

Unsupervised Methods

(mostly clustering)

Unsupervised Methods

- All this time in class we have seen *supervised* methods:
 - Data have outcomes: (x_i, y_i)
- In this section, we will look at data *without* outcomes
- Previously, we cared about $P(X, Y)$, but concentrated on $P(Y | X)$ since that's what matters for *prediction*
- Now, we want $P(X)$ since there is no outcome Y

Outline

- **Hierarchical Clustering**
- **K-means (and K-medoids) clustering**
- **Model-Based clustering (Gaussian Mixture Models)**
 - **EM algorithm**

© 2002 Nature Publishing Group

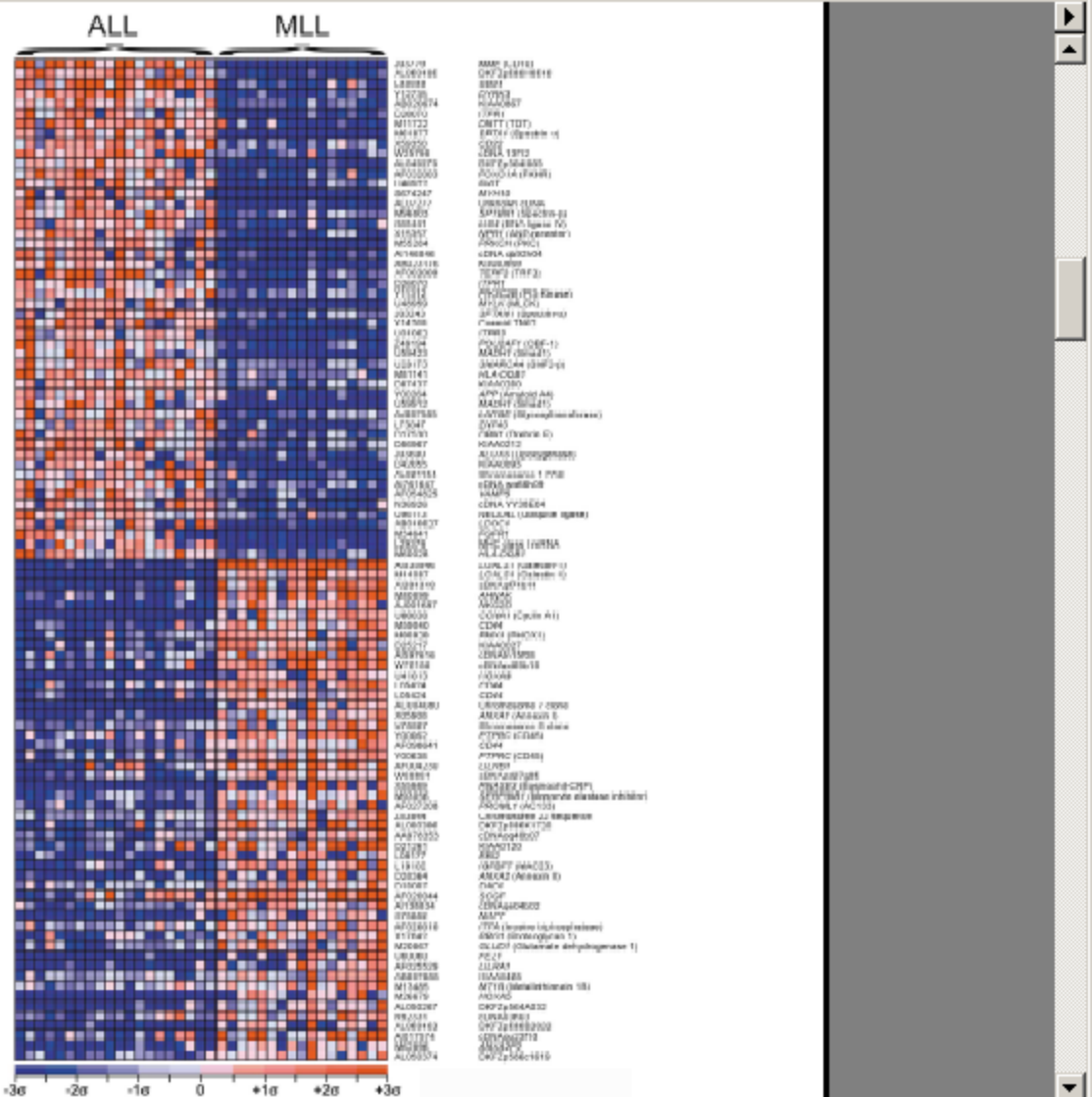
precursor ALL bearing an ALL translocation against those from individuals diagnosed with conventional B-precursor ALL that lack this translocation. Initially, we collected samples from 20 individuals with conventional childhood ALL (denoted ALL), 10 of which had a TEL/AML1 translocation. In addition, we collected samples from 17 individuals affected with the MLL translocation (denoted MLL). Details of the affected individuals and expression data are available online (Methods).

First, we determined whether there were genes among the 12,600 tested whose expression pattern correlated with the presence of an MLL translocation. We sorted the genes by their degree of correlation with the MLL/ALL distinction (Fig. 1) and used permutation testing to assess the statistical significance of the observed differences in gene expression¹³. For the 37 samples tested, roughly 1,000 genes are underexpressed in MLL as compared with conventional ALL, and about 200 genes are relatively highly expressed (data not shown). Thus, MLL shows a gene expression profile markedly different from that of conventional ALL.

MLL shows multilineage gene expression

Inspection of the genes differentially expressed between MLL and ALL is instructive (Fig. 1). Many genes underexpressed in MLL have a function in early B-cell development. These include genes expressed in early B cells^{14,15}, MME, CD24, CD22

Fig. 1 Genes that distinguish ALL from MLL. The 100 genes that are most highly correlated with the class distinction are shown. Each column represents a leukemia sample, and each row represents an individual gene. Expression levels are normalized for each gene, where the mean is 0, expression levels greater than the mean are shown in red and levels less than the mean are in blue. Increasing distance from the mean is represented by increasing color intensity. The top 50 genes are relatively underexpressed and the bottom 50 genes relatively overexpressed in MLL. Gene accession numbers and the gene symbols or DNA sequence names are labeled on the right. Individual samples are arranged such that column 1 corresponds to ALL patient 1, column 2 corresponds to ALL patient 2, and so on. Information about the samples along with the top 200 genes that make the ALL/MLL distinction and their accession numbers can be found on our web site



exactly the germinal centre phenotype *in vitro*, as determined by the failure of a variety of activation conditions to induce the expression of BCL-6 protein, a highly specific marker for germinal centre B signature of germinal centre B cells was reproduced virtually unchanged in FL, supporting the view that this lymphoma arises from this stage of B-cell differentiation (Fig. 2).

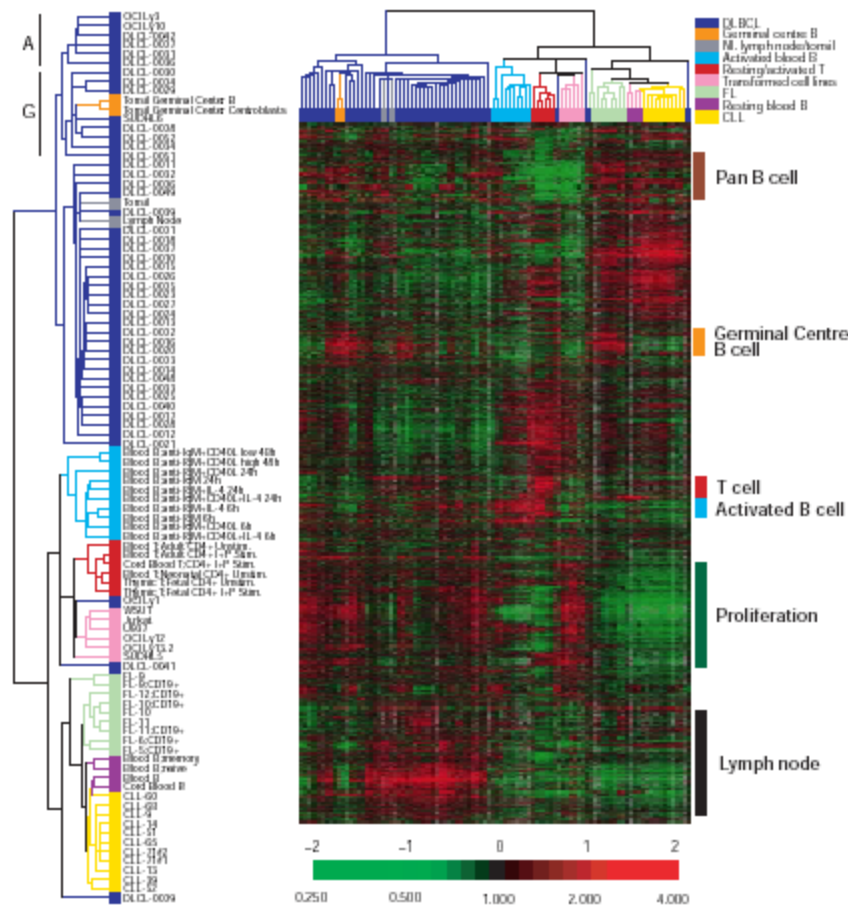


Figure 1 Hierarchical clustering of gene expression data. Depicted are the ~1.8 million measurements of gene expression from 126 microarray analyses of 96 samples of normal and malignant lymphocytes. The dendrogram at the left lists the samples studied and provides a measure of the relatedness of gene expression in each sample. The dendrogram is colour coded according to the category of mRNA sample studied (see

