

CS6375: Machine Learning

Gautam Kunapuli

Clustering



THE UNIVERSITY OF TEXAS AT DALLAS

Erik Jonsson School of Engineering and Computer Science

Supervised vs Unsupervised Learning

Supervised learning: Given **labeled** data $(x_i, y_i), i = 1, \dots, n$, learn a function $f : x \rightarrow y$

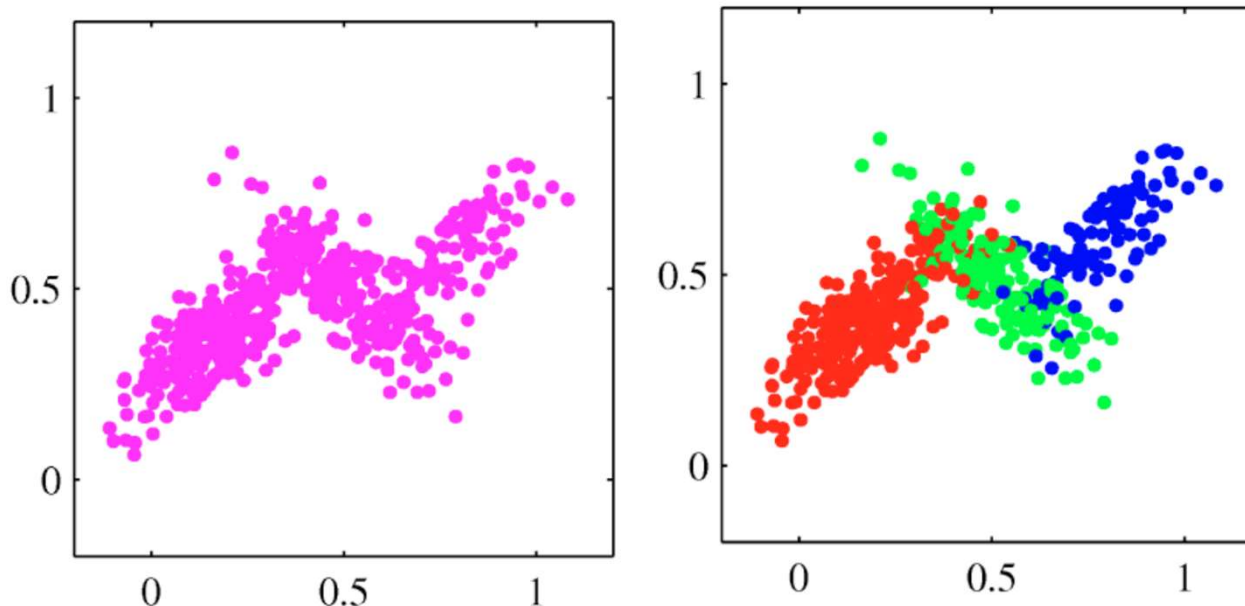
- Categorical y : **classification**
- Continuous y : **regression**

Unsupervised learning: Given **unlabeled** data $x_i, i = 1, \dots, n$, can we infer the underlying structure of X ?

Why do unsupervised learning?

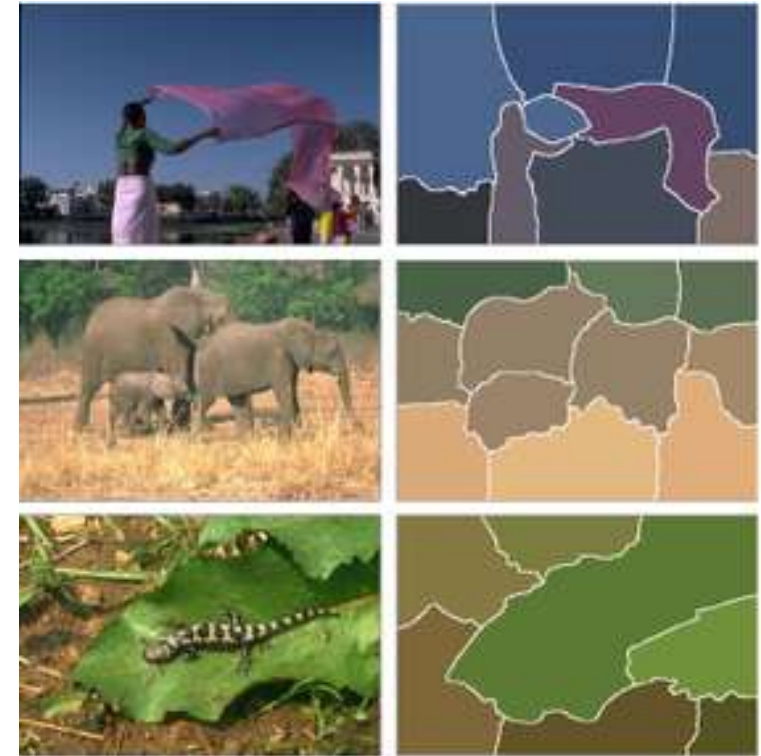
- raw data cheap; labeled data expensive
- save memory/computation.
- **reduce noise** in high-dimensional data
- useful in **exploratory data analysis**
- **pre-processing** for supervised learning
 - e.g., pca for dimensionality reduction

Basic Idea for Clustering: discover groups such that samples **within a group are more similar** to each other than samples across groups

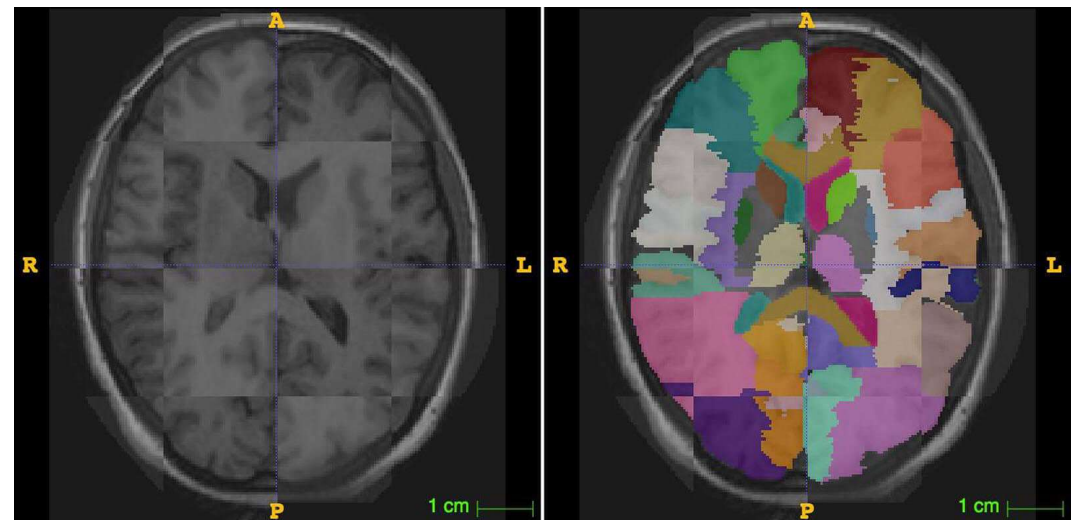


Example: Image Segmentation

Example: Partition a digital image of pixels into **segments** (also known as super-pixels). The goal of **segmentation** is to extract a **higher-order representation** that is more meaningful and (semantically) easier to analyze.



Medical image segmentation partitions a medical image into different meaningful segments. Segments often correspond to different tissues, organs, diseases, pathologies. Medical image segmentation is challenging due to low **contrast**, **noise**, **differences in individuals** etc.,



k-Means Clustering

Ingredients of k-means

- a **distance function** to identify the “closest” cluster centers
- a **loss function** to evaluate clusters
- an **algorithm** that optimizes this loss function

k-means Clustering

Given: Unlabeled data, $x_i, i = 1, \dots, n$

Initialize: Pick k random points as cluster centers

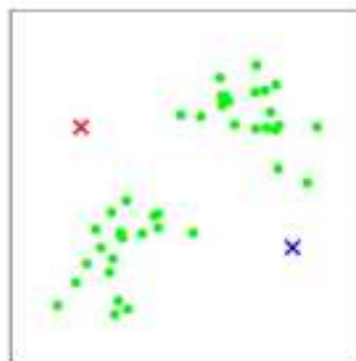
$\mu_j, j = 1, \dots, k$

while not converged do

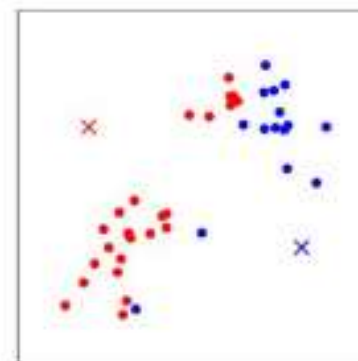
- **Assign** data points x_i to **closest** cluster center μ_j
- **Update** the cluster centers μ_j to the mean (average) of the points assigned to that cluster
- if the assignments no longer change, **converged = true**



