

Expectation-Maximization (EM)

Aarti Singh

Machine Learning 10-315

Nov 23, 2020

Some slides courtesy of Eric Xing, Carlos Guestrin



MACHINE LEARNING DEPARTMENT

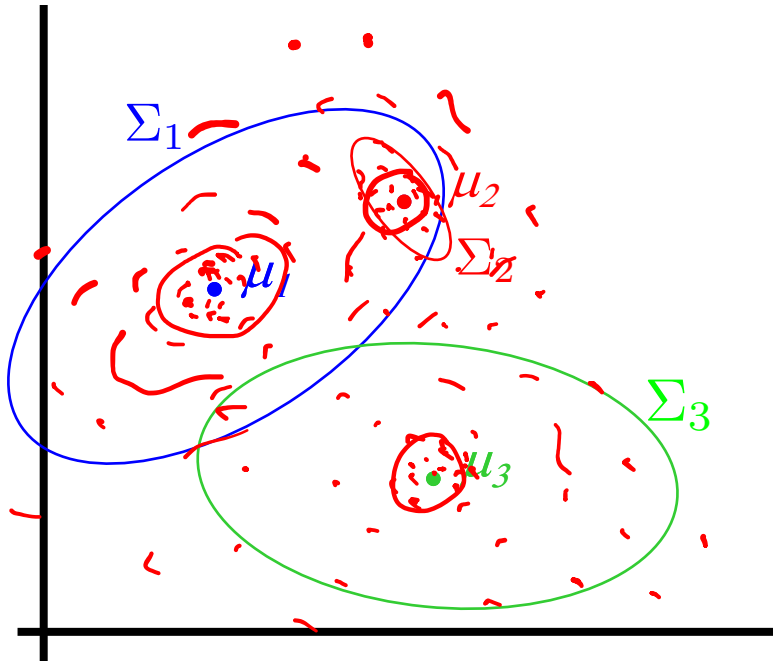
Carnegie Mellon.
School of Computer Science

Mixture models (Gaussian)

~~...~~

$$\rightarrow x_1, \dots, x_m \sim p(x) = \sum_{i=1}^k p(x|Y = i) P(Y = i)$$

$x_i \in \mathbb{R}^d$ (handwritten)
 $\mathcal{N}(\mu_i, \Sigma_i)$ (handwritten)
 Mixture component (handwritten)
 Mixture proportion, p_i (handwritten)



Gaussian mixture model


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


Parameters: $\{p_i, \mu_i, \Sigma_i\}_{i=1}^K$

- How to estimate parameters? Max ^{Marginal} Likelihood $p(x)$
 But don't know labels Y (recall Gaussian Bayes classifier)

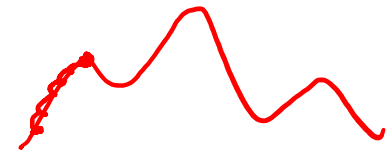
Expectation-Maximization (EM)

A general algorithm to deal with hidden data, but we will study it in the context of unsupervised learning (hidden labels)

- 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 marginal likelihood (or leaves it unchanged). Thus it always converges to a local optimum of the likelihood.



EM for spherical, same variance GMMs

same mixture proportions $\rightarrow \Sigma_i$

$\rightarrow P_i = P(Y=i)$

\rightarrow Initialize: $\mu_1, \mu_2, \dots, \mu_K$ randomly

E-step

Compute "expected" classes of all datapoints for each class

$\forall i, j$

$$P(y=i | x_j, \mu_1, \dots, \mu_K) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y=i)$$

$=$

In K-means "E-step"
we do hard assignment

EM does soft assignment

$$\Sigma_i = \sigma^2 I$$

$$P_i = \frac{1}{K}$$

$$\{P_i, \mu_i, \Sigma_i\}_{i=1}^K$$

$$p(y=i | x) \propto \frac{p(x | y=i) p(y=i)}{\sim N(\mu_i, \Sigma_i)}$$

Bayes rule

$$p(y | x) = \frac{p(x | y) p(y)}{p(x)}$$

$$p(x) = \sum_{i=1}^K p(x | y=i) p(y=i)$$

$\sim N(\mu_i, \Sigma_i)$ P_i

EM for spherical, same variance GMMs same mixture proportions

Initialize: $\mu_1, \mu_2, \dots, \mu_K$ randomly

E-step

Compute “expected” classes of all datapoints for each class

$$P(y = i | x_j, \mu_1, \dots, \mu_k) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y = i)$$

In K-means “E-step”
we do hard assignment

EM does soft assignment

M-step

Compute Max. like μ given our data’s class membership distributions (weights)

$$\mu_i = \frac{\sum_{j=1}^m P(y = i | x_j) x_j}{\sum_{j=1}^m P(y = i | x_j)}$$

$$\frac{1}{m} \sum_{j=1}^m x_j \mathbb{1}_{x_j \in C_i} \leftarrow \text{K-means}$$

Exactly same as MLE with
weighted data

Iterate.

EM for general GMMs

Iterate. On iteration t let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \dots p_k^{(t)} \}$$

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

Initialize: $\lambda_0 \leftarrow$ random
E-step

Compute "expected" classes of all datapoints for each class

$$\rightarrow P(y = i | x_j, \lambda_t) \propto p_i^{(t)} p(x_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

$$\equiv \overline{P(x_j | y=i)}$$

Just evaluate a Gaussian at x_j

M-step

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

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

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

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

$\sum_j \mathbb{1}_{x_j \in C_i}$
 $\equiv \frac{m_i}{m}$ ← equal wts.

