

EM Algorithm

Aarti Singh

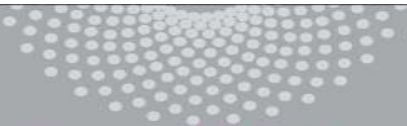
Slides courtesy: Eric Xing, Carlos Guetrin

Machine Learning 10-701/15-781

Feb 17, 2010

The logo consists of the letters 'ML' in a bold, black, sans-serif font. A thick red horizontal line is positioned directly beneath the letters.

MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

K-means Recap ...

- Randomly initialize k centers
 - $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$
- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:
 - $C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$
- **Recenter:** μ_i becomes centroid of its point:
 - $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C(j)=i} \|\mu - x_j\|^2$
 - Equivalent to $\mu_i \leftarrow$ average of its points!

What is K-means optimizing?

- Potential function $F(\mu, C)$ of centers μ and point allocations C :

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

- Optimal K-means:
 - $\min_{\mu} \min_C F(\mu, C)$

K-means algorithm

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$$

- **K-means algorithm:**

(1) Fix μ , optimize C

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

$$= \sum_{j=1}^m \underbrace{\min_{C(j)} \|\mu_{C(j)} - x_j\|^2}_{\text{Exactly first step -- assign each point to the nearest cluster center}}$$

Exactly first step – assign each point to the nearest cluster center

K-means algorithm

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$$

- **K-means algorithm:**

(2) Fix C, optimize μ

$$\begin{aligned} & \min_{\mu_1, \mu_2, \dots, \mu_K} \sum_{i=1}^K \sum_{j:C(j)=i} \|\mu_i - x_j\|^2 \\ &= \sum_{i=1}^K \underbrace{\min_{\mu_i} \sum_{j:C(j)=i} \|\mu_i - x_j\|^2}_{\text{Solution: average of points in cluster } i} \end{aligned}$$

Solution: average of points in cluster i
Exactly second step (re-center)

K-means algorithm

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$$

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

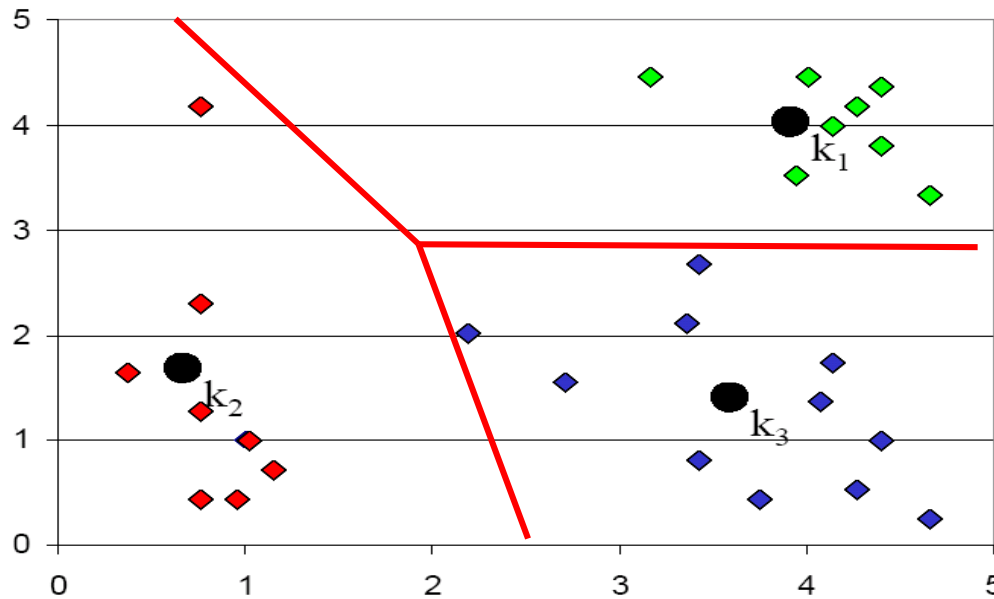
(1) Fix μ , optimize C **Expectation step**

(2) Fix C, optimize μ **Maximization step**

Today, we will see a generalization of this approach:

EM algorithm

K-means Decision boundaries



**“Linear”
Decision
Boundaries**

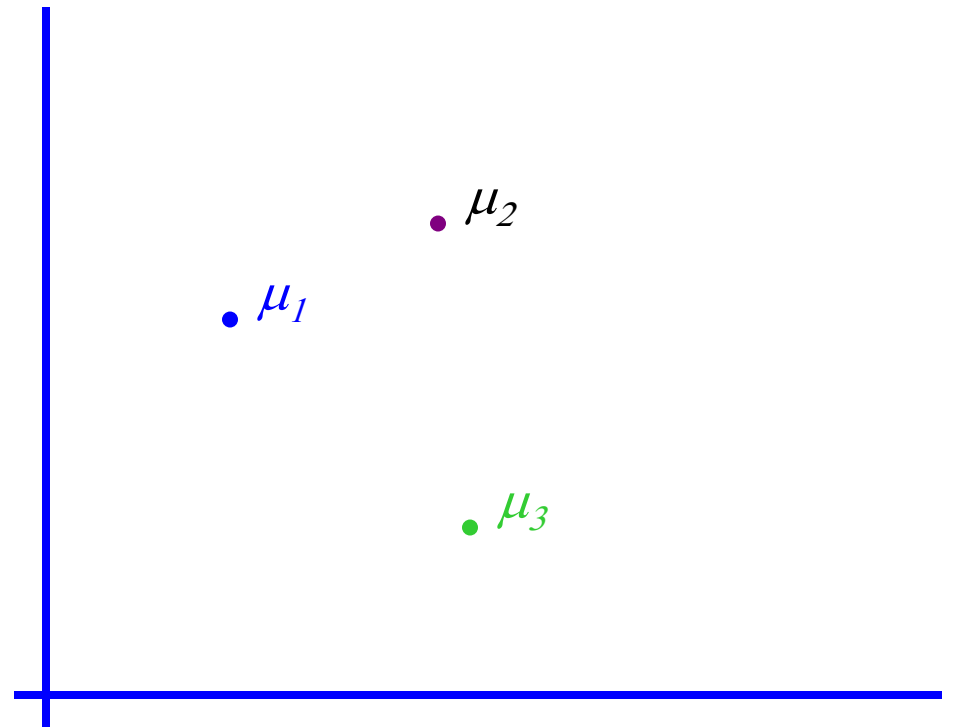
Generative Model:

Assume data comes from a mixture of K Gaussians distributions with same variance

K-means: Generative model

Mixture of K Gaussians distributions: (Multi-modal distribution)

- There are k components
- Component i has an associated mean vector μ_i

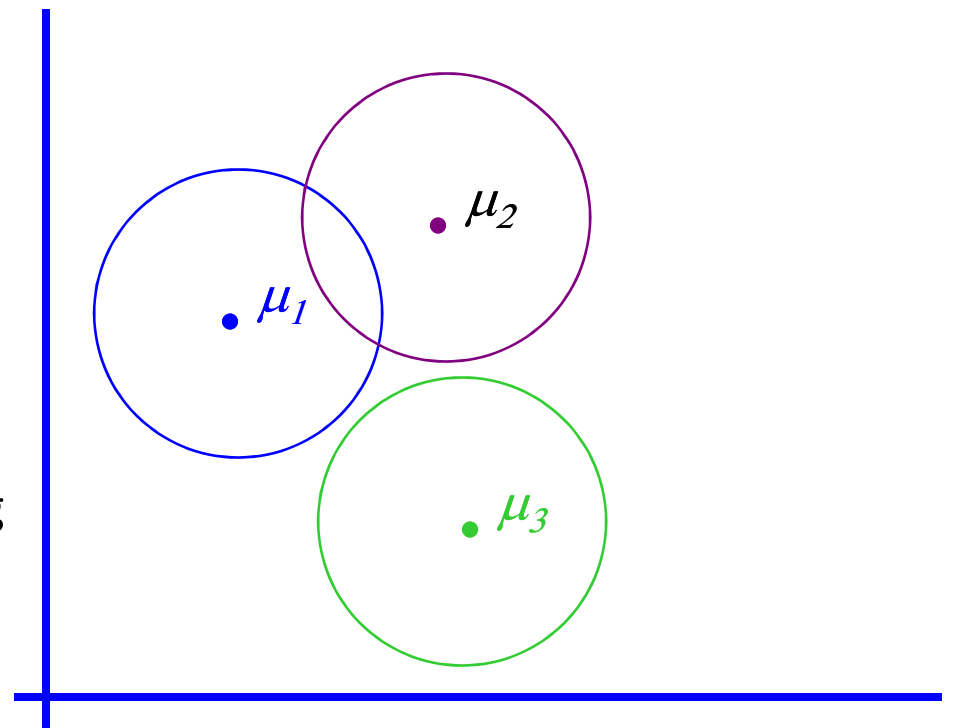


K-means: Generative model

Mixture of K Gaussians distributions: (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 $\sigma^2 I$

Each data point is generated according to the following recipe:



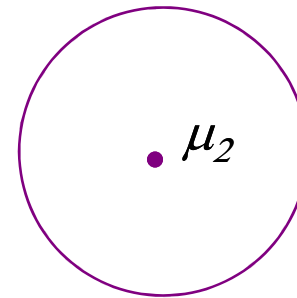
K-means: Generative model

Mixture of K Gaussians distributions: (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 $\sigma^2 I$

Each data point is generated according to the following recipe:

- 1) Pick a component at random:
Choose component i with probability $P(y=i)$



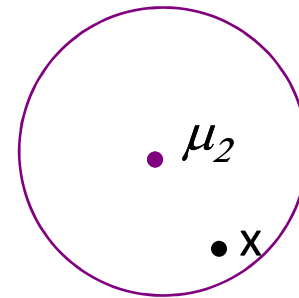
K-means: Generative model

Mixture of K Gaussians distributions: (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 $\sigma^2 I$

Each data point is generated according to the following recipe:

- 1) Pick a component at random:
Choose component i with probability $P(y=i)$
- 2) Datapoint $x \sim N(\mu_i, \sigma^2 I)$



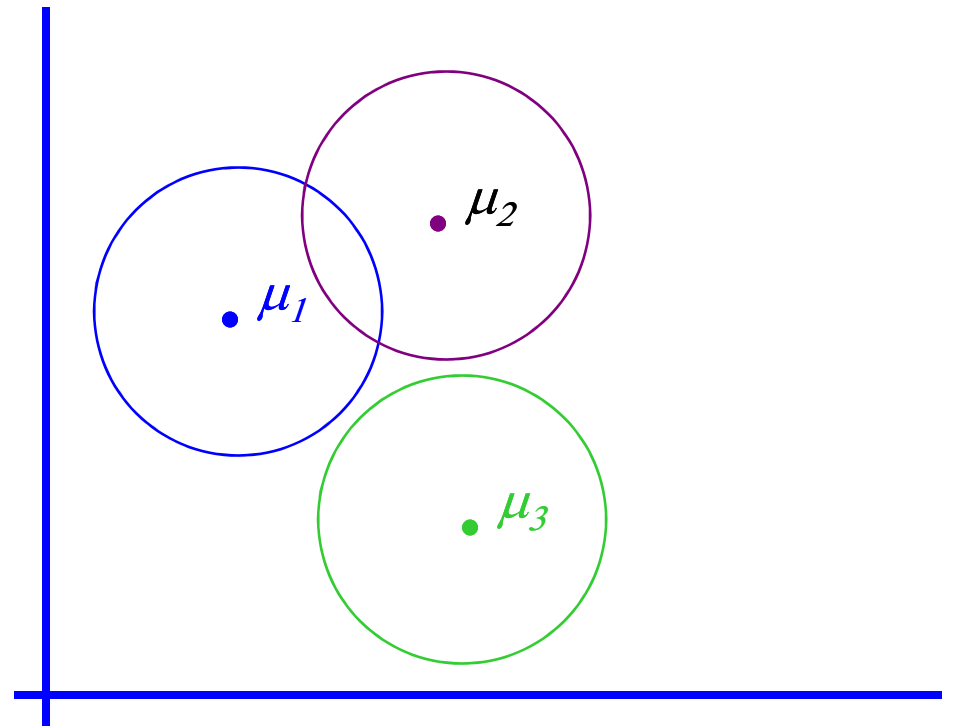
K-means: Generative model

Mixture of K Gaussians distributions: (Multi-modal distribution)

$$p(x|y=i) \sim N(\mu_i, \sigma^2 I)$$

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

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



K-means: Generative model

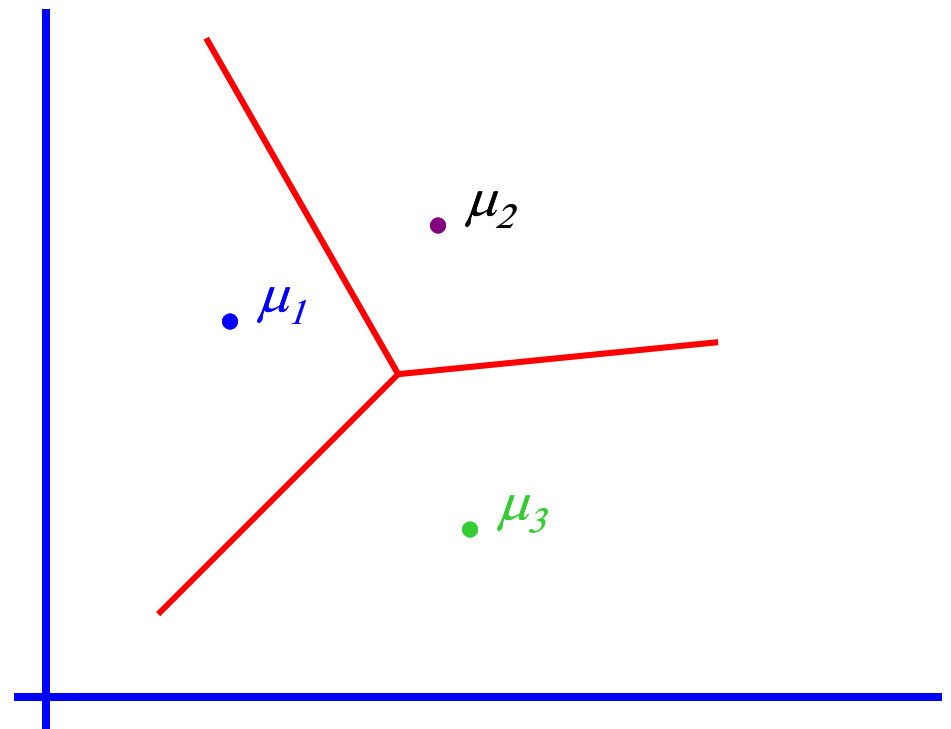
Mixture of K Gaussians distributions: (Multi-modal distribution)

$$p(x|y=i) \sim N(\mu_i, \sigma^2 I)$$

Gaussian Bayes Classifier:

$$\begin{aligned} \log \frac{P(y = i | \mathbf{x})}{P(y = j | \mathbf{x})} \\ &= \log \frac{p(\mathbf{x} | y = i)P(y = i)}{p(\mathbf{x} | y = j)P(y = j)} \\ &= \mathbf{W}^T \mathbf{x} \end{aligned}$$

Depends on $\mu_1, \mu_2, \dots, \mu_K, \sigma^2, P(y=1), \dots, P(y=K)$



“Linear Decision boundary” – Recall that second-order terms cancel out

K-means: Generative model

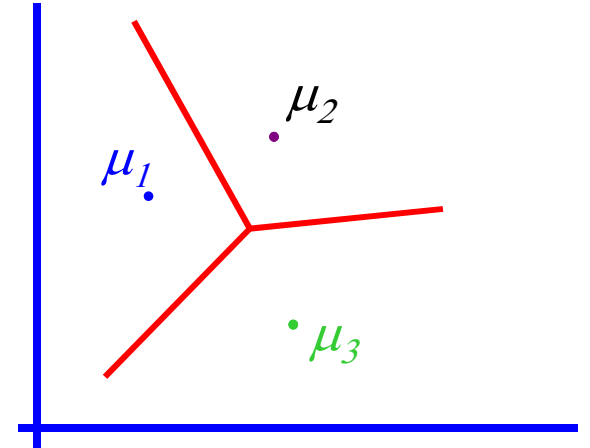
Maximum Likelihood Estimate (MLE)

$$\operatorname{argmax}_{\mu_1, \mu_2, \dots, \mu_K, \sigma^2, