6.047 / 6.878 Computational Biology: Genomes, Networks, Evolution
Fall 2008

# Clustering

# Structure in High-Dimensional Data
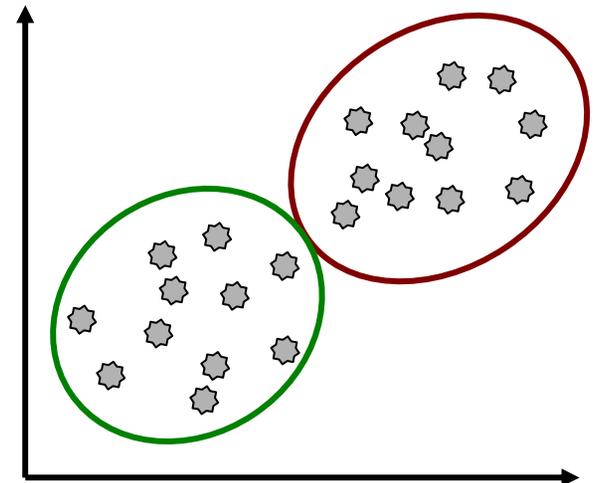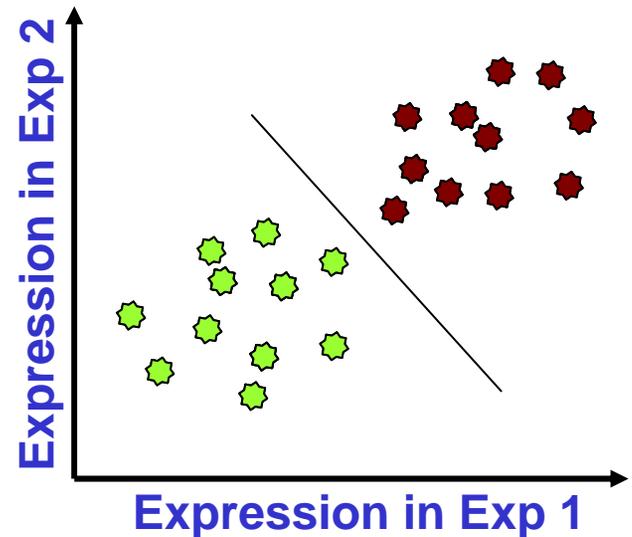
Gyulassy, Atilla, et al. "Topologically Clean Distance Fields." *IEEE Transactions on Visualization and Computer Graphics* 13, no. 6 (2007): 1432-1439.

- Structure can be used to reduce dimensionality of data
- Structure can tell us something useful about the underlying phenomena
- Structure can be used to make inferences about new data

# Clustering vs Classification

- **Objects** characterized by one or more features

- **Classification**
  - Have <u>labels</u> for some points
  - Want a "rule" that will accurately assign labels to new points
  - Supervised learning

- **Clustering**
  - No labels
  - Group points into clusters based on how "near" they are to one another
  - Identify <u>structure</u> in data
  - Unsupervised learning

**Expression in Exp 2**

**Expression in Exp 1**

# Today

- Microarray Data

- K-means clustering

- Expectation Maximization

- Hierarchical Clustering

# Central Dogma

# Expression Microarrays

- A way to measure the levels of mRNA in every gene

- Two basic types
  - Affymetrix gene chips
  - Spotted oligonucleotides

- Both work on same principle
  - Put DNA probe on slide
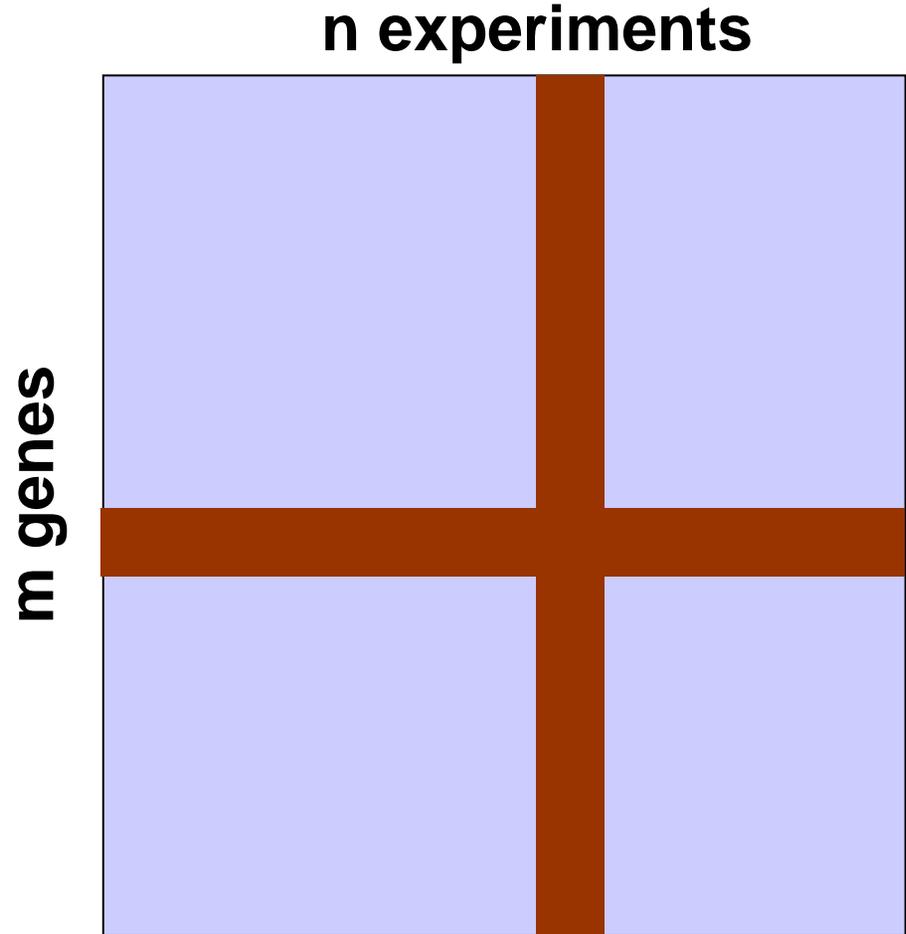  - Complementary hybridization