hybridization of fluorescent cDNA probes prepared from each experimental mRNA samples to a reference mRNA sample. These ratios are a measure of relative gene expression in each experimental sample and were depicted according to the colour scale shown at the bottom. As indicated, the scale extends from fluorescence ratios of 0.25 to 4 (-2 to +2 in log base 2 units). Grey indicates missing or excluded data. See

the pathogen remains localized to the tissue granuloma. The difference in the gene expression profiles between the two leprosy subclasses suggested that their differences in disease manifestation may reflect two opposing gene expression programs that influence the type of host response (10). The most pronounced differences in T-lep and L-lep gene expression profiles were among genes within the immune response family (Fig. 3). The gene expression profiles were consistent with previous data showing that type 1 cytokines associated with cell-mediated immunity predominate in T-lep lesions, whereas type 2 cytokines predominate in L-lep lesions (2, 11, 12). For example, genes encoding the type 1 cytokines lymphotoxin- α , interleukin (IL)-7, and IL-15 were comparatively up-regulated in T-lep lesions, as well as genes encoding CD1b and signaling lymphocytic activation molecule (SLAM), two molecules previously linked to cell-mediated immunity and type-1 cytokine production in these patients (13-15). In contrast, L-lep lesions differentially expressed the type 2 cytokines transforming growth factor- β and IL-5, as well as IL-4 and IL-10, although the differences in gene expression for the latter two genes between the patient groups were not statistically significant (P not ≤ 0.05) and were instead confirmed by quantitative polymerase chain reaction (qPCR) (8). As part of the type 2 pattern, L-lep lesions also exhibited marked up-regulation of genes related to humoral immunity, including immunoglobulin (Ig) heavy and light chains and molecules involved in B cell activation.

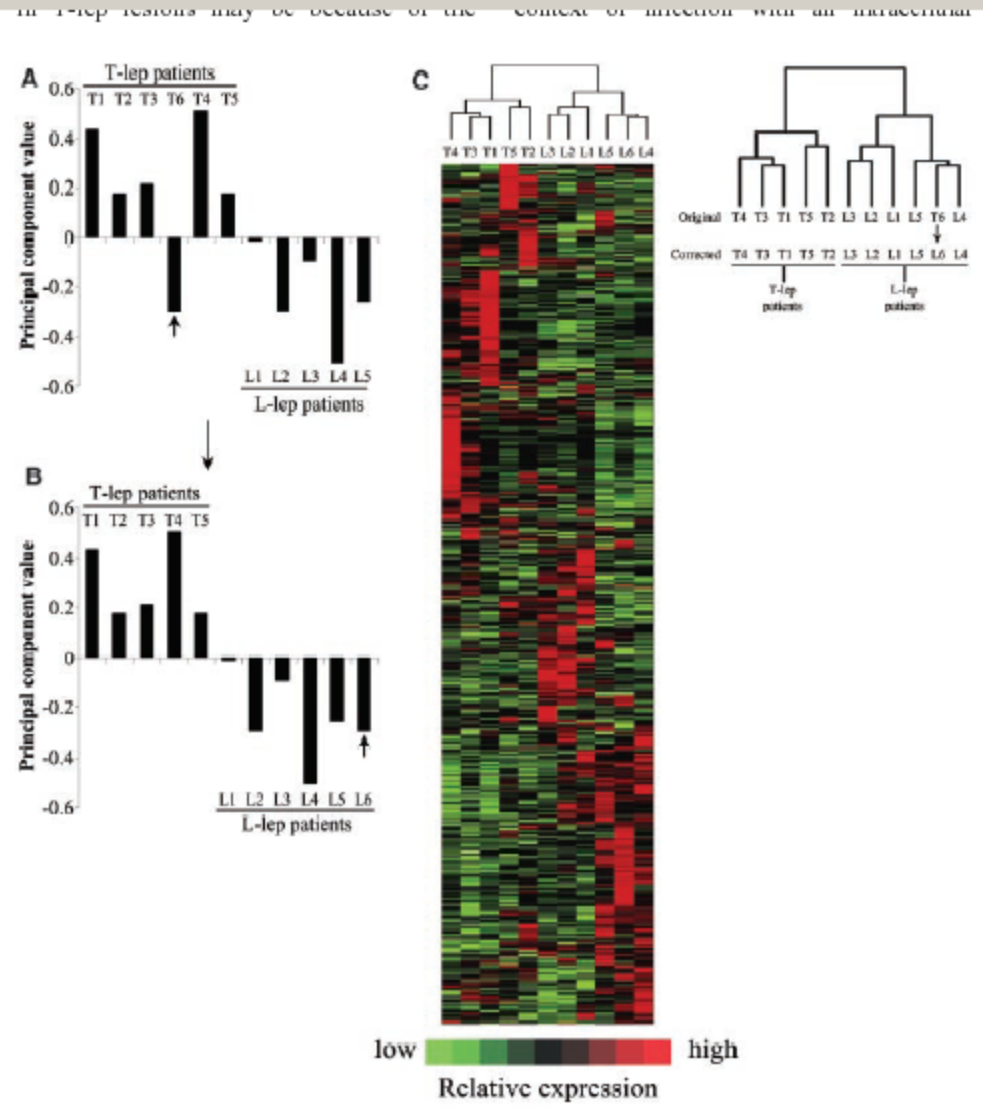


Fig. 1. Two unsupervised data analyses separate leprosy patients into clinically relevant subclasses.

Distance

- **Clustering organizes things that are *close* into groups**
- **What does it mean for two genes to be close?**
- **What does it mean for two samples to be close?**
- **Once we know this, how do we define groups?**

Distance

- **We need a mathematical definition of distance between two points**
- **What are points?**
- **If each gene is a point, what is the mathematical definition of a point?**

Points

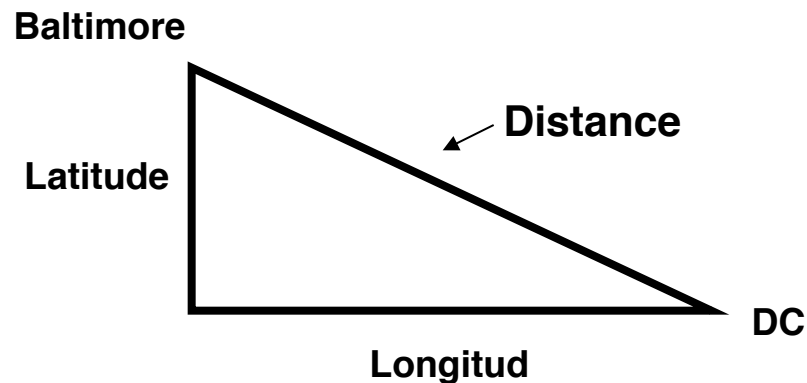
- **Gene1 = $(E_{11}, E_{12}, \dots, E_{1N})'$**
- **Gene2 = $(E_{21}, E_{22}, \dots, E_{2N})'$**

- **Sample1 = $(E_{11}, E_{21}, \dots, E_{G1})'$**
- **Sample2 = $(E_{12}, E_{22}, \dots, E_{G2})'$**

- **E_{gi} = expression gene g , sample i**

Most Famous Distance

- **Euclidean distance**
 - Example distance between gene 1 and 2:
 - Sqrt of Sum of $(E_{1i}-E_{2i})^2, i=1, \dots, N$
- **When N is 2, this is distance as we know it:**

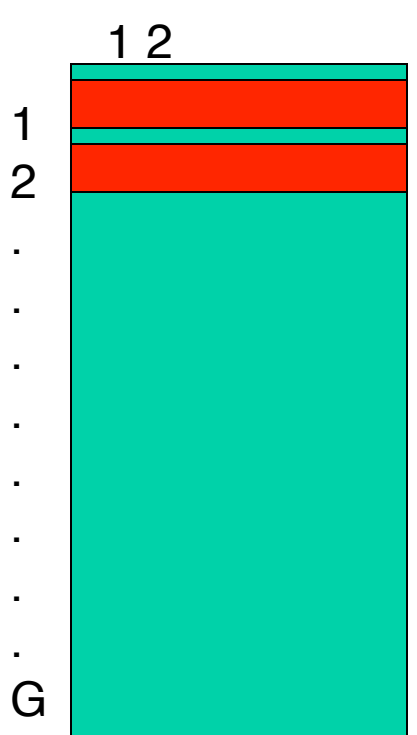


When N is 20,000 you have to think abstractly

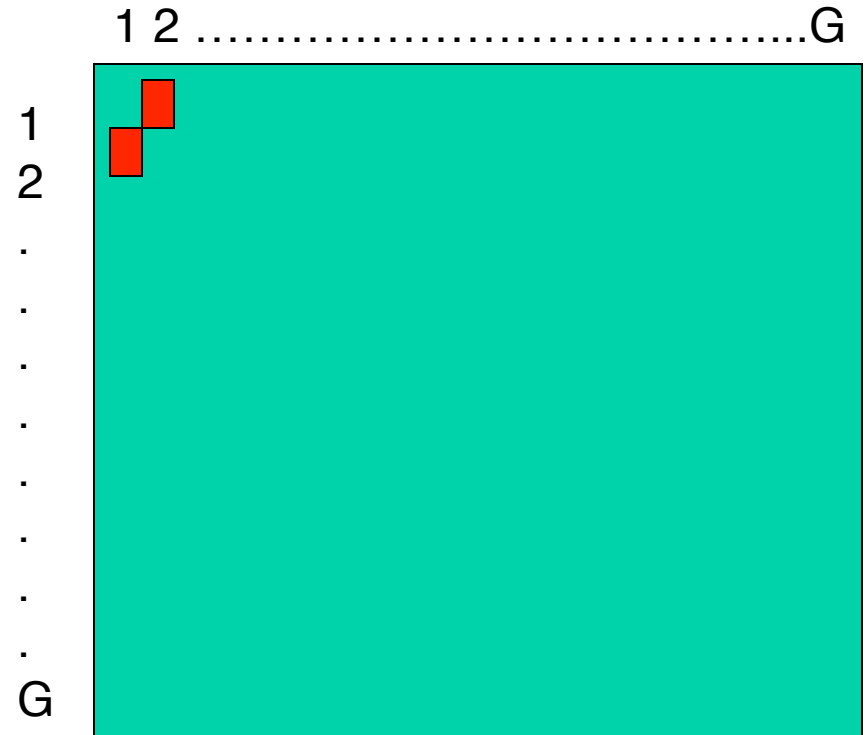
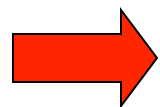
Similarity

