

Clustering contd...

Aarti Singh

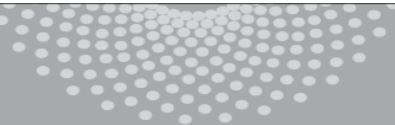
Machine Learning 10-701

Apr 26, 2023

Some slides courtesy of Eric Xing, Carlos Guestrin



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

K-means Recap ...



given

- Randomly initialize k centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

Iterate $t = 0, 1, 2, \dots$

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $C^{(t)}(j) \leftarrow \arg \min_{i=1, \dots, k} \|\mu_i^{(t)} - x_j\|^2$

- **Recenter:** μ_i becomes centroid of its points:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2$ $i \in \{1, \dots, k\}$

$\equiv \mu_i^{(t+1)} = \frac{\sum_{j \in C^{(t)}(j)=i} x_j}{\# \dots}$

- Equivalent to $\mu_i \leftarrow$ average of its points!

K-means algorithm

- Optimize potential function: $\sum_{j=1}^m \|\mu_i - x_j\|^2$
 $\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$
 $\mu_1 \dots \mu_k$

- K-means algorithm:** (coordinate descent on F)

(1) Fix μ , optimize C

Expected cluster assignment

(2) Fix C , optimize μ

Maximum likelihood for center

Similar to EM/Baum Welch algorithm for learning HMM parameters

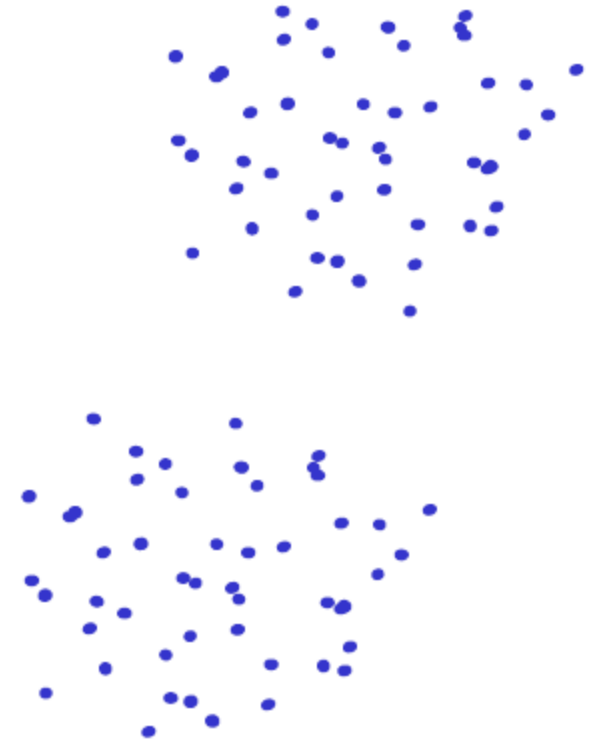
$\text{latent states} \equiv \text{latent cluster assignment}$

Computational Complexity

- At each iteration,
 - Computing distance between each of the n objects and the K cluster centers is $O(Kn)$. ✓
 - Computing cluster centers: Each object gets added once to some cluster: $O(n)$.
- Assume these two steps are each done once for l iterations: $O(lKn)$.

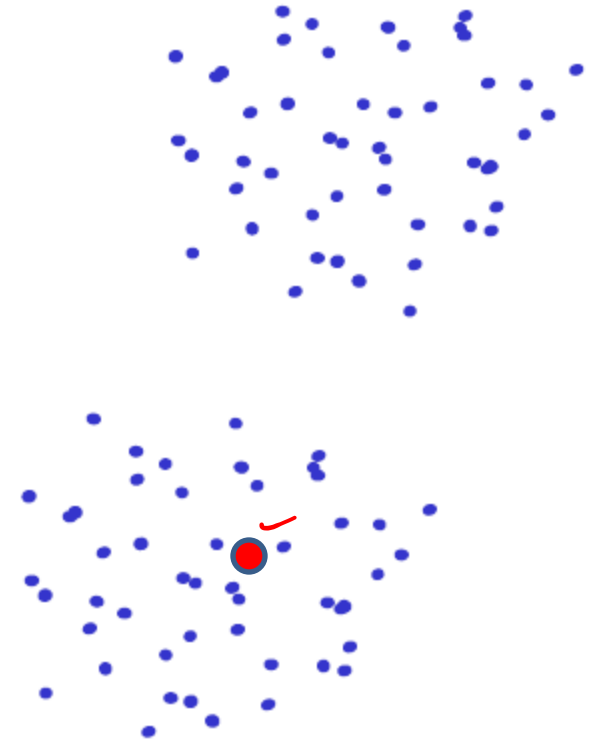
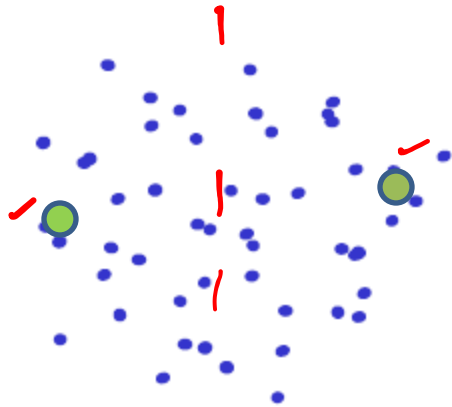
Seed Choice

- Results are quite sensitive to seed selection.



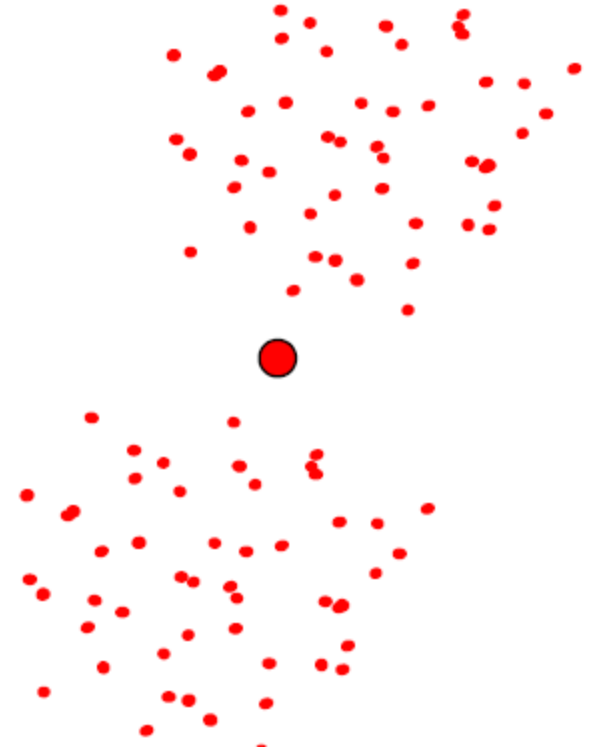
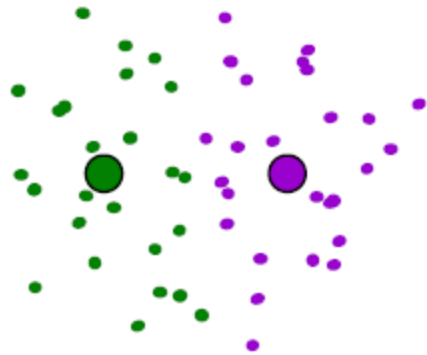
Seed Choice

- Results are quite sensitive to seed selection.



Seed Choice

- Results are quite sensitive to seed selection.



Seed Choice

- Results can vary based on random seed selection.
 - Some seeds can result in poor convergence rate, or convergence to sub-optimal clustering.
 - Try out multiple starting points (very important!!!) ✓
 - k-means ++ algorithm of Arthur and Vassilvitskii
- key idea: choose centers that are far apart
- (probability of picking a point as cluster center \propto distance from nearest center picked so far)

Other Issues

- Number of clusters K

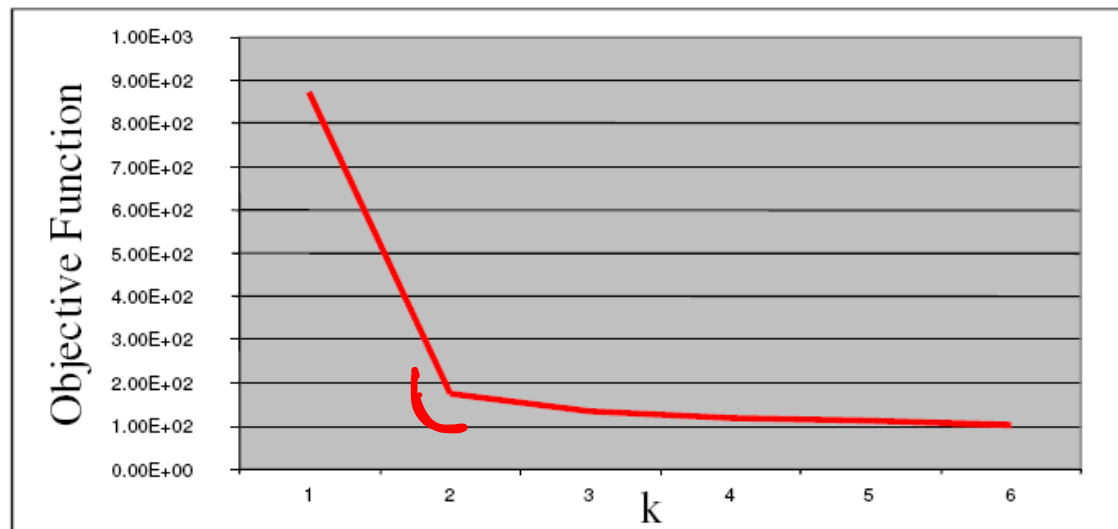
- Objective function

$$\min_K \min_{\mu, c} \sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2$$

(Note: A red arrow points to the sum term in the original image.)