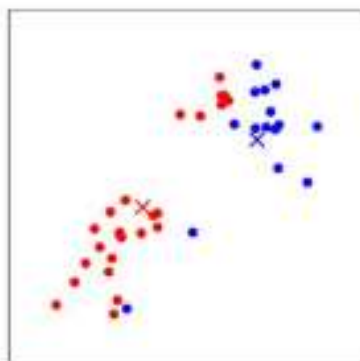
(a)



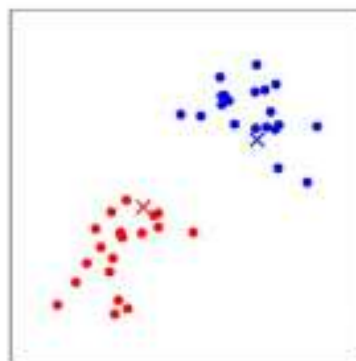
(b)



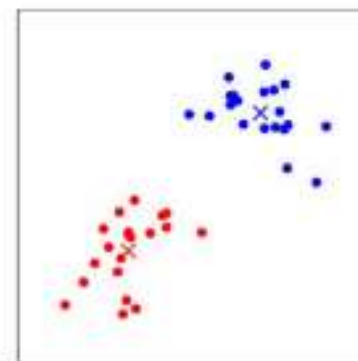
(c)



(d)



(e)

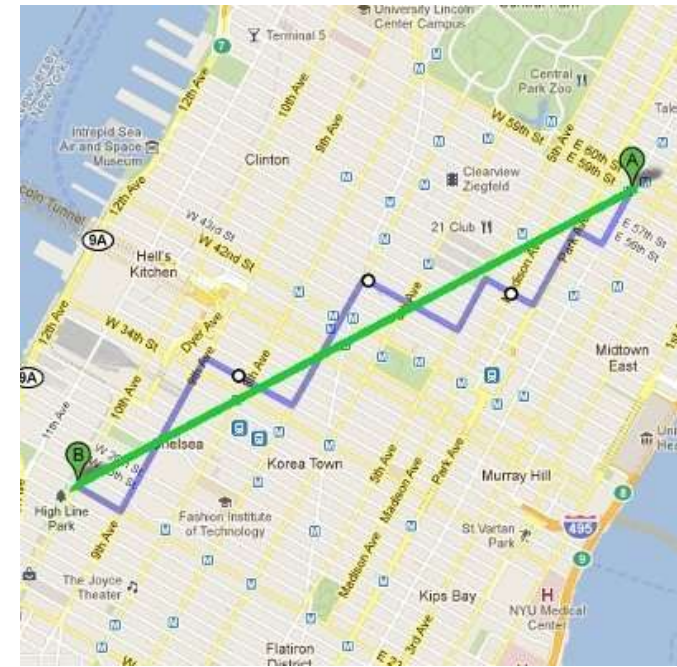


(f)

k-Means: Distance Functions

Properties of a distance function (also applies to *k*-nearest neighbors)

- **symmetry** $d(x, z) = d(z, x)$
 - if symmetry does not hold, we can say that x looks like z , but z does not look like x , which is not meaningful
 - Euclidean distance is symmetric, but KL divergence is not!
- **positivity**, $d(x, z) \geq 0$; and self-similarity $d(x, z) = 0$ if and only if $x = z$
 - if these do not hold, the distance function cannot tell apart two different objects, which is not useful
- **triangle inequality**: $d(a, b) + d(b, c) \geq d(a, c)$
 - if the triangle inequality does not hold, we can say a is like b and b is like c but a is not like c , which is not meaningful



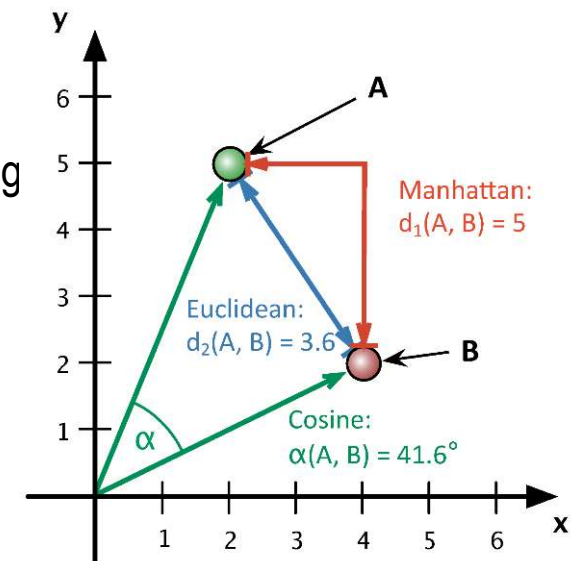
Two commonly-used distance measures produce two commonly used clustering algorithms. Consider a d -dimensional data set (d features) with points x, z

k-Means uses Euclidean distance:

$$d_2(x, z) = \|x - z\|_2 = \sqrt{(x_1 - z_1)^2 + \dots + (x_d - z_d)^2}$$

k-Medoids uses Manhattan distance:

$$d_1(x, z) = \|x - z\|_1 = |x_1 - z_1| + \dots + |x_d - z_d|$$



k-Means: Loss Function

The **key idea** behind **k-means clustering** is to find k clusters each described by a **prototype** (cluster centers) $\mu_j, j = 1, \dots, k$

- Assignment of training example x_i to clusters can be represented by **responsibilities** $r_{ij} \in \{0, 1\}$
 - $r_{ij} = 1$ if example x_i is assigned to the j -th cluster
 - need to ensure that $\sum_{j=1}^k r_{ij} = 1$ to ensure that the example x_i is assigned to one and only one cluster
 - cluster assignments for each data point can be read off the **responsibility matrix**
 - columns sum give us the **size of each cluster**

responsibilities for a data set with 6 examples and 3 clusters

| | cl ₁ | cl ₂ | cl ₃ |
|-------|-----------------|-----------------|-----------------|
| x_1 | 0 | 0 | 1 |
| x_2 | 0 | 1 | 0 |
| x_3 | 1 | 0 | 0 |
| x_4 | 0 | 0 | 1 |
| x_5 | 1 | 0 | 0 |
| x_6 | 0 | 1 | 0 |

The responsibilities can be used to define a **loss function**:

$$J = \sum_{i=1}^n \sum_{j=1}^k r_{ij} \cdot d(x_i, \mu_j)$$

- if a training example is assigned to a cluster that is not closest to it, the loss function will **increase**
- e.g., x_i is closer to cluster **3** rather than to cluster **5**; it should be assigned to cluster **3**, otherwise the loss function will be higher since $d(x_i, \mu_5) > d(x_i, \mu_3)$

the loss function depends on the choice of **distance function**, which is **application dependent**

- need to consider the type of features
 - **Categorical, ordinal or continuous**
- can learn distance function from data
 - ML research area called **metric learning**

k-Means: Algorithm

How do we minimize loss with respect to assignments and cluster centers (r_{ij}, μ_j) ?

- Chicken and egg problem
- If prototypes (μ_j) known, can assign data points to clusters and get responsibilities (r_{ij})
- If responsibilities (r_{ij}) known, can compute prototypes (μ_j) by averaging data points in each cluster

Objective

$$\min_{\mu} \min_C \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2$$

1. Fix μ , optimize C :

$$\min_C \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2 = \min_c \sum_i |x_i - \mu_{x_i}|^2$$

Step 1 of kmeans

- **Implementing Step 1:** fix prototypes, update cluster assignment of data points to closest prototype

2. Fix C , optimize μ :

$$\min_{\mu} \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2$$

- Take partial derivative of μ_i and set to zero, we have

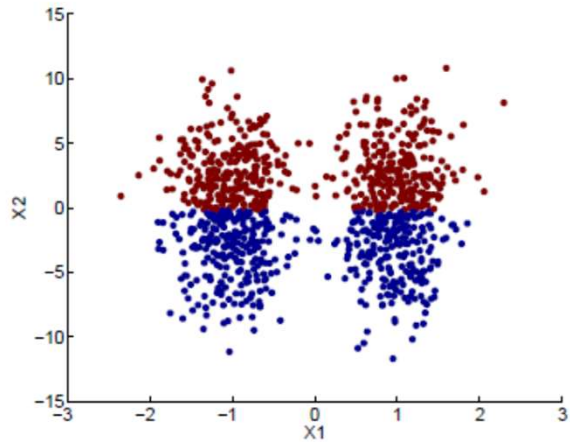
$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

Step 2 of kmeans

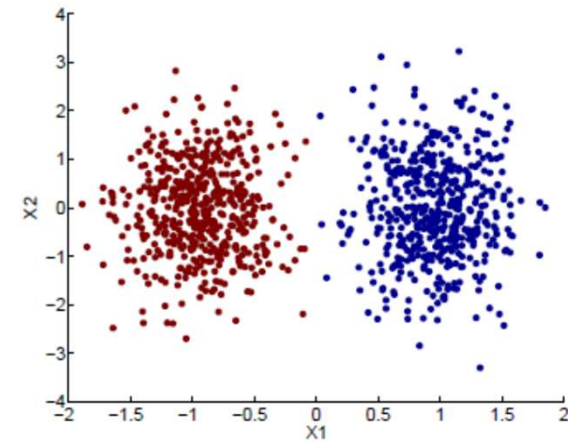
- **Implementing Step 2:** fix cluster assignments, update prototypes using data points

