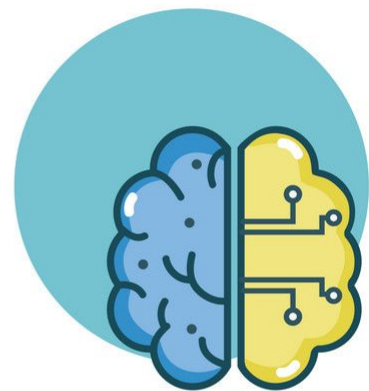


INTRODUCTION TO MACHINE LEARNING

CLUSTERING

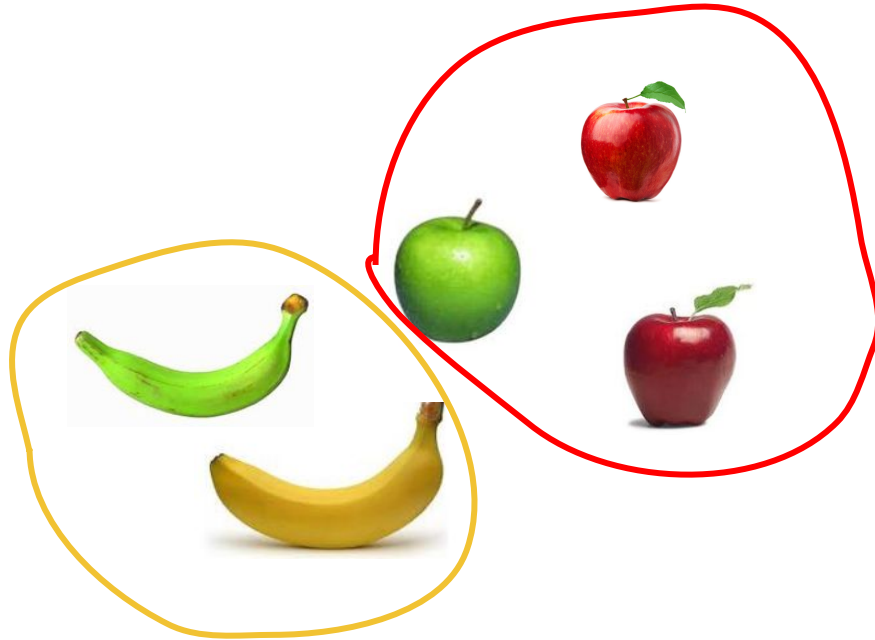


Elisa Ricci



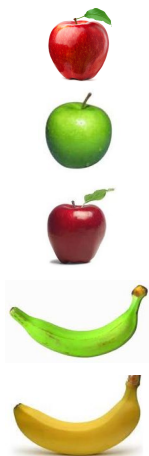
UNSUPERVISED LEARNING: CLUSTERING

Unsupervised learning setting: we are given data, i.e. examples, but no labels



UNSUPERVISED LEARNING: CLUSTERING

Raw data



features

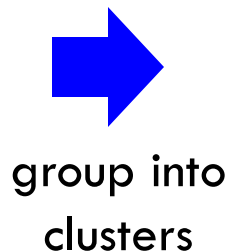
$f_1, f_2, f_3, \dots, f_n$

$f_1, f_2, f_3, \dots, f_n$

$f_1, f_2, f_3, \dots, f_n$

$f_1, f_2, f_3, \dots, f_n$

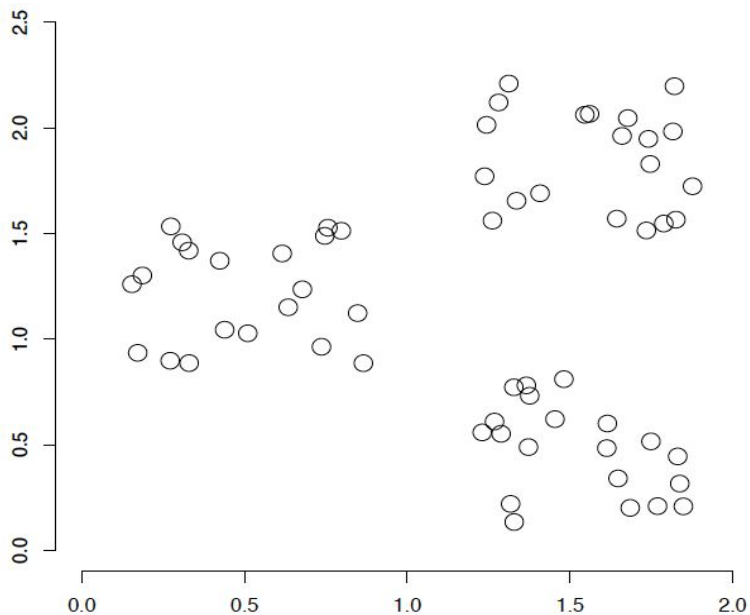
$f_1, f_2, f_3, \dots, f_n$



No supervision: we are only given data and want to find groupings

UNSUPERVISED LEARNING: CLUSTERING

Clustering: the process of grouping a set of objects into classes of similar objects

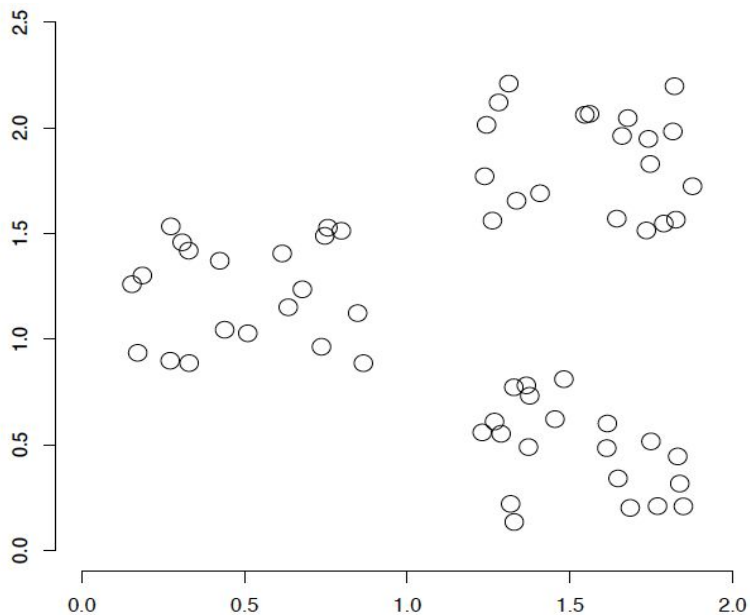


Which clustering algorithms can we use?

What are some of the issues for clustering?

UNSUPERVISED LEARNING: CLUSTERING

Clustering: the process of grouping a set of objects into classes of similar objects