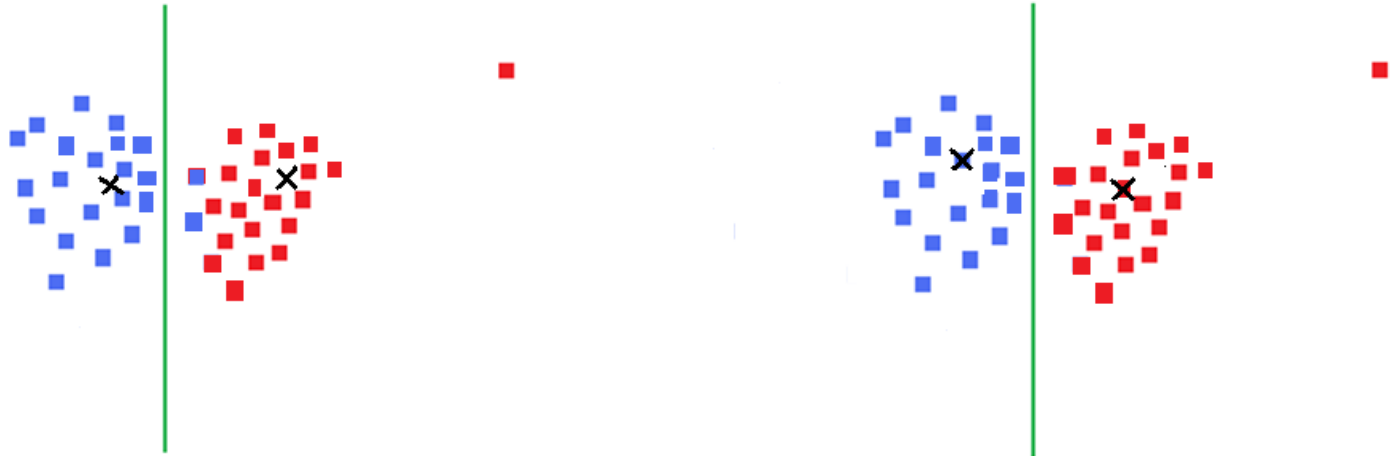
- Can you pick K by minimizing the objective over K?

- Look for “Knee” in objective function



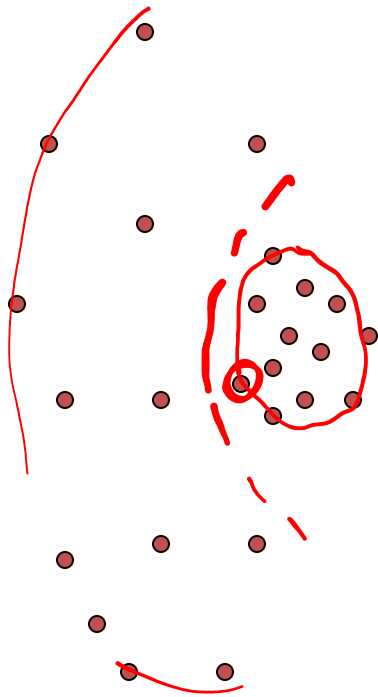
Other Issues

- Sensitive to Outliers
 - use K-medoids



- Shape of clusters
 - Assumes isotropic, equal variance, convex clusters

K-means limitations



- Clusters may overlap *"soft" assignment*
- Some clusters may be "wider" than others
- Clusters may not be linearly separable ✓

Partitioning Algorithms

- K-means
 - **hard assignment**: each object belongs to only one cluster
- Mixture modeling
 - **soft assignment**: probability that an object belongs to a cluster

$p(x)$

Generative approach

Mixture models

GMM – Gaussian Mixture Model (Multi-modal distribution)

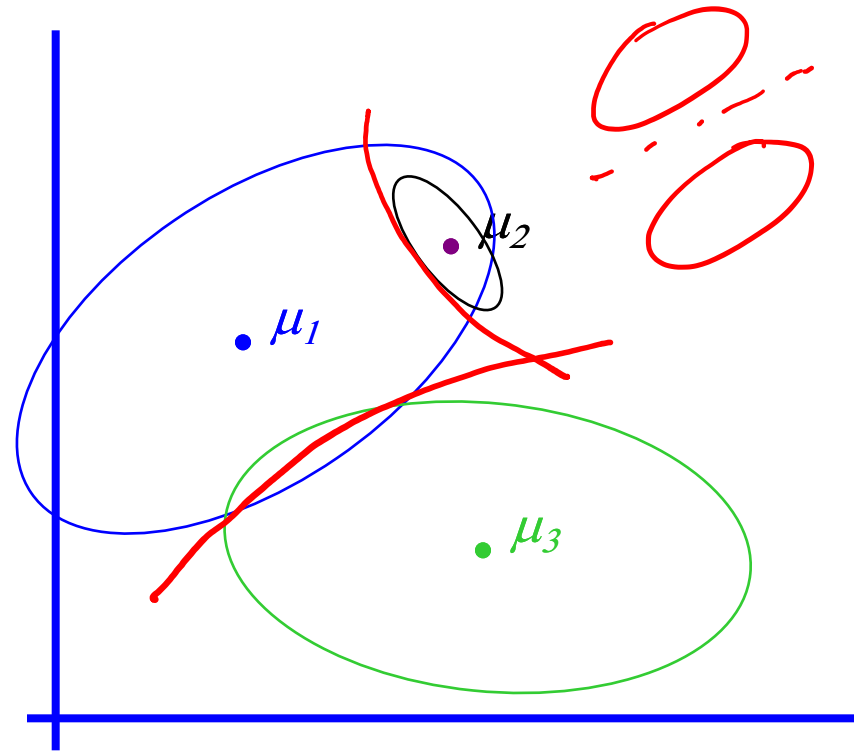
$$p(x/y=i) \sim \mathcal{N}(\mu_i, \Sigma_i)$$

$$p(x) = \sum_i p(x/y=i) P(y=i)$$

Mixture
component

Mixture
proportion

$$p(y=i|x) \propto p(y=i) p(x|y=i)$$



Mixture models

GMM – Gaussian Mixture Model (Multi-modal distribution)

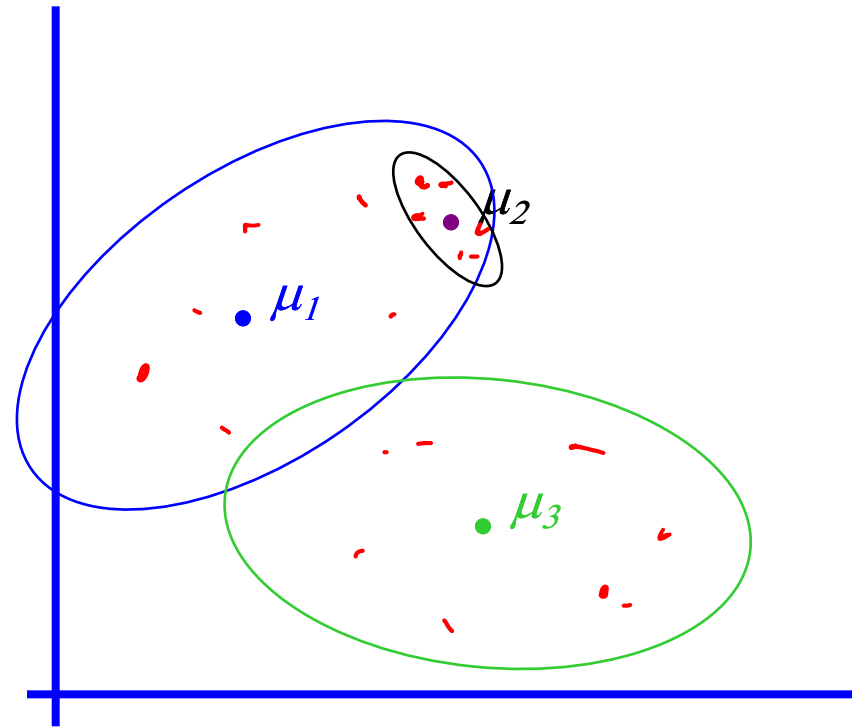
- There are k components
- Component i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix Σ_i

↕ Each data point is generated according to the following recipe:

1) Pick a component at random:
Choose component i with

probability $P(y=i) \checkmark \equiv p_i$

2) Datapoint $x \sim N(\checkmark \mu_i, \checkmark \Sigma_i)$



Learning GMMs via EM algorithm

Iterate. On iteration t let our estimates be $P(y=i) \equiv p_i$

$$\lambda_t = \{ \underline{\mu}_1^{(t)}, \underline{\mu}_2^{(t)} \dots \underline{\mu}_k^{(t)}, \underline{\Sigma}_1^{(t)}, \underline{\Sigma}_2^{(t)} \dots \underline{\Sigma}_k^{(t)}, \underline{p}_1^{(t)}, \underline{p}_2^{(t)} \dots \underline{p}_k^{(t)} \}$$

$p_i^{(t)}$ is shorthand for estimate of $P(y=i)$ on t 'th iteration

E-step

Compute "expected" classes of all datapoints for each class

$$\Rightarrow \underline{w}_{ji} P(\underline{y} = i | x_j, \lambda_t) \propto p_i^{(t)} p(x_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

Just evaluate a Gaussian at x_j

M-step

Compute MLEs given our data's class membership distributions (weights)

$$\mu_i^{(t+1)} = \frac{\sum_j \omega_{ji} P(\underline{y} = i | x_j, \lambda_t) x_j}{\sum_j P(\underline{y} = i | x_j, \lambda_t)}$$

$$\Sigma_i^{(t+1)} = \frac{\sum_j P(\underline{y} = i | x_j, \lambda_t) (x_j - \mu_i^{(t+1)})(x_j - \mu_i^{(t+1)})^T}{\sum_j P(\underline{y} = i | x_j, \lambda_t)}$$