$m = \#$ data points

EM for general GMMs: Example

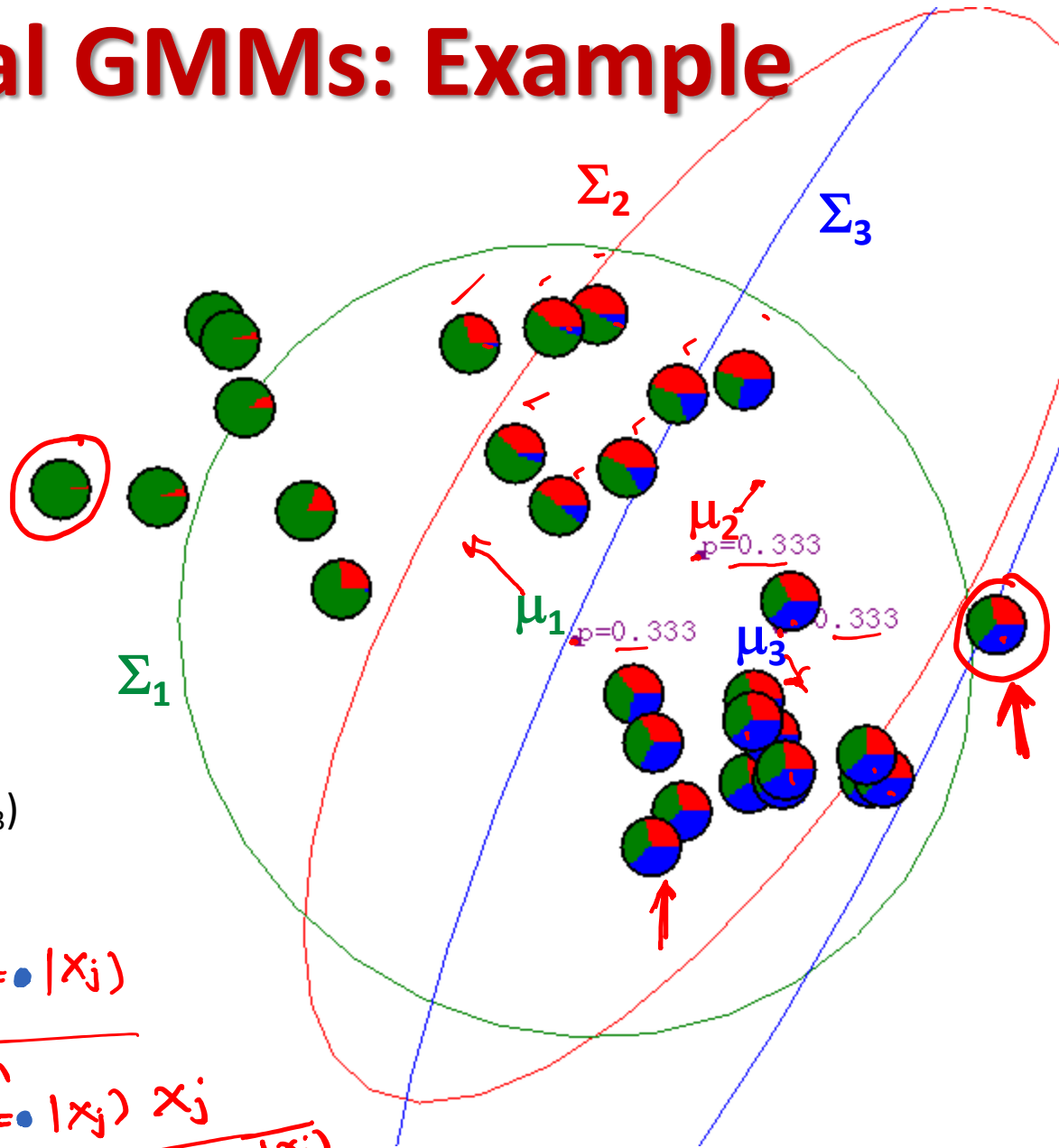
$K=3$

① Random initialization

μ_1, μ_2, μ_3

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

$\Sigma_1, \Sigma_2, \Sigma_3$



② E-step

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



$p_3 \leftarrow$

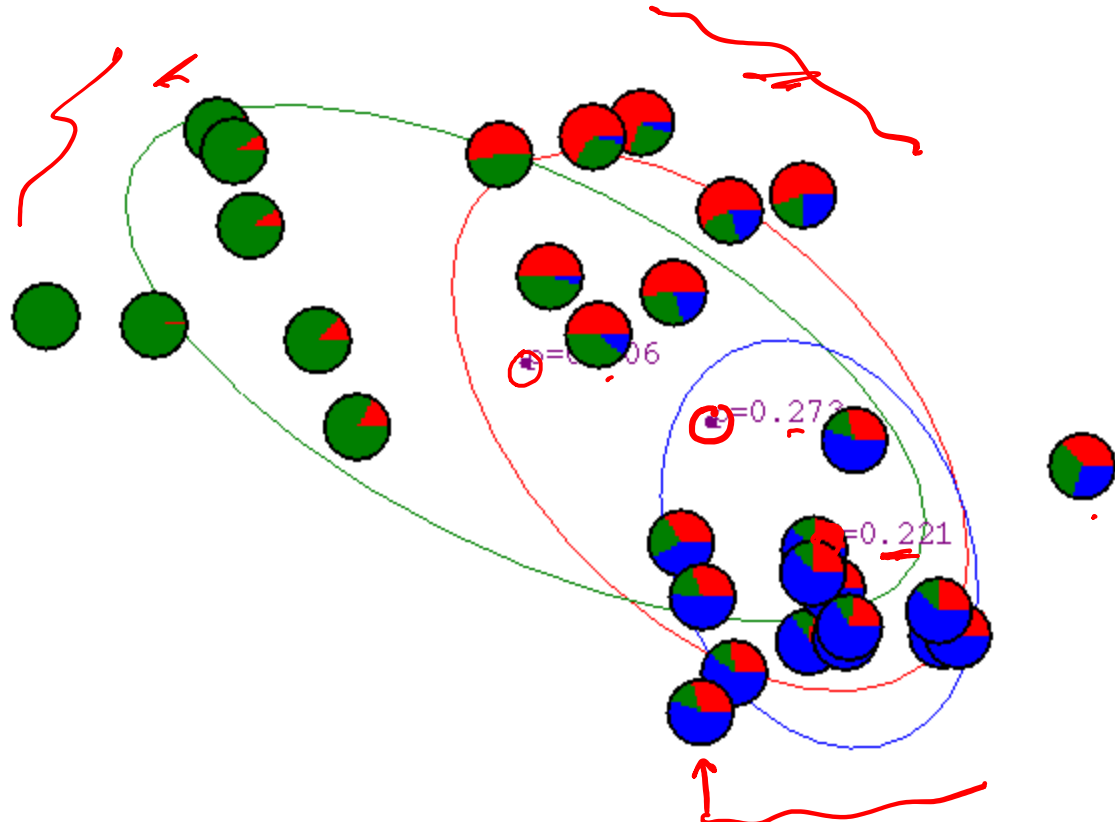
$$\frac{\sum_{j=1}^m p(y = \bullet | x_j)}{m}$$

③ M-step

$\mu_3 \leftarrow$

$$\frac{\sum_{j=1}^m p(y = \bullet | x_j) x_j}{\sum_{j=1}^m p(y = \bullet | x_j)}$$