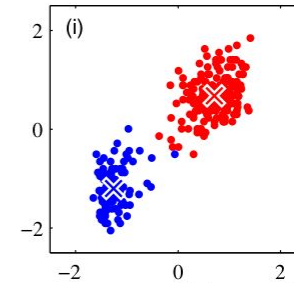
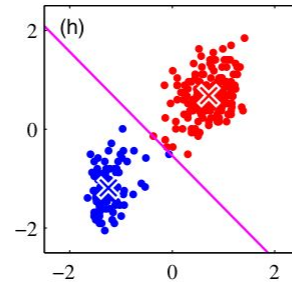
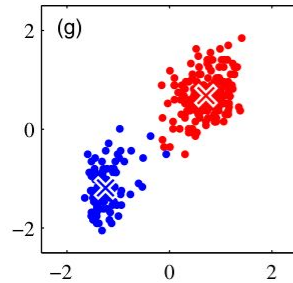
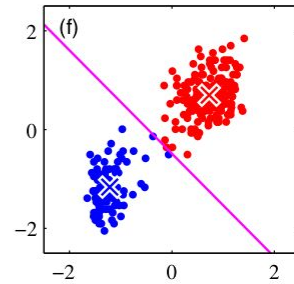
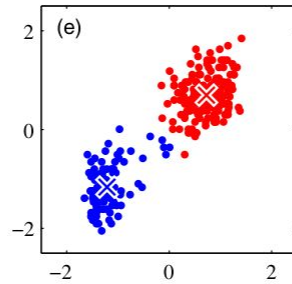
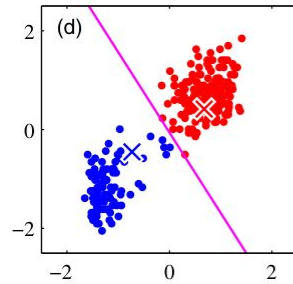
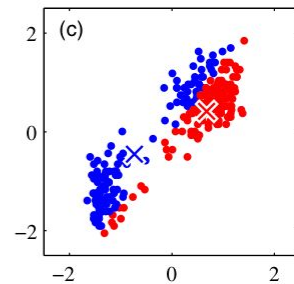
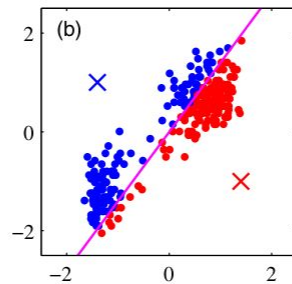
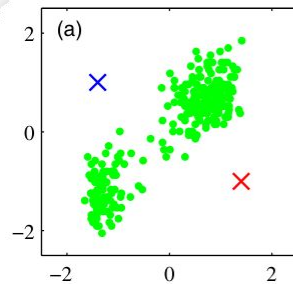
Which clustering
algorithms can we use?

K-MEANS

Most well-known and popular clustering algorithm.

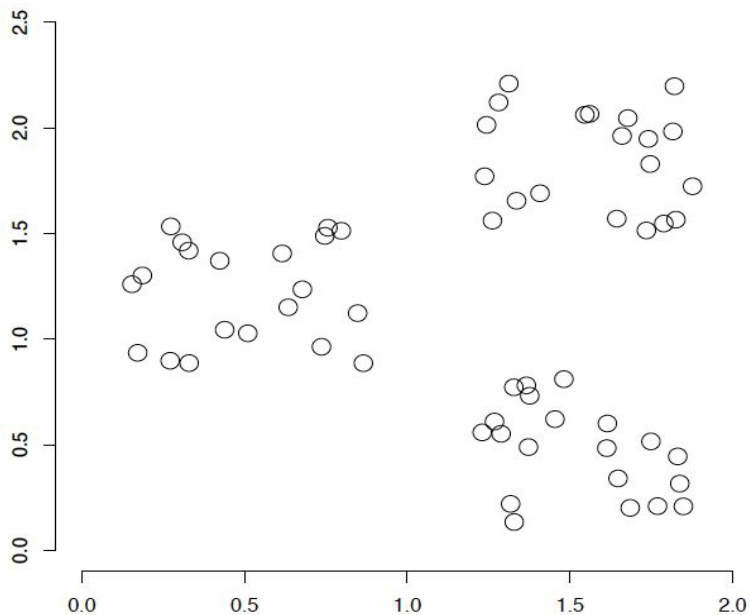
- Start with some initial cluster centroids
- Iterate:
 - Assign each example to closest centroid
 - Recalculate centers as the mean of the points in a cluster

K-MEANS



UNSUPERVISED LEARNING: CLUSTERING

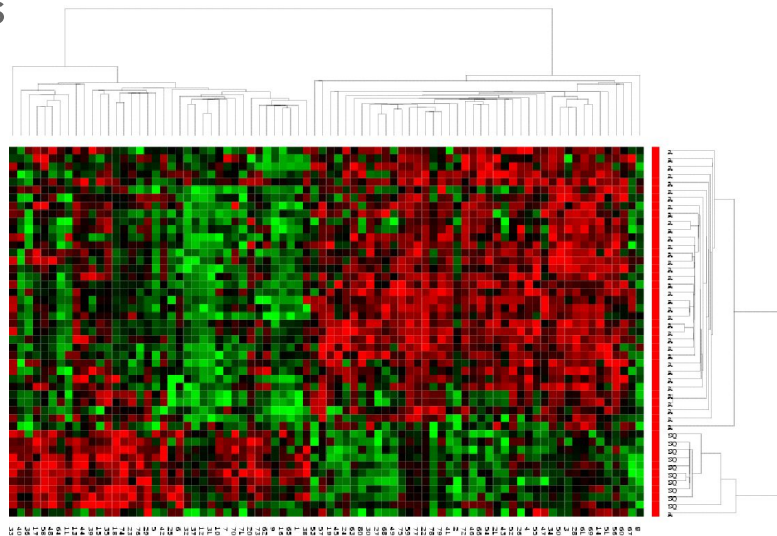
Clustering: the process of grouping a set of objects into classes of similar objects



What are some of the issues for clustering?

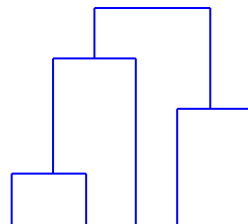
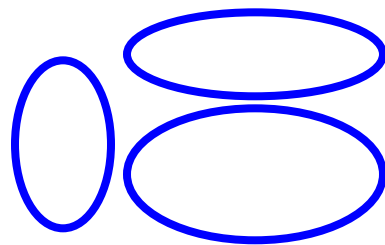
ISSUES FOR CLUSTERING

- Representation: how do we represent examples?
 - features?
- Similarity/distance between examples
- Flat clustering or hierarchical
- Number of clusters
 - Fixed a priori
 - Data driven



CLUSTERING ALGORITHMS

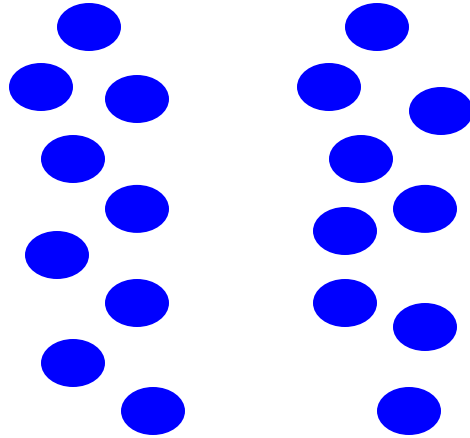
- **Flat algorithms**
 - K-means clustering (usually start with a random partitioning and refine it iteratively)
 - Spectral clustering
- **Hierarchical algorithms**
 - Bottom-up, agglomerative
 - Top-down, divisive



HARD VS. SOFT CLUSTERING

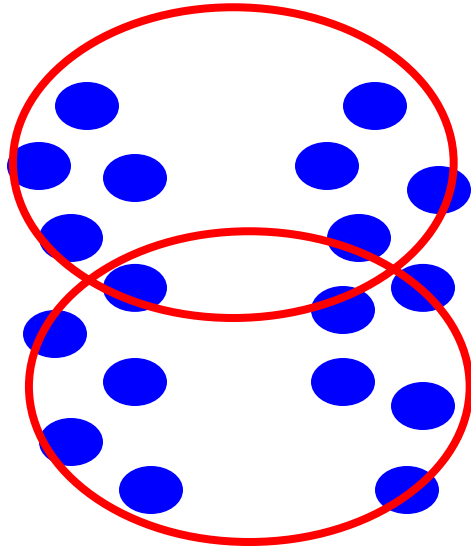
- **Hard clustering:** Each example belongs to exactly one cluster
 - K-Means
- **Soft clustering:** An example can belong to more than one cluster (probabilistic)
 - Makes more sense for applications like creating browsable hierarchies, e.g. you may want to put a pair of sneakers in two clusters: (i) sports apparel and (ii) shoes