$$p_i^{(t+1)} = \frac{\sum_j P(\underline{y} = i | x_j, \lambda_t)}{m}$$

$m = \#$ data points

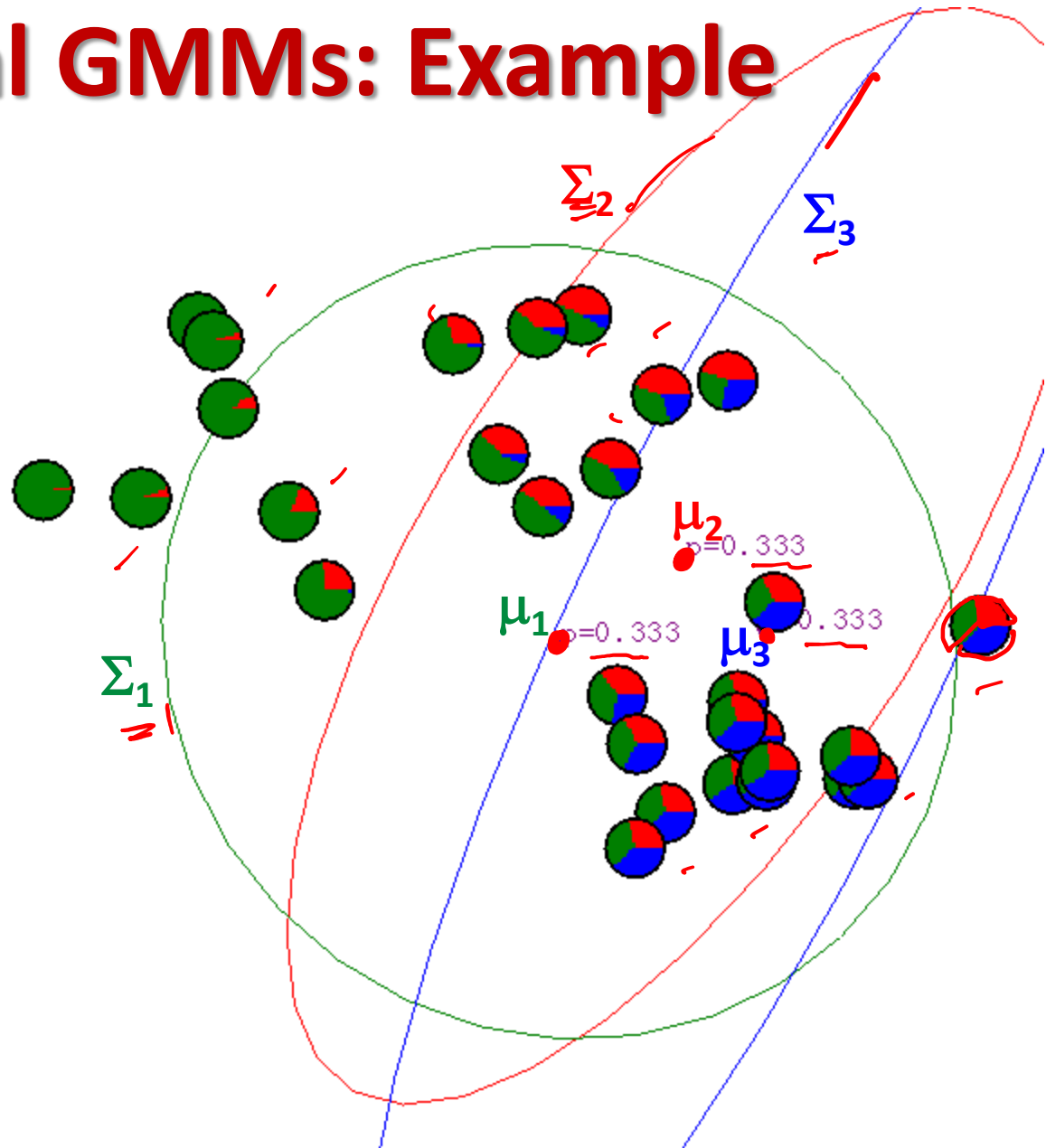
EM for general GMMs: Example

$$k=3$$


$$\mu_1, \mu_2, \mu_3$$

$$\Sigma_1, \Sigma_2, \Sigma_3$$

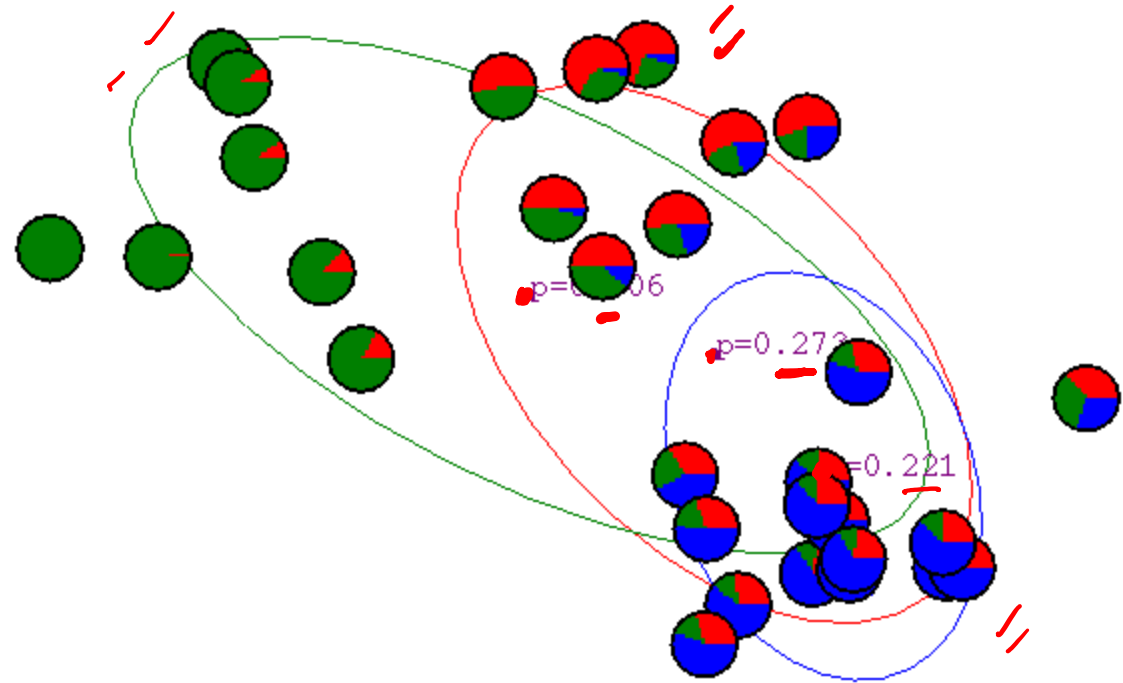
$$p_1, p_2, p_3 = \frac{1}{3}$$



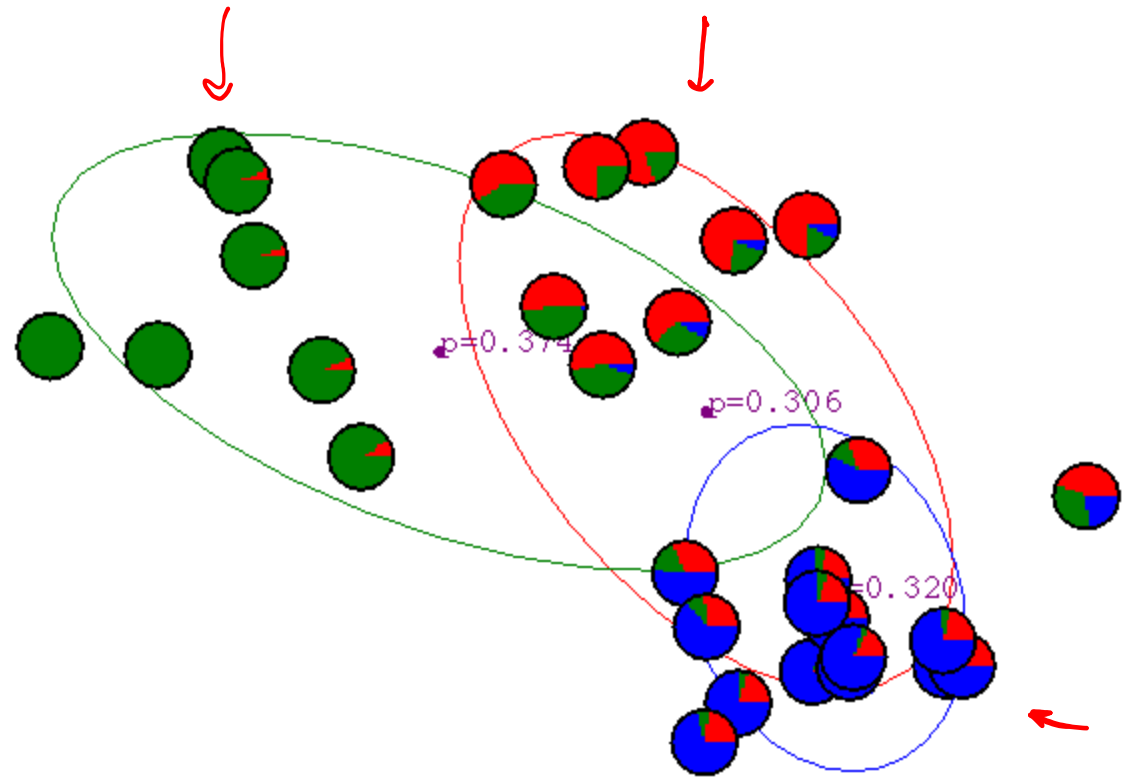
$$P(y = \bullet | x_j, \mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3, p_1, p_2, p_3)$$

