification' tree, also called a <u>dendrogram</u>

☐ A more intuitive example is probably obtained from molecular sequence data (an early example of clustering applications)

☐ We have a collection of aligned nucleotide sequences from different species, and wish to construct their evolutionary hierarchy/history – a phylogeny.

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# Hierarchical clustering

☐ Consider the following distance matrix on 5 nucleotide (partial mitochondrial genome) sequences. The values are *p-distances* defined as the number of nucleotide differences normalized by the length of the sequence.
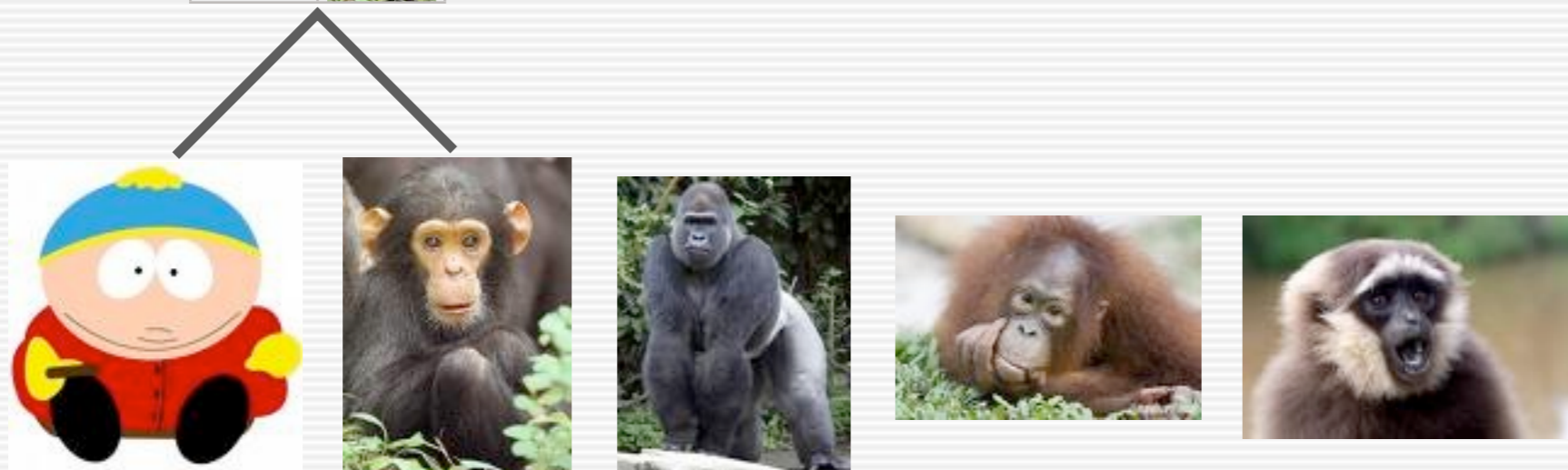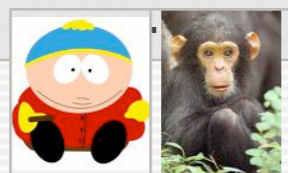
| | Human | Chimpanzee | Gorilla | Orangutan | Gibbon |
|---|---|---|---|---|---|
| Human | - | 0.0882682 | 0.102793 | 0.159598 | 0.179688 |
| Chimpanzee | - | - | 0.106145 | 0.170759 | 0.1875 |
| Gorilla | - | - | - | 0.166295 | 0.1875 |
| Orangutan | - | - | - | - | 0.188616 |
| Gibbon | - | - | - | - | - |

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# Clustering procedure

- ☐ At each step, we select the two closest sequences and **join** them to form a clade.

- ☐ We then replace the two just joined sequences with their ancestor

- ☐ This reduces the size of the data matrix by one

- ☐ We need to compute the distances from the new ancestor to the remaining sequences

|  | Human | Chimpanzee | Gorilla | Orangutan | Gibbon |
|---|---|---|---|---|---|
| Human | - | **0.0882682** | 0.102793 | 0.159598 | 0.179688 |
| Chimpanzee | - | - | 0.106145 | 0.170759 | 0.1875 |
| Gorilla | - | - | - | 0.166295 | 0.1875 |
| Orangutan | - | - | - | - | 0.188616 |
| Gibbon | - | - | - | - | - |