# Expression Microarrays

- Measure the level of mRNA messages in a cell

# Expression Microarray Data Matrix

- Genes are typically given as rows
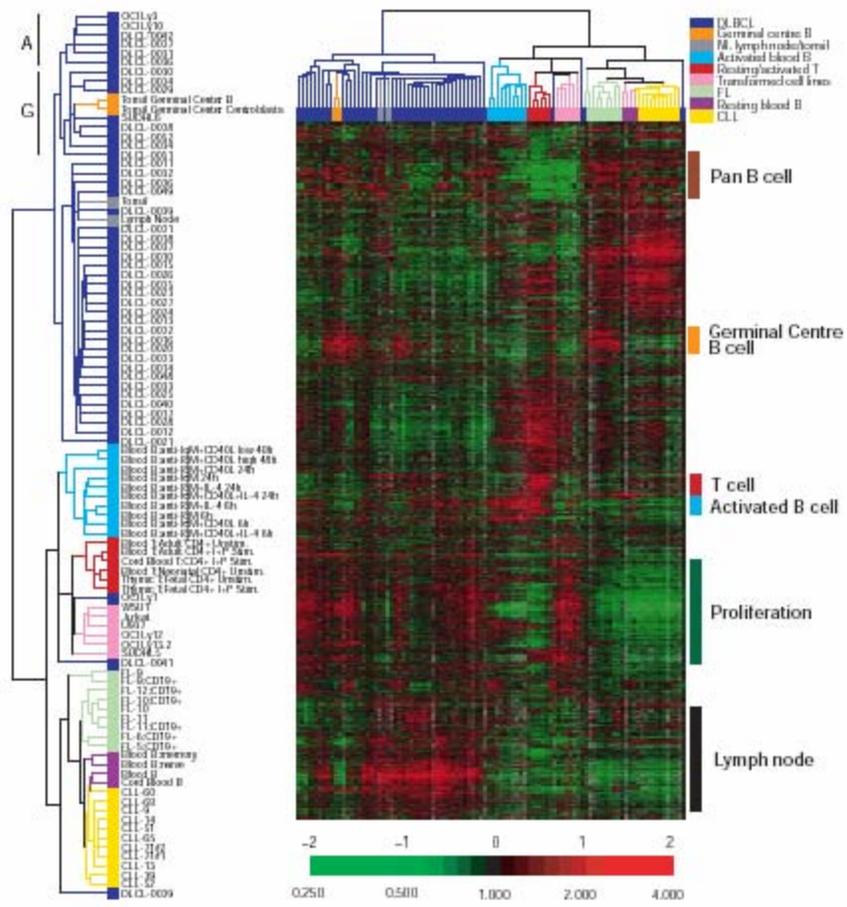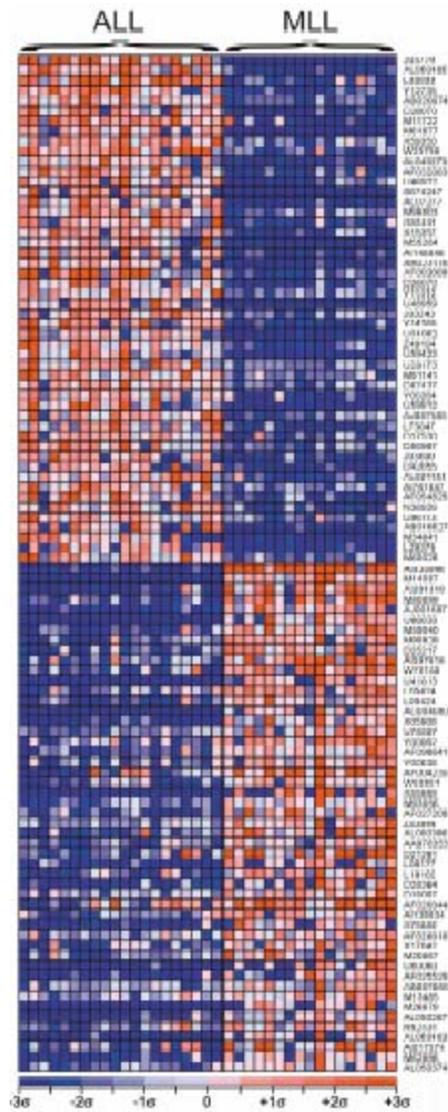
- Experiment are given by columns

**n experiments**

**m genes**

# Clustering and Classification in Genomics

- ## Classification

  - ➢ <u>Microarray data</u>: classify cell state (i.e. AML vs ALL) using expression data

  - ➢ <u>Protein/gene sequences</u>:  predict function, localization, etc.

- ## Clustering

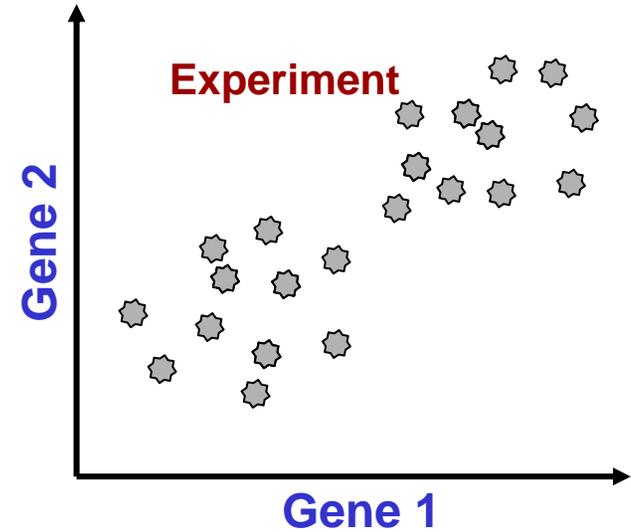  - ➢ <u>Microarray data</u>:  groups of genes that share similar function have similar expression patterns – identify regulons

  - ➢ <u>Protein sequence</u>:  group related proteins to infer function

  - ➢ <u>EST data</u>:  collapse redundant sequences
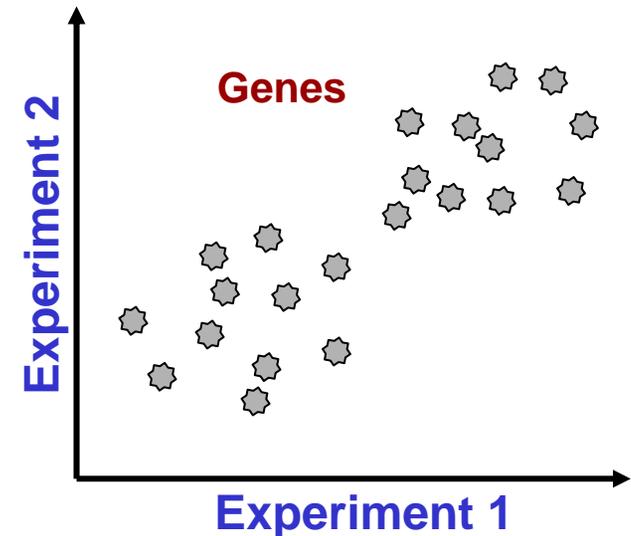
# Clustering Expression Data

- ## Cluster Experiments
  - Group by similar expression profiles

- ## Cluster Genes
  - Group by similar expression in different conditions

# Why Cluster Genes by Expression?

- ## Data Exploration
  - Summarize data
  - Explore without getting lost in each data point
  - Enhance visualization

- ## Co-regulated Genes
  - Common expression may imply common regulation
  - Predict *cis*-regulatory promoter sequences

- ## Functional Annotation
  - Similar function from similar expression

**GCN4**

**His2**          **Amino Acids**
**His3**          **Amino Acids**
**Unknown**

# Clustering Algorithms

- Partitioning
  - Divides objects into non-overlapping clusters such that each data object is in exactly one subset