- **Instead of distance, clustering can use *similarity***
- **If we standardize points then Euclidean distance is equivalent to using absolute value of correlation as a similarity index**
- **Other examples:**
 - Spearman correlation
 - Categorical measures

The similarity/distance matrices

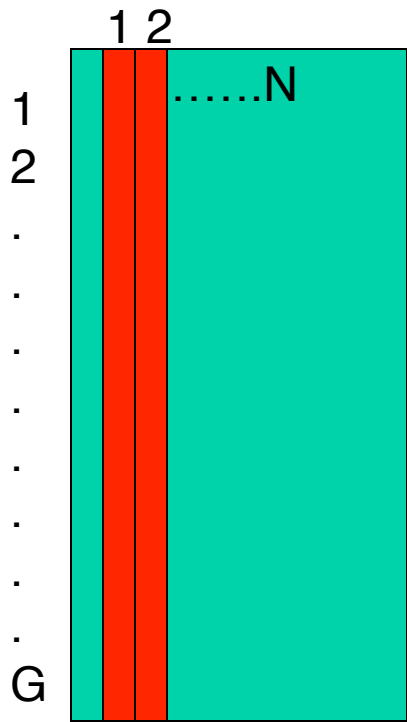


DATA MATRIX

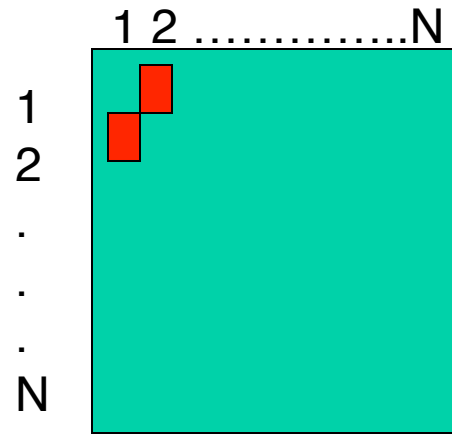
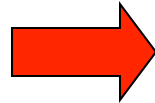


GENE SIMILARITY MATRIX

The similarity/distance matrices



DATA MATRIX



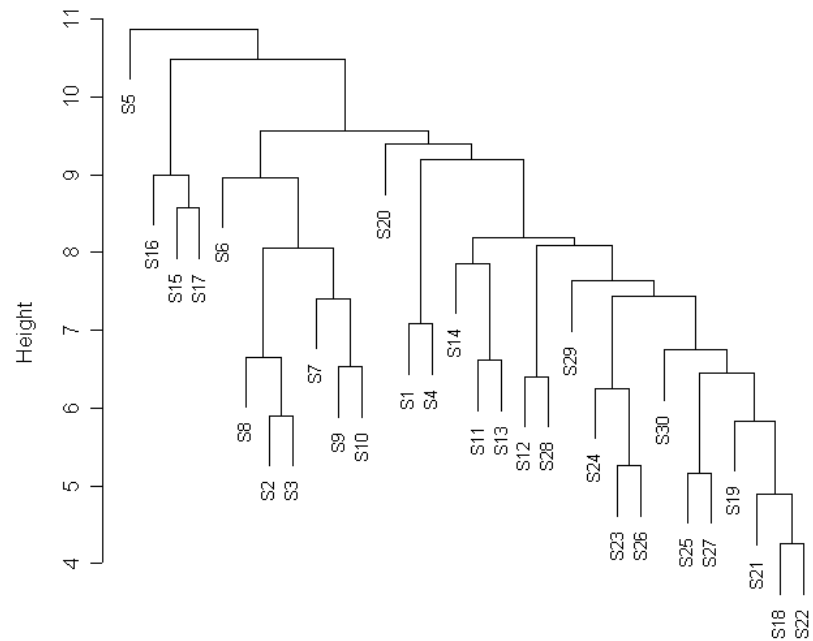
SAMPLE SIMILARITY MATRIX

Hierarchical

- **Divide all points into 2. Then divide each group into 2. Keep going until you have groups of 1 and can not divide further.**
- **This is divisive or top-down hierarchical clustering. There is also agglomerative clustering or bottom-up**

Dendrograms

- We can then make dendrograms showing divisions
- The y-axis represents the distance between the groups divided at that point

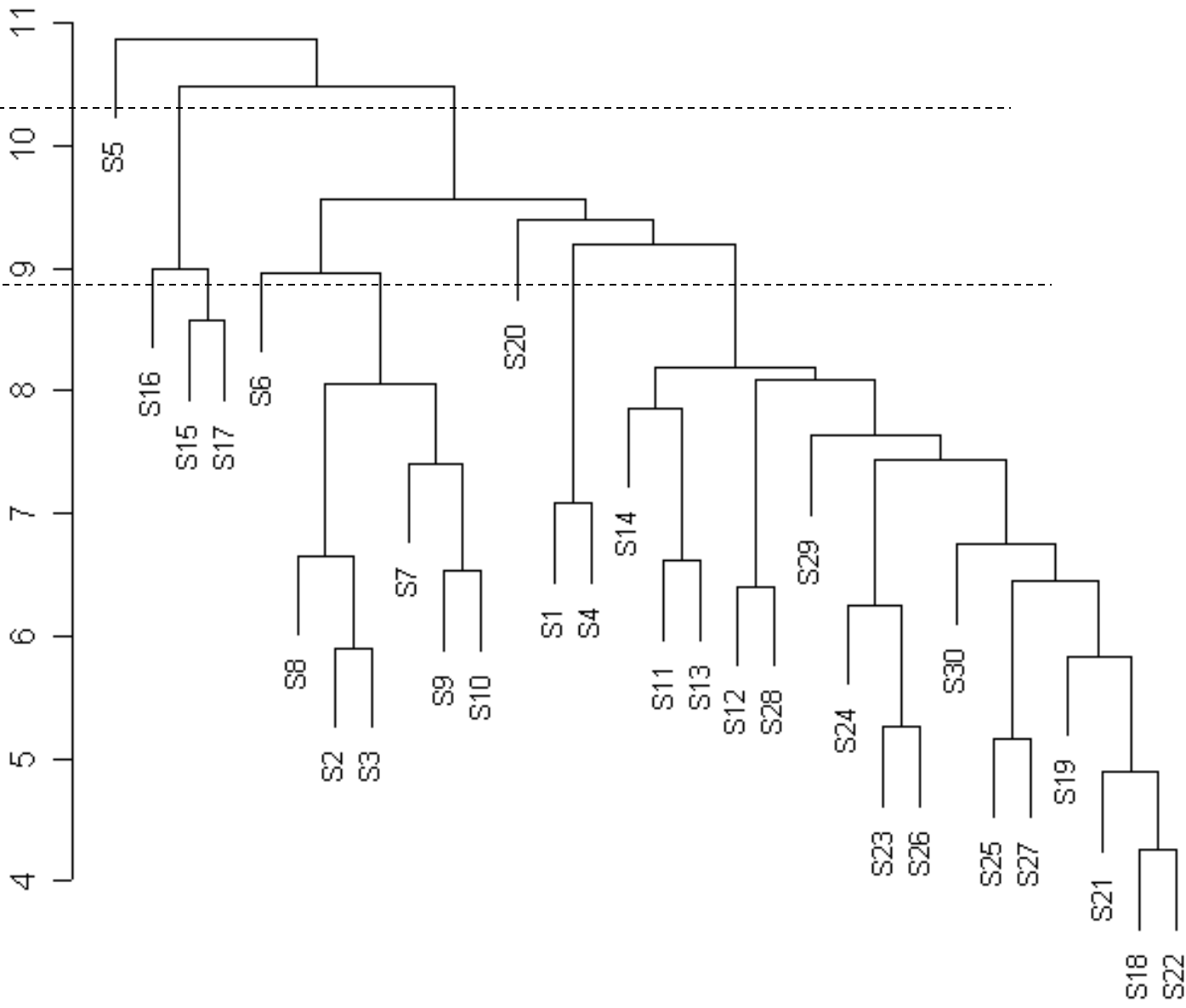


**Note: Left and right is assigned arbitrarily.
Look at the height of division to find out distance.
For example, S5 and S16 are very far.**

But how do we form actual clusters?