After 1st iteration



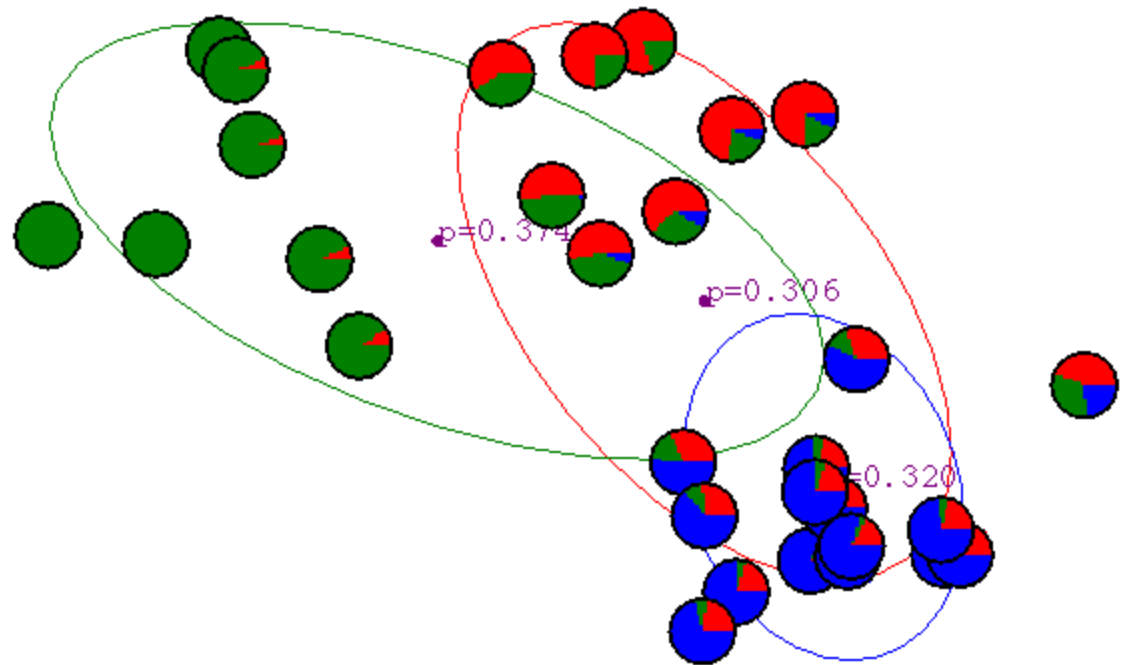
① E-step

$$P(y = \square | X_j, \dots)$$

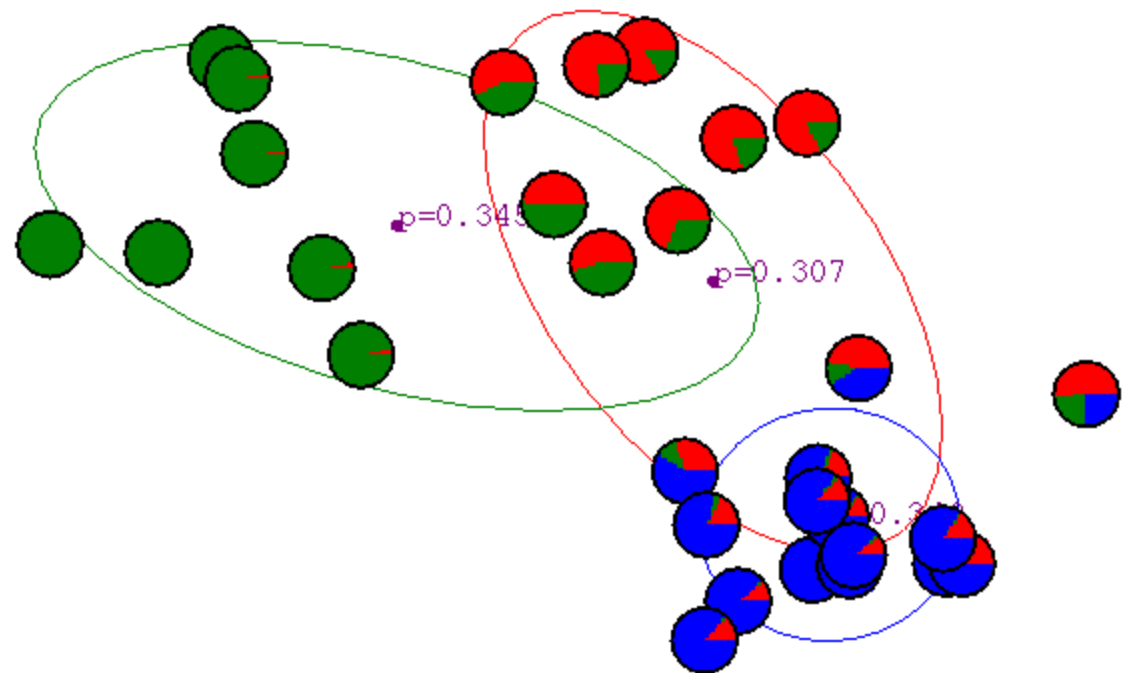
② M-step

$$\mu_i, \Sigma_i, P_i$$

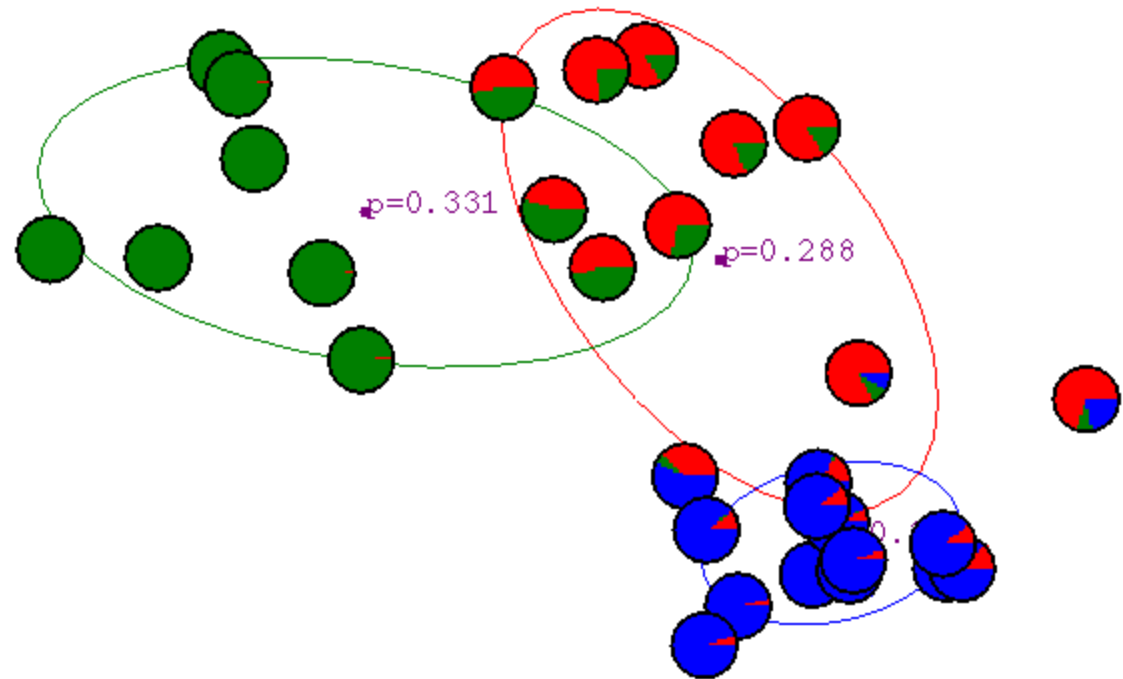
After 2nd iteration



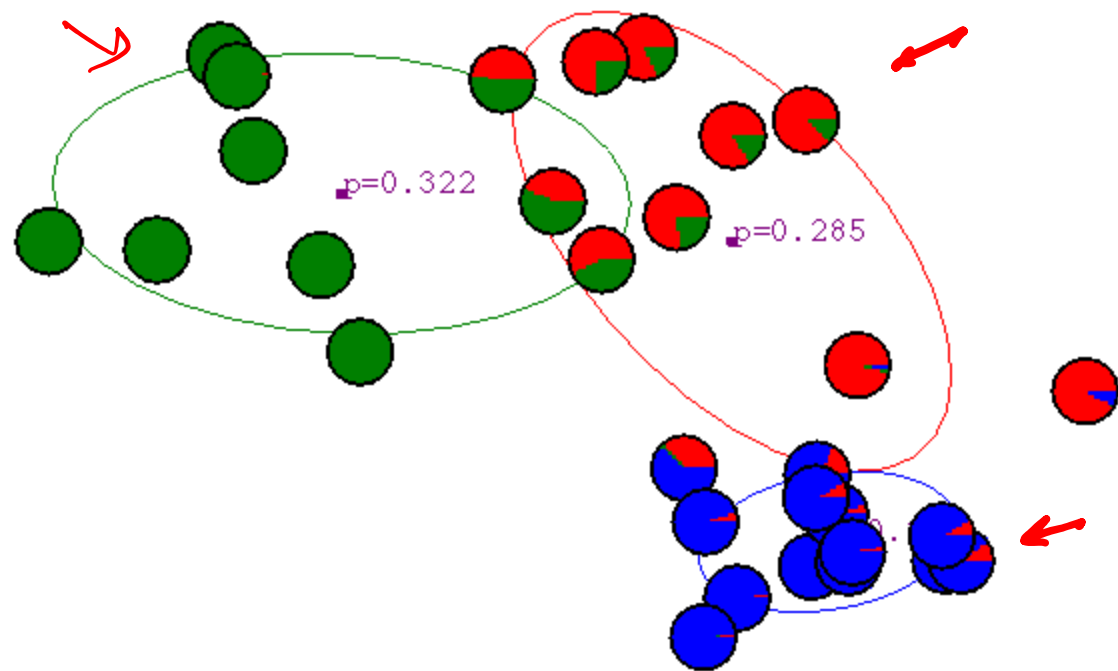