- Agglomerative
  - A set of nested clusters organized as a hierarchy

# K-Means Clustering

The Basic Idea

- Assume a fixed number of clusters, K

- Goal: create "compact" clusters

# More Formally

1. **Initialize** K centers $\mathbf{u}_k$

**For each iteration** n until convergence

2. Assign each $\mathbf{x}_i$ the label of the nearest center, where the distance between $\mathbf{x}_i$ and $\mathbf{u}_k$ is

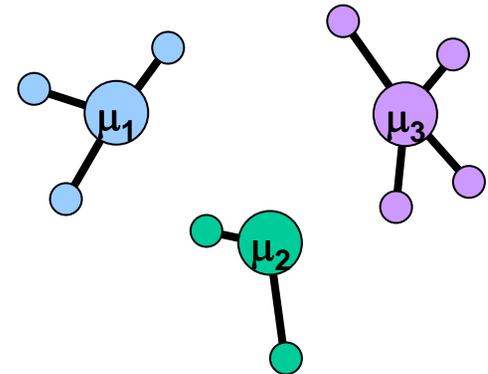$$d_{i,k} = \left( \mathbf{x}_i - \boldsymbol{\mu}_k \right)^2$$

3. Move the position of each $\mathbf{u}_k$ to the centroid of the points with that label

$$\boldsymbol{\mu}_k(n+1) = \sum_{x_i \text{ with label } j} \frac{\mathbf{x}_i}{\left| \mathbf{x}^k \right|} \;,\; \left| \mathbf{x}^k \right| = \#\mathbf{x}_i \text{ with label k}$$

# Cost Criterion

**We can think of K-means as trying to create clusters that <span style="color:navy">minimize a cost criterion</span> associated with the size of the cluster**

$$\mathrm{COST}\left(\mathbf{x}_1,\mathbf{x}_2,\mathbf{x}_3,...,\mathbf{x}_n\right)=\sum_{\boldsymbol{\mu}_k}\ \sum_{\mathbf{x}_i\text{ with label k}}\left(\mathbf{x}_i-\boldsymbol{\mu}_k\right)^2$$

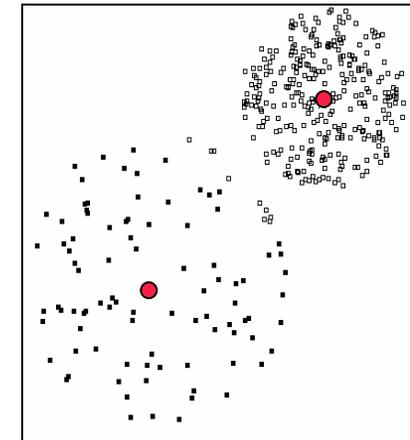**Minimizing this means minimizing each cluster term separately:**

$$\sum_{\mathbf{x}_i\text{ with label k}}\left(\mathbf{x}_i-\boldsymbol{\mu}_k\right)^2$$
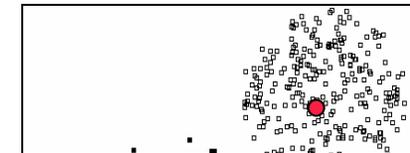
# Fuzzy K-Means

- Initialize K centers $\mathbf{u}_k$

- For each point calculate the probability of membership for each category

$$P(\text{label K} \mid \mathbf{x}_i, \boldsymbol{\mu}_k)$$
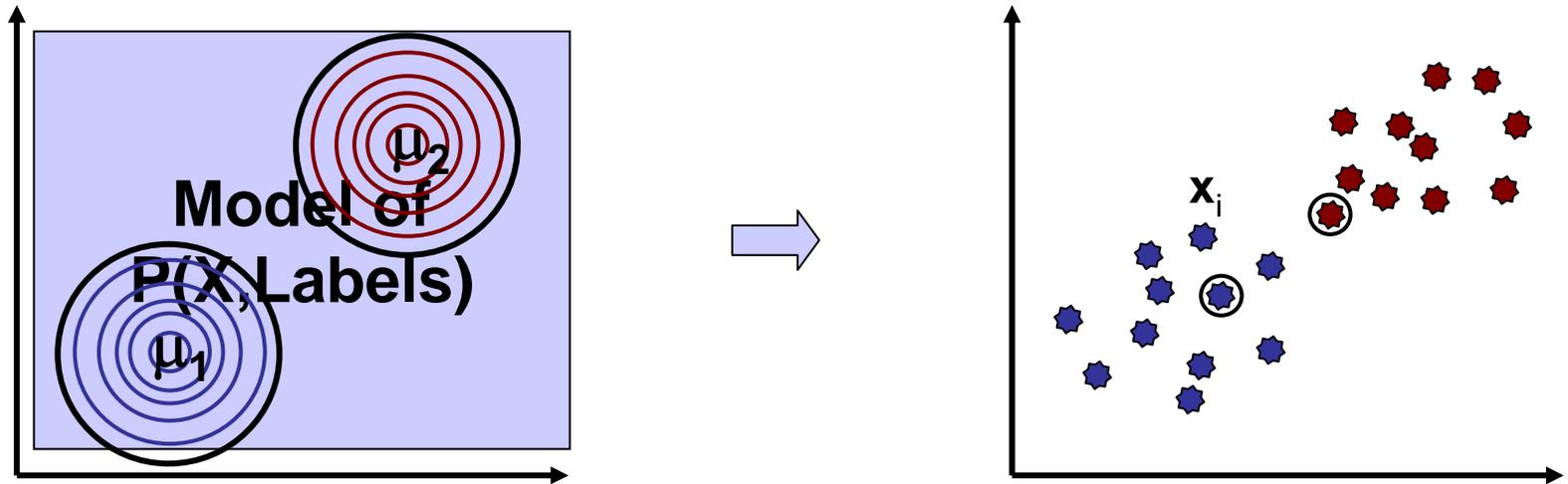
**K-means**

- Move the position of each $\mathbf{u}_k$ to the weighted centroid :

**Of course, K-Means just special case where**

$$P(\text{label K} \mid \mathbf{x}_i, \boldsymbol{\mu}_k) = \begin{cases} 1 & \text{if } \mathbf{x}_i \text{ is closest to } \boldsymbol{\mu}_k \\ 0 & \text{otherwise} \end{cases}$$

# K-Means as a Generative Model



**Model of P(X,Labels)**
$\mu_2$
$\mu_1$

$\mathbf{x}_i$

**Samples drawn from two equally normal distributions with unit variance - a *Gaussian Mixture Model***

$$P\left(\mathbf{x}_i \mid \mathbf{u}_j\right) = \frac{1}{\sqrt{2\pi}} \exp\left\{ -\frac{\left(\mathbf{x}_i - \mathbf{u}_j\right)^2}{2} \right\}$$