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# UPDATING DISTANCES

□ There are multiple strategies for computing the distances to the new 'ancestral' sequence **a** that joins sequences **m** and **n**

| | |
|---|---|
| Single Linkage | $d(x,a) = \min\left[d(x,m), d(x,n)\right]$ |
| Complete Linkage | $d(x,a) = \max\left[d(x,m), d(x,n)\right]$ |
| UPGMA <small>Unweighted Pair Group Method with Arithmetic Mean</small> | $d(x,a) = \dfrac{d(x,m) + d(x,n)}{2}$ |
| WPGMA <small>Weighted Pair Group Method with Arithmetic Mean</small> | $d(x,a) = \dfrac{s(m)d(x,m) + s(n)d(x,n)}{s(m) + s(n)}$ |

s(**n**) counts the number of actual sequences represented by node **n**.

Sergei L Kosakovsky Pond [spond@ucsd.edu]

# EXAMPLE CONTINUED

☐ Use complete linkage. Joining human and chimp...

| | Human | Chimpanzee | Gorilla | Orangutan | Gibbon |
|---|---|---|---|---|---|
| Human | - | 0.0882682 | 0.102793 | 0.159598 | 0.179688 |
| Chimpanzee | - | - | 0.106145 | 0.170759 | 0.1875 |
| Gorilla | - | - | - | 0.166295 | 0.1875 |
| Orangutan | - | - | - | - | 0.188616 |
| Gibbon | - | - | - | - | - |

⬇

| | Human-Chimpanzee | Gorilla | Orangutan | Gibbon |
|---|---|---|---|---|
| Human-Chimpanzee | - | **0.106145** | **0.170759** | **0.1875** |
| Gorilla | - | - | 0.166295 | 0.1875 |
| Orangutan | - | - | - | 0.188616 |
| Gibbon | - | - | - | - |

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

| | Human-Chimpanzee | Gorilla | Orangutan | Gibbon |
|---|---|---|---|---|
| Human-Chimpanzee | - | 0.106145 | 0.170759 | 0.1875 |
| Gorilla | - | - | 0.166295 | 0.1875 |
| Orangutan | - | - | - | 0.188616 |
| Gibbon | - | - | - | - |



| | Human-Chimpanzee-Gorilla | Orangutan | Gibbon |
|---|---|---|---|
| Human-Chimpanzee-Gorilla | - | 0.170759 | 0.1875 |
| Orangutan | - | - | 0.188616 |
| Gibbon | - | - | - |

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

| | Orangutan | Gibbon |
|---|---|---|
| Human-Chimpanzee-Gorilla | 0.170759 | 0.1875 |
| Orangutan | - | 0.188616 |
| Gibbon | - | - |

| | Gibbon |
|---|---|
| Hum-Chimp-Gor-Orang | 0.188616 |
| Gibbon | - |

SERGEI L KOSAKOVSKY POND [SPOND@UCSD.EDU]

# A note on WPGMA

|  | Gorilla | Orangutan | Gibbon |
|---|---|---|---|
| Human-Chimpanzee | 0.104469 | 0.165179 | 0.183594 |
| Gorilla | - | 0.166295 | 0.1875 |
| Orangutan | - | - | 0.188616 |
| Gibbon | - | - | - |

|  | Orangutan | Gibbon |
|---|---|---|
| Human-Chimpanzee-Gorilla | 0.165551 | 0.184896 |
| Orangutan | - | 0.188616 |
| Gibbon | - | - |

☐ **d(HCG-Orang)** = 1/3 [2 d(HC-Orang) + d (Gor-Orang)]

# BACK TO MICROARRAYS...

☐ Clustering plots can be interpreted as gene/condition hierarchy



HTTP://UPLOAD.WIKIMEDIA.ORG/WIKIPEDIA/COMMONS/4/48/HEATMAP.PNG

# A few other applications

## Clustering of highly homologous sequences to reduce the size of large protein databases

Weizhong Li[1], Lukasz Jaroszewski[2] and Adam Godzik[2,*]