PROBLEMS WITH K-MEANS



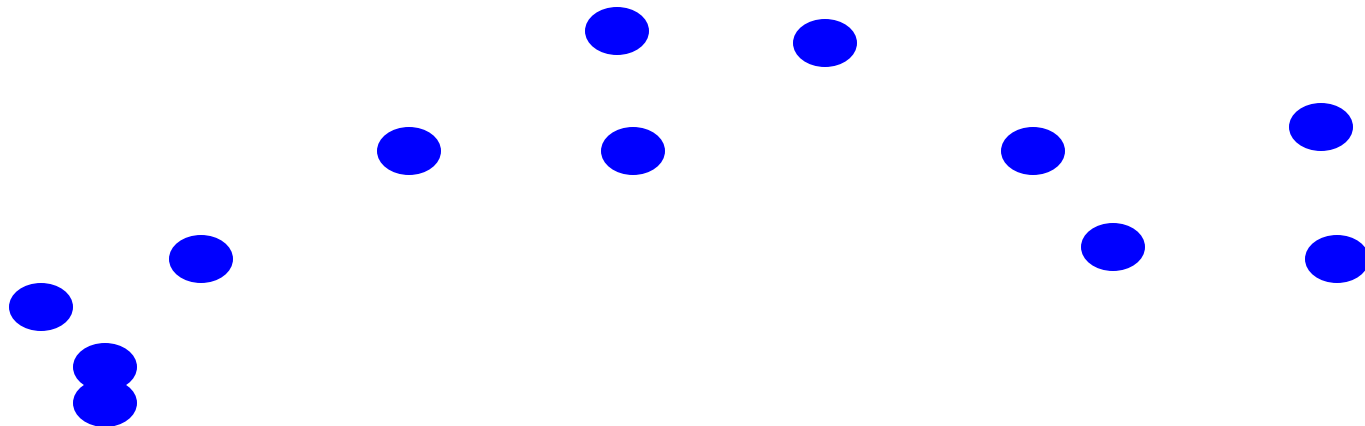
What would K-means give us here?

PROBLEMS WITH K-MEANS

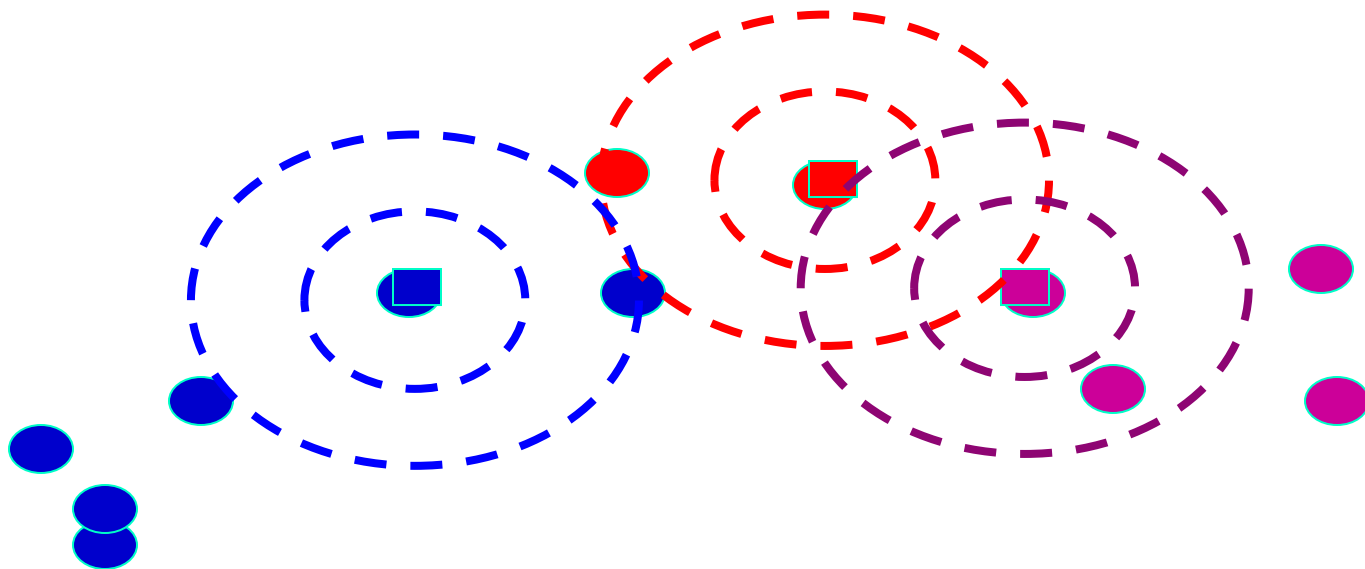


k-means assumes spherical clusters!