- **guaranteed to converge** in a **finite** number of iterations
- **alternating minimization algorithm** is a variant of **expectation-maximization** (EM)
- Running time per iteration
 - E-step: fix prototypes, update assignment of data points to closest prototype is $O(kn)$
 - M-step: fix cluster assignments, update prototypes using data points is $O(n)$

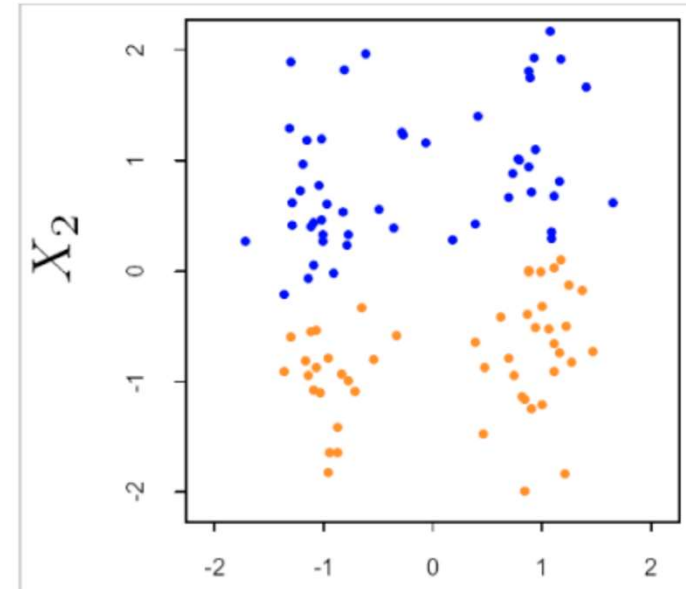
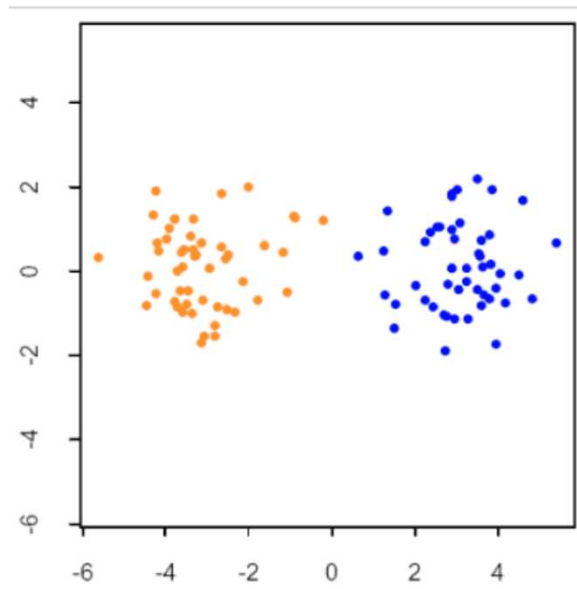
To Standardize or To Not Standardize?



Without standardization



With standardization



How To Select k ?

- **Cross-validation:** Partition data into two sets. Estimate prototypes on train and use these to compute the loss function on validation
- **Stability of clusters:** Measure the change in the clusters obtained by resampling or splitting the data.
- **Non-parametric approach:** Place a prior on k
- **Gap statistic:** select a k such that the “compactness of clustering” is best compared to a reference distribution (that has no obvious clustering) see [\[Tibshirani et al., 2001\]](#) for more details

K=2



K=3



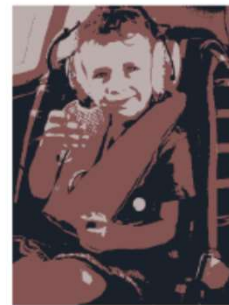
K=10



Original



4%



8%

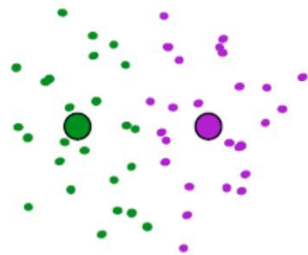


17%

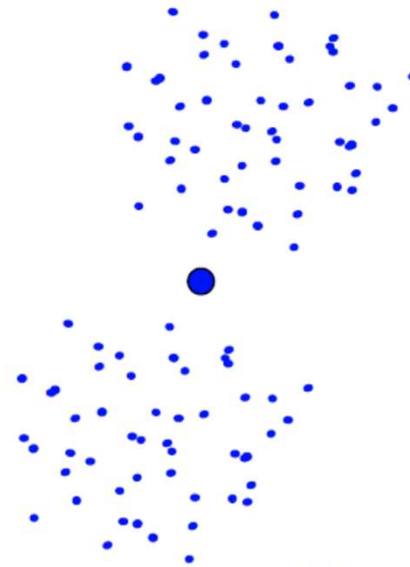


Limitations of k-Means

- k-means will **converge to a local minima**
- **different initializations** can lead to very **different results**
- **run k-means several times** with **random** starting points, pick clustering with **smallest loss**



Would be better to have
one cluster here



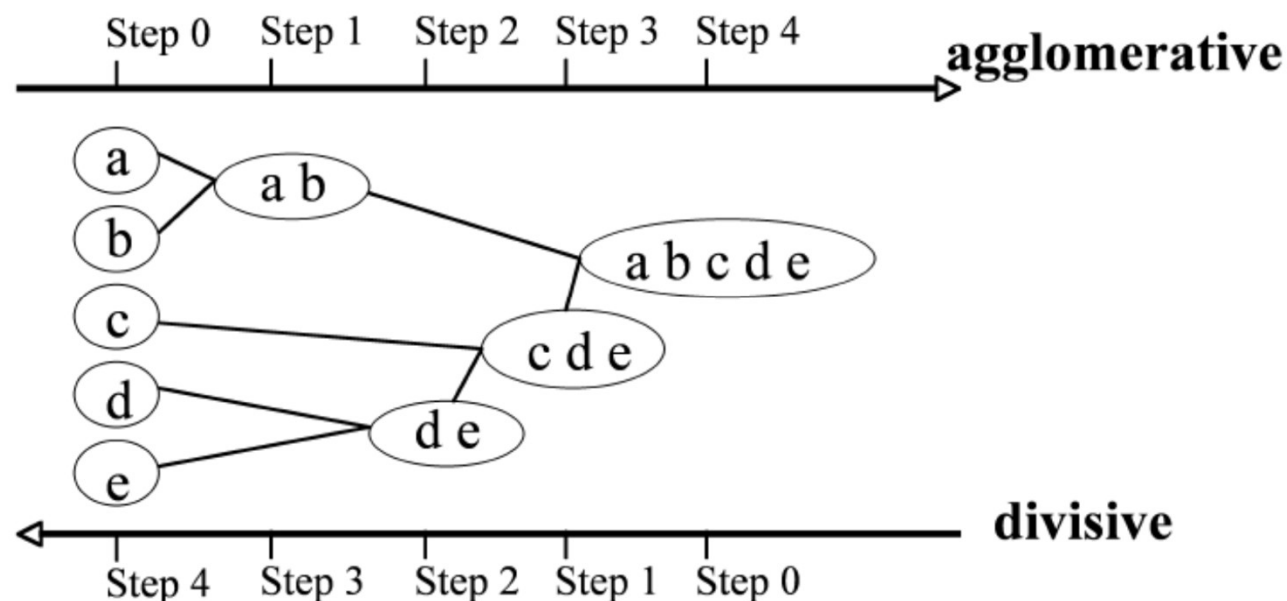
... and two clusters here

- results can change dramatically when k changes

Hierarchical Clustering

Another limitation of k-means: clusters change arbitrarily for different k

Solution: organize clusters in a hierarchical way by grouping similar clusters



- **Bottom-up (agglomerative):** Recursively merge two groups with the smallest **between-cluster similarity**
- **Top-down (divisive):** Recursively split a **least-coherent** (e.g. largest diameter) cluster
- Users can then choose a cut through the hierarchy to represent the most natural division into clusters (e.g. where intergroup similarity exceeds some threshold).

Agglomerative Clustering

How do we define “closest” for clusters?