# Unsupervised Learning



**Learn?**

$x_i$

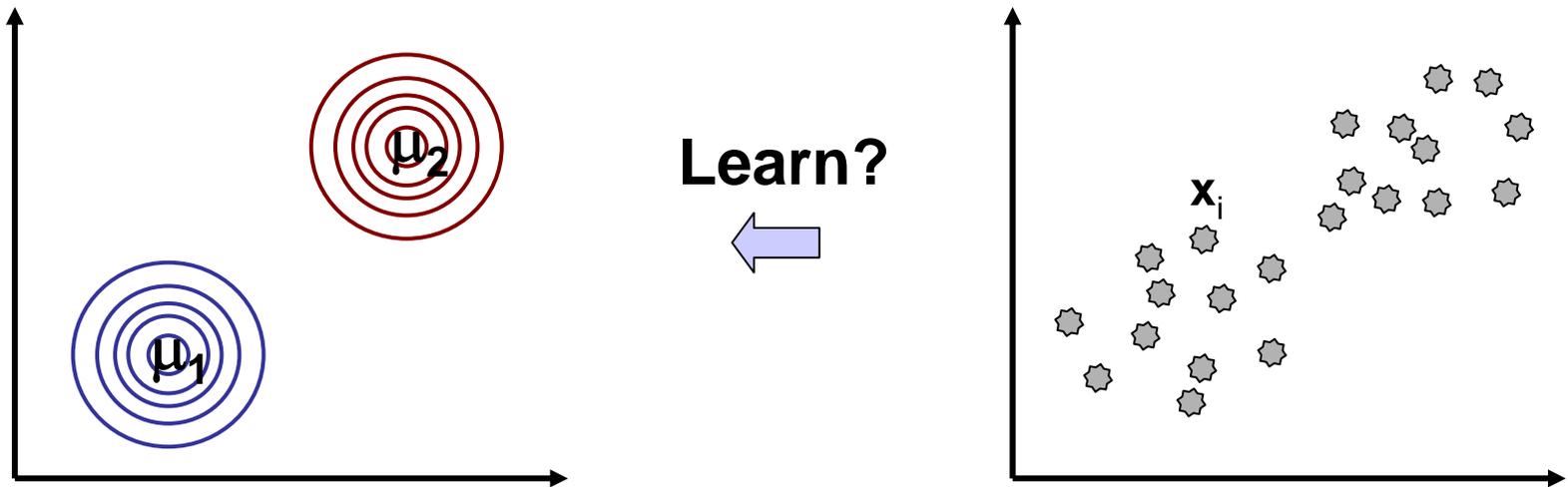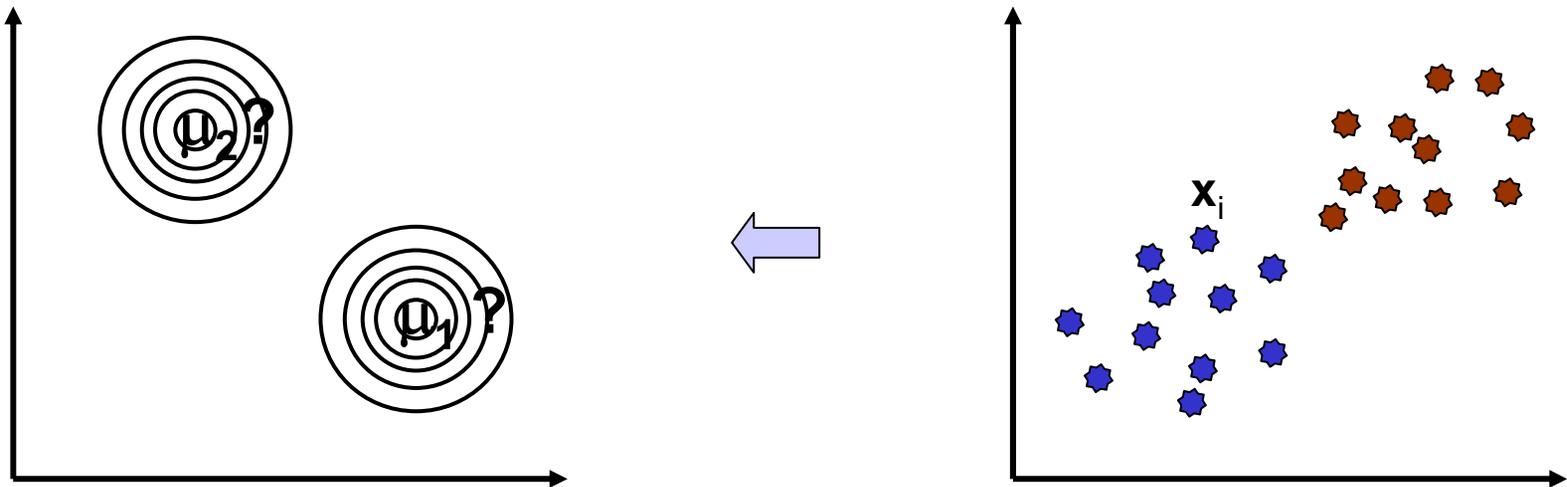**Samples drawn from two equally normal distributions with unit variance - a _Gaussian Mixture Model_**

$$P\left(\mathbf{x}_i \mid \mathbf{u}_j\right) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{\left(\mathbf{x}_i - \mathbf{u}_j\right)^2}{2}\right\}$$