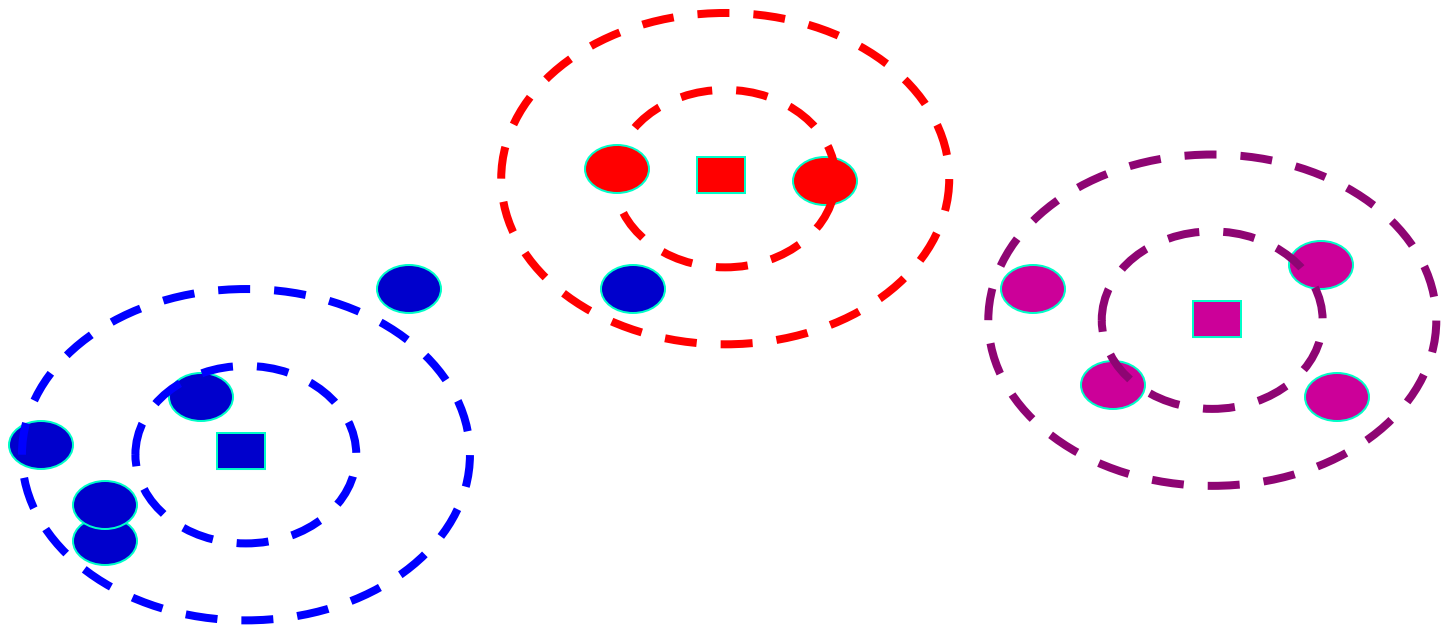
SPHERICAL CLUSTERS



SPHERICAL CLUSTERS



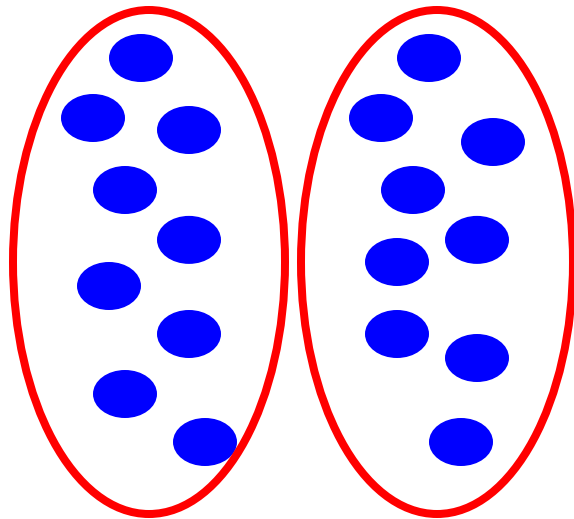
SPHERICAL CLUSTERS



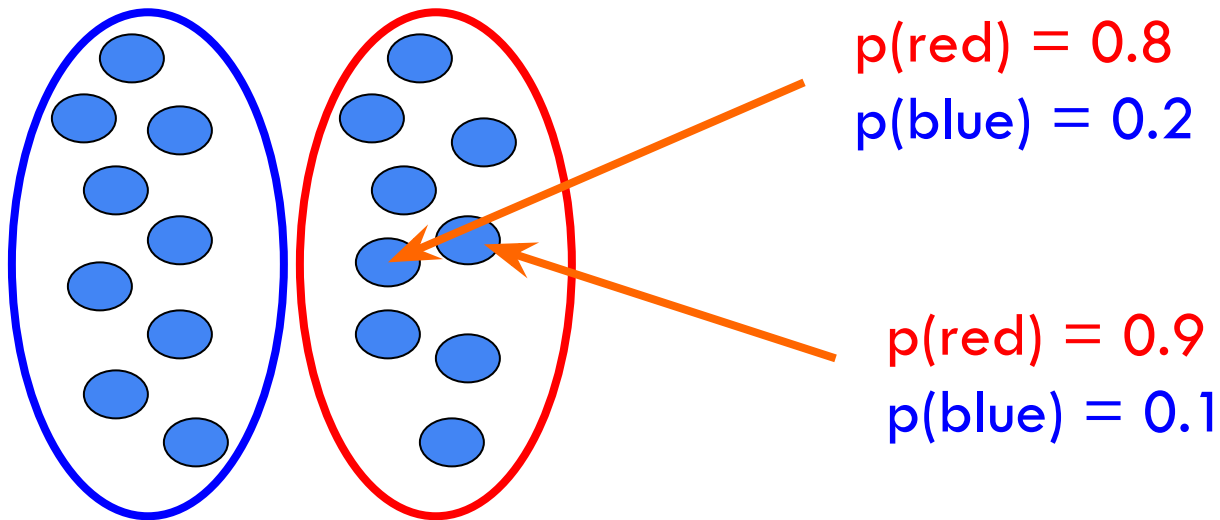
Iteratively learning a collection of spherical clusters

A BETTER ALGORITHM: EM CLUSTERING

Expectation Maximization (EM): assume data came from a **mixture of Gaussians** (elliptical data), assign data to cluster with a certain probability (**soft clustering**)



SOFT CLUSTERING



EM CLUSTERING

- Very similar at a high-level to K-means
 - Iterate between assigning points and recalculating cluster centers
- Two main differences between K-means and EM clustering
 - We assume **elliptical clusters** (instead of spherical)
 - It is a **soft clustering** algorithm

EM CLUSTERING

- Start with some initial cluster centers
- Iterate:
 - Soft assign points to each cluster

Calculate: $p(\theta_c | x)$

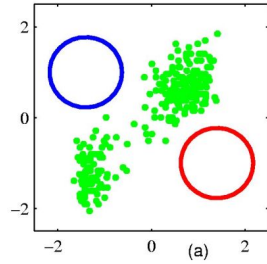
the probability of each point belonging to each cluster

- Recalculate the cluster centers

Calculate new cluster parameters θ_c

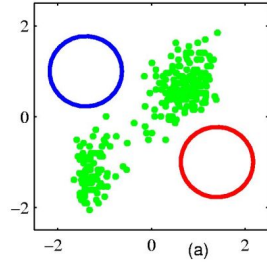
maximum likelihood cluster centers given the current soft clustering

EXAMPLE



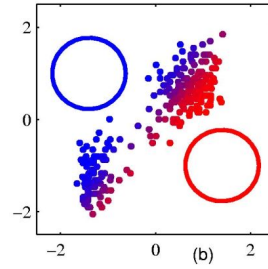
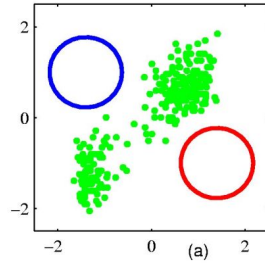
Start with some initial cluster centers

EXAMPLE



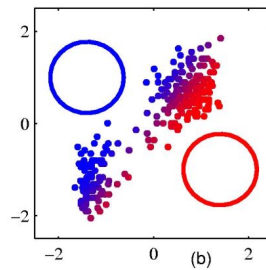
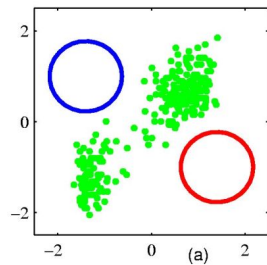
Which points belong to which clusters (soft)?

EXAMPLE



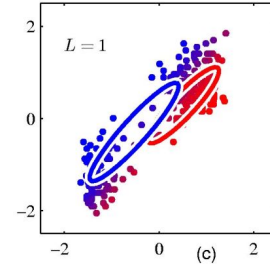
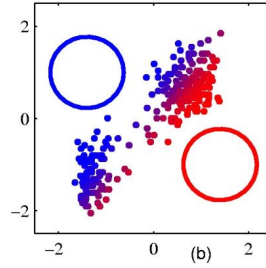
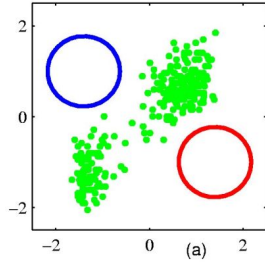
It is a soft (probabilistic) assignment

EXAMPLE



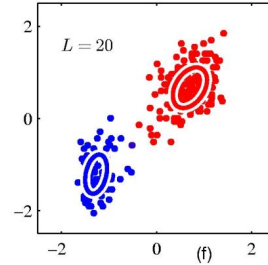
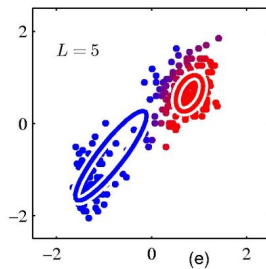
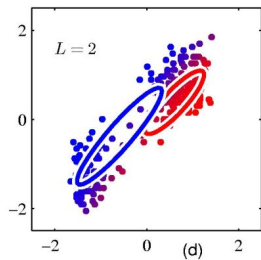
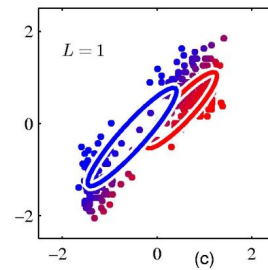
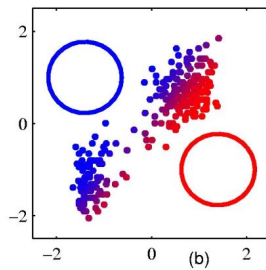
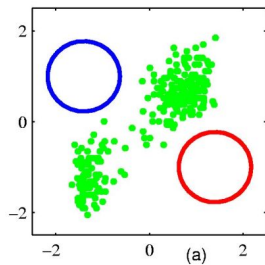
What do the new centers look like?