 \rightarrow w_j red
datapoint

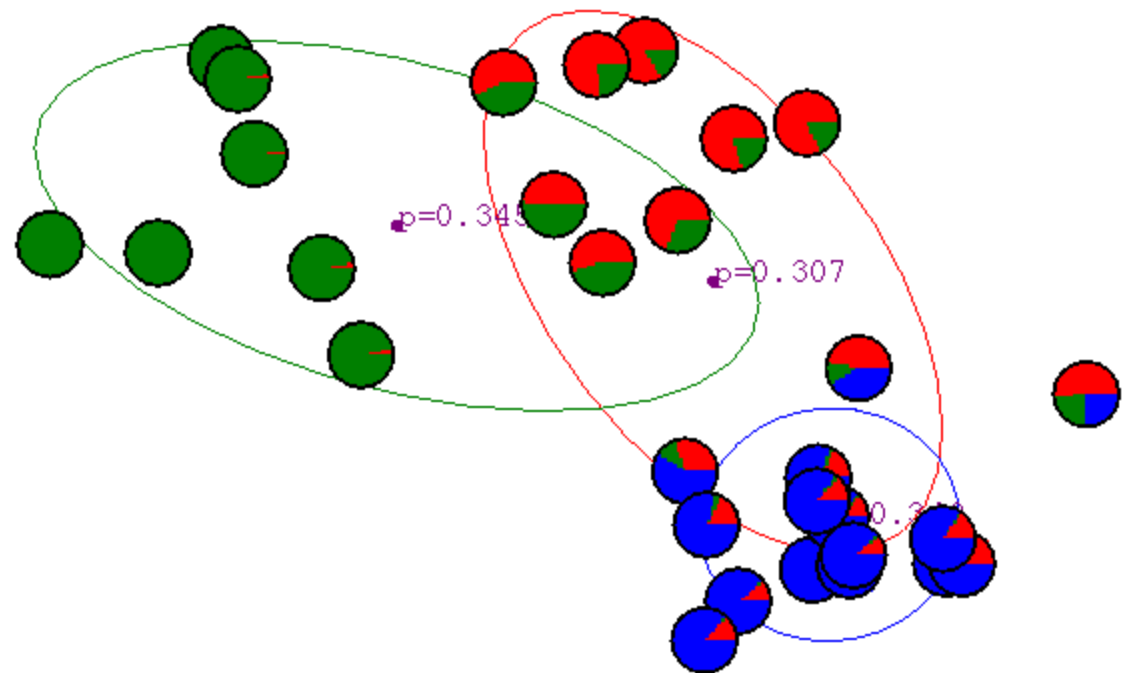
After 1st iteration



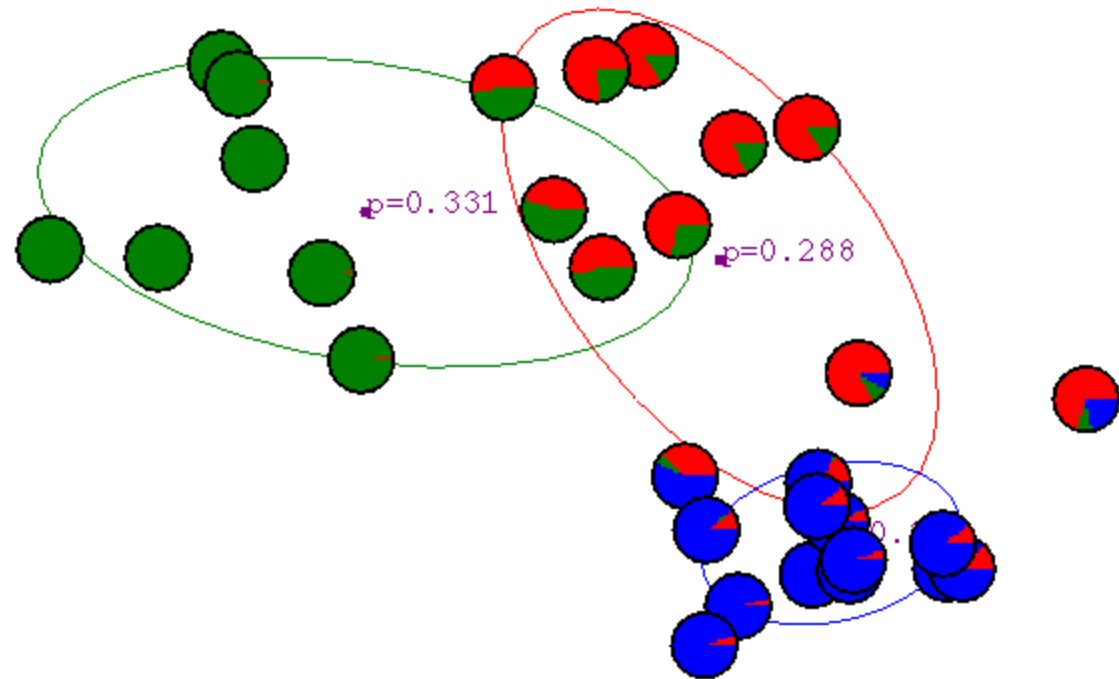
After 2nd iteration



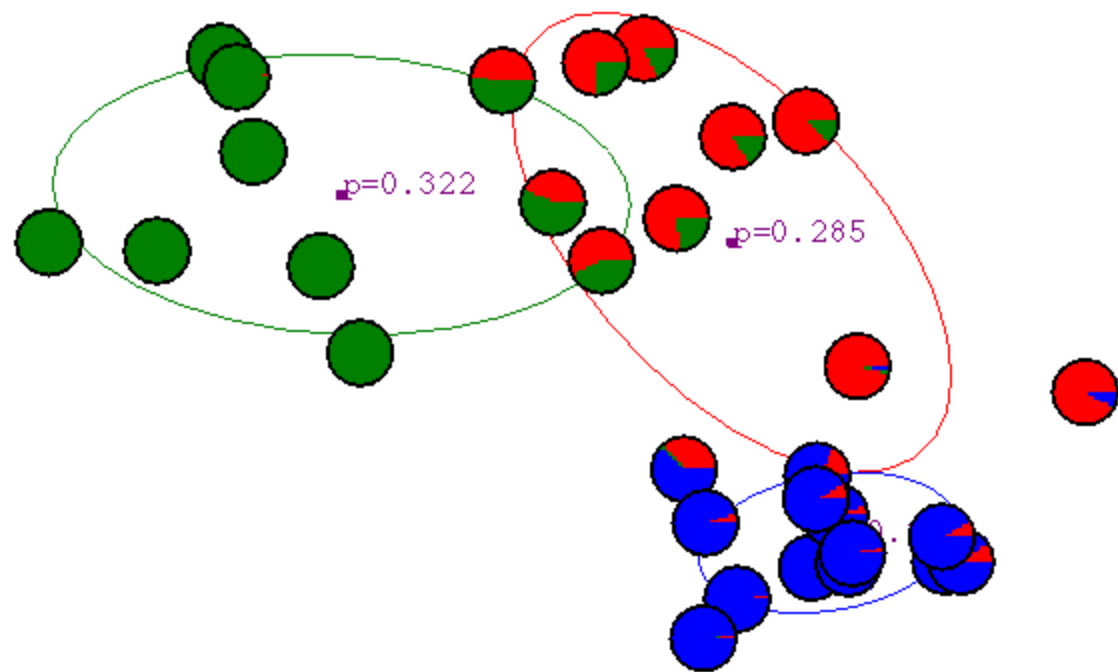
After 3rd iteration



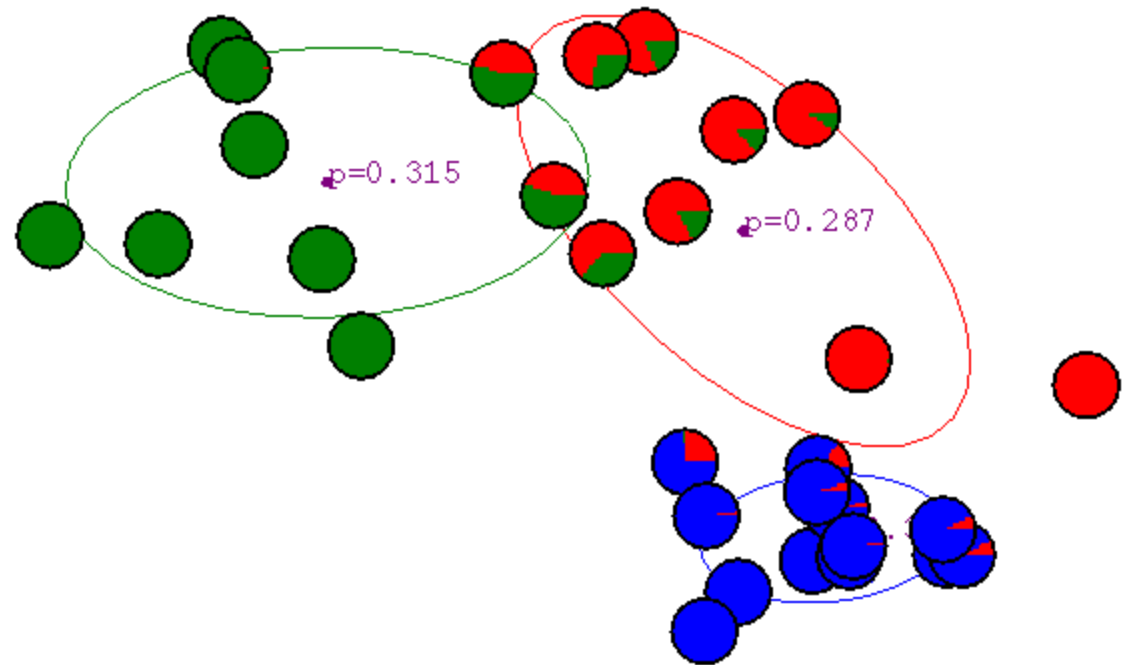
After 4th iteration



After 5th iteration

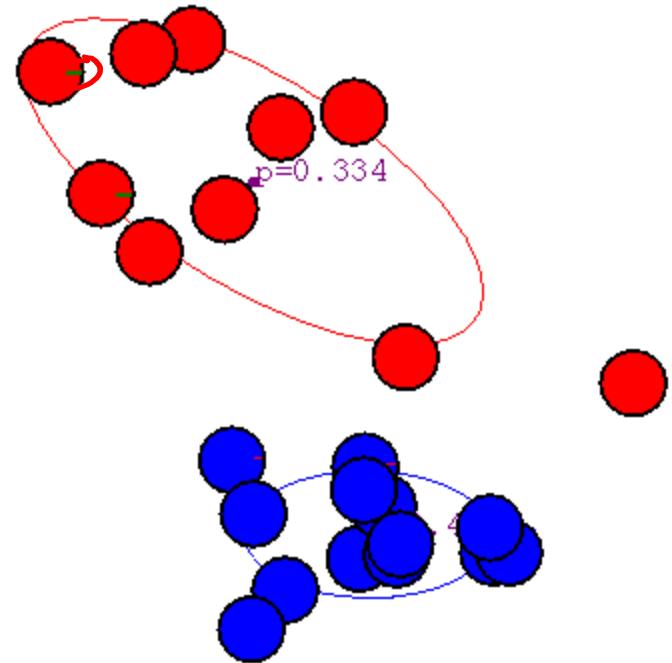
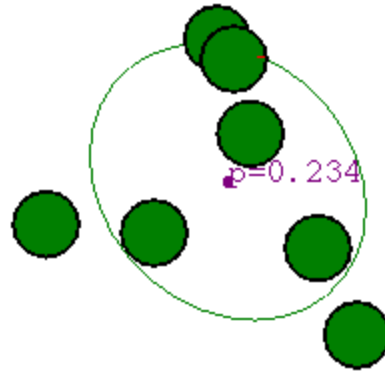