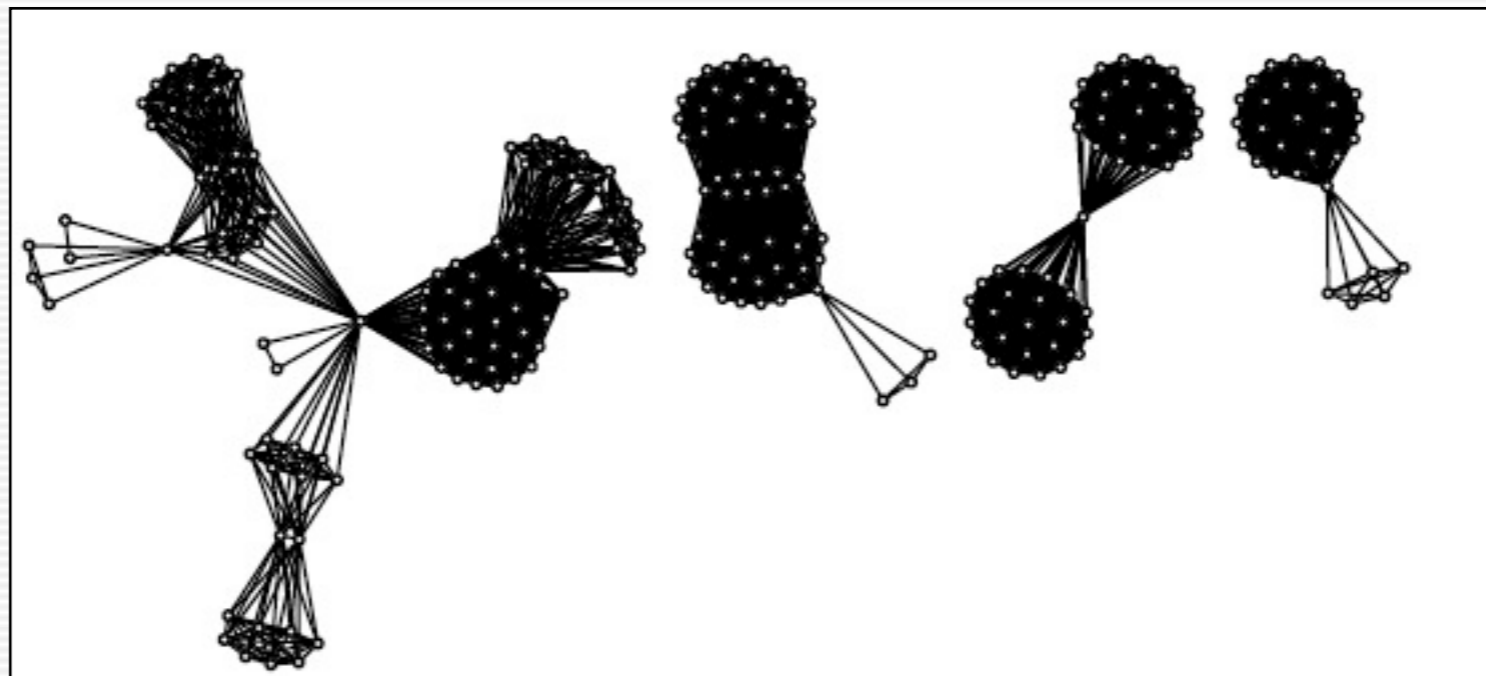[1]San Diego Supercomputer Center, La Jolla, CA 92093, USA and [2]The Burnham Institute, La Jolla, CA 92037, USA

☐ Use clustering of similar sequences in protein databases to reduce complexity and speed up comparisons. Each cluster of similar sequences is represented by a single sequence.

☐ Complexity reduction is an important application of clustering

| Database | Threshold (%) | Word length | Clusters | Time (minute) | Time of nrdb90 |
|---|---|---|---|---|---|
| PDB | 90 | 5 | 4 342 | 0.25 | 8.5 |
| 8732 sequences | 80 | 4 | 3 907 | 0.26 | NA |
| 1 850 235 letters | 75 | 3 | 3 767 | 0.44 | NA |
|  | 65 | 2 | 3 535 | 3.96 | NA |
| SWISS-PROT | 90 | 5 | 71 180 | 6.4 | 208 |
| 88 780 sequences | 80 | 4 | 62 272 | 15.6 | NA |
| 31 984 247 letters | 75 | 3 | 58 651 | 96.4 | NA |
| NR | 90 | 5 | 316 436 | 117 | 2176 |
| 563 276 sequences | 80 | 4 | 264 495 | 325 | NA |
| 177 028 588 letters | 75 | 3 | 247 074 | 1597 | NA |