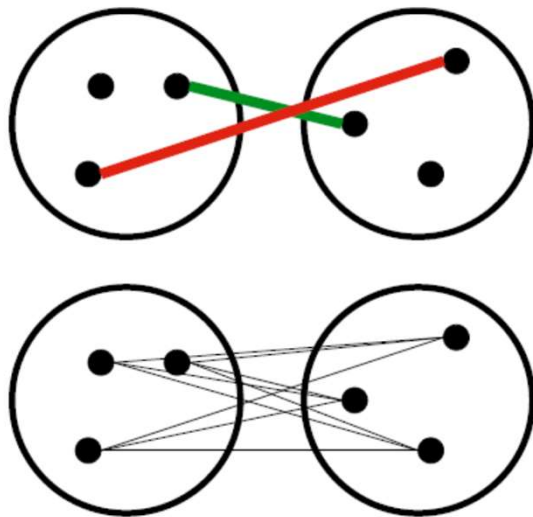
Closest pair (single-link clustering) tends to yield elongated clusters

Farthest pair (complete-link clustering) tends to yield rounder, more spherical clusters

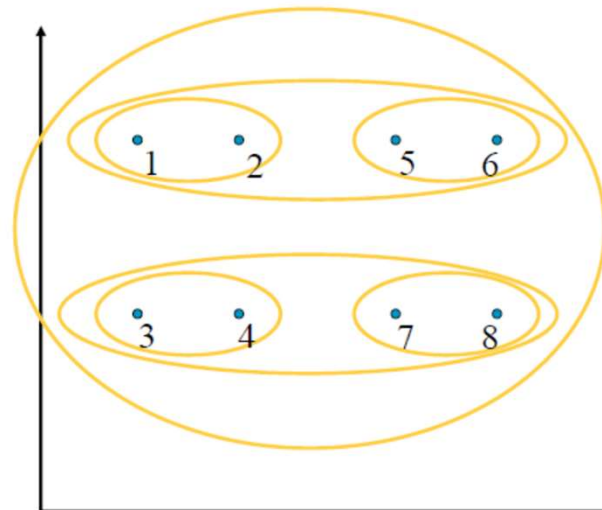
Average of all pairs trades-off between single and complete link

Algorithm:

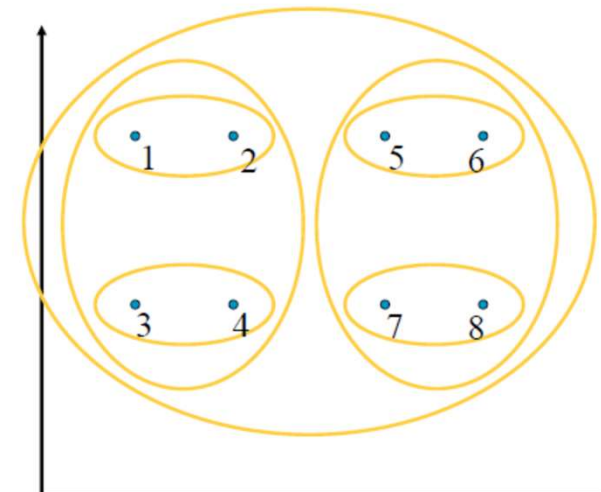
- **Initialize:** each training example is its own cluster
- **repeat until** only one cluster left
 - pick two “closest” clusters
 - merge them into a new cluster



Closest pair
(single-link clustering)



Farthest pair
(complete-link clustering)



Agglomerative Clustering

How do we define “closest” for clusters?

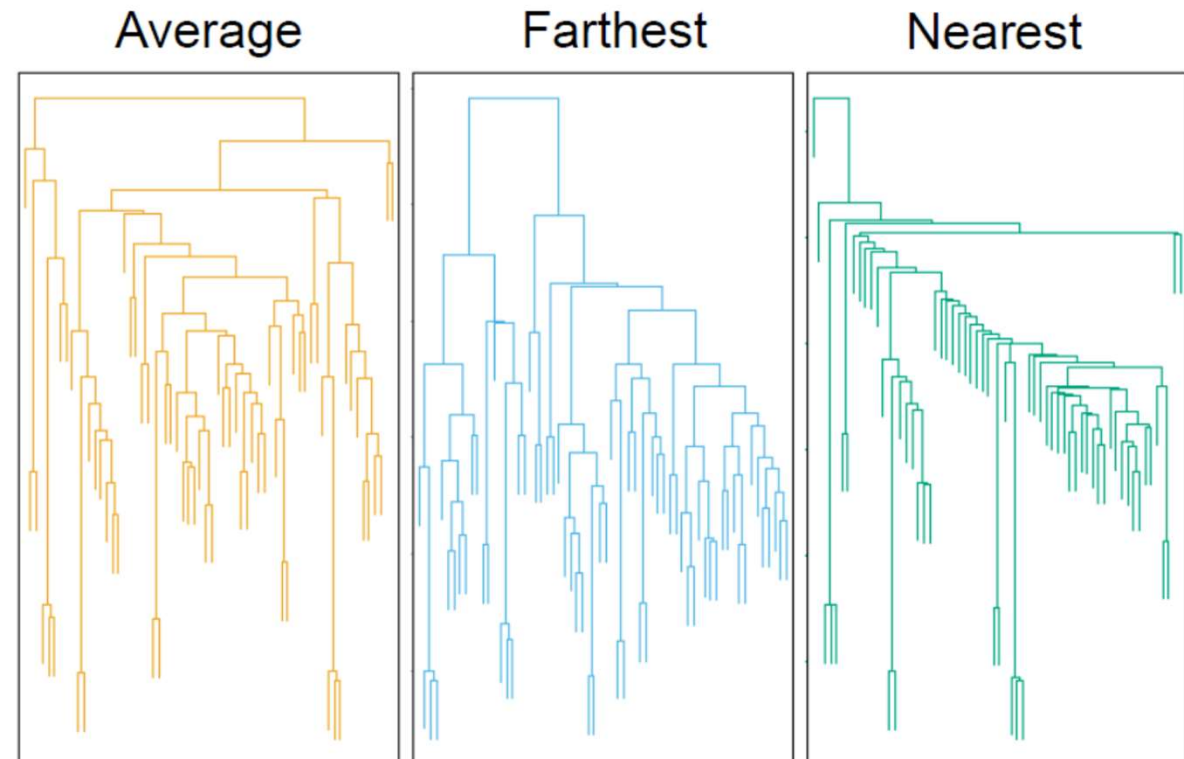
Closest pair (single-link clustering) tends to yield elongated clusters

Farthest pair (complete-link clustering) tends to yield rounder, more spherical clusters

Average of all pairs trades-off between single and complete link

Algorithm:

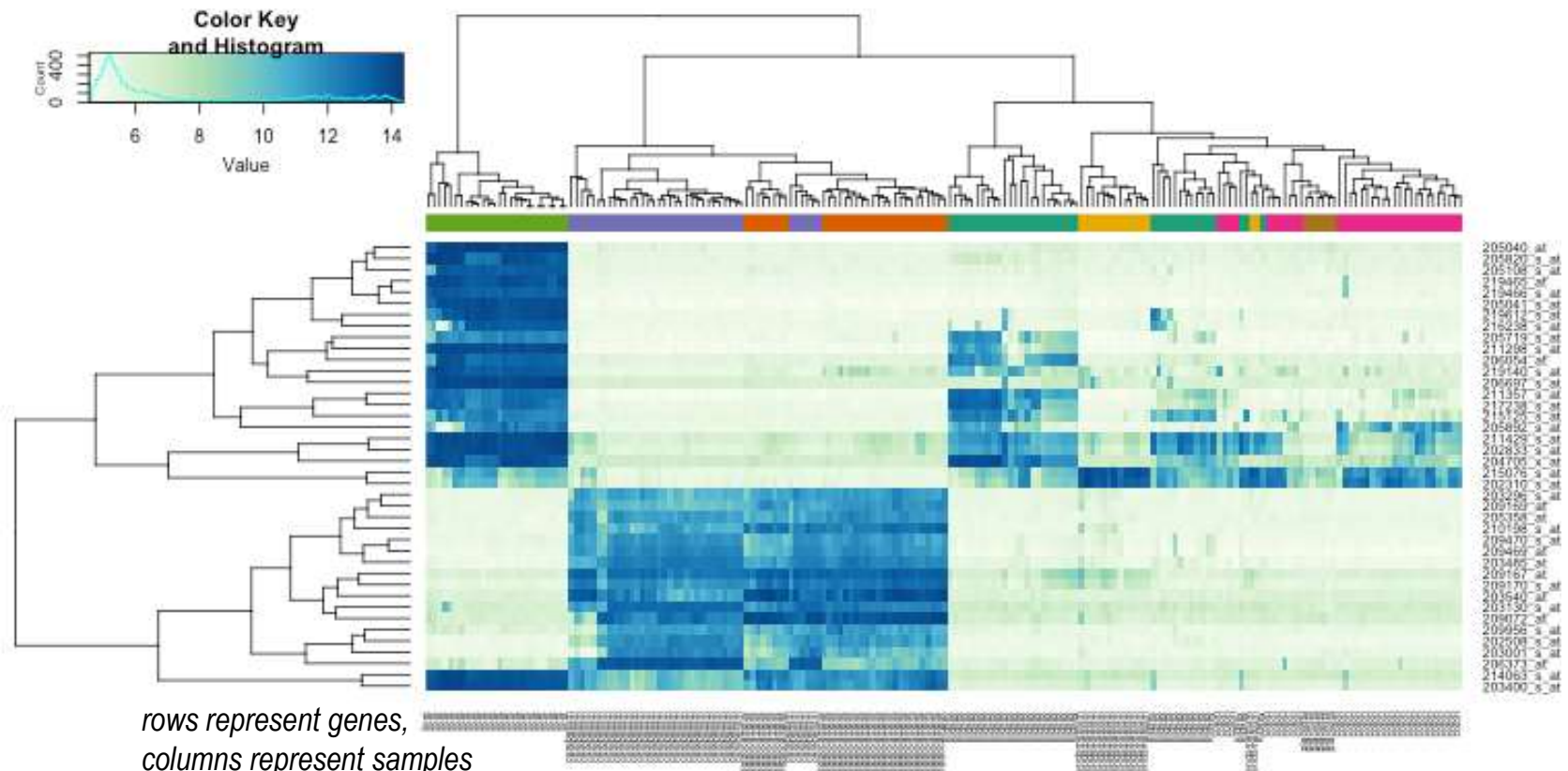
- **Initialize:** each training example is its own cluster
- **repeat until** only one cluster left
 - pick two “closest” clusters
 - merge them into a new cluster