We need to pick a height

Height



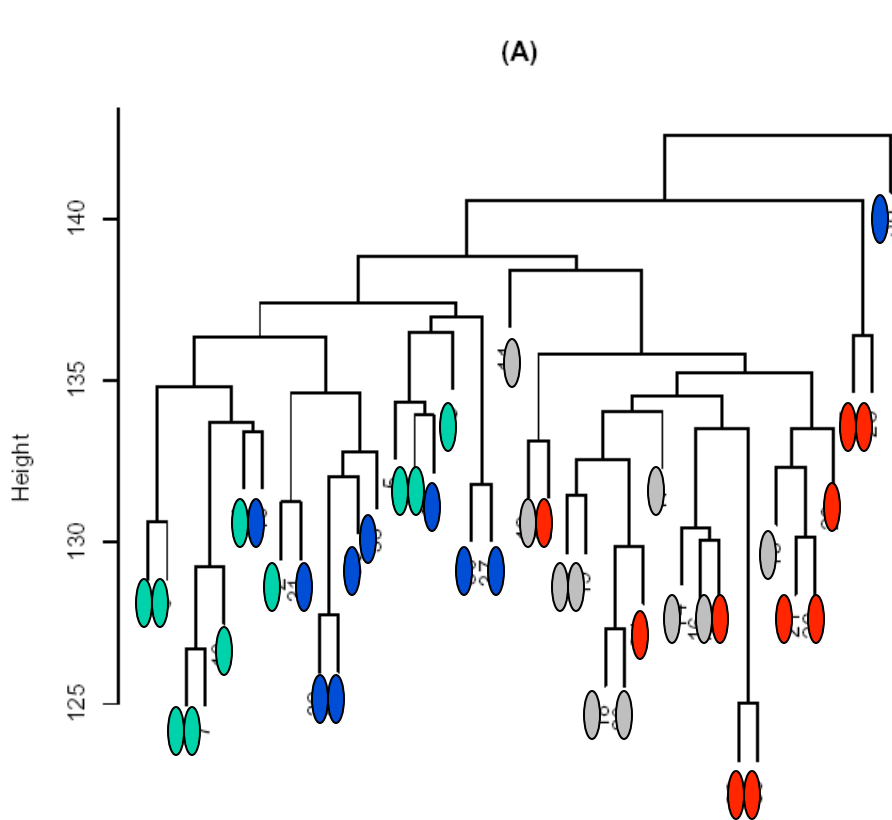
How to make a hierarchical clustering

- 1. Choose samples and genes to include in cluster analysis**
- 2. Choose similarity/distance metric**
- 3. Choose clustering direction (top-down or bottom-up)**
- 4. Choose linkage method (if bottom-up)**
- 5. Calculate dendrogram**
- 6. Choose height/number of clusters for interpretation**
- 7. Assess cluster fit and stability**
- 8. Interpret resulting cluster structure**

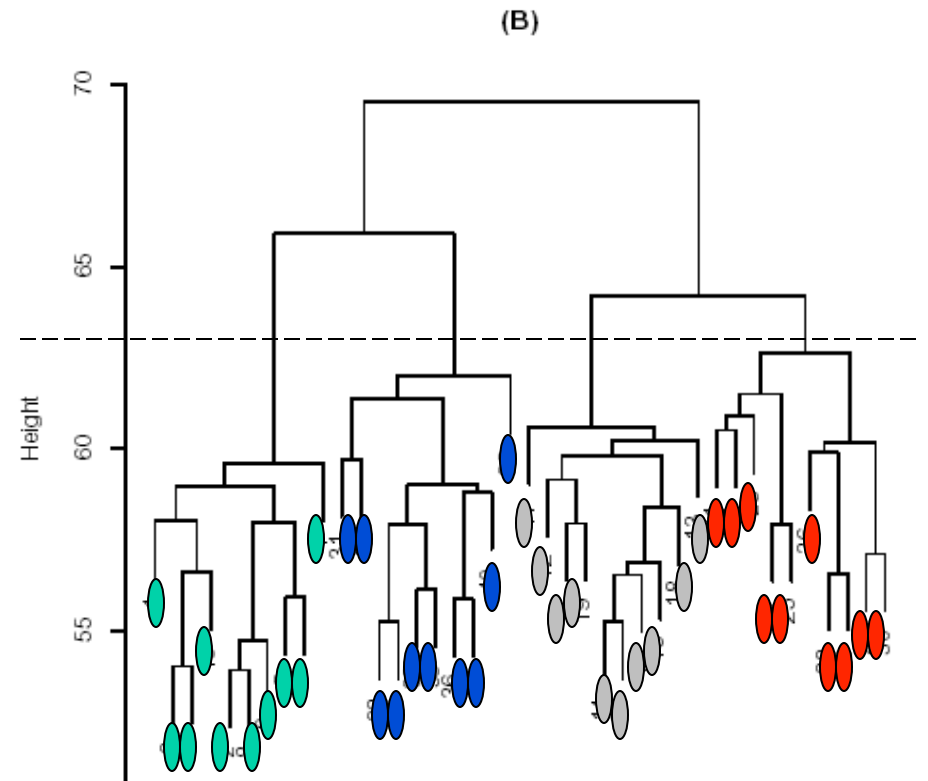
1. Choose samples and genes to include

- **Important step!**
- **Do you want housekeeping genes included?**
- **What to do about replicates from the same individual/tumor?**
- **Genes that contribute noise will affect your results.**
- **Including all genes: dendrogram can't all be seen at the same time.**
- **Perhaps screen the genes?**

Simulated Data with 4 clusters: 1-10, 11-20, 21-30, 31-40



A: 450 relevant genes plus 450 “noise” genes.



B: 450 relevant genes.

2. Choose similarity/distance matrix

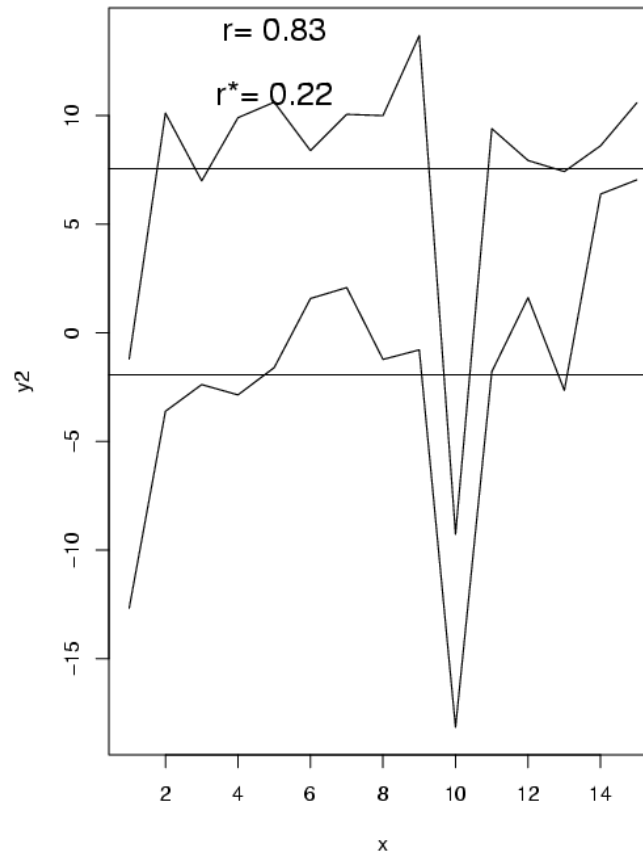
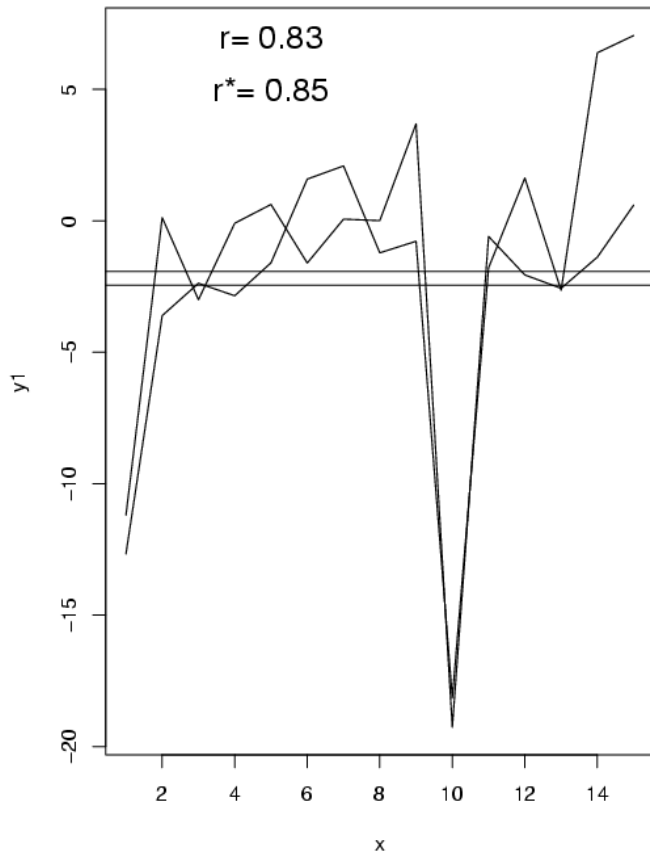
- **Think hard about this step!**
- **Remember: garbage in → garbage out**
- **The metric that you pick should be a valid measure of the distance/similarity of genes.**
- **Examples:**
 - **Applying correlation to highly skewed data will provide misleading results.**
 - **Applying Euclidean distance to data measured on categorical scale will be invalid.**
- **Not just “wrong”, but which makes most sense**

Some correlations to choose from

- **Pearson Correlation:**
$$s(x_1, x_2) = \frac{\sum_{k=1}^K (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2)}{\sqrt{\sum_{k=1}^K (x_{1k} - \bar{x}_1)^2 \sum_{k=1}^K (x_{2k} - \bar{x}_2)^2}}$$

- **Uncentered Correlation:**
$$s(x_1, x_2) = \frac{\sum_{k=1}^K x_{1k} x_{2k}}{\sqrt{\sum_{k=1}^K x_{1k}^2 \sum_{k=1}^K x_{2k}^2}}$$

- **Absolute Value of Correlation:**
$$s(x_1, x_2) = \left| \frac{\sum_{k=1}^K (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2)}{\sqrt{\sum_{k=1}^K (x_{1k} - \bar{x}_1)^2 \sum_{k=1}^K (x_{2k} - \bar{x}_2)^2}} \right|$$



The difference is that, if you have two vectors X and Y with identical shape, but which are offset relative to each other by a fixed value, they will have a standard Pearson correlation (centered correlation) of 1 but will not have an uncentered correlation of 1.