- The structure of proteins interactions can be represented by a graph

    - Node = proteins, Edges = interactions

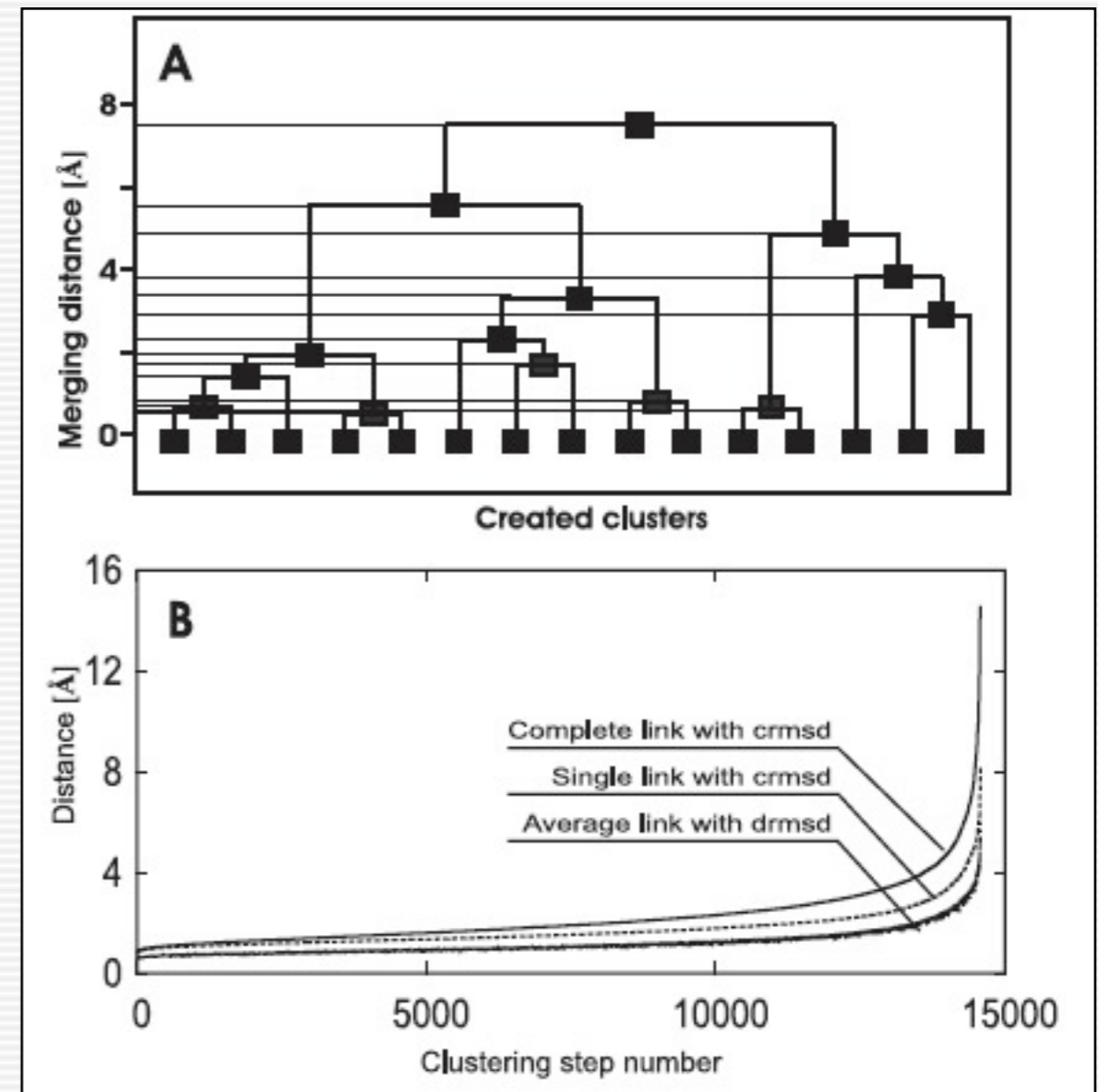    - Look for clusters (densely connected components) in graphs