Example: Gene Expression Analysis

Example: Discover patterns in **gene expression data**; e.g., new types of cancer from gene expression profiles, drug responses from genotypes etc.,

Gene expression analysis makes extensive use of **hierarchical clustering**, where **each example is assigned its own cluster**, and then clusters are grouped together **bottom-up**. Each level of the hierarchy can be interpreted as a clustering of different granularity and is visualized as a **dendrogram**.

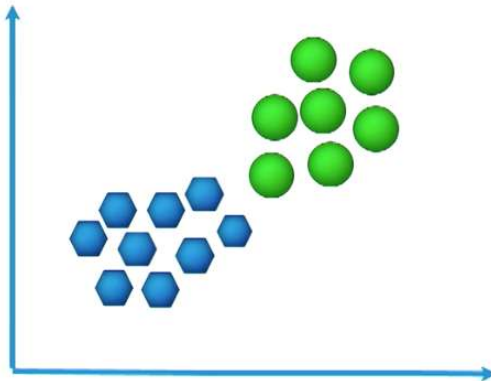


Limitations of k-Means

- k-means performs **hard clustering**
 - cluster assignments using **responsibilities** $r_{ij} \in \{0, 1\}$
 - a small perturbation to a data point (noise) to flip it to another cluster (instability)
 - assumes spherical clusters and equal probabilities for each cluster
- **soft clustering** can **soften** assignments to a range $r_{ij} \in [0, 1]$
 - interpretation: training example can now belong to more than one cluster with a probability r_{ij}
 - approaches: fuzzy clustering, **Gaussian mixture models**
 - compare with perceptron vs. logistic regression

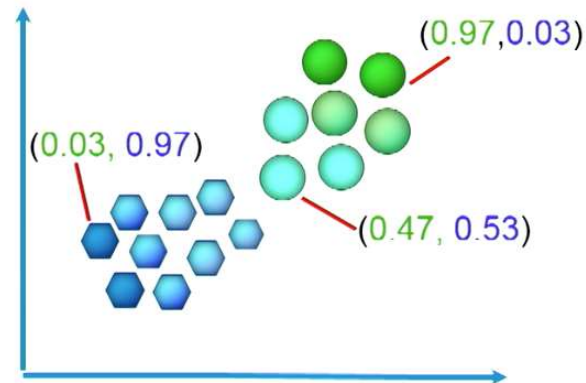
Hard clustering

Each observation belongs to exactly one cluster



Soft clustering

An observation can belong to more than one cluster to a certain degree (e.g. likelihood of belonging to the cluster)

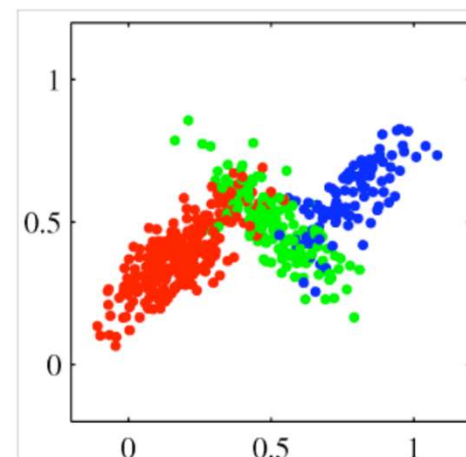
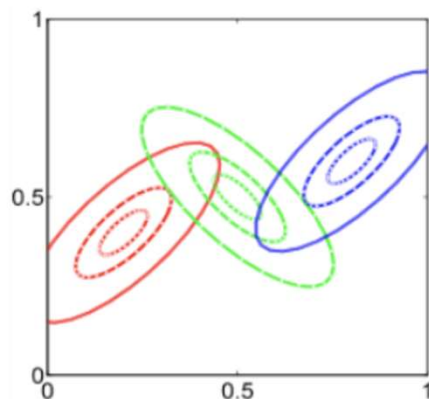


Gaussian Mixture Models

Probabilistic clustering (generative model):

Each cluster is associated with a **Gaussian distribution**. To **generate** data,

- randomly choose a cluster j with probability $P(y = j)$
 - *distribution over the clusters is modeled as a multinomial distribution*
- generate from the distribution of the j -th cluster: $P(\mathbf{x}_i | y = j) = N(\mathbf{x} | \boldsymbol{\mu}_j, \Sigma_j)$
 - *distribution over each cluster is modeled as a multivariate Gaussian distribution*

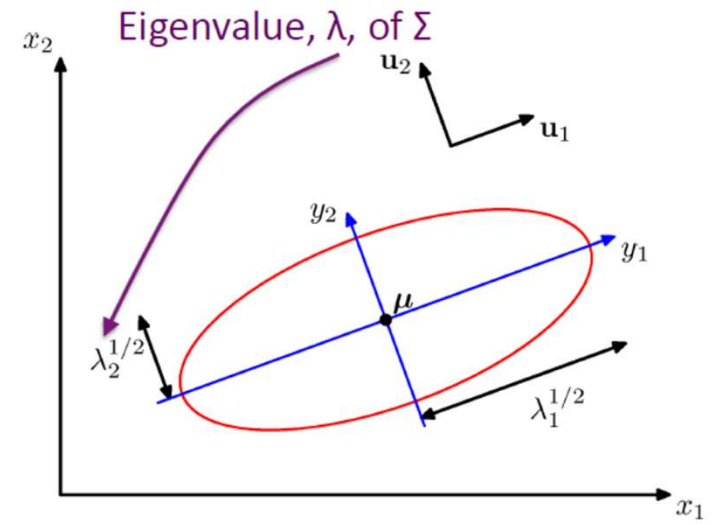


$$P(y = j) = \pi_j, \sum_{j=1}^k \pi_j = 1 \text{ (multinomial distribution)}$$

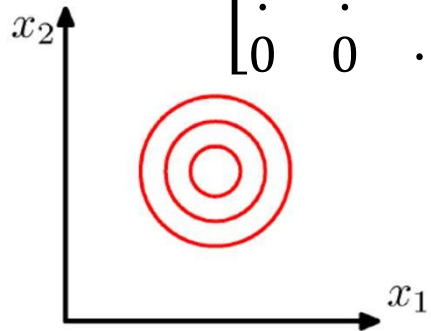
$$P(\mathbf{x}_i | y = j) = N(\mathbf{x} | \boldsymbol{\mu}_j, \Sigma_j) \text{ (multivariate Gaussian)}$$

The Multivariate Gaussian Distribution

$$P(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

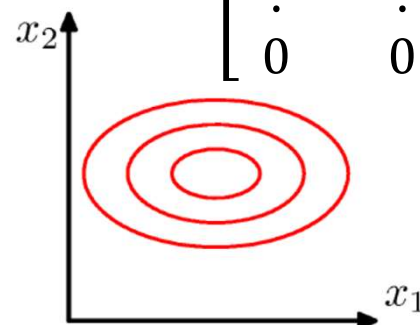


$$\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$



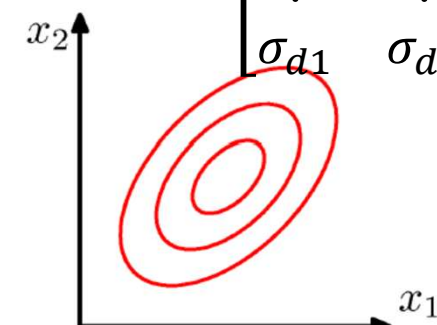
identity matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{11} & 0 & \cdots & 0 \\ 0 & \sigma_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{dd} \end{bmatrix}$$



diagonal matrix:
Gaussian Naïve Bayes

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_{dd} \end{bmatrix}$$



arbitrary positive semi-definite matrix:
eigenvectors specify rotation,
eigenvalues specify (relative) elongation

Mixtures of Gaussians

For a single data point \mathbf{x}_i

$$P(\mathbf{x}_i) = \sum_{j=1}^k \pi_j \cdot N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