EXAMPLE



Cluster centers get a **weighted** contribution from points

EXAMPLE



... after several iterations

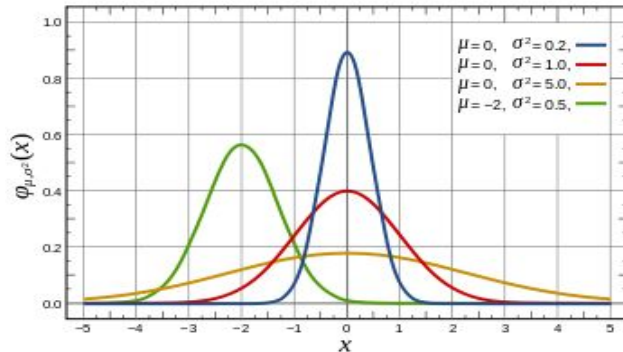
RECAP: MIXTURE OF GAUSSIANS

How do you define a Gaussian (i.e. ellipse)?

In 1D

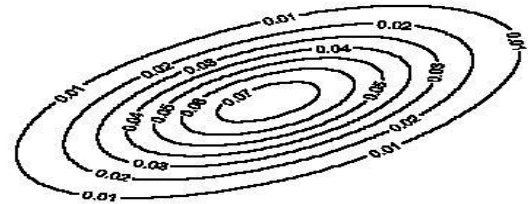
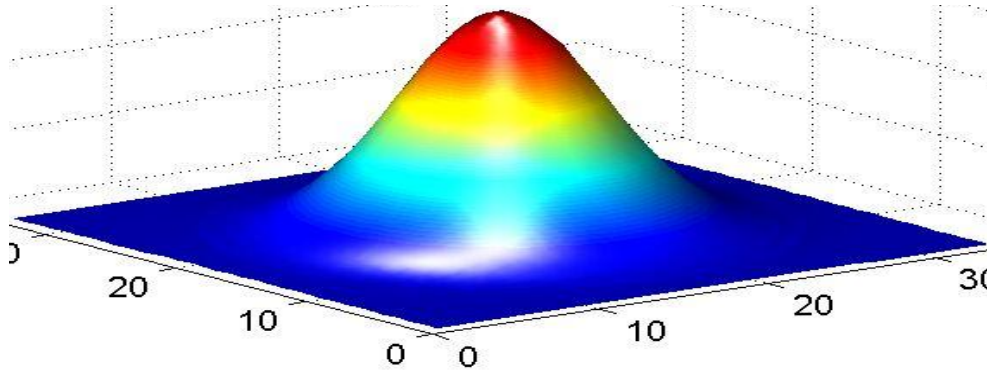
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

parameterized by the mean and the standard deviation



RECAP: MIXTURE OF GAUSSIANS

In 2D



Covariance determines
the shape of these contours

RECAP: MIXTURE OF GAUSSIANS

In general...

$$N[\mathbf{x}; \mu, \Sigma] = \frac{1}{(2\pi)^{d/2} \sqrt{\det(\Sigma)}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]$$

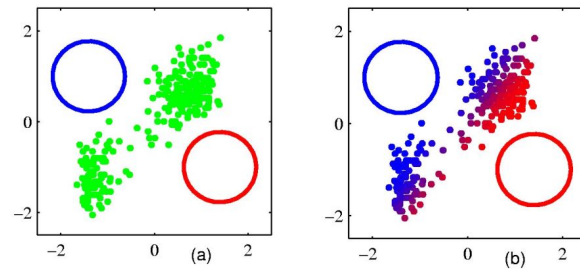
We learn the **means of each cluster** (i.e. the center) and the **covariance matrix** (i.e. how spread out it is in any given direction)

STEP 1: SOFT CLUSTER POINTS

Soft assign points to each cluster

- Calculate: $p(\theta_c | \mathbf{x})$

(the probability of each point belonging to each cluster)

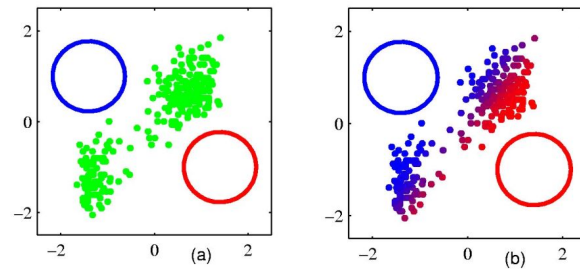


STEP 1: SOFT CLUSTER POINTS

Soft assign points to each cluster

- Calculate: $p(\theta_c | \mathbf{x})$

(the probability of each point belonging to each cluster)



How do we calculate these probabilities?

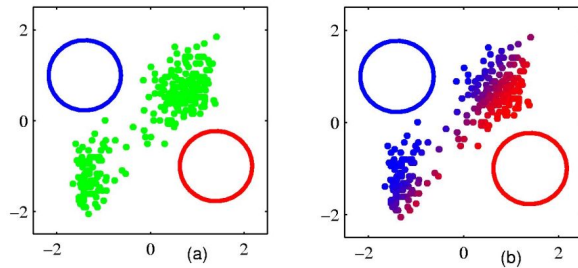
STEP 1: SOFT CLUSTER POINTS

Soft assign points to each cluster

- Calculate: $p(\theta_c|x)$

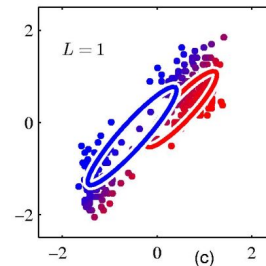
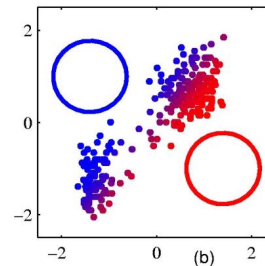
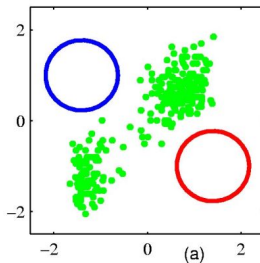
(the probability of each point belonging to each cluster)

Just plug into the probability the Gaussian equation for each cluster
(normalized to make a probability)



STEP 2: RECALCULATE CENTERS

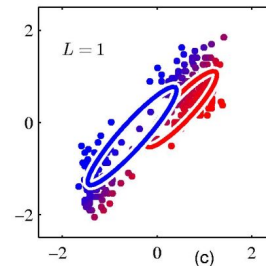
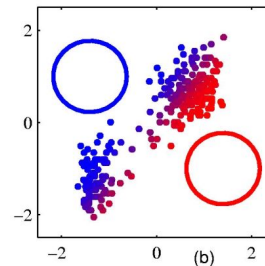
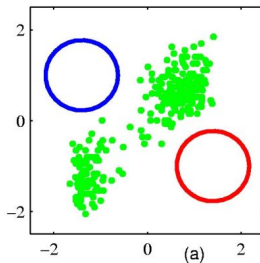
- Recalculate centers:
 - Calculate new cluster parameters θ_c
 - Maximum likelihood cluster centers given the current soft clustering



STEP 2: RECALCULATE CENTERS

- Recalculate centers:
 - Calculate new cluster parameters θ_c
 - Maximum likelihood cluster centers given the current soft clustering

How do we do that?



FITTING A GAUSSIAN: IDEA

- What is the “best”-fit Gaussian for this data?

10, 10, 10, 9, 9, 8, 11, 7, 6, ...

- Recall this is the 1-D Gaussian equation:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The Maximum Likelihood Estimation (MLE) is just the mean and variance of the data!