After 6th iteration



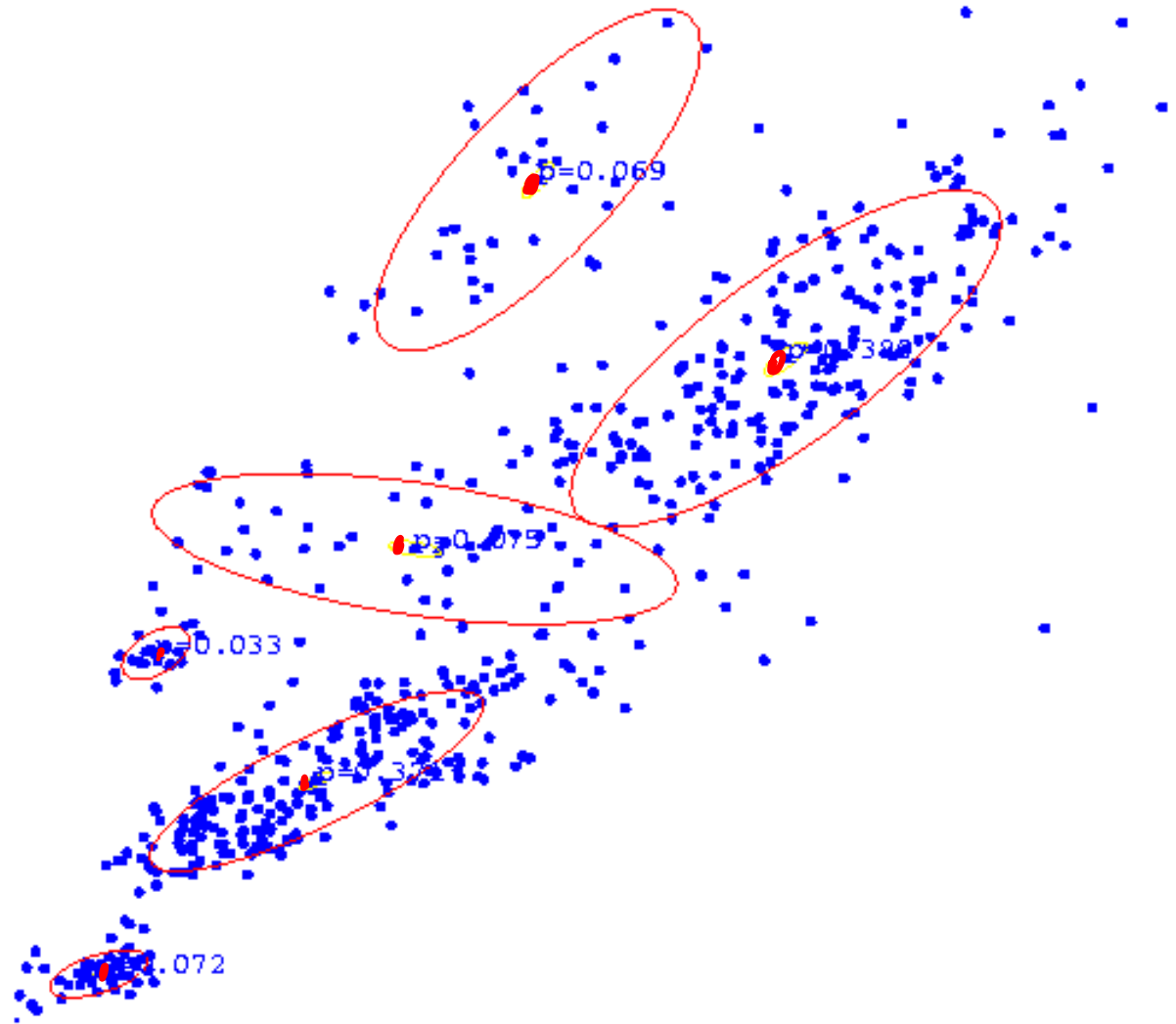
After 20th iteration

$p(x)$
↑
 $p(y=i|x)$

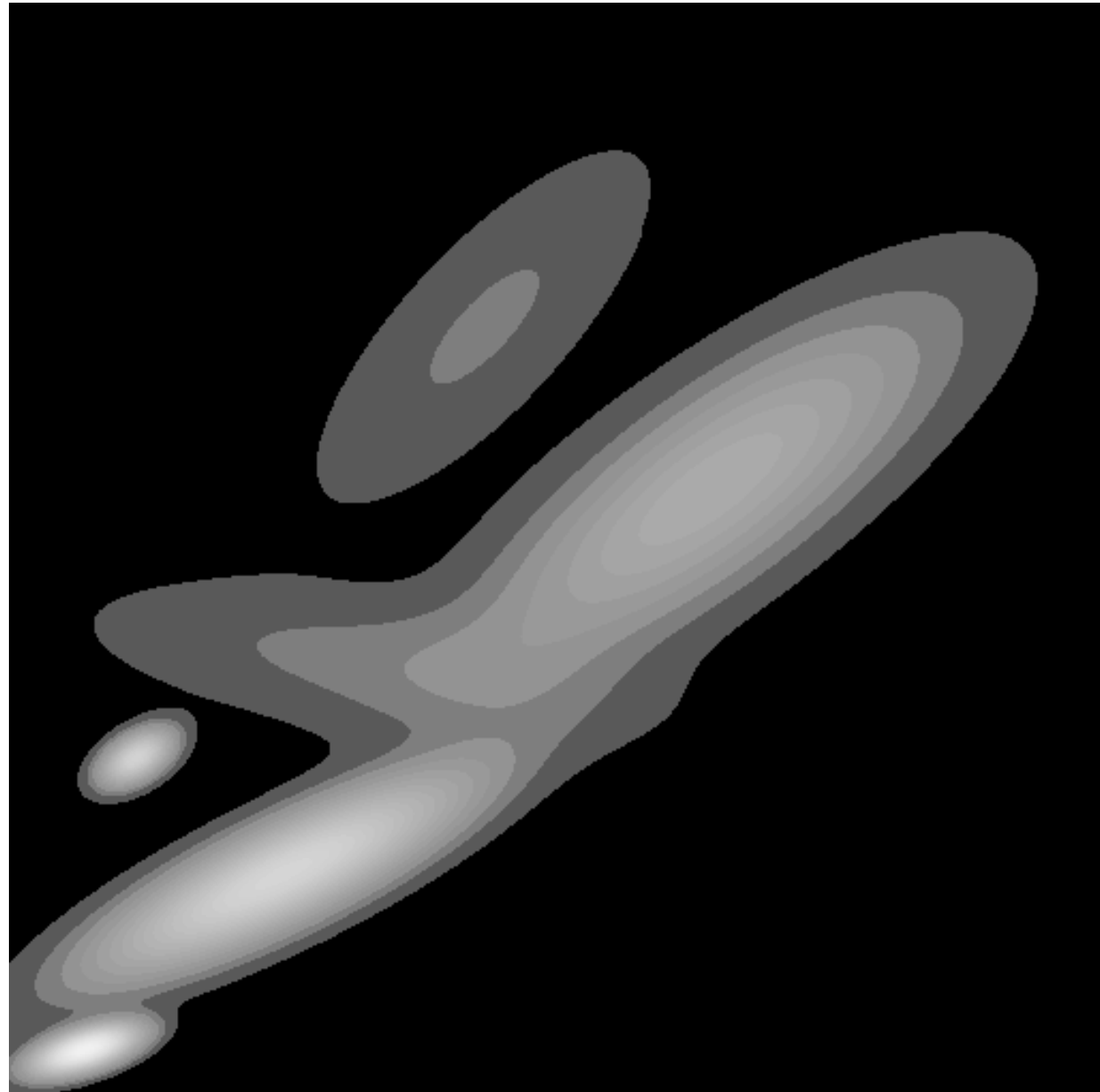


GMM clustering of assay data

$k=6$



Resulting Density Estimator

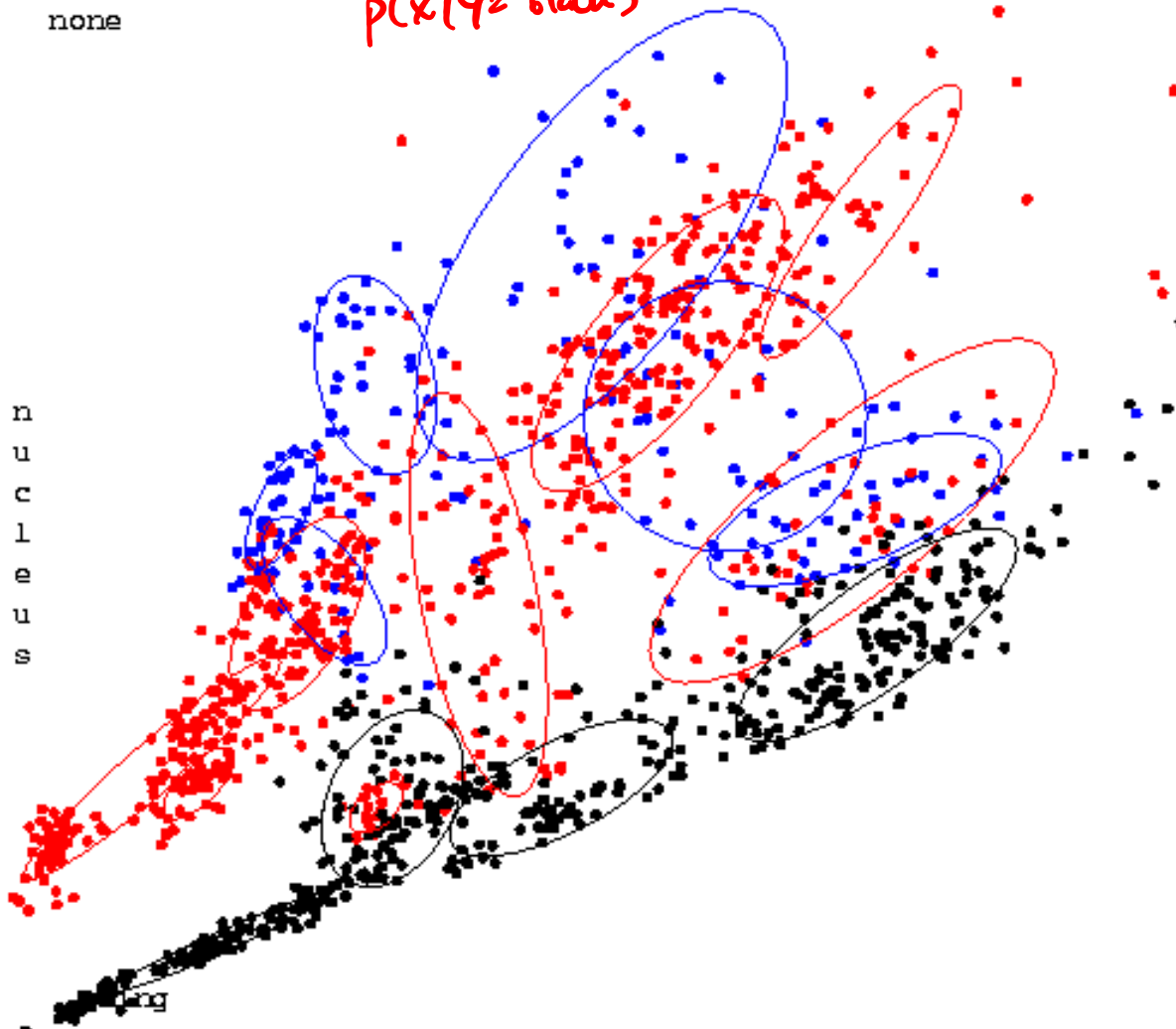


Three classes of assay

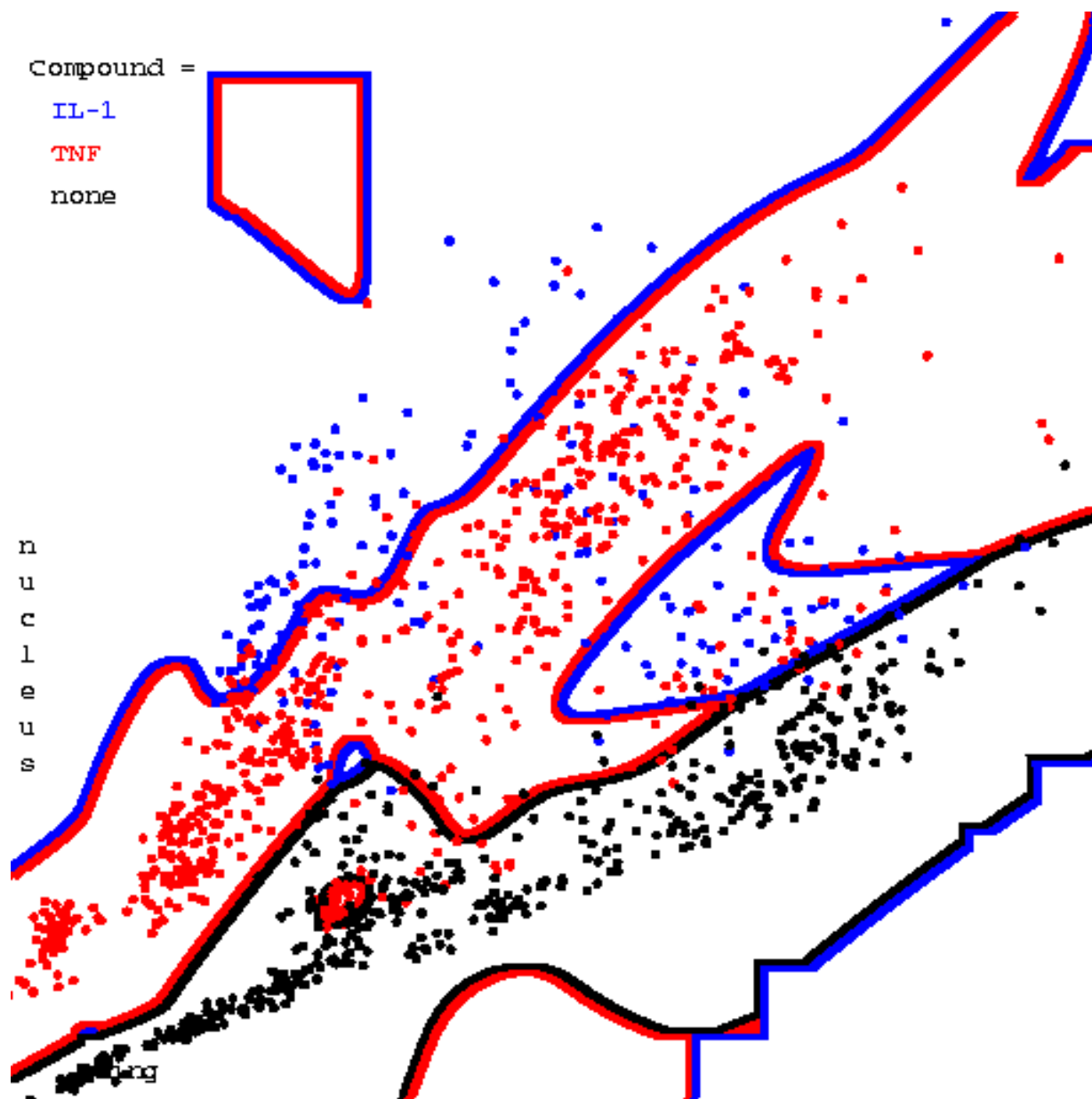
(each learned with its own mixture model)

Compound =
IL-1
TNF
none

$$p(x|Y=blue) \checkmark \leftarrow \mu_y^{(1)}, \Sigma_y^{(1)}$$
$$p(x|Y=red)$$
$$p(x|Y=black)$$
$$\mu_y^{(k)}, \Sigma_y^{(k)}$$



Resulting Bayes Classifier



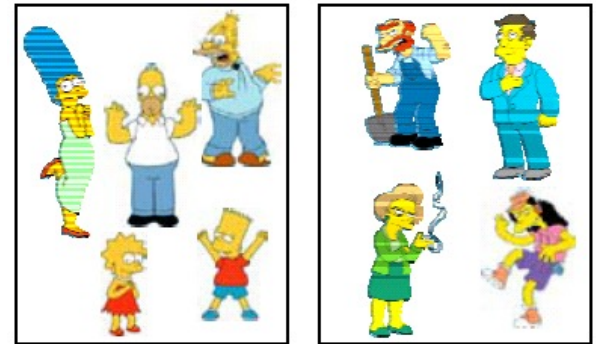
Expectation-Maximization (EM)

A general algorithm to deal with hidden data

- No need to choose step size as in Gradient methods.
- EM is an Iterative algorithm with two linked steps:
 - E-step: fill-in hidden data (Y) using inference ←
 - M-step: apply standard MLE/MAP method to estimate parameters ←
 $\{\rho_i, \mu_i, \Sigma_i\}_{i=1}^k$
- This procedure monotonically improves the likelihood (or leaves it unchanged). Thus it always converges to a local optimum of the likelihood.

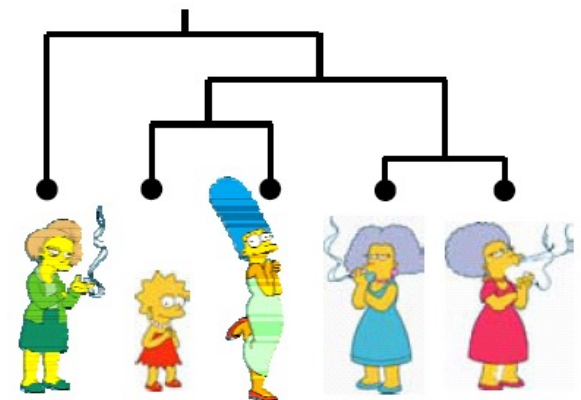
Clustering Algorithms

- Partition algorithms
 - K means clustering
 - Mixture-Model based clustering



- Hierarchical algorithms

- Single-linkage
- Average-linkage