# If We Have Labeled Points

**Need to estimate unknown gaussian centers from data**

**In general, how could we do this?**
**How could we "estimate" the "best" $u_k$?**



*Choose $u_k$ to maximize probability of model*

# If We Have Labeled Points

**Need to estimate unknown gaussian centers from data**

**In general, how could we do this?**
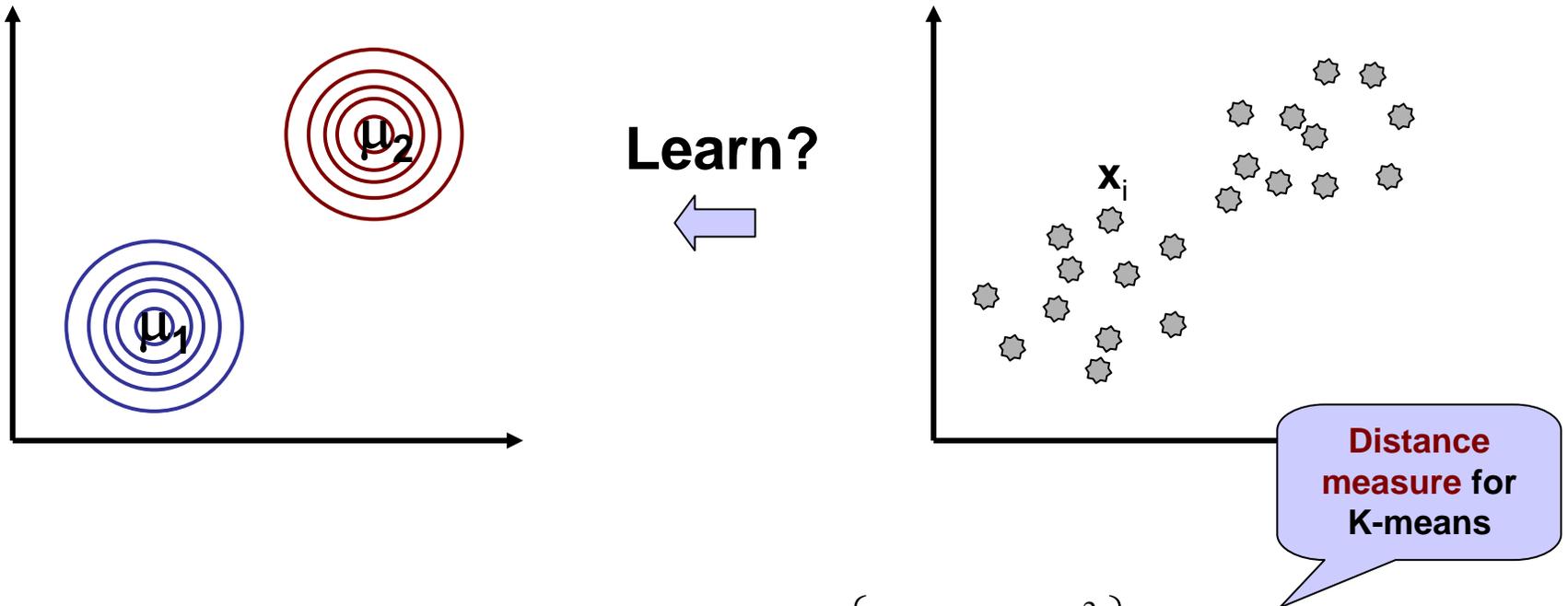**How could we "estimate" the "best" u$_k$?**

Given a set of **x**$_i$, all with label k, we can find the
<u>maximum likelihood</u> $\mu_k$ from

$$\arg\max_{\mu} \left\{ \log \prod_i P\left(\mathbf{x}_i \mid \mu\right) \right\} = \arg\max_{\mu} \sum_i \left\{ -\frac{1}{2}(\mathbf{x}_i - \mathbf{u})^2 + \log\left(\frac{1}{\sqrt{2\pi}}\right) \right\}$$

$$= \arg\min_{\mu} \sum_i (\mathbf{x}_i - \mathbf{u})^2$$

**Solution is
the centroid
of the x$_i$**

# If We Know Cluster Centers

**Need to estimate labels for the data**

$\mu_2$

**Learn?**

$\mathbf{x}_i$

**Distance measure for K-means**

$$\arg\max_k \ \mathrm{P}_k\left(\mathbf{x}_i \mid \boldsymbol{\mu}_i\right) = \arg\max_k \ \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{\left(\mathbf{x}_i - \mathbf{u}_k\right)^2}{2}\right\} = \arg\min_k \left(\mathbf{x}_i - \mathbf{u}_k\right)^2$$

$\mu_1$

# What If We Have Neither?

An idea:

1. Imagine we start with some $u_k^0$
2. We *could* calculate the most likely labels for $x_i^0$ given these $u_k^0$
3. We *could* then use these labels to choose $u_k^1$
4. And iterate (to convergence)

# Expectation Maximization (EM)

1. **Initialize parameters**

2. **E Step** **Estimate** <u>probability of hidden labels</u> **, Q, given parameters and sequence**

$$Q = P(label_i | x, u_k{}^{t-1})$$

3. **M Step** **Choose new parameters to** **maximize** <u>expected likelihood</u> **of parameters given Q**

$$u_k{}^t = \arg\max_u E_Q \left[ \log P(labels | x, u_k{}^{t-1}) \right]$$

4. **Iterate**

**P(x|Model)** *guaranteed* **to increase each iteration**

# Expectation Maximization (EM)

***Remember the basic idea!***

**1.Use model to estimate (distribution of) missing data**
**2.Use estimate to update model**
**3.Repeat until convergence**

Model is the <u>gaussian distributions</u>

Missing data are the <u>data point labels</u>

# Revisiting K-Means

1. Initialize K centers $\mathbf{u}_k$

2. Assign each $\mathbf{x}_i$ the label of the nearest center, where the distance between $\mathbf{x}_i$ and $\mathbf{u}_k$ is

$$d_{i,k} = \left( \mathbf{x}_i - \boldsymbol{\mu}_k \right)^2$$

⟹ **The most likely label k for a point $x_i$**

3. Move the position of each $\mathbf{u}_k$ to the centroid of the points with that label

⟹ **Maximum likelihood parameter $\mu_k$ given most likely label**

4. Iterate

# Revisiting K-Means

**Generative Model Perspective**

1. Initialize K centers $\mathbf{u}_k$

2. Assign each $\mathbf{x}_i$ the label of the nearest center, where the distance between $\mathbf{x}_i$ and $\mathbf{u}_k$ is

$$d_{i,k} = \left( \mathbf{x}_i - \mathbf{\mu}_k \right)^2$$

3. Move the position of each $\mathbf{u}_k$ to the centroid of the points with that label

4. Iterate

1. Initialize parameters

2. E Step Estimate most likely missing label given previous parameter

3. M Step Choose new parameters to maximize likelihood of parameters given estimated labels

4. Iterate

# Revisiting K-Means

**This is analogous to <u>Viterbi Learning from HMMs</u>**

1. Initialize K centers $\mathbf{u}_k$

2. Assign each $\mathbf{x}_i$ the label of the nearest center, where the

   **Analogy with HMM is to use <u>Viterbi</u> to find most likely <u>missing *path*</u> labels**

   **(see Durbin book)**

3. that label

4. Iterate

1. Initialize parameters

2. E Step Estimate <u>most likely missing label</u> given previous parameter

3. M Step Choose new parameters to <u>maximize likelihood</u> of parameters given estimated labels
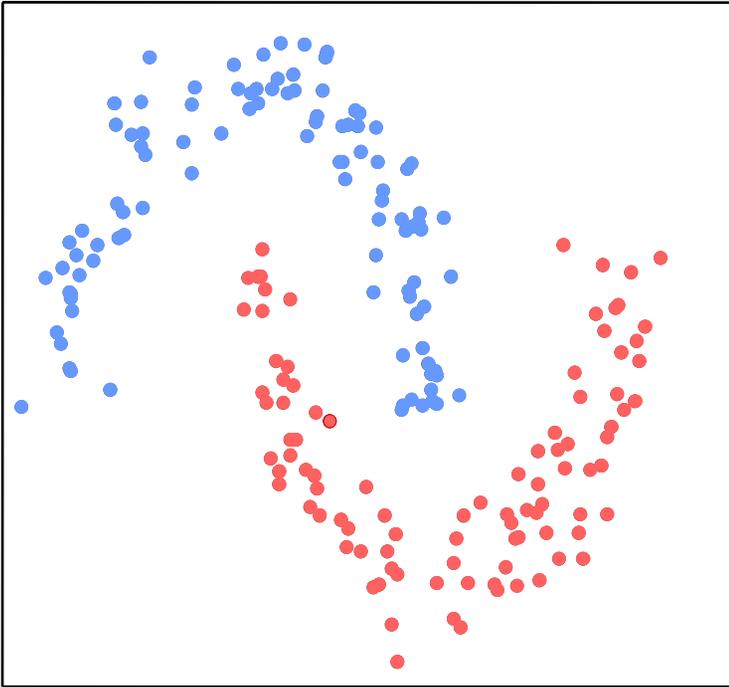
4. Iterate