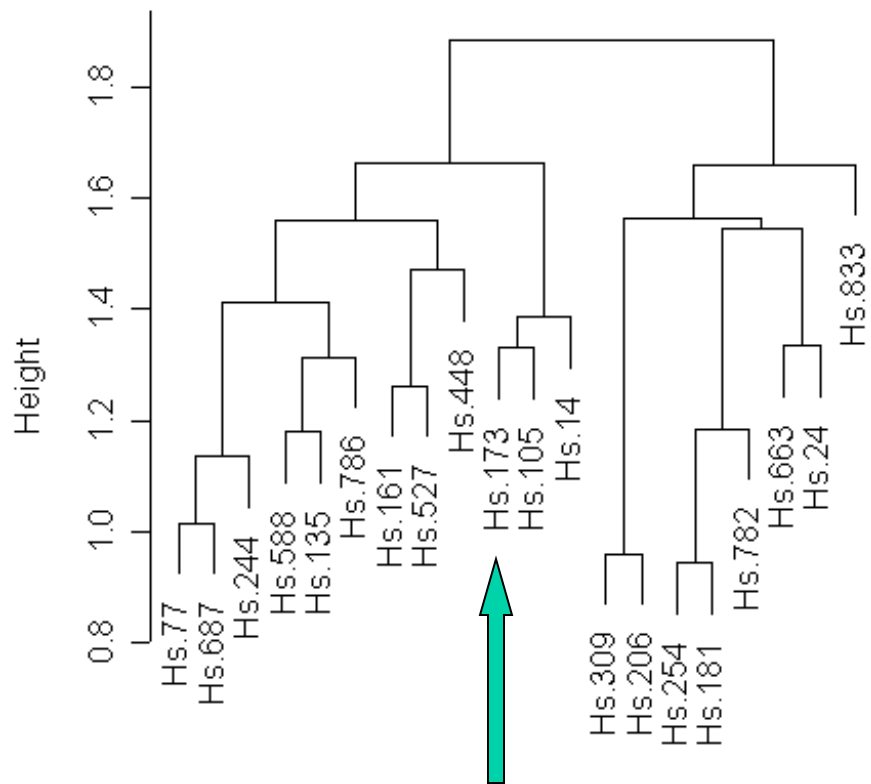
3. Choose clustering direction (top-down or bottom-up)

- **Agglomerative clustering (bottom-up)**
 - Starts with as each gene in its own cluster
 - Joins the two most similar clusters
 - Then, joins next two most similar clusters
 - Continues until all genes are in one cluster
- **Divisive clustering (top-down)**
 - Starts with all genes in one cluster
 - Choose split so that genes in the two clusters are most similar (maximize “distance” between clusters)
 - Find next split in same manner
 - Continue until all genes are in single gene clusters

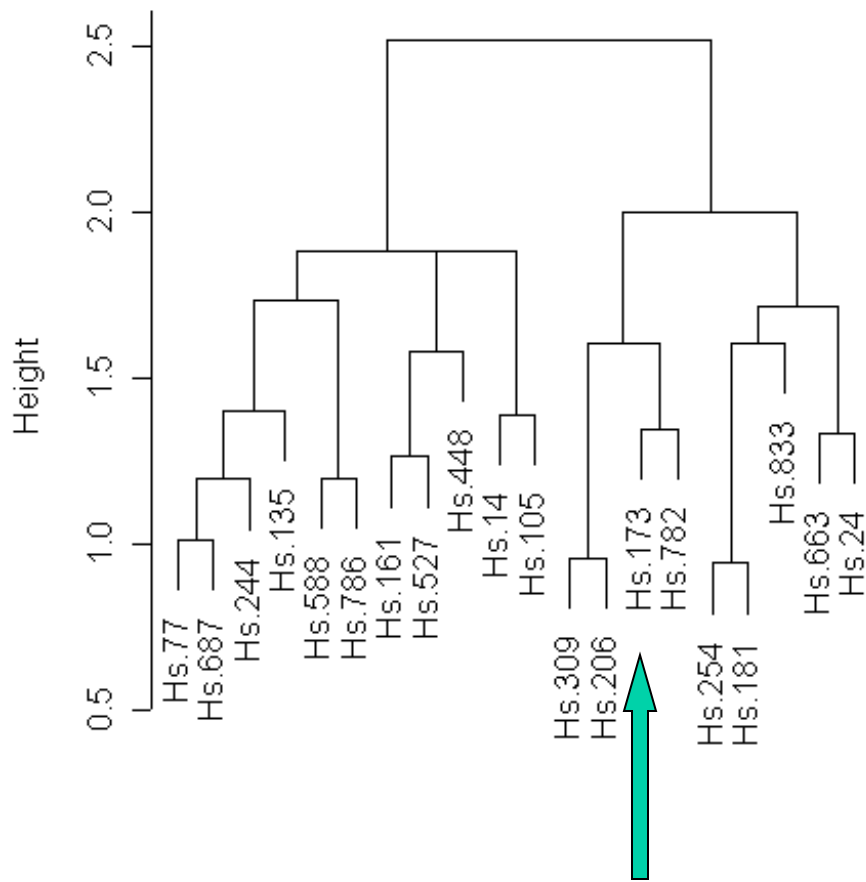
Which to use?

- **Both are only ‘step-wise’ optimal: at each step the optimal split or merge is performed**
- **This does not imply that the final cluster structure is optimal!**
- **Agglomerative/Bottom-Up**
 - **Computationally simpler, and more available.**
 - **More “precision” at bottom of tree**
 - **When looking for small clusters and/or many clusters, use agglomerative**
- **Divisive/Top-Down**
 - **More “precision” at top of tree.**
 - **When looking for large and/or few clusters, use divisive**
- **In gene expression applications, divisive makes more sense.**
- **Results ARE sensitive to choice!**

C: Agglom,Cor,Average

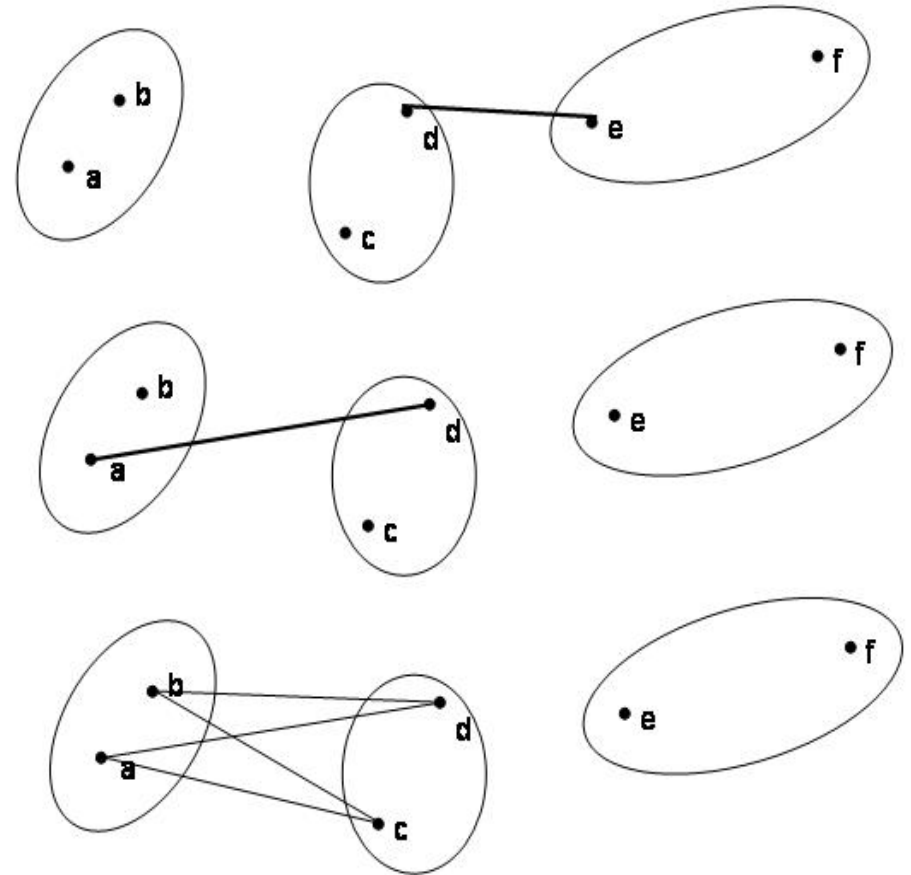


G: Div,Cor



4. Choose linkage method (if bottom-up)

- **Single Linkage:** join clusters whose distance between closest genes is smallest (elliptical)
- **Complete Linkage:** join clusters whose distance between furthest genes is smallest (spherical)
- **Average Linkage:** join clusters whose average distance is the smallest.