## HCPM—program for hierarchical clustering of protein models

Dominik Gront* and Andrzej Kolinski

Faculty of Chemistry, Warsaw University, Pasteura 1, 02-093 Warsaw, Poland

☐ Hierarchical clustering to improve protein structure prediction by merging the predictions made by a large number of alternative conformation models
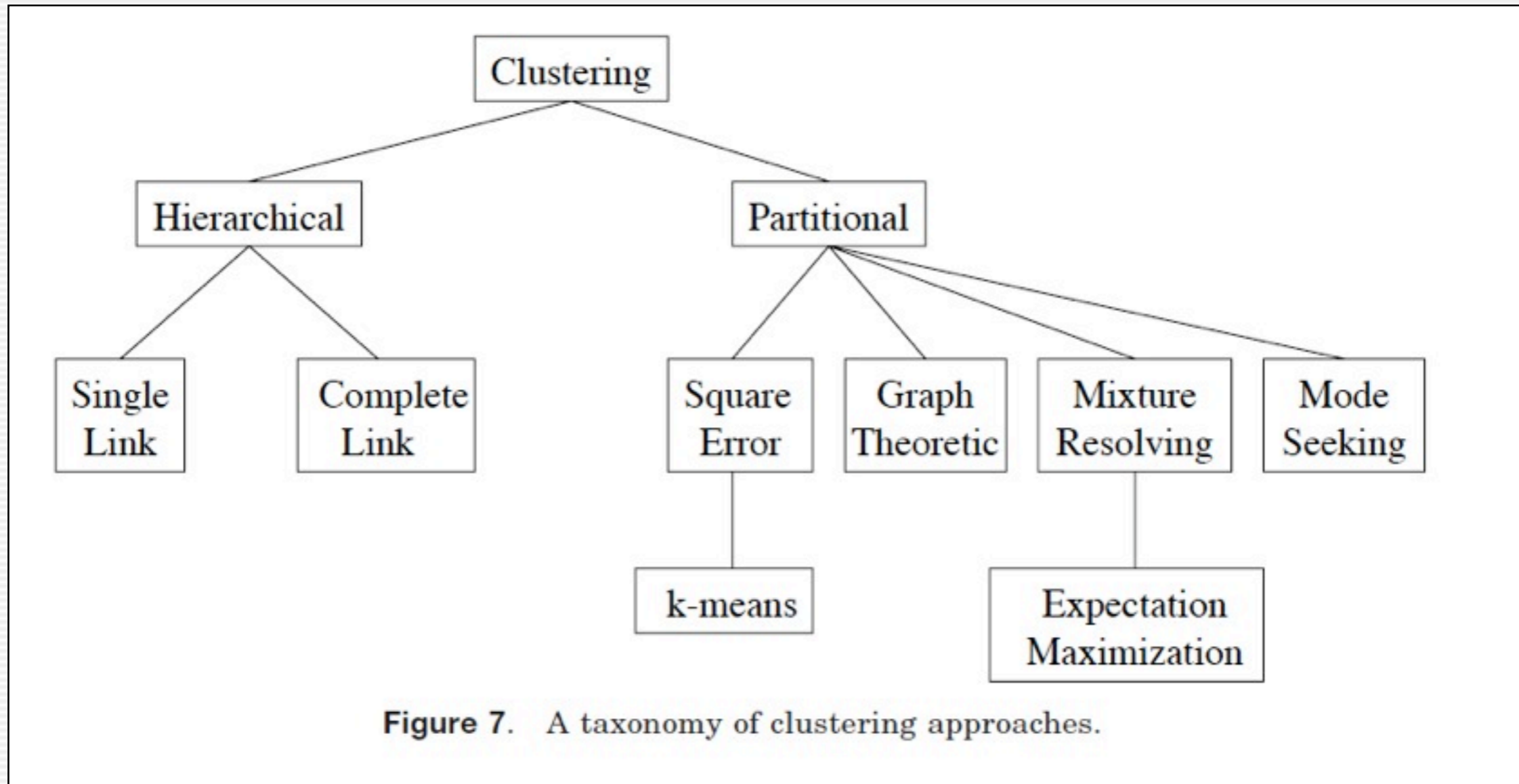
# FURTHER READING...



Figure 7. A taxonomy of clustering approaches.

**Figure 15.** Using the minimal spanning tree to form clusters.

# DEFINES THE CONCEPT OF 'ELEMENT BELONGS TO A PARTITION WITH A PROBABILITY'

# BUILD A MINIMUM SPANNING TREE AND DELETE LONGEST EDGES TO CREATE PARTITIONS



{(1,0.9), (2,0.8), (3,0.7), (4,0.6), (5,0.55), (6,0.2), (7,0.2), (8,0.0), (9,0.0)}

{(1,0.0), (2,0.0), (3,0.0), (4,0.1), (5,0.15), (6,0.4), (7,0.35), (8,1.0), (9,0.9)}

**Figure 16.** Fuzzy clusters.