# Revisting Fuzzy K-Means

**Recall that instead of assigning each point $x_i$ to a label k, we calculate the probability of each label for that point (fuzzy membership):**

$$P(\text{label K} \mid \mathbf{x}_i, \boldsymbol{\mu}_k)$$

**Recall that given a set of $x_i$, all with label k, we select a new $\mu_k$ with the update:**

<span style="color:darkred">**Looking at case b=1**</span>

$$\boldsymbol{\mu}_k(n+1) = \sum_{x_i \text{ with label j}} \mathbf{x}_i \, P(\boldsymbol{\mu}_k \mid \mathbf{x}_i) \Bigg/ \sum_{x_i \text{ with label j}} P(\boldsymbol{\mu}_k \mid \mathbf{x}_i)^b$$

**It can be shown that this update rule follows from assuming the gaussian mixture generative models and performing *Expectation-Maximization***

# Revisiting Fuzzy K-Means

## This is analogous to __Baum Welch from HMMs__

1. Initialize K centers $\mathbf{u}_k$

2. For each point calculate the probability of membership for each category

$$P(\text{label K} \mid \mathbf{x}_i, \boldsymbol{\mu}_k)$$

3. Move the position of each $\mathbf{u}_k$ to the weighted centroid :

$$\boldsymbol{\mu}_k(n+1) = \sum_{\mathbf{x}_i \text{ with label j}} \mathbf{x}_i \, P(\boldsymbol{\mu}_k \mid \mathbf{x}_i)^b \Bigg/ \sum_{\mathbf{x}_i \text{ with label j}} P(\boldsymbol{\mu}_k \mid \mathbf{x}_i)^b$$

4. Iterate

---

1. Initialize parameters

2. E Step Estimate probability over missing labels given previous parameter

3. M Step Choose new parameters to maximize expected likelihood of parameters given estimated labels

4. Iterate

# K-Means, Viterbi learning & EM

**K-Means** and **Fuzzy K-means** are two related methods that can be seen performing unsupervised learning on a **gaussian mixture model**

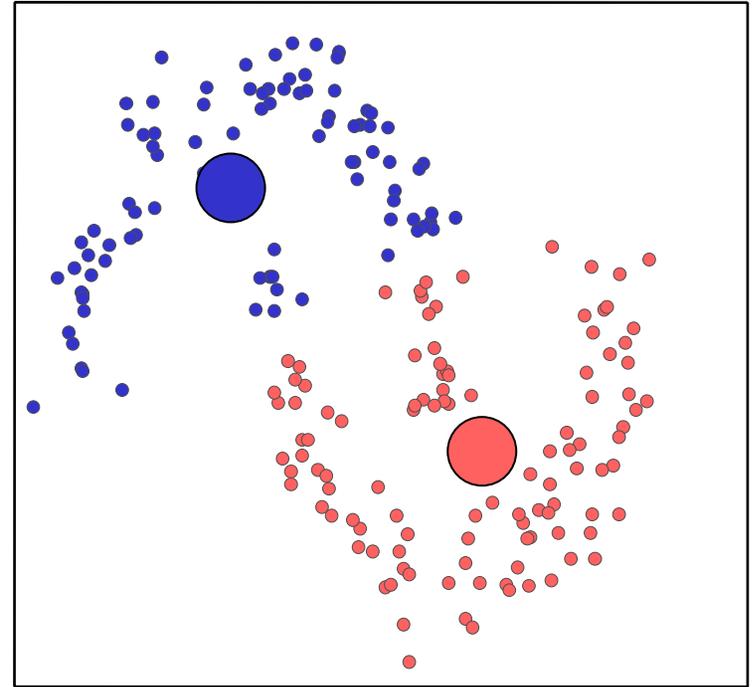**Reveal assumptions about underlying data model**

**Can relax assumptions by relaxing constraints on model**
- **Including explicit covariance matrix**
- **Relaxing assumption that all gaussians are equally likely**

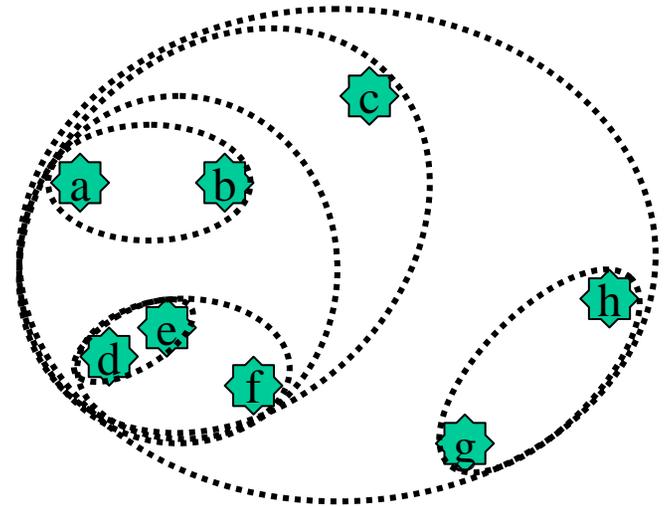# Implications: Non-globular Clusters



**Actual Clustering**

**K-means (K = 2)**

# But How Many clusters?

- How do we select K?
  - We can always make clusters "more compact" by increasing K
  - e.g. What happens is if K=number of data points?
  - What is a meaningful improvement?
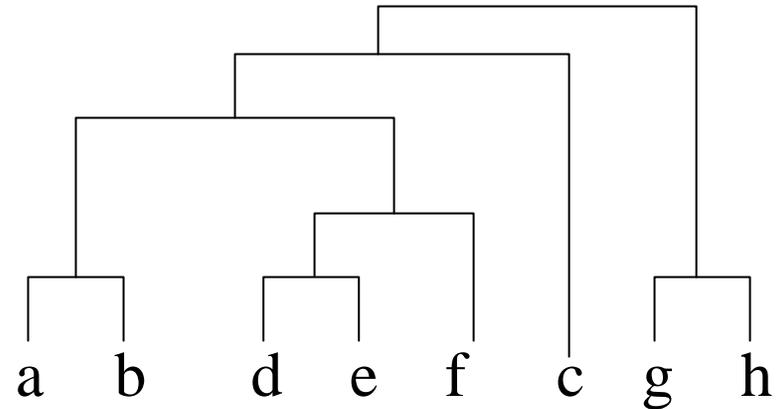- Hierarchical clustering side-steps this issue