- Complete-linkage
- Centroid-based



Hierarchical Clustering

- Bottom-Up Agglomerative Clustering

Starts with each object in a separate cluster, and repeat:

- Joins the most similar pair of clusters,
 - Update the similarity of the new cluster to others
- until there is only one cluster.



Greedy - less accurate but simple to implement

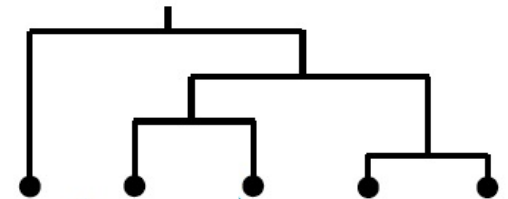
- Top-Down divisive

Starts with all the data in a single cluster, and repeat:

- Split each cluster into two using a partition algorithm
- Until each object is a separate cluster.



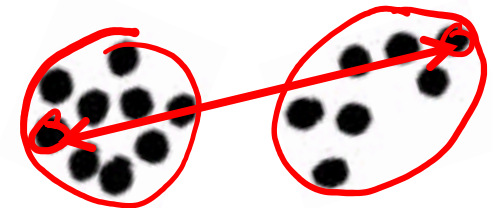
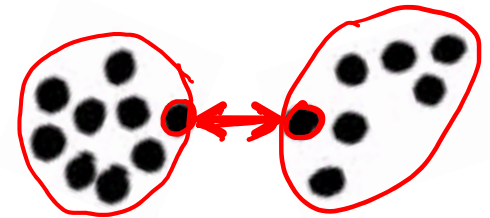
More accurate but complex to implement



Bottom-up Agglomerative clustering

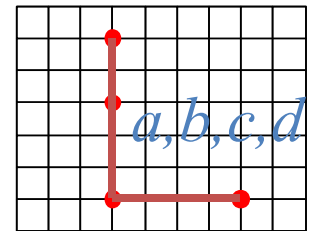
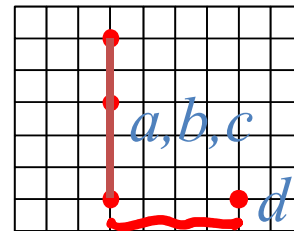
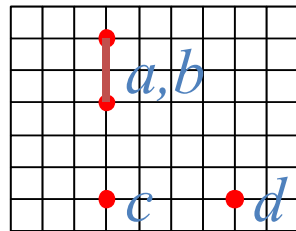
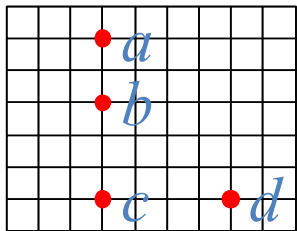
Different algorithms differ in how the similarities are defined (and hence updated) between two clusters

- Single-Linkage
 - Nearest Neighbor: similarity between their closest members.
- Complete-Linkage
 - Furthest Neighbor: similarity between their furthest members.
- Centroid
 - Similarity between the centers of gravity
- Average-Linkage
 - Average similarity of all cross-cluster pairs.