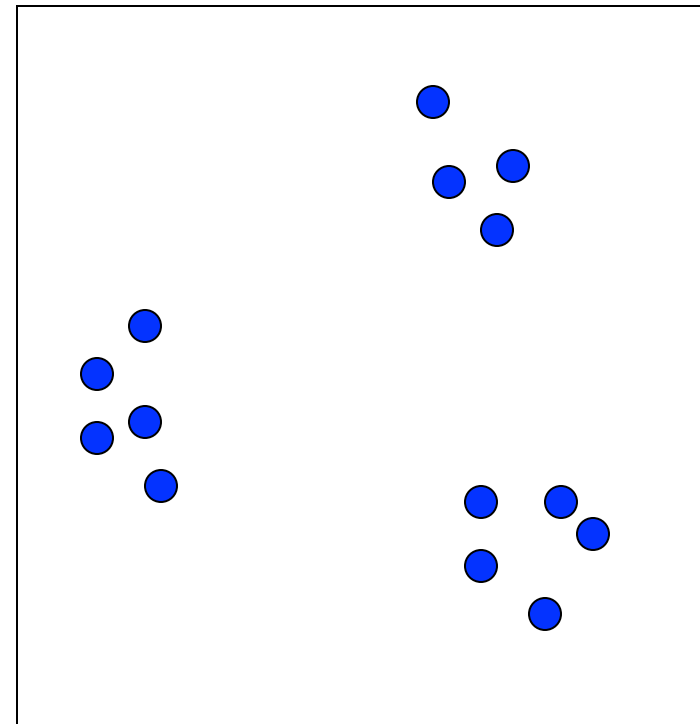
5. Calculate dendrogram

6. Choose height/number of clusters for interpretation

- **In gene expression, we don't see "rule-based" approach to choosing cutoff very often.**
- **Tend to look for what makes a good story.**
- **There are more rigorous methods. (more later)**
- **"Homogeneity" and "Separation" of clusters can be considered. (Chen et al. Statistica Sinica, 2002)**
- **Other methods for assessing cluster fit can help determine a reasonable way to "cut" your tree.**

K-means

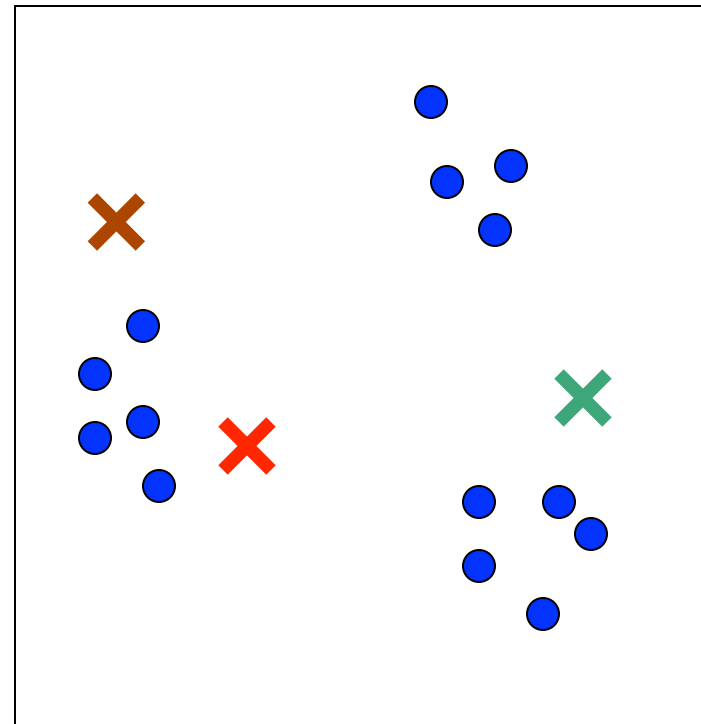
- **We start with some data**
- **Interpretation:**
 - We are showing expression for two samples for 14 genes
 - We are showing expression for two genes for 14 samples
- **This is simplification**



Iteration = 0

K-means

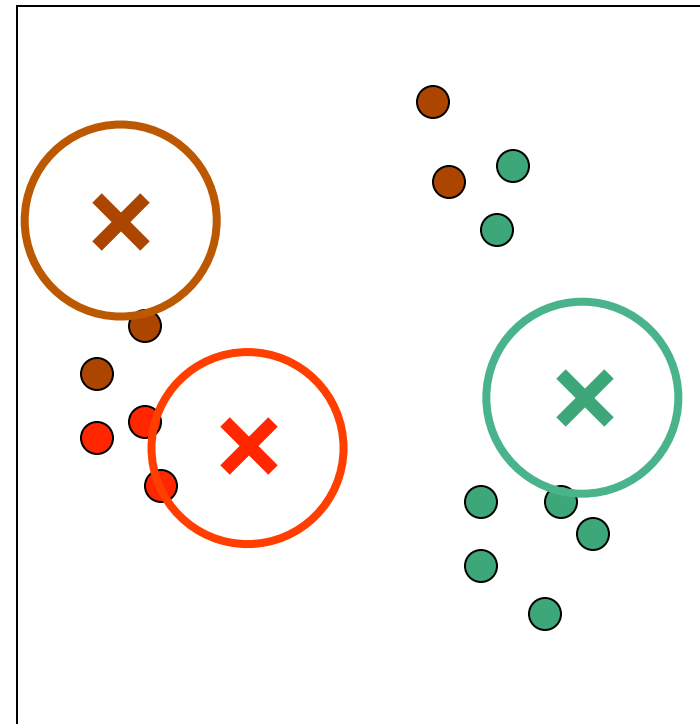
- Choose K *centroids*
- These are starting values that the user picks.
- There are some data driven ways to do it



Iteration = 0

K-means

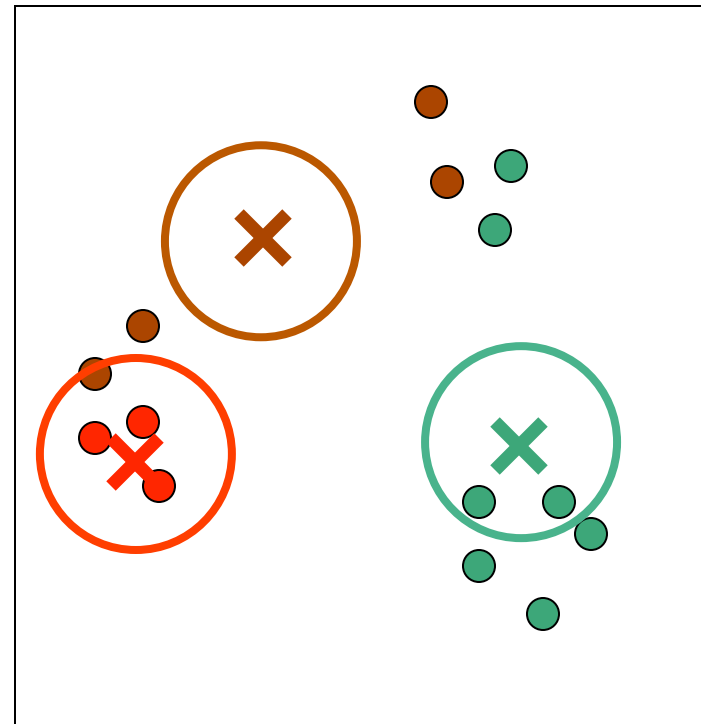
- Make first *partition* by finding the closest centroid for each point
- This is where distance is used



Iteration = 1

K-means

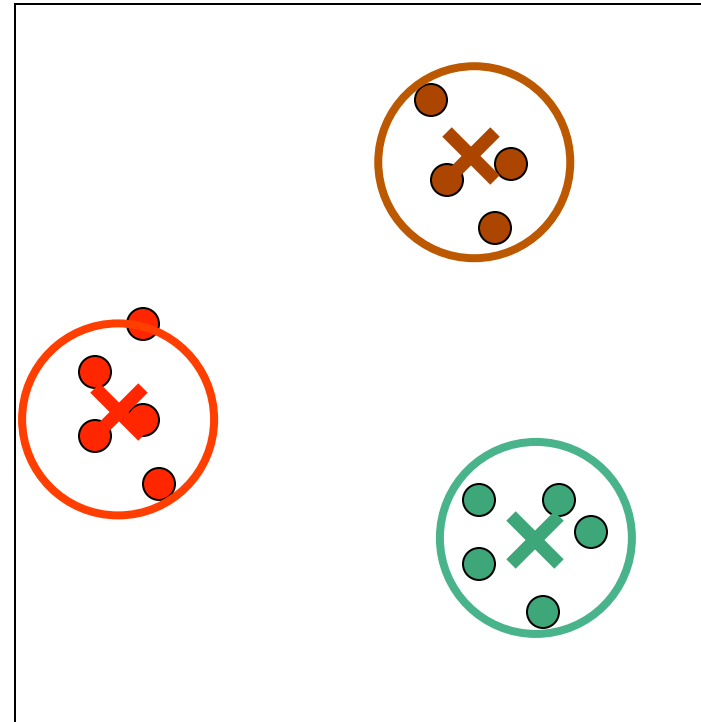
- Now re-compute the centroids by taking the *middle* of each cluster