After 3rd iteration



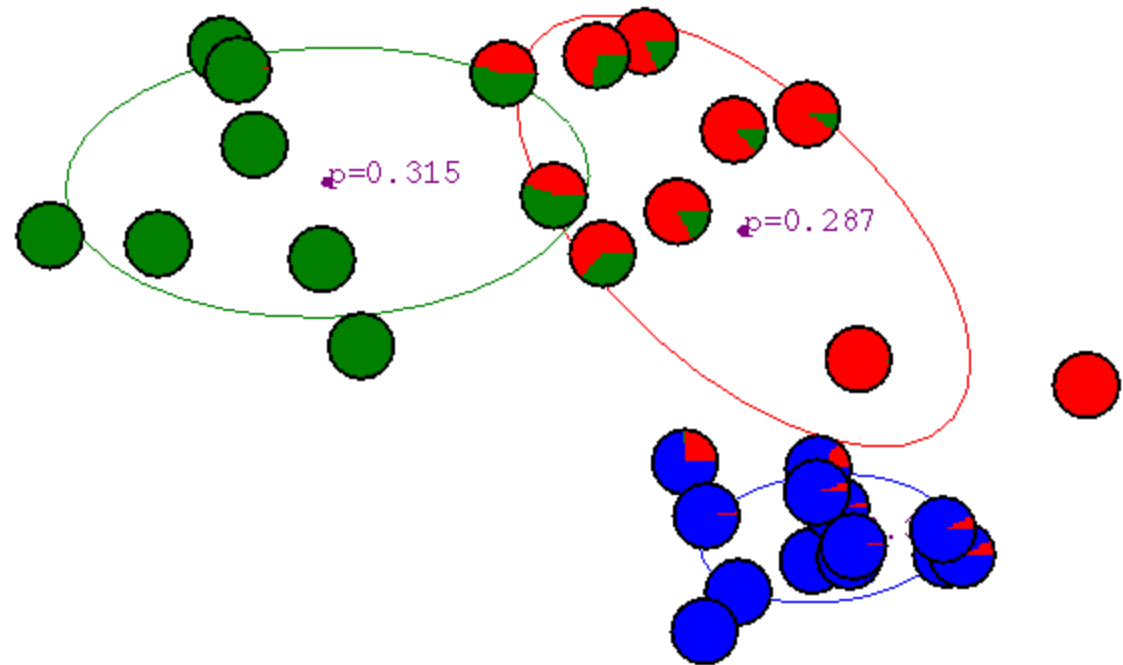
After 4th iteration



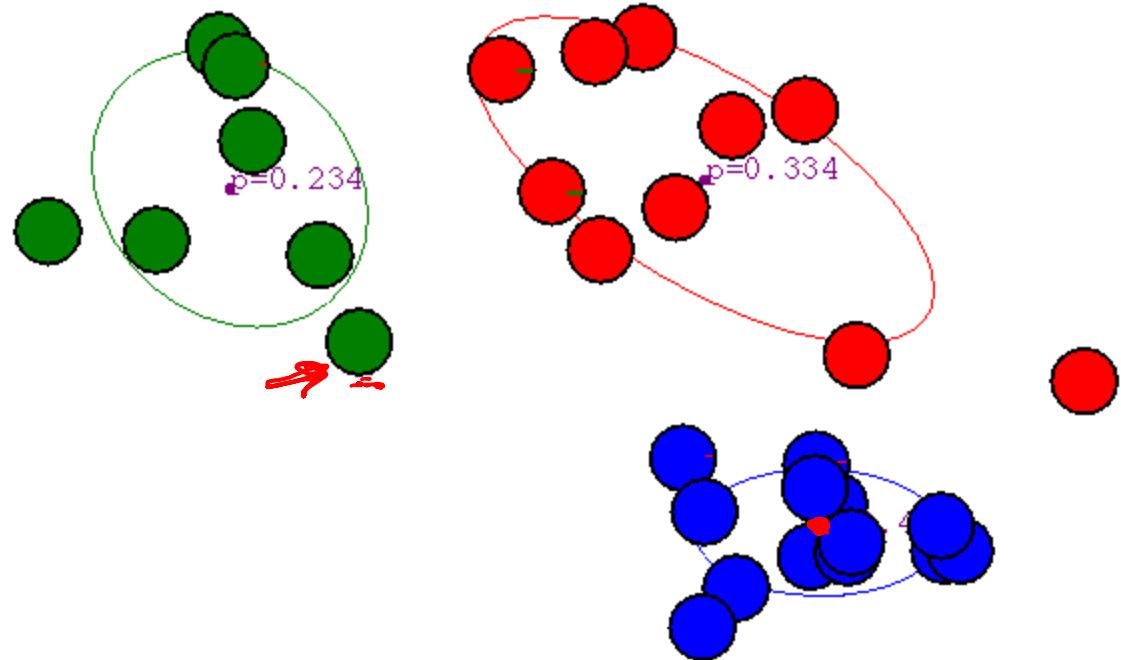
After 5th iteration



After 6th iteration

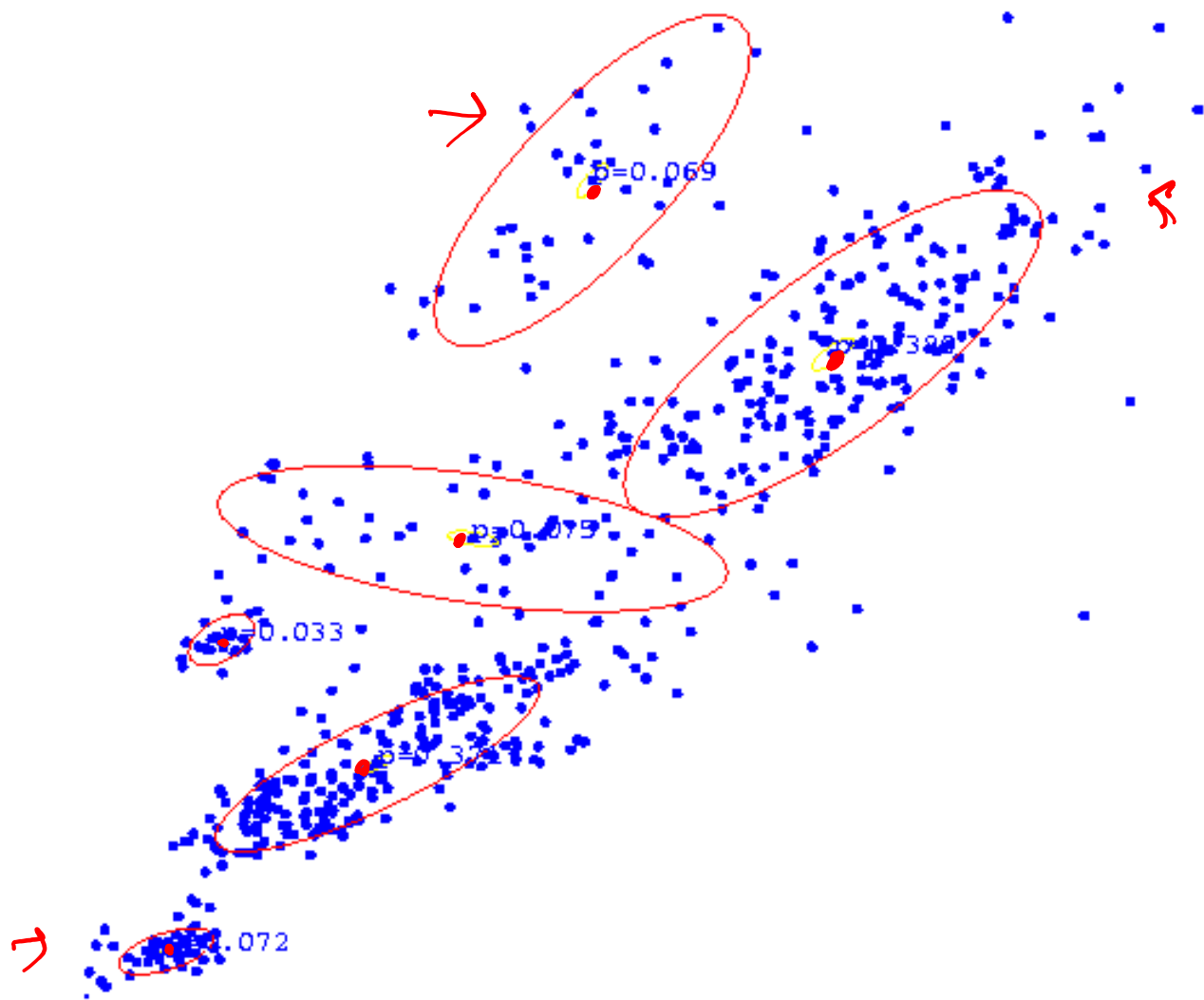


After 20th iteration



GMM clustering of assay data

$p(x)$
K choice
"6"