# Hierarchical clustering

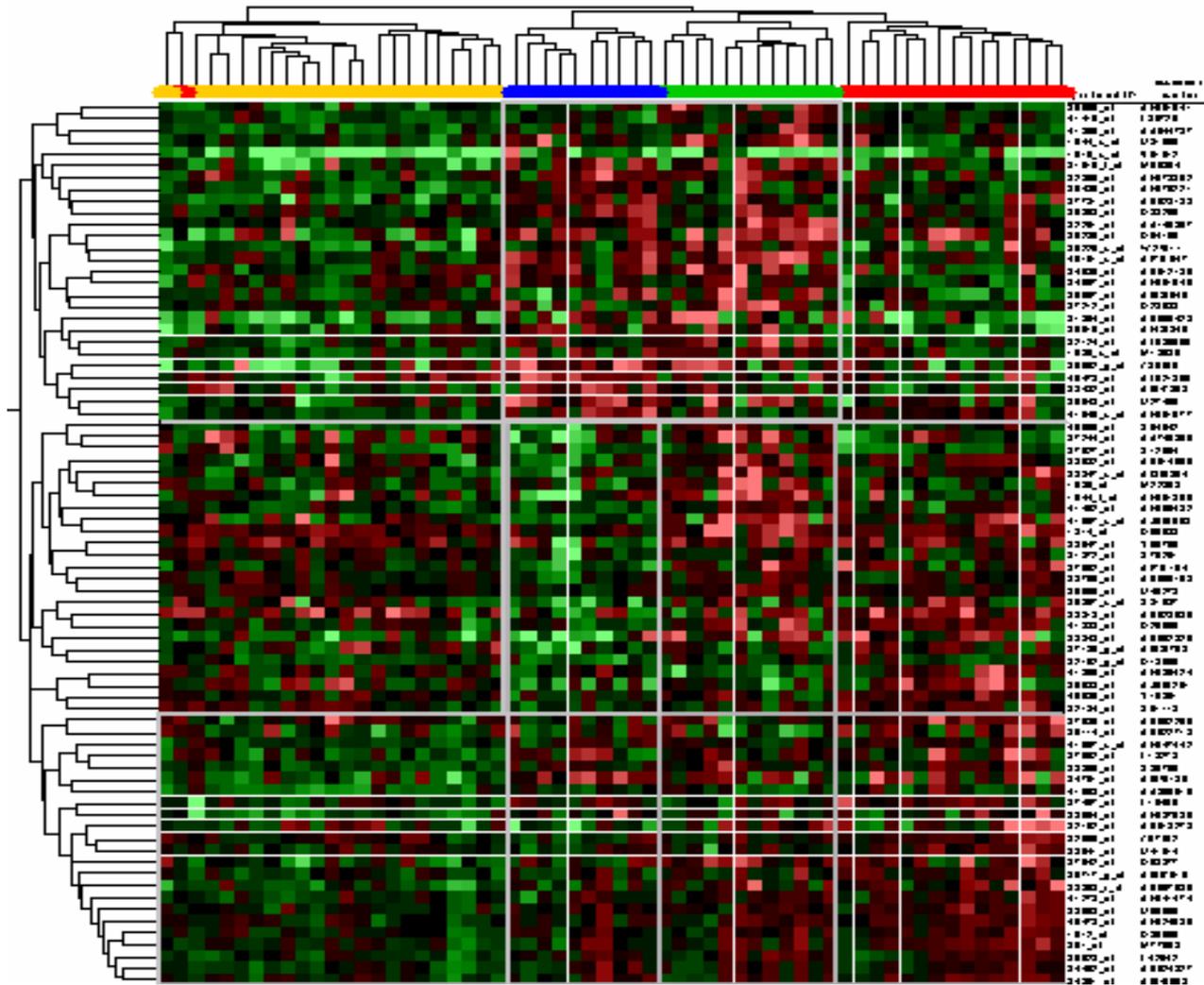Most widely used algorithm for expression data

- Start with each point in a separate cluster

- At each step:
  - Choose the pair of **closest clusters**
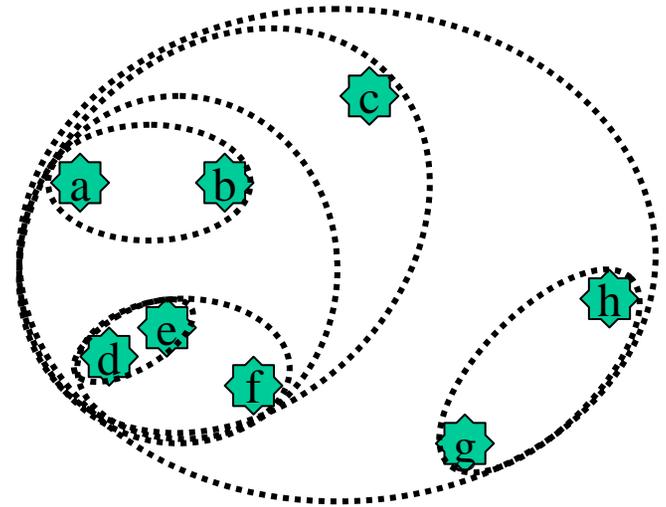  - Merge

➡ Phylogeny (UMPGA)

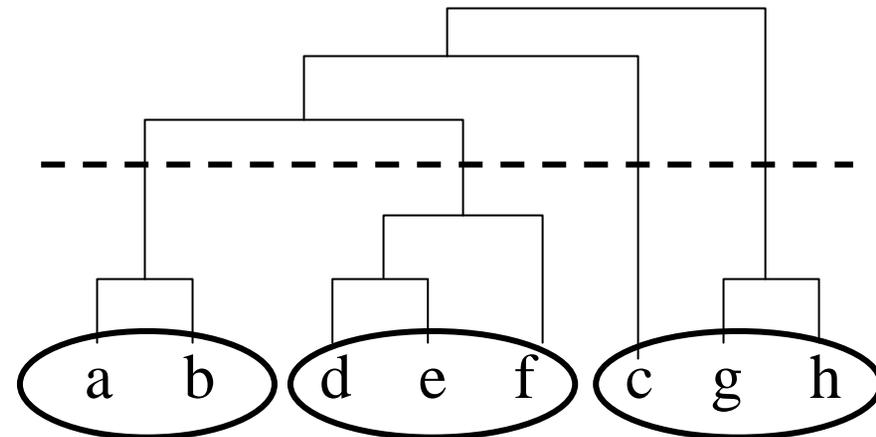# Visualization of results

# Hierarchical clustering

*Avoid needing to select number of clusters*

Produces clusters at all levels

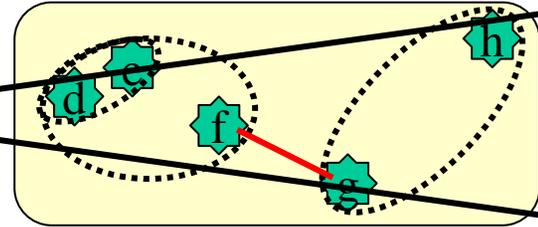We can always select a "cut level" to create disjoint clusters