EM CLUSTERING

- EM stands for Expectation Maximization
 - **Expectation:** Given the current model, figure out the expected probabilities of the data points to each cluster

$$p(\theta_c|x)$$

What is the probability of each point belonging to each cluster?

- **Maximization:** Given the probabilistic assignment of all the points, estimate a new model θ_c

Maximum likelihood estimation

EM CLUSTERING

- Similar to K-Means
 - **Assign/cluster each point to closest center**

$p(\theta_c|x)$ Soft assignment

- **Recalculate centers as the mean of the points in a cluster**

Estimate a new model θ_c

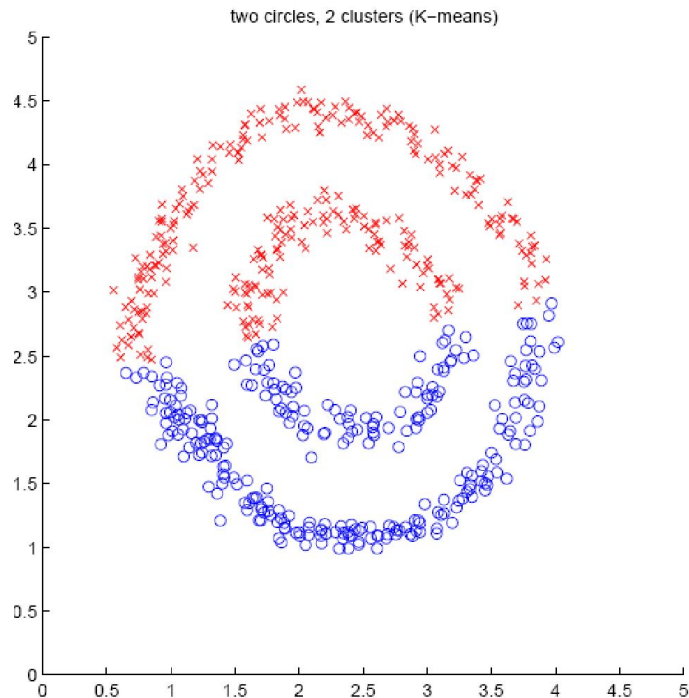
EXPECTATION MAXIMIZATION: SUMMARY

- Similar to K-Means
- Each iterations increases the likelihood of the data and is guaranteed to converge (though to a local optimum)
- Not just for clustering while K-means is just for clustering
- EM is a general purpose approach for training a model when you do not have labels

OTHER CLUSTERING ALGORITHMS

- K-means and EM-clustering are by far the most popular for clustering
- However, they cannot handle all clustering tasks
- What types of clustering problems cannot they handle?

NON-GAUSSIAN DATA



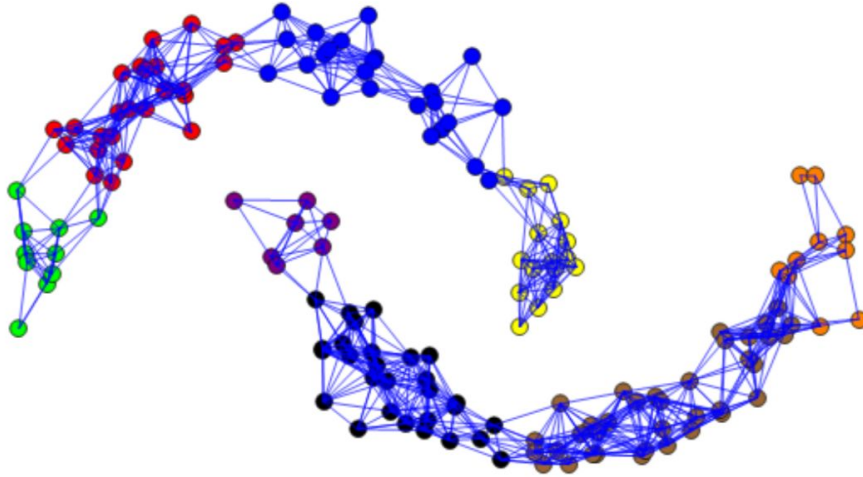
What is the problem?

Similar to classification:
global decision (linear model)
vs. local decision (K-NN)

Spectral clustering

SPECTRAL CLUSTERING

Group points based on links in a graph



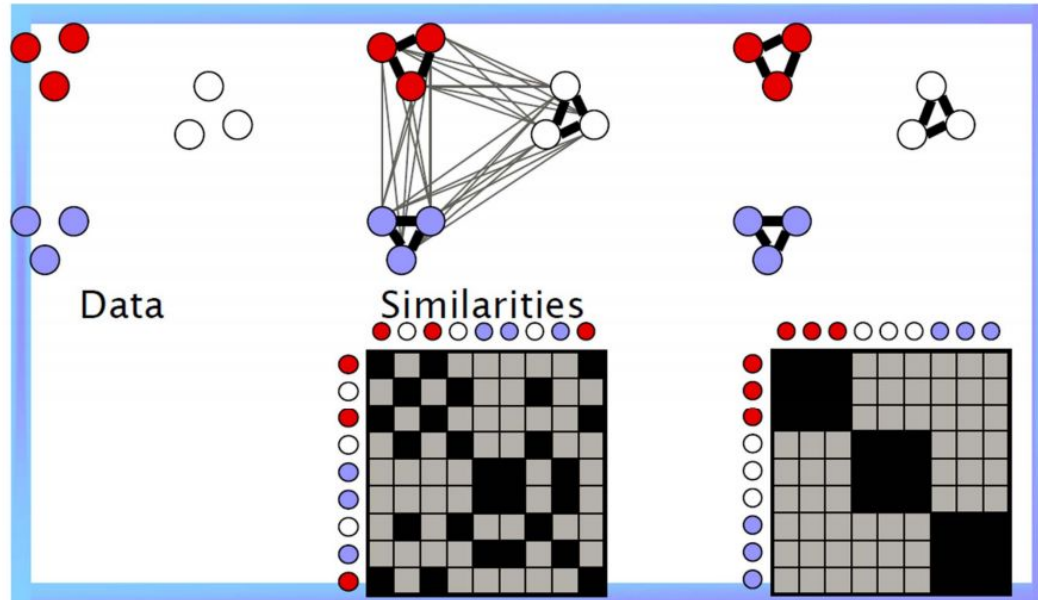
CREATE A GRAPH

- We can create a fully connected graph or a K-nearest neighbor graph (each node is only connected to its K closest neighbors)
- We can create a graph with some notion of similarity among nodes. It is common to use a Gaussian Kernel

$$W(i, j) = \exp \frac{-|x_i - x_j|^2}{\sigma^2}$$

GRAPH PARTITIONING

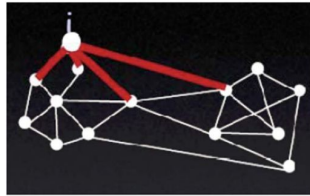
Define a similarity matrix and cut the graph