Single-Linkage Method

Euclidean Distance



(1)

(2)

(3)

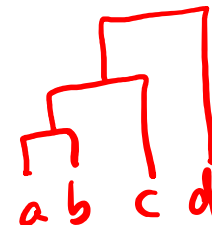
	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

	<i>c</i>	<i>d</i>
<i>a, b</i>	3	5
<i>c</i>		4

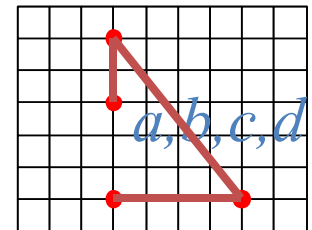
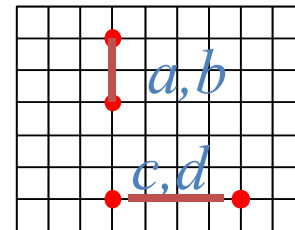
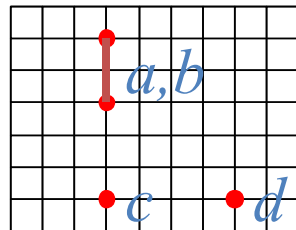
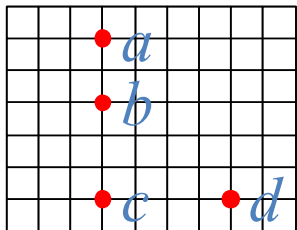
	<i>d</i>
<i>a, b, c</i>	4

Distance Matrix



Complete-Linkage Method

Euclidean Distance



(1)

(2)

(3)

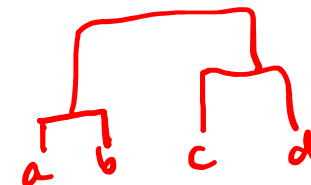
	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

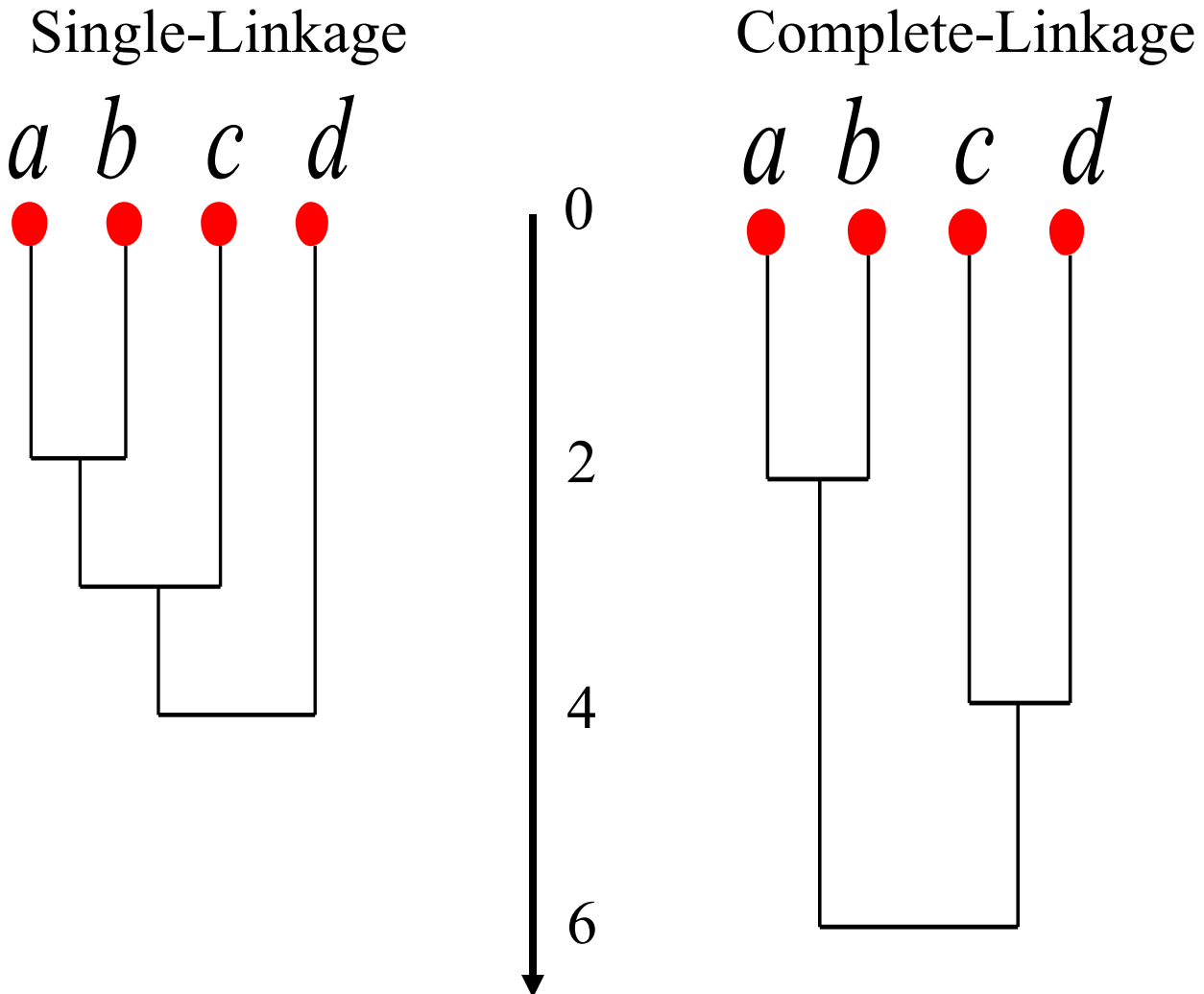
	<i>c</i>	<i>d</i>
<i>a, b</i>	5	6
<i>c</i>		4

	<i>c, d</i>
<i>a, b</i>	6

Distance Matrix

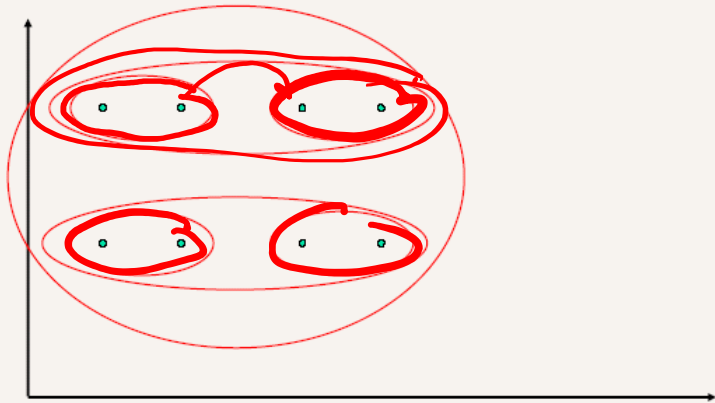


Dendrograms

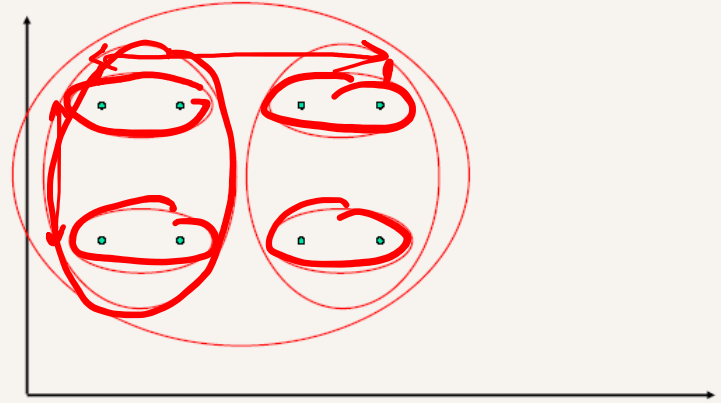


Another Example

Single Link Example



Complete Link Example



Single vs. Complete Linkage

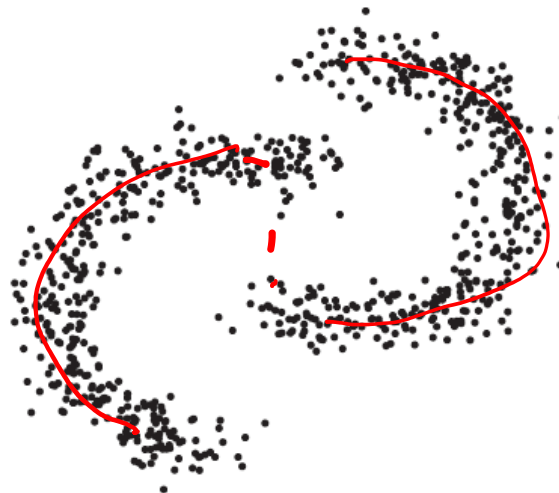
Shape of clusters

Single-linkage

allows anisotropic and non-convex shapes

Complete-linkage

assumes isotropic, convex shapes



Computational Complexity

- All hierarchical clustering methods need to compute similarity of all pairs of n individual instances which is $O(n^2)$.
- At each iteration,
 - Sort similarities to find largest one $O(n^2 \log n)$.
 - Update similarity between merged cluster and other clusters.
Computing similarity to each other cluster can be done in constant time.
- So we get $O(n^2 \log n)$ or $O(n^3)$ (if naively implemented)

$nK \times \text{iteration}$

What you need to know...

- Partition based clustering algorithms
 - K-means ✓
 - Coordinate descent
 - Seeding
 - Choosing K
 - Mixture models ✓
EM algorithm
- Hierarchical clustering algorithms
 - Single-linkage
 - Complete-linkage
 - Centroid-linkage
 - Average-linkage