*But how do we define distances between clusters?*
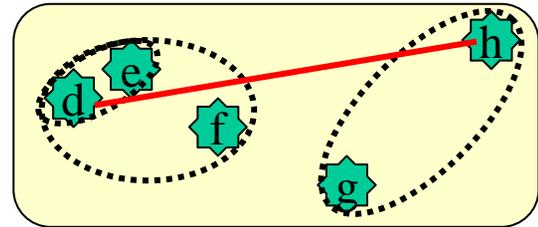
# Distance between clusters

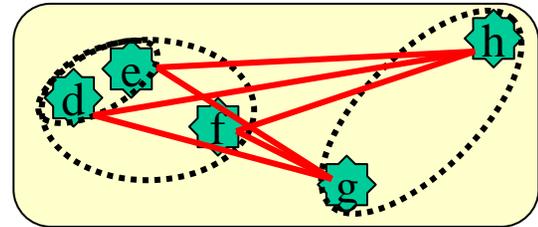- CD(X,Y)=min$_{x \in X, y \in Y}$ D(x,y)

  *Single-link method*

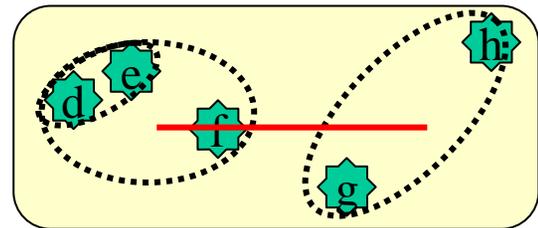- CD(X,Y)=max$_{x \in X, y \in Y}$ D(x,y)

  *Complete-link method*

- CD(X,Y)=avg$_{x \in X, y \in Y}$ D(x,y)

  *Average-link method*

- CD(X,Y)=D( avg(X) , avg(Y) )

  *Centroid method*

# (Dis)Similarity Measures
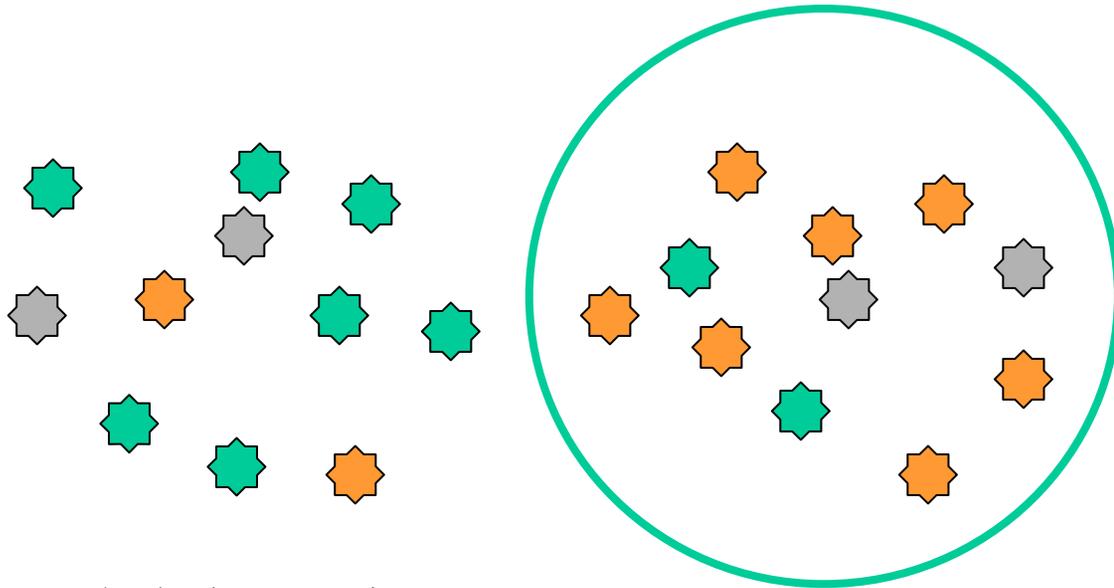
Image removed due to copyright restrictions.

Table 1, Gene expression similarity measures. D'haeseleer, Patrik. "How Does Gene Expression Clustering Work?" *Nature Biotechnology* 23 (2005): 1499-1501.

# Evaluating Cluster Performance

**In general, it depends on your goals in clustering**

- Robustness
  - Select random samples from data set and cluster
  - Repeat
  - Robust clusters show up in all clusters

- Category Enrichment
  - Look for categories of genes "over-represented" in particular clusters
  - Also used in Motif Discovery
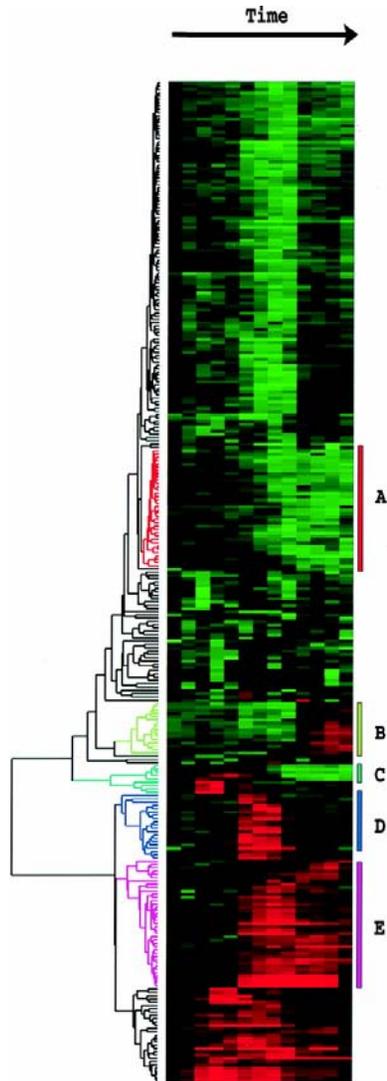
# Evaluating clusters – Hypergeometric Distribution



$$P(pos \geq r) = \sum_{m \geq r} \frac{\binom{p}{m}\binom{N-p}{k-m}}{\binom{N}{k}}$$

P-value of uniformity
in computed cluster

Prob that a randomly chosen
set of k experiments would
result in m positive and k-m
negative

- N experiments, p labeled **+**, (N-p) **–**
- Cluster: k elements, m labeled **+**
- P-value of *single* cluster containing k elements of which at least r are **+**

# Similar Genes Can Cluster



Eisen, Michael et al. "Cluster Analysis and Display of Genome-wide Expression Patterns." *PNAS* 95, no. 25 (1998): 14863-14868. Copyright (1998) National Academy of Sciences, U.S.A.

**Clustered 8600 human genes using expression time course in fibroblasts**

**(A) Cholesterol biosynthesis**

**(B) Cell cycle**

**(C) Immediate early response**

**(D) Signalling and angiogenesis**

**(E) Wound healing**

**(Eisen (1998) PNAS)**

# Clusters and Motif Discovery

Expression from
15 time points
during yeast
cell cycle



Figure by MIT OpenCourseWare.

**Tavazoie & Church (1999)**

# Next Lecture

**The other side of the coin… Classification**