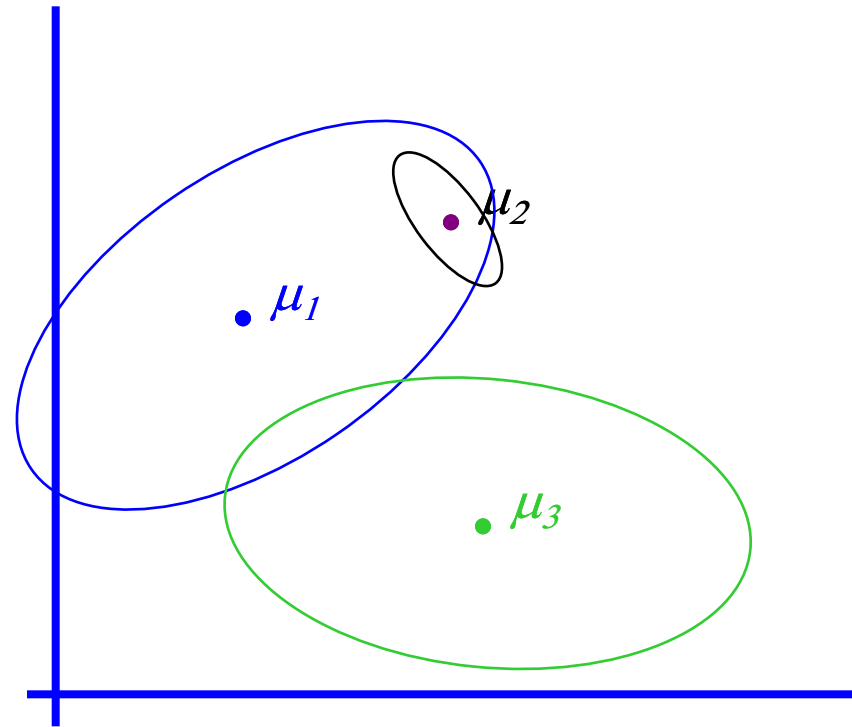
General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

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

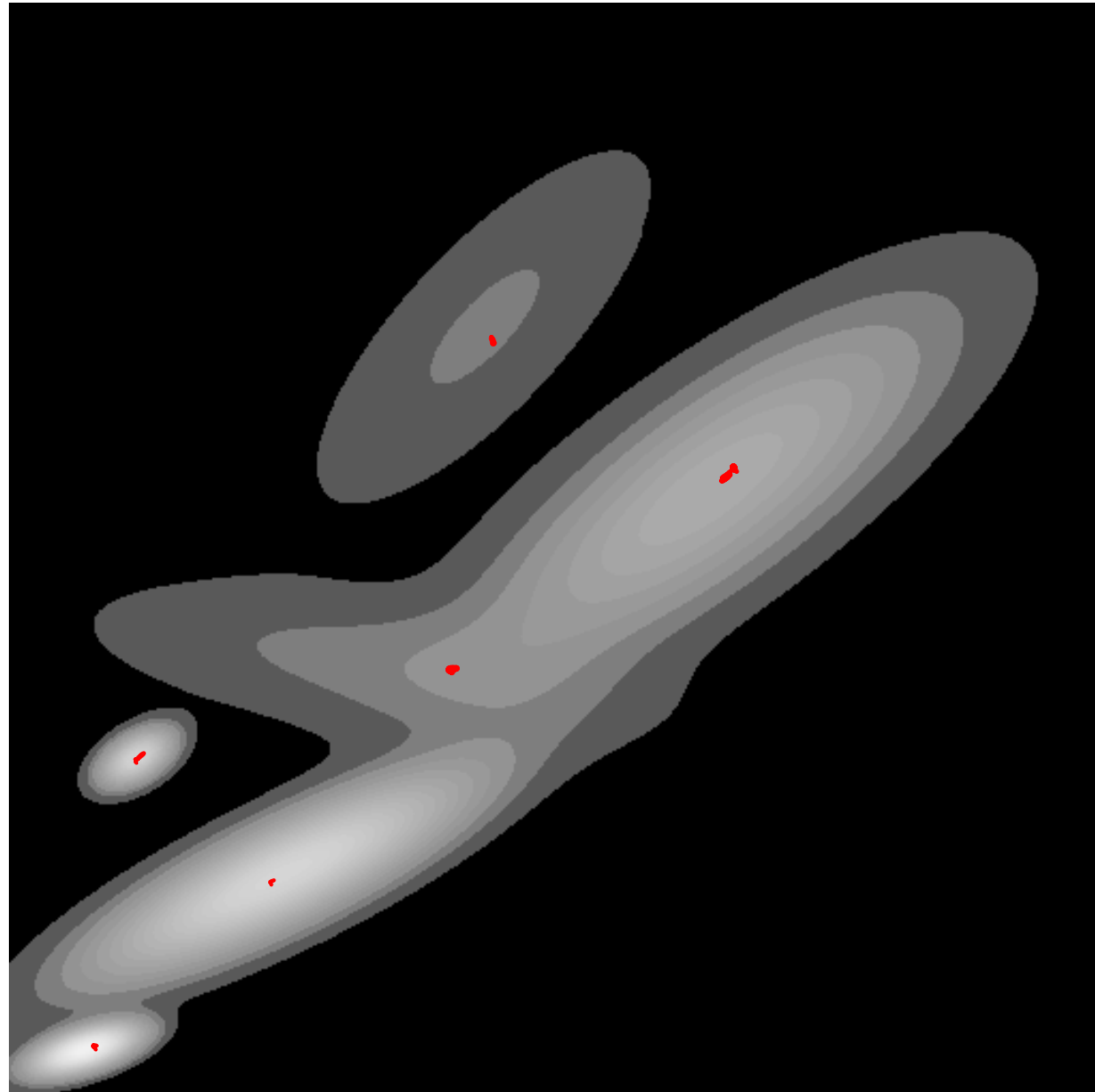
↓ ↓
Mixture **Mixture**
component **proportion**