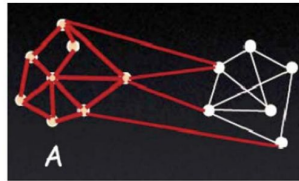
GRAPH TERMINOLOGY

Degree of a node and volume of a set

$$d_i = \sum_j w_{i,j}$$

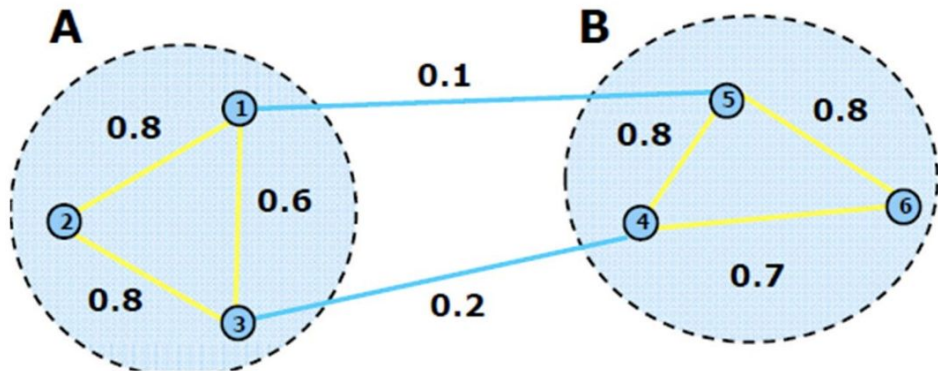


$$\text{vol}(A) = \sum_{i \in A} d_i, A \subseteq V$$



GRAPH CUT

- Consider a partition of the graph into two parts A and B
- $Cut(A, B)$: sum of the weights of the set of edges that connect the two groups
- An intuitive goal is **find the partition that minimizes the cut**



$$Cut(A, B) = 0.3$$

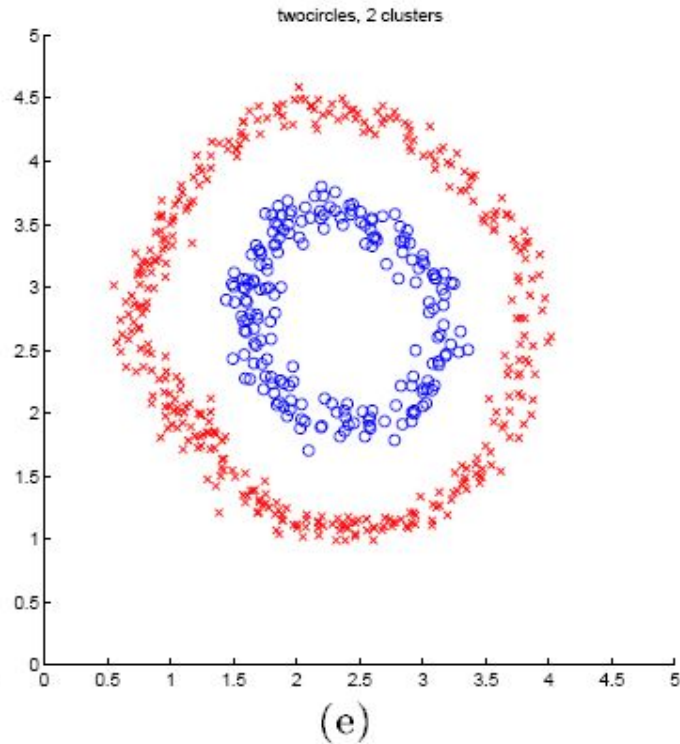
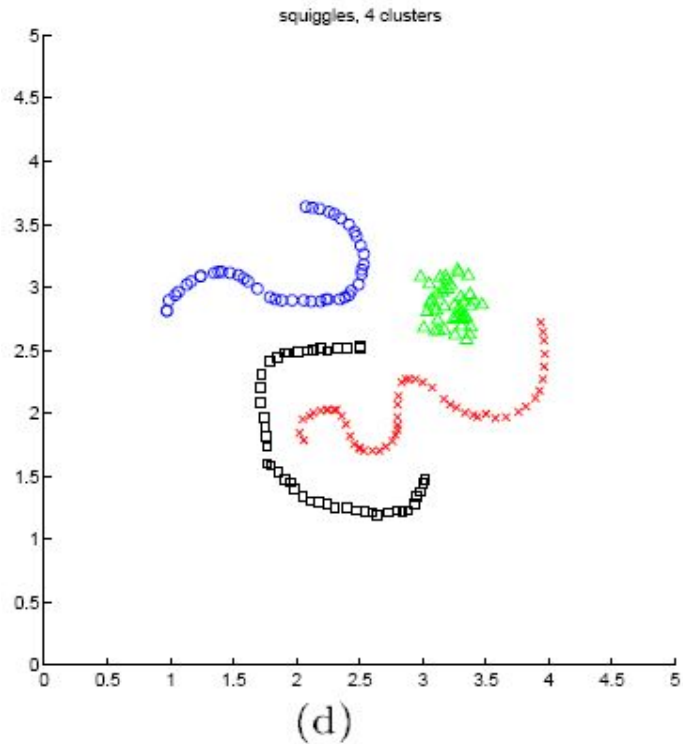
NORMALIZED CUT

- Minimize

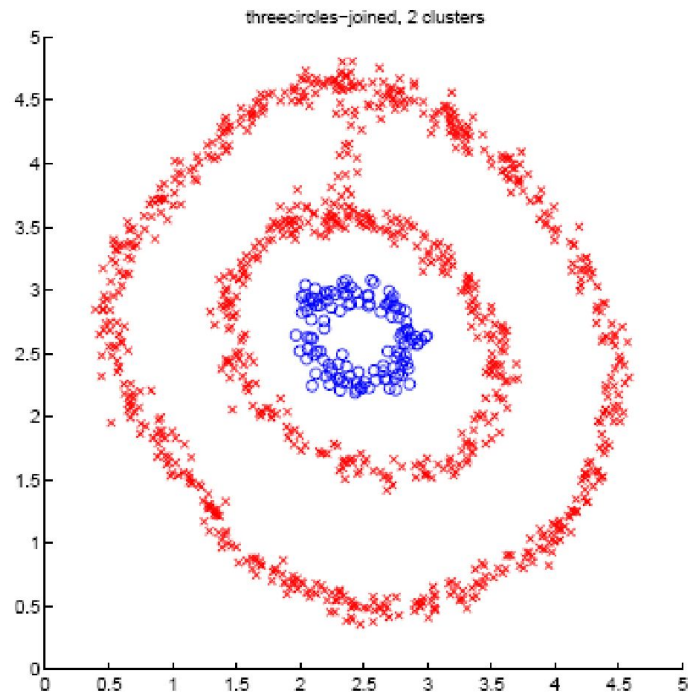
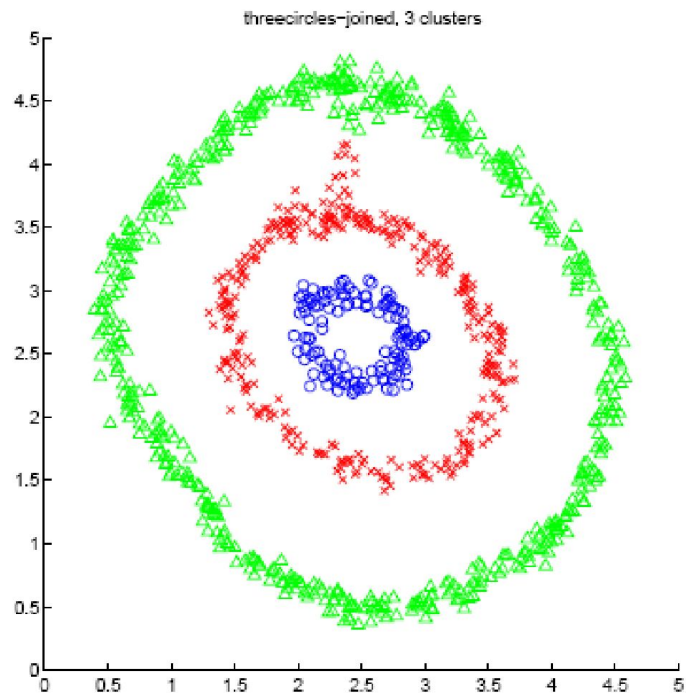
$$NCut(A, B) = Cut(A, B)/Vol(A) + Cut(A, B)/Vol(B)$$

- This problem can formalized as a discrete optimization problem and then relaxed in the continuous domain and become a generalized eigenvalue problem