Iteration = 2

K-means

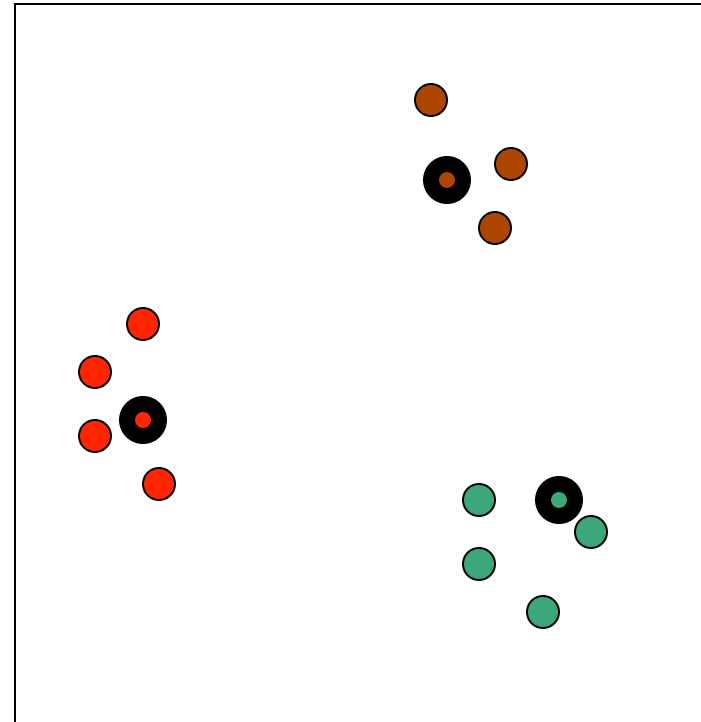
- Repeat until the centroids stop moving or until you get tired of waiting



Iteration = 3

K-medoids

- **A little different**
- **Centroid: The average of the samples within a cluster**
- **Medoid: The “representative object” within a cluster.**
- **Initializing requires choosing medoids at random.**

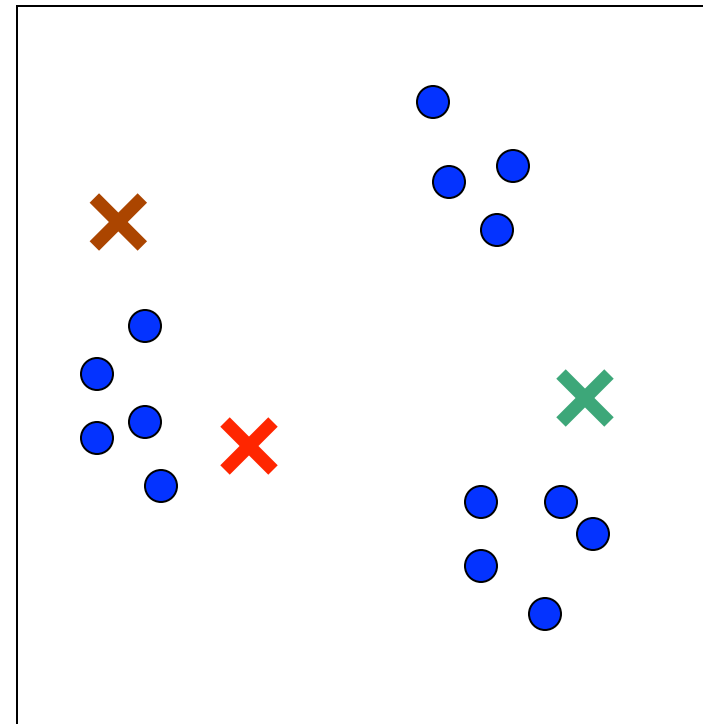


K-means Limitations

- **Final results depend on starting values**
- **How do we chose K? There are methods but not much theory saying what is best.**
- **Where are the pretty pictures?**

Model-Based Clustering

- Choose K *centroids*
- These are starting values that the user picks.
- There are some data driven ways to do it

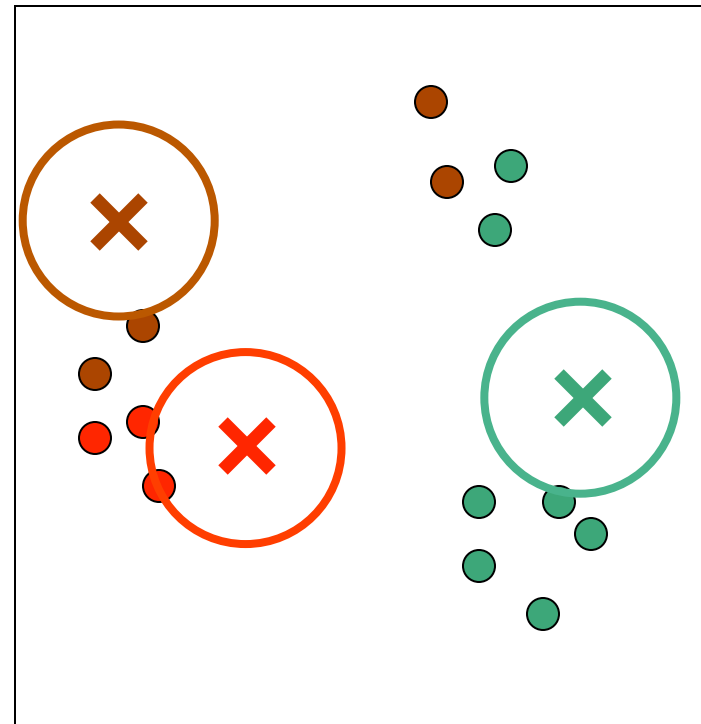


Iteration = 0

Model-Based Clustering

- No *partitions* now
- Assumption:
 - Each cluster can be modeled by a parametric distribution