↖ mixing coefficient
↘ individual Gaussian component

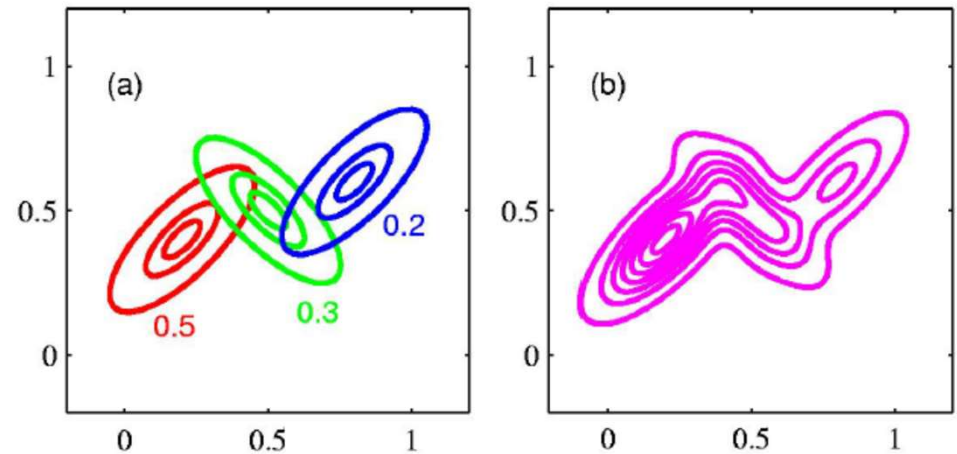
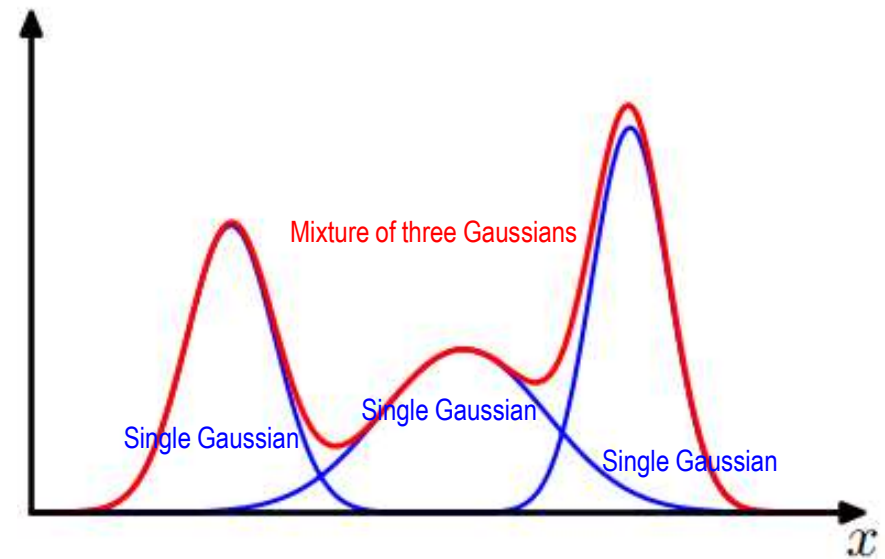
$$\pi_j \geq 0, \sum_{j=1}^k \pi_j = 1$$

For a data set $\mathbf{x}_i, i = 1, \dots, n$

$$P(\mathbf{x}) = \prod_{i=1}^n \sum_{j=1}^k \pi_j \cdot N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

The log-likelihood simplifies the problem, but the parameters of the components are still coupled

$$\log P(\mathbf{x}) = \sum_{i=1}^n \log \sum_{j=1}^k \pi_j \cdot N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$



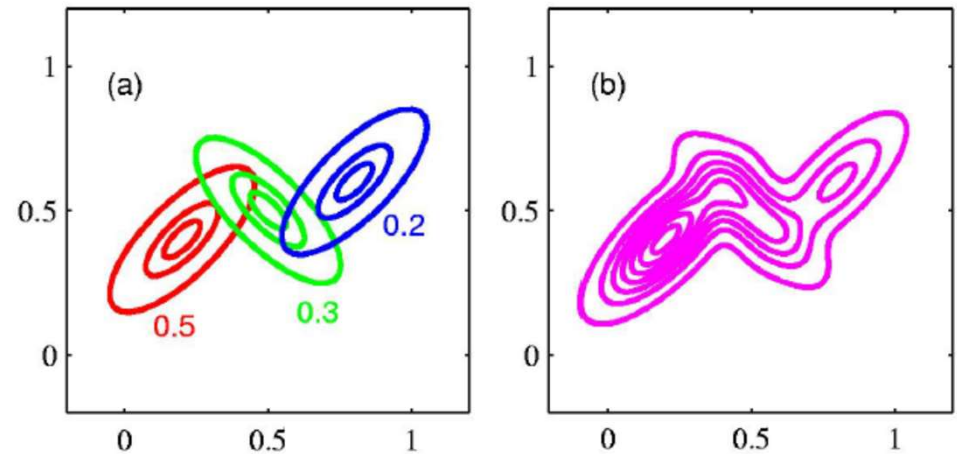
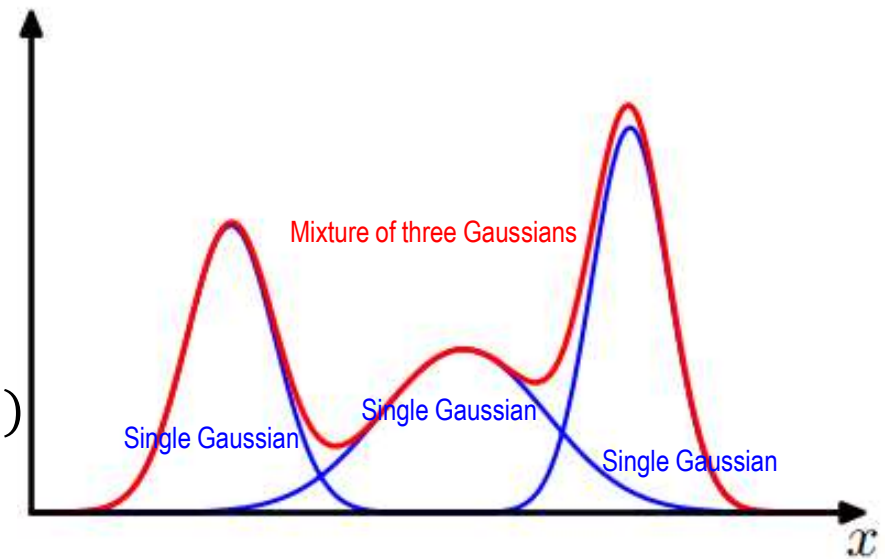
Mixtures of Gaussians

Introduce **latent variables** for each training example (these latent variables are essentially cluster responsibilities)

$$\mathbf{z}_i = [z_{i1}, z_{i2}, \dots, z_{ij}, \dots, z_{ik}]$$

$$\log P(\mathbf{x}, \mathbf{z}) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} \cdot (\log \pi_j + \log N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j))$$

The introduction of the latent variables causes the parameters π_j and $(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ to **decouple!**



Expectation Maximization for GMMs

E-step: Fix parameters, compute responsibilities

$$r_{ij} = E(z_{ij}) = \frac{\pi_j N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{j=1}^k \pi_j N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

M-step: Fix responsibilities, compute parameters by maximizing expected complete log-likelihood

$$\text{maximize } \log P(\mathbf{x}, \mathbf{z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^n \sum_{j=1}^k r_{ij} \cdot (\log \pi_j + \log N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j))$$

$$\pi_j = \frac{\sum_{i=1}^n r_{ij}}{n}$$

$$\boldsymbol{\mu}_j = \frac{\sum_{i=1}^n r_{ij} \mathbf{x}_i}{\sum_{i=1}^n r_{ij}}$$

$$\boldsymbol{\Sigma}_j = \frac{\sum_{i=1}^n r_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)^T (\mathbf{x}_i - \boldsymbol{\mu}_j)}{\sum_{i=1}^n r_{ij}}$$

Loss function:

minimize sum of squared distance.

- Hard assignment of points to clusters.
- Assumes spherical clusters with equal probability of a cluster.