SPECTRAL CLUSTERING

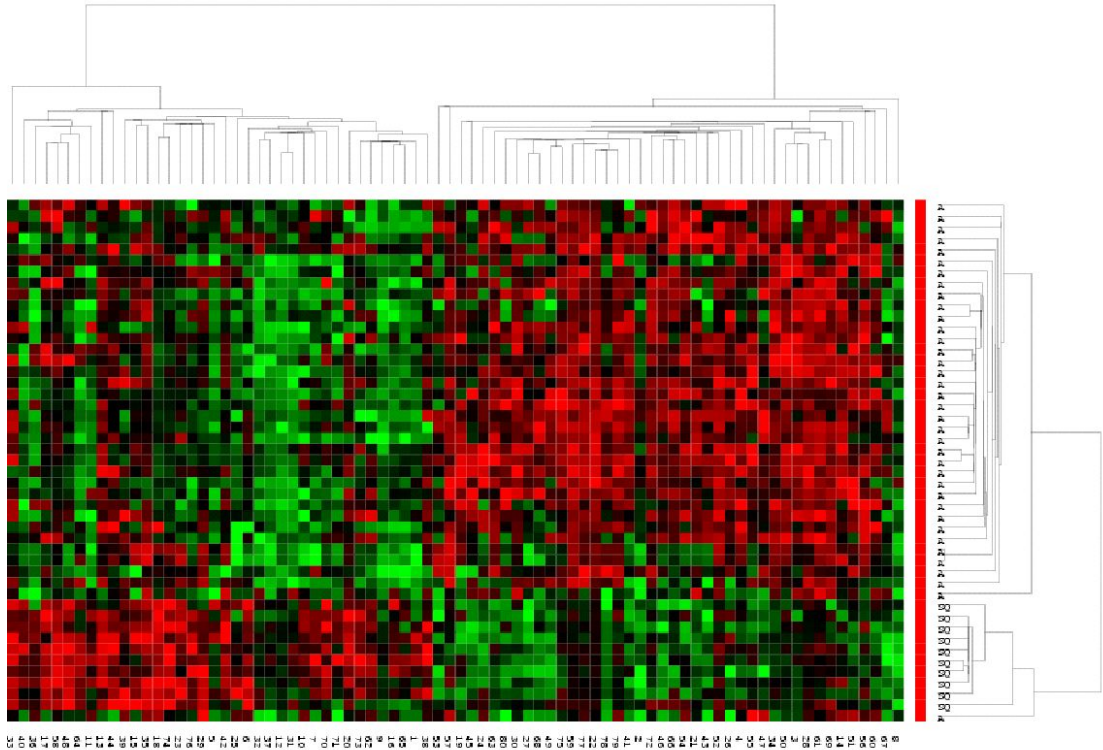


SPECTRAL CLUSTERING



HIERARCHICAL CLUSTERING

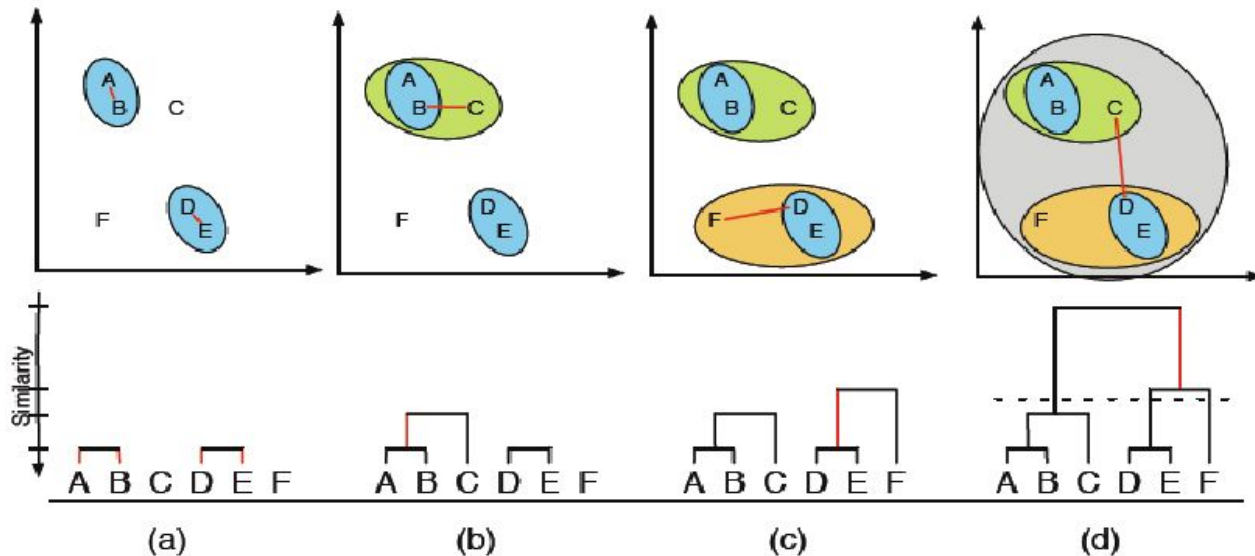
Gene profiles



HIERARCHICAL CLUSTERING

Produces a set of nested clusters organized as a hierarchical tree

The tree is called **dendrogram**

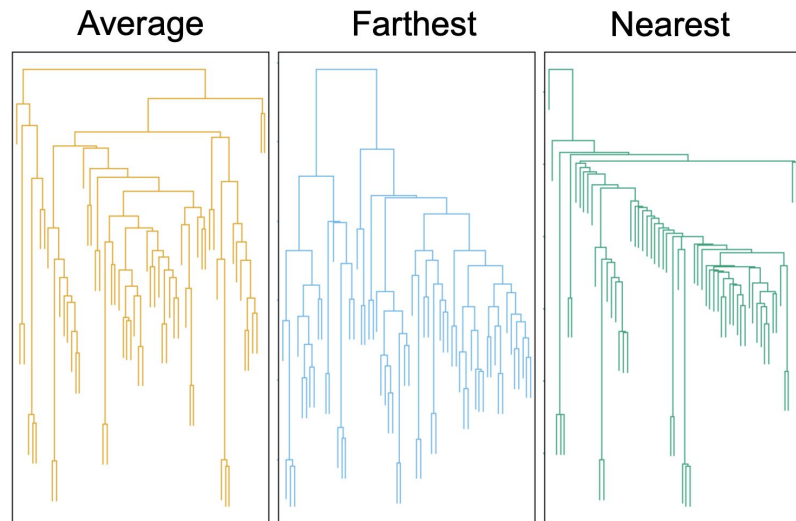
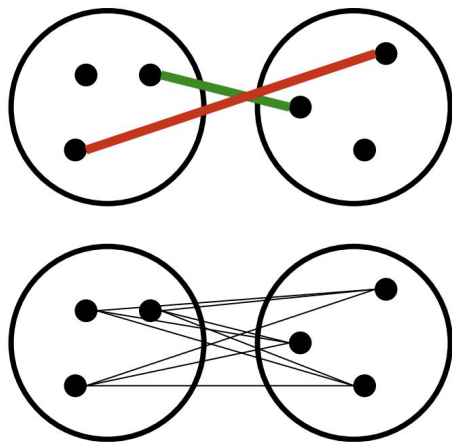


AGGLOMERATIVE CLUSTERING

- **Idea:** first merge very similar instances and then incrementally build larger clusters out of smaller clusters
- **Algorithm:**
 - Initially, each instance in its own cluster
 - Repeat:
 - Pick the two closest clusters
 - Merge them into a new cluster
 - Stop when there is only one cluster left
- Produces not one clustering, but **a family of clusterings** represented by a dendrogram

AGGLOMERATIVE CLUSTERING

- How should we define “closest” for clusters with multiple elements?
- Different choices (closest pair, farthest pair, average among all distances, etc.) creates different solutions



QUESTIONS?