$$p(x|y=i) \sim N(\mu_i, \Sigma_i)$$



Resulting Density Estimator

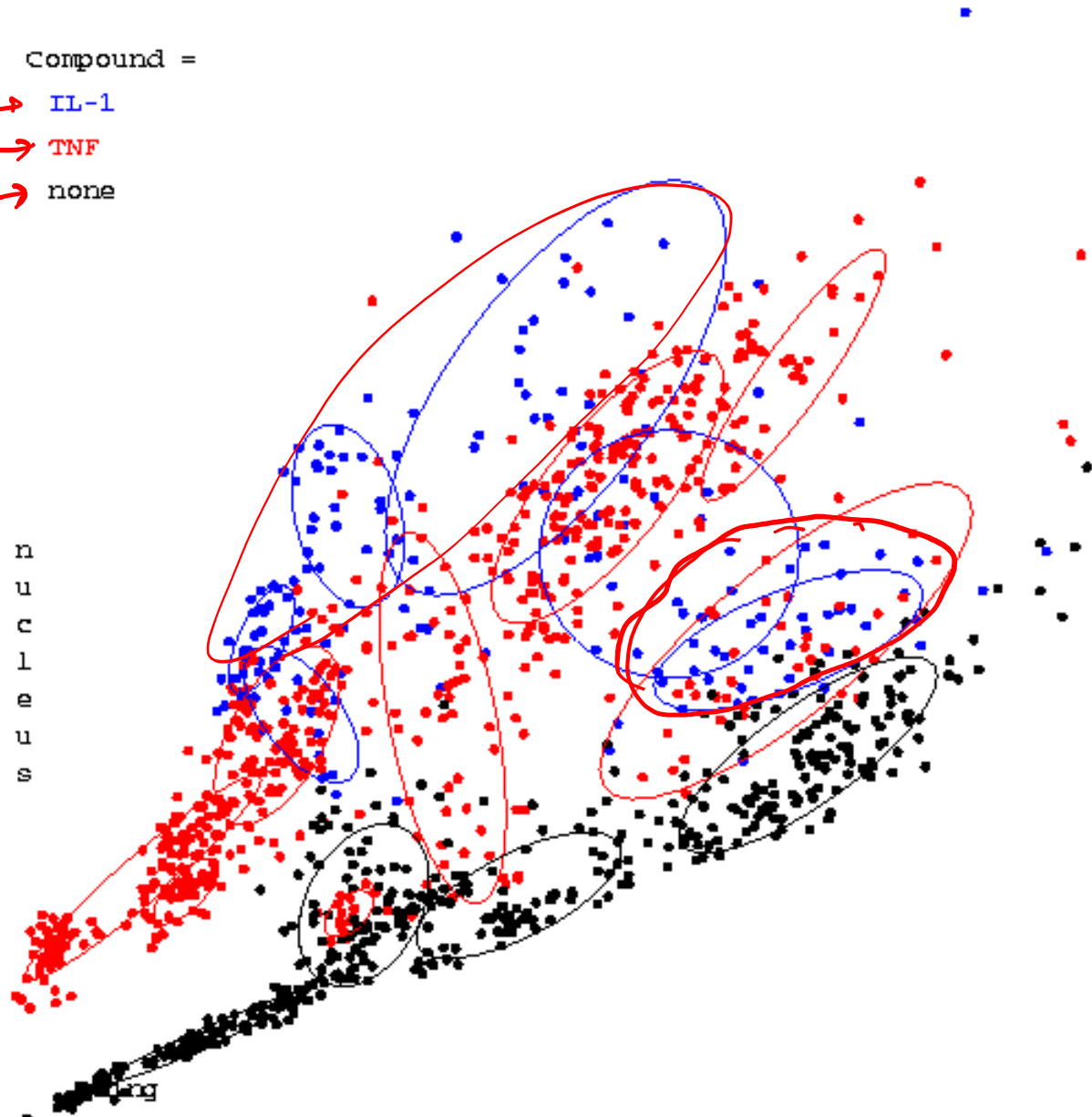
Contour plots
of $p(x)$
 $p(x|y=\bullet) \sim \text{GMM}$
- multimodal



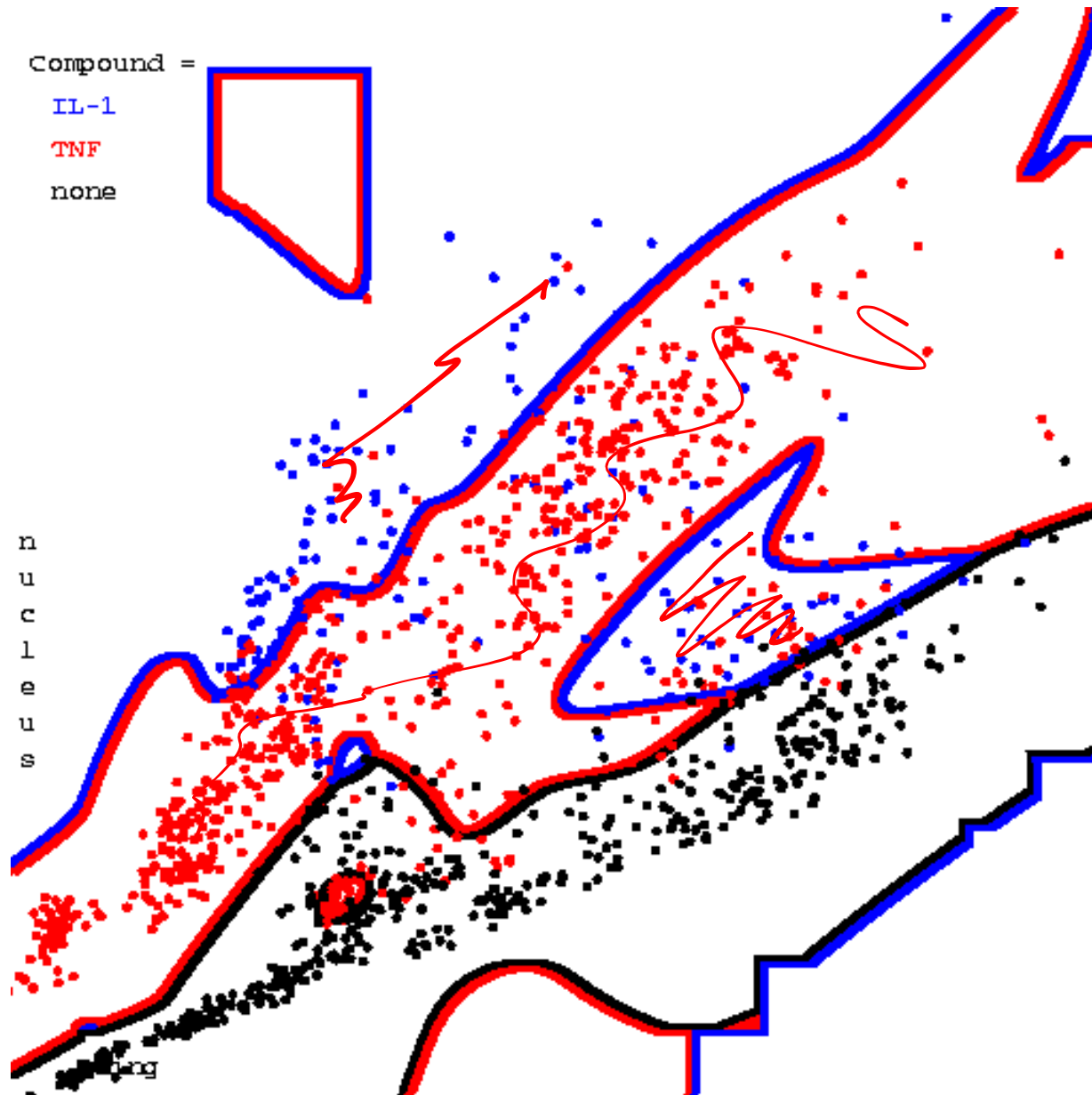
Three classes of assay

(each learned with its own mixture model)