$$f_k(x) \sim N(\mu_k, \sigma^2 \mathbf{I})$$

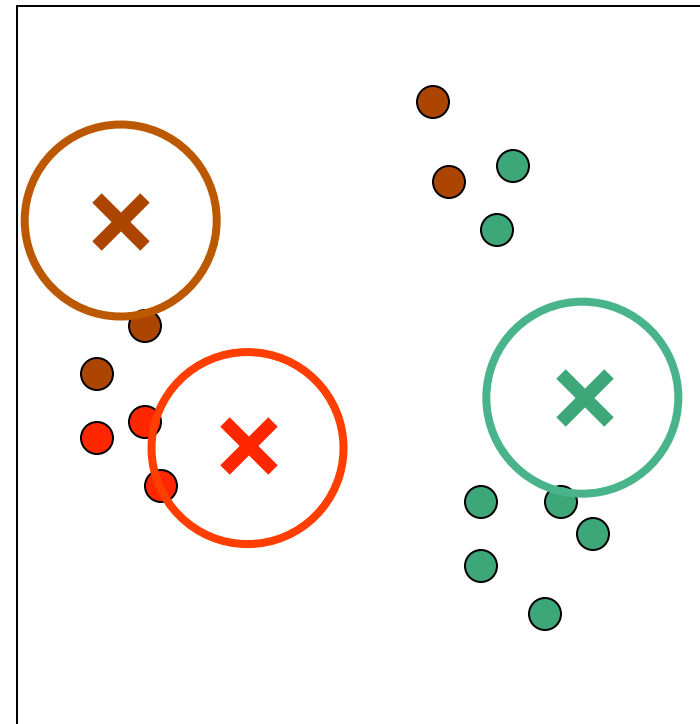


Iteration = 1

Model-Based Clustering

- No *partitions* now
- Points can be assigned to clusters with a *probability*

$$P(\text{cl}(x) = k | \Theta) = \frac{f_k(x)\pi_k}{\sum_l f_l(x)\pi_l}$$



Iteration = 1

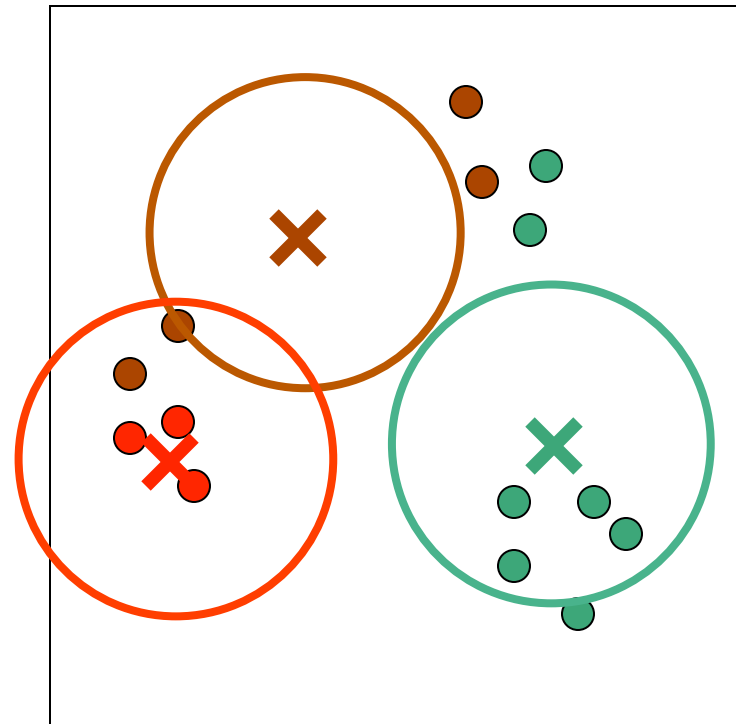
Model-Based Clustering

- Now re-compute the centroids by taking the *weighted mean* of each cluster

$$\hat{\mu}_k = \frac{\sum_i z_{ik} x_i}{\sum_i z_{ik}}$$

$$z_{ik} = P(\text{cl}(x_i) = k | \Theta)$$

- ~~re~~ **New: re-compute scale (σ^2) from a weighted variance**

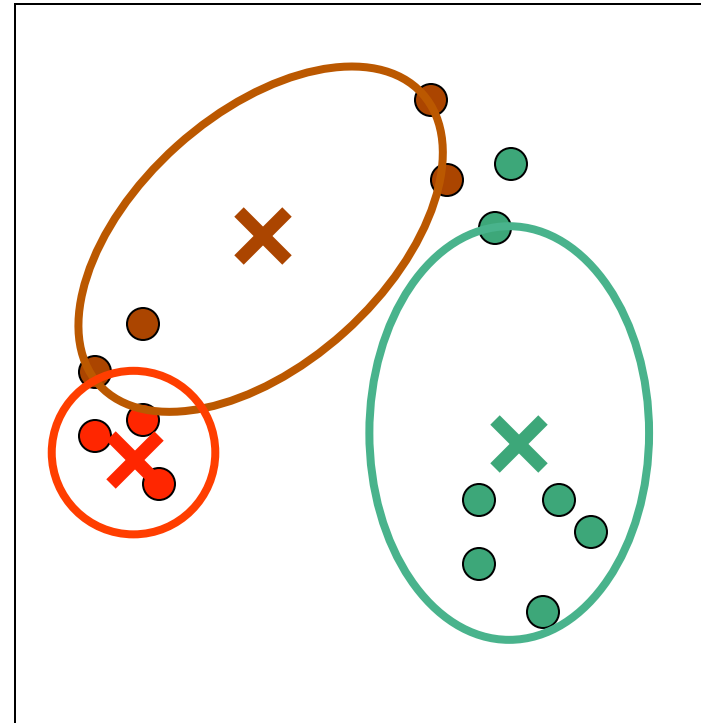


Iteration = 2

Model-Based Clustering

- **The general case:**

$$f_k(x) \sim N(\mu_k, \Sigma_k)$$



Iteration = 2

Model-Based Clustering

- Another way to look at it:

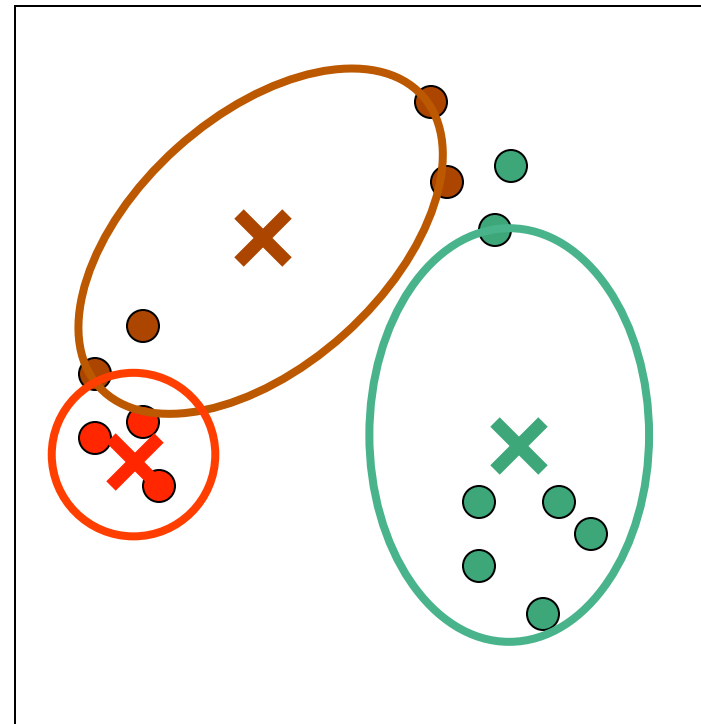
- we have data points

$$(x_i, \Delta_i)$$

- Δ_i is cluster assignment

- Which we *don't* observe

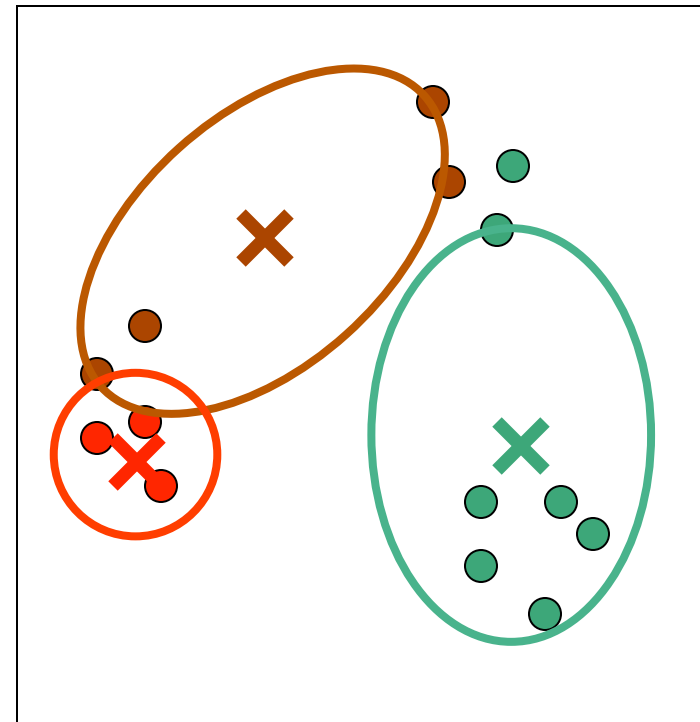
- i.e. *missing data*



Iteration = 2

Model-Based Clustering

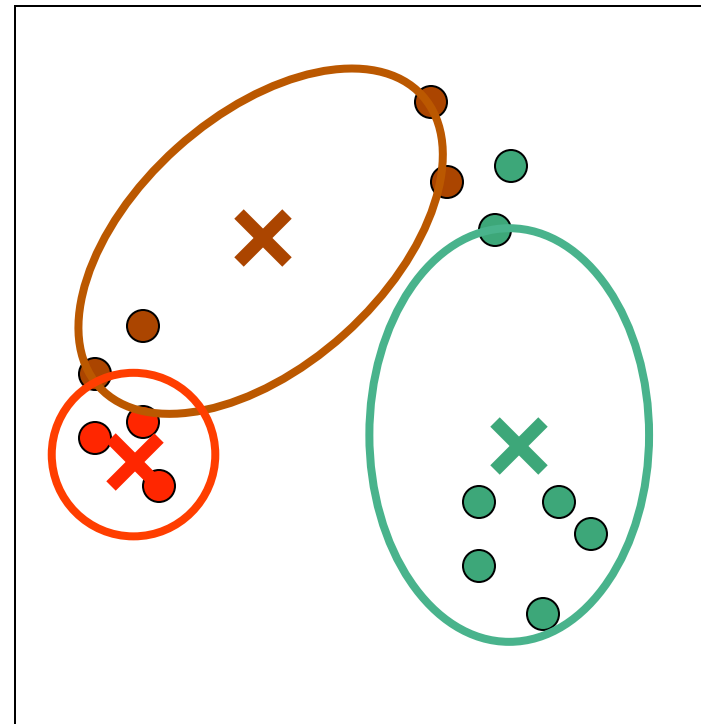
- Another way to look at it:
 - we have data points
 - (x_i, Δ_i)
 - Δ_i is cluster assignment
 - We just described the *EM* algorithm to get MLE's of means and variances



Iteration = 2

Model-Based Clustering

- Yet another (but similar) way to look at it:
 - $f(X)$ is a *mixture* of normals
- $$f(x) = \sum_k \pi_k f_k(x)$$
- **EM** is *one* algorithm to get MLEs of the mixture components



Iteration = 2

7. Assess cluster fit and stability

- **PART OF THE MISUNDERSTOOD!**
- **Most often ignored.**
- **Cluster structure is treated as reliable and precise**
- **BUT! Usually the structure is rather unstable, at least at the bottom.**
- **Can be VERY sensitive to noise and to outliers**
- **Homogeneity and Separation**
- **Cluster Silhouettes and Silhouette coefficient: how similar genes within a cluster are to genes in other clusters (composite separation and homogeneity) (more later with K-medoids) (Rousseeuw Journal of Computation and Applied Mathematics, 1987)**

Assess cluster fit and stability (continued)

- **WADP: Weighted Average Discrepant Pairs**
 - Bittner et al. Nature, 2000
 - Fit cluster analysis using a dataset
 - Add random noise to the original dataset
 - Fit cluster analysis to the noise-added dataset
 - Repeat many times.
 - Compare the clusters across the noise-added datasets.
- **Consensus Trees**
 - Zhang and Zhao Functional and Integrative Genomics, 2000.
 - Use parametric bootstrap approach to sample new data using original dataset
 - Proceed similarly to WADP.
 - Look for nodes that are in a “majority” of the bootstrapped trees.
- **More not mentioned.....**

Careful though....

- **Some validation approaches are more suited to some clustering approaches than others.**
- **Most of the methods require us to define number of clusters, even for hierarchical clustering.**
 - **Requires choosing a cut-point**
 - **If true structure is hierarchical, a cut tree won't appear as good as it might truly be.**

Final Thoughts

- **The most overused statistical method in gene expression analysis**
- **Gives us pretty red-green picture with patterns**
- **But, pretty picture tends to be pretty unstable.**
- **Many different ways to perform hierarchical clustering**
- **Tend to be sensitive to small changes in the data**
- **Provided with clusters of every size: where to “cut” the dendrogram is user-determined**