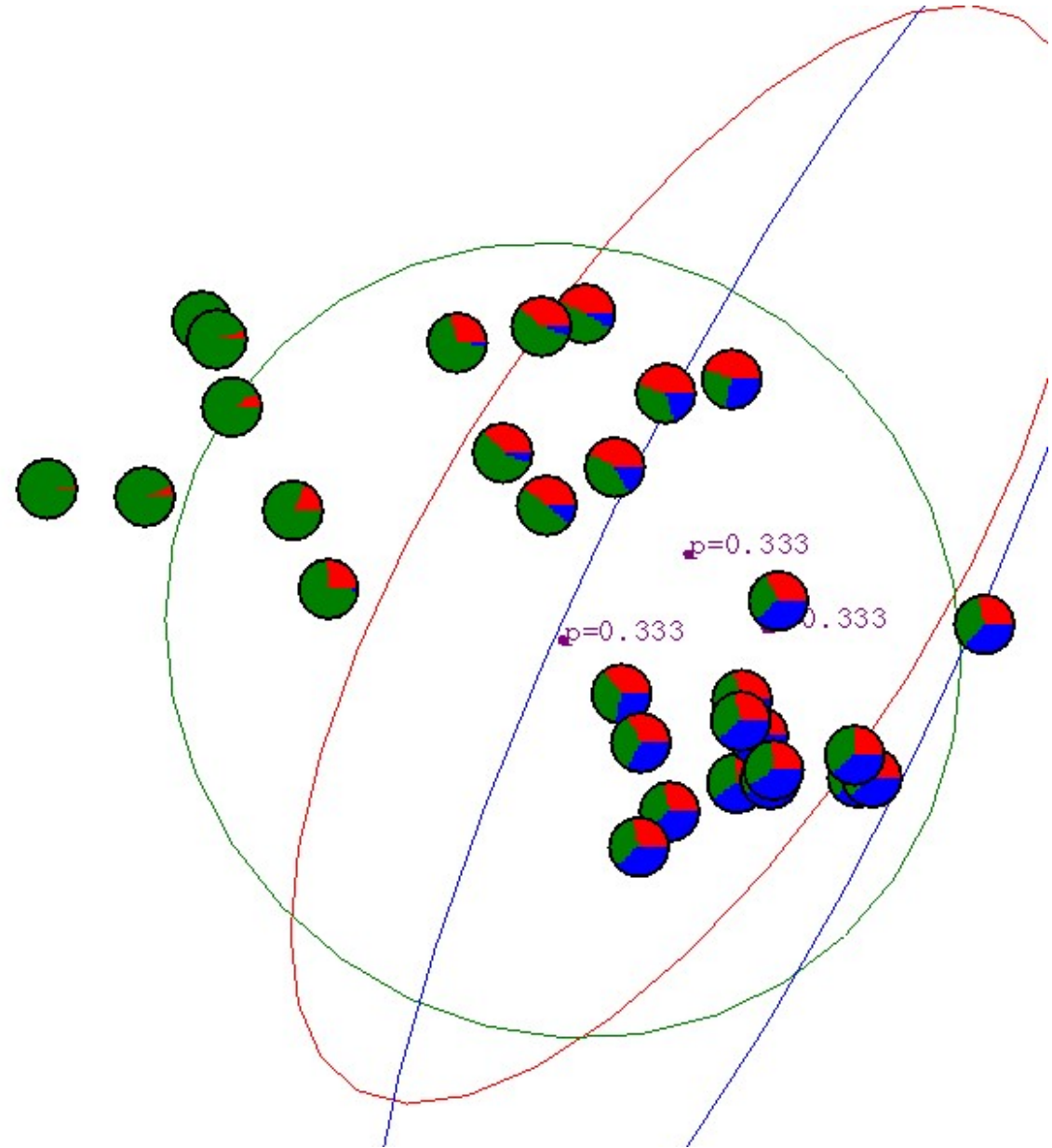
k-Means

- **hard** assignment of points to clusters
- minimize sum of squared distance
- assumes spherical clusters with equal probability of clusters

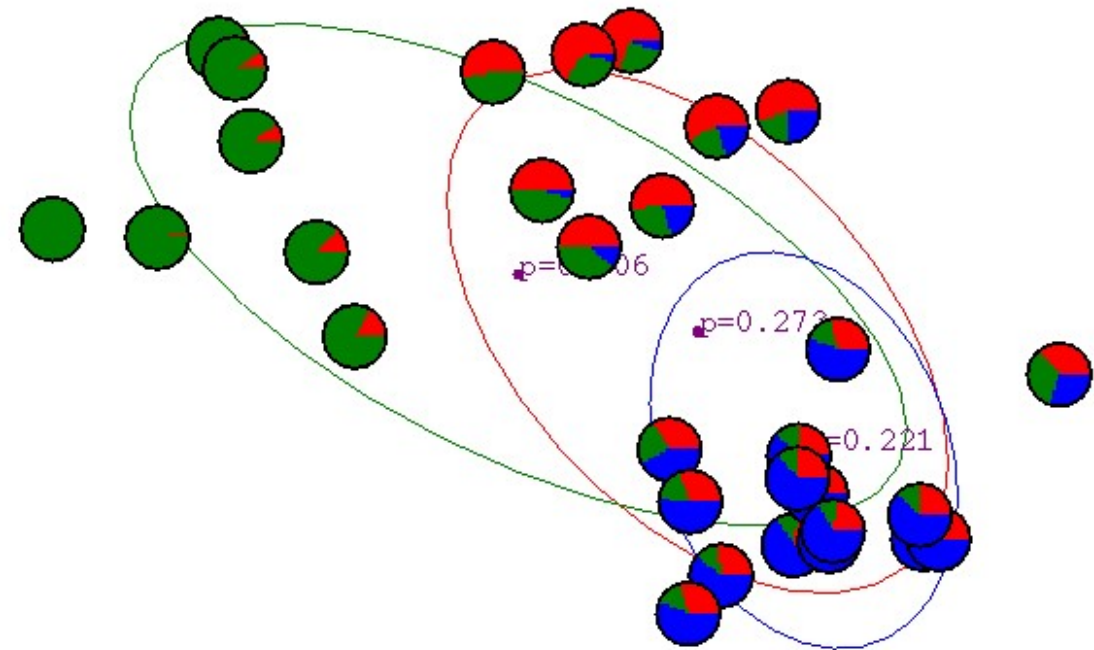
Gaussian Mixture Models

- **soft** assignment of points to clusters
- maximize log likelihood.
- Can be used for non-spherical clusters with different probabilities

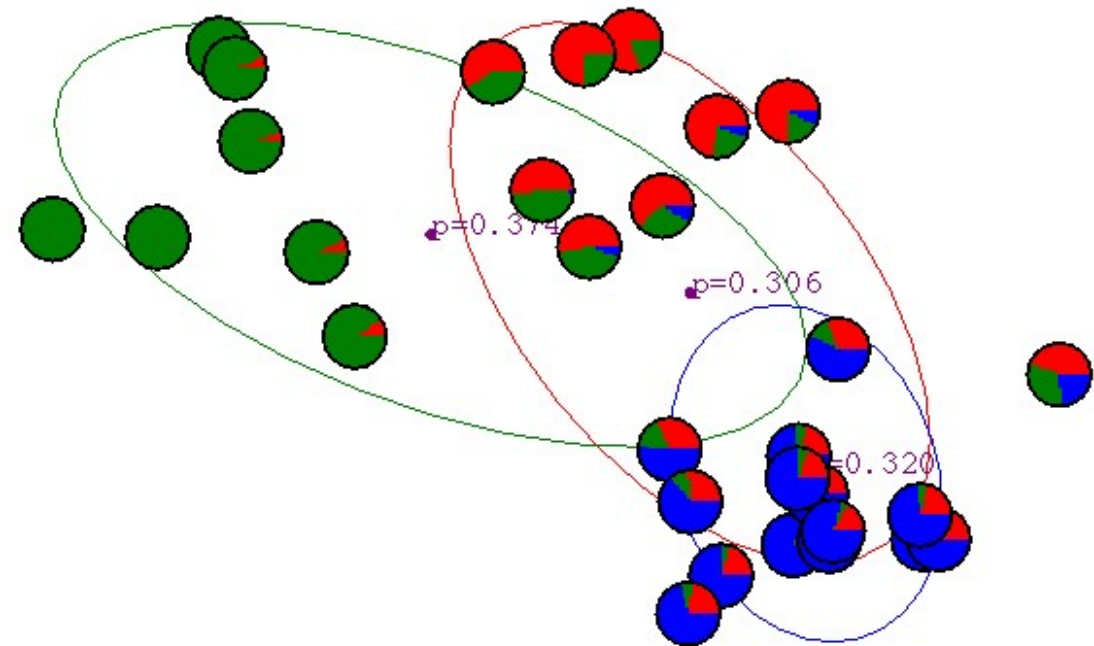
Expectation Maximization for GMMs



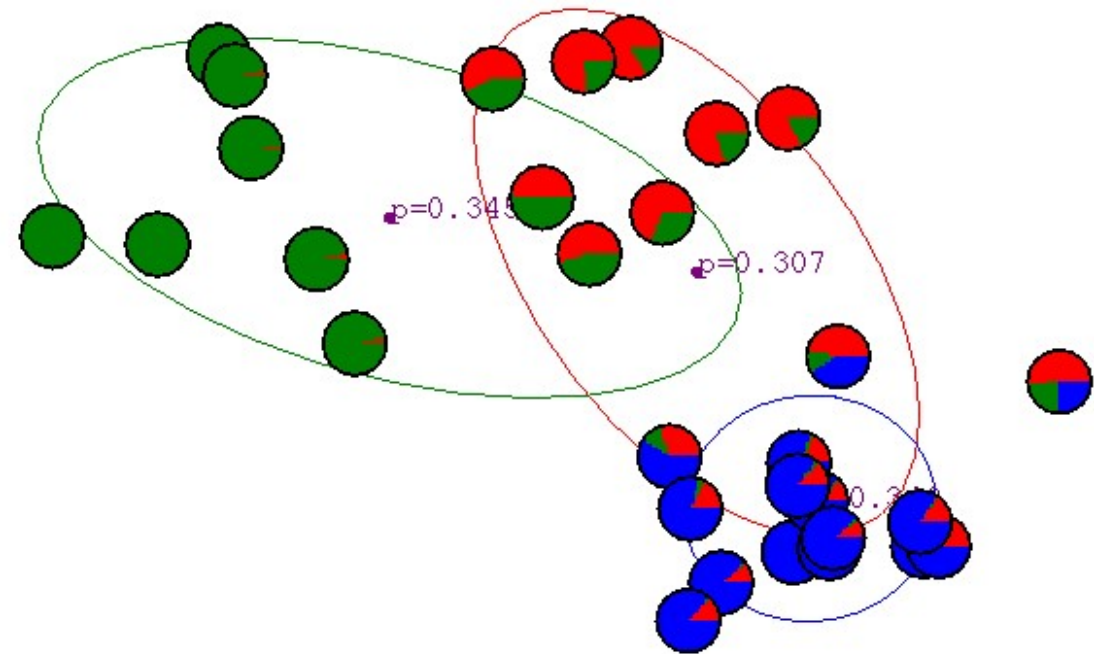
Expectation Maximization for GMMs



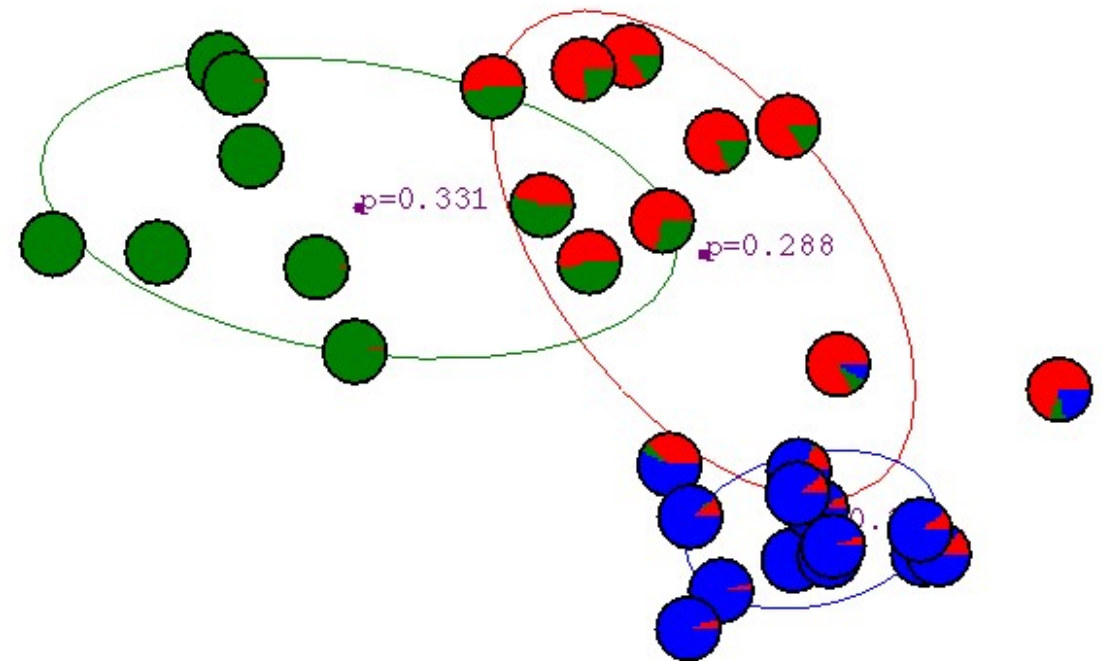
Expectation Maximization for GMMs



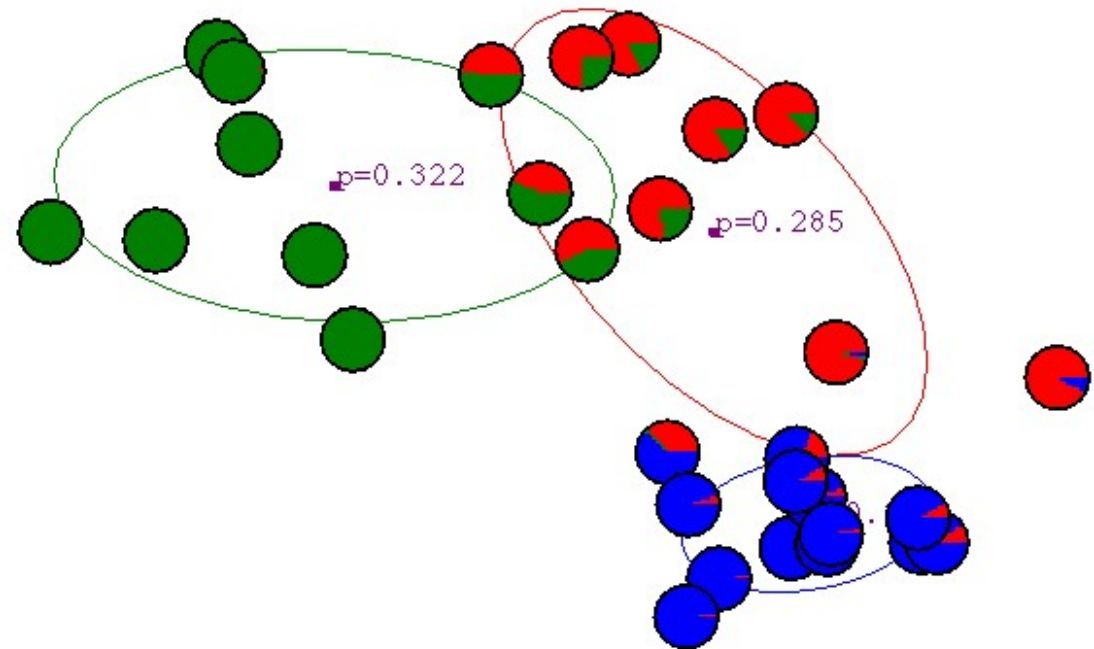
Expectation Maximization for GMMs



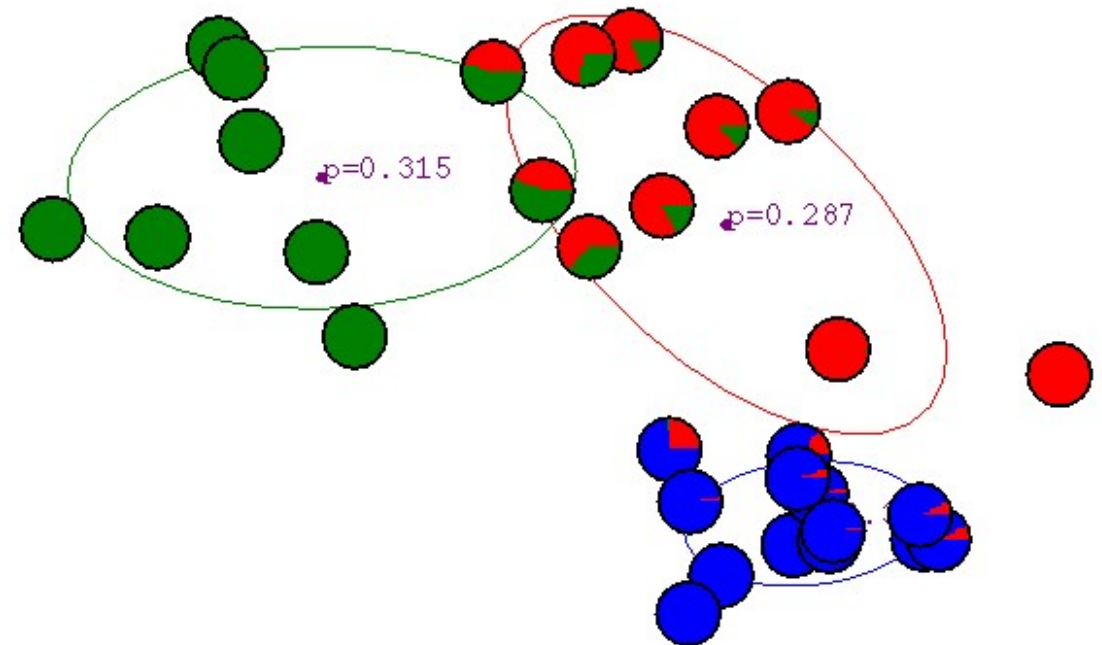
Expectation Maximization for GMMs



Expectation Maximization for GMMs



Expectation Maximization for GMMs



Expectation Maximization for GMMs