Compound =
→ IL-1
→ TNF
→ none



Resulting Bayes Classifier

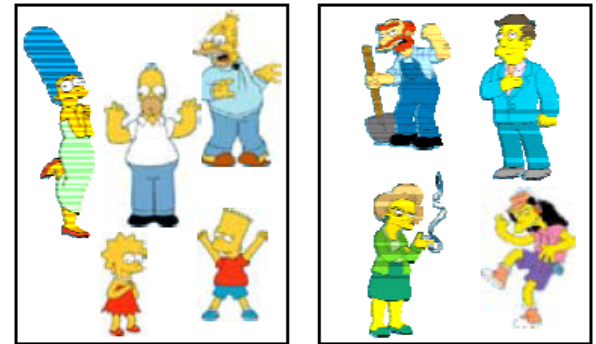


Summary: EM Algorithm

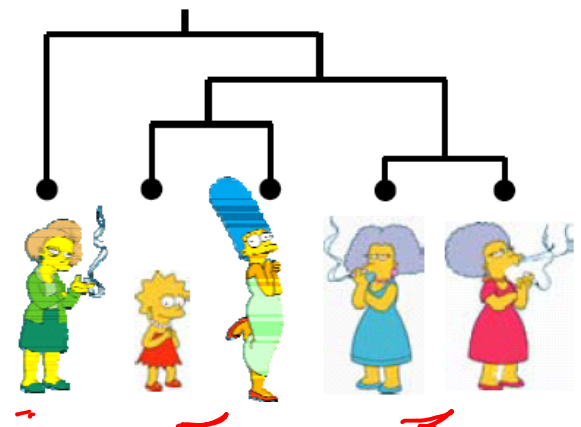
- A way of maximizing likelihood function for hidden variable models. Finds MLE of parameters when the original (hard) problem can be broken up into two (easy) pieces:
 1. Estimate some “missing” or “unobserved” data from observed data and current parameters.
 2. Using this “complete” data, find the maximum likelihood parameter estimates.
- Alternate between filling in the latent variables using the best guess (posterior) and updating the parameters based on this guess:
 1. E-step: soft cluster assignment for each data point
 2. M-step: update parameters of each mixture component
- EM can get stuck in local minima. *though guaranteed to converge.*
- BUT Extremely popular in practice.

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

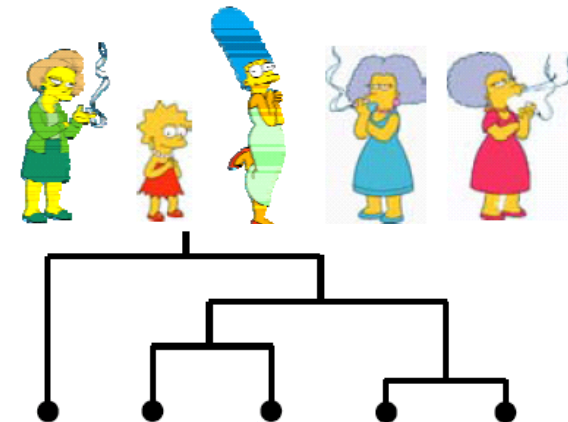


- Top-Down divisive

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

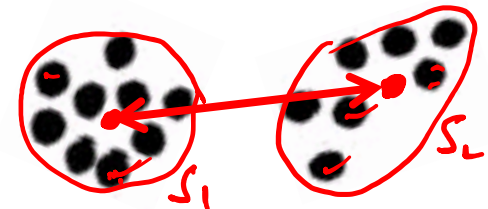
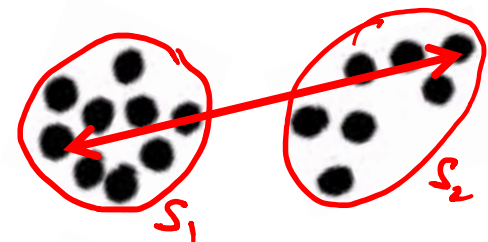
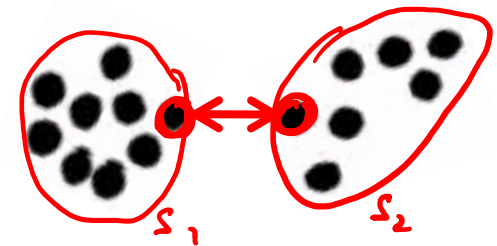




Bottom-up Agglomerative clustering

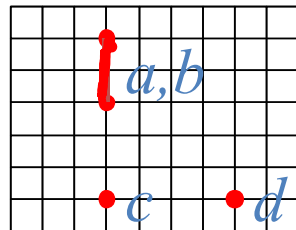
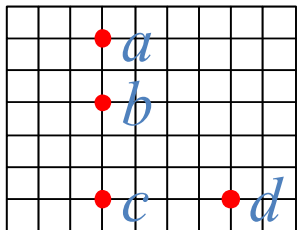
Different algorithms differ in how the similarities are defined (and hence updated) between two clusters

- 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.

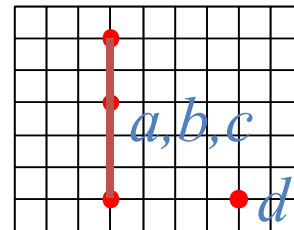


Single-Linkage Method

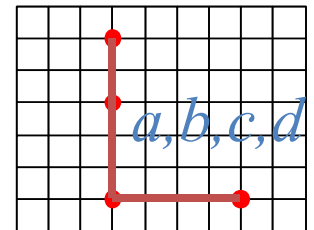
Euclidean Distance



(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

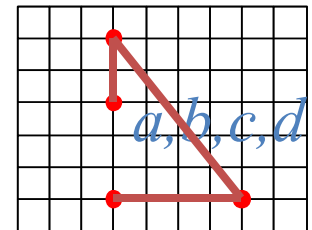
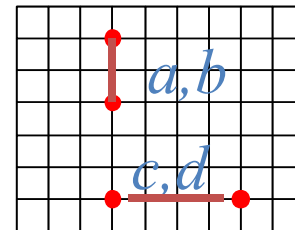
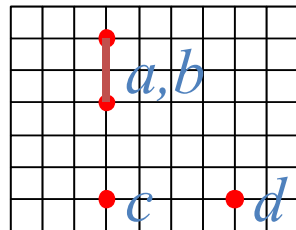
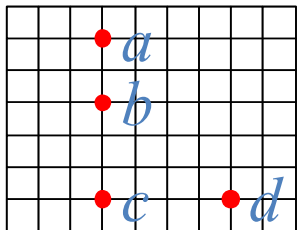
	<i>c</i>	<i>d</i>
<i>a, b</i>	3	5
<i>c</i>		4

	<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

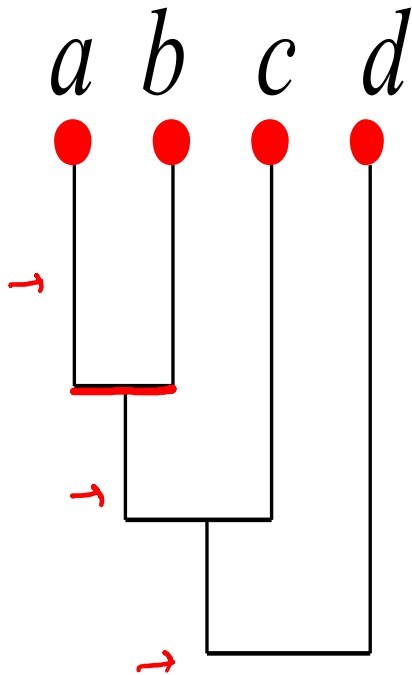
	<i>c</i>	<i>d</i>
<i>a, b</i>	5	6
<i>c</i>		4

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

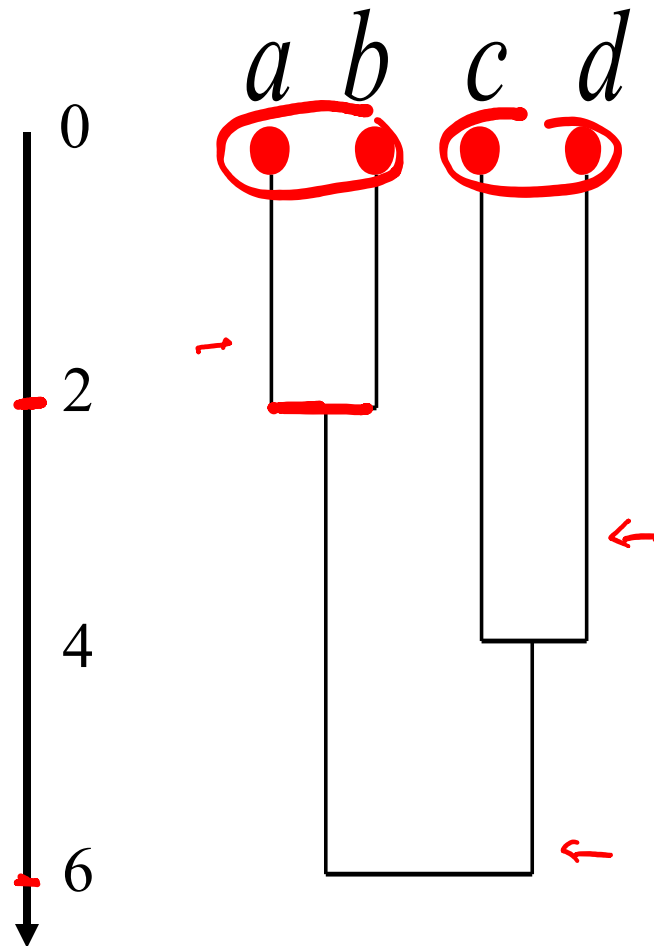
Distance Matrix

Dendrograms

Single-Linkage

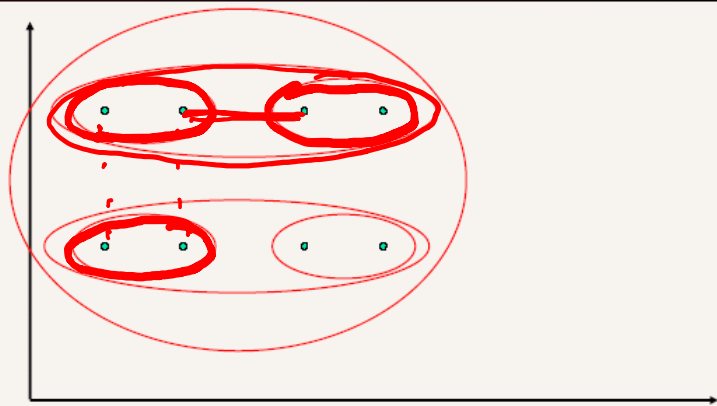


Complete-Linkage

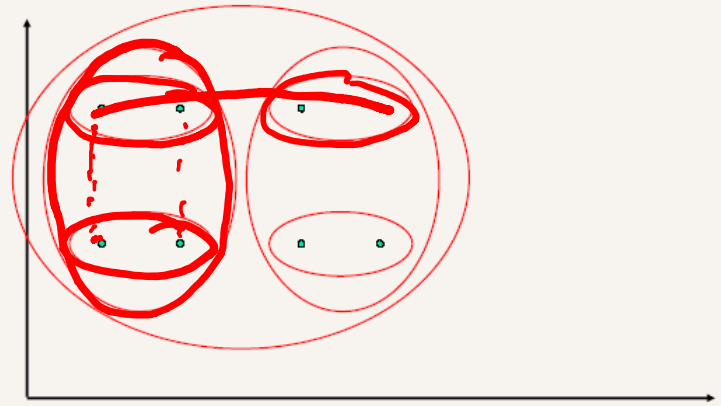


Another Example

Single Link Example



Complete Link Example



Single vs. Complete Linkage

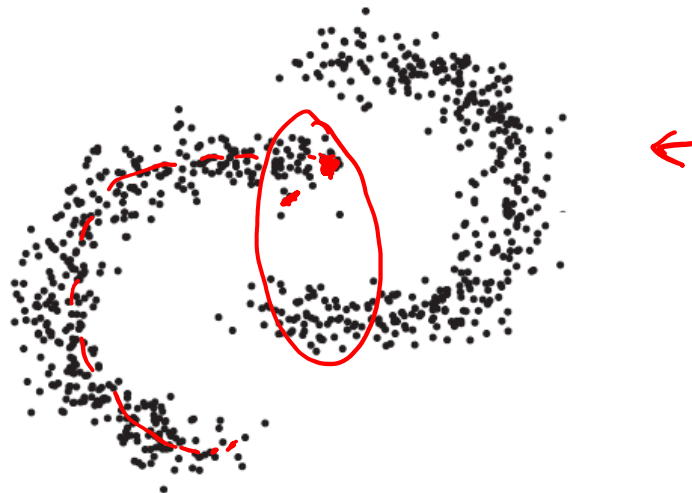
Shape of clusters

Single-linkage

allows anisotropic and non-convex shapes ←

Complete-linkage

assumes isotropic, convex shapes



Computational Complexity

bottom-up (linkage)

- 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)

Computational Complexity (K-means)

- 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: $O(lKn)$.

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 .

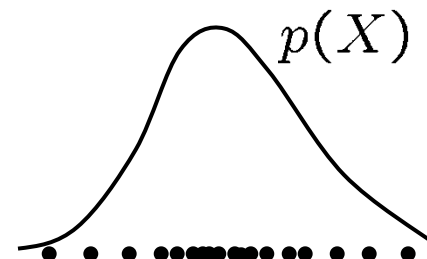
Unsupervised Learning

“Learning from unlabeled/unannotated data” (without supervision)



What can we predict from unlabeled data?

- Density estimation



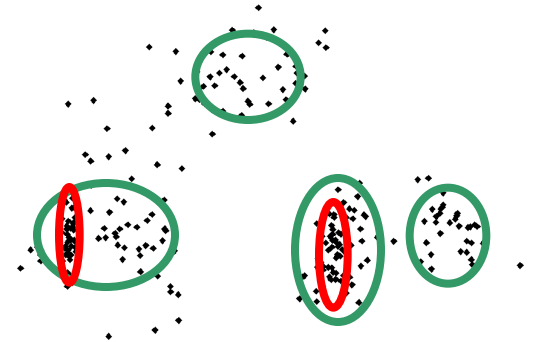
Unsupervised Learning

“Learning from unlabeled/unannotated data” (without supervision)



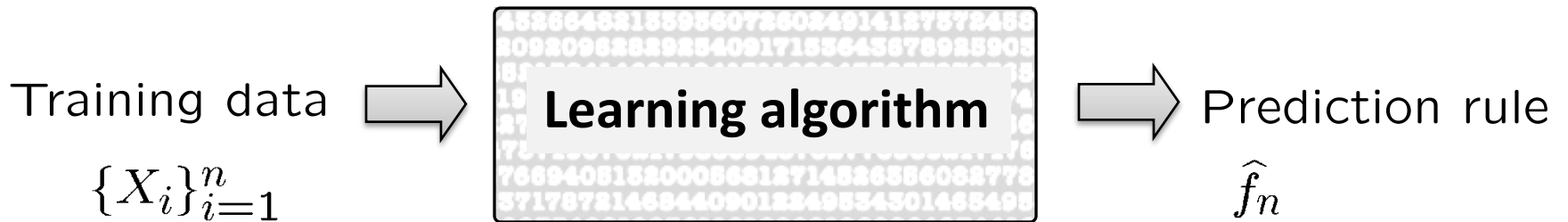
What can we predict from unlabeled data?

- Density estimation
- Groups or clusters in the data





Unsupervised Learning

“Learning from unlabeled/unannotated data” (without supervision)



What can we predict from unlabeled data?

- Density estimation 
- Groups or clusters in the data 
- Dimensionality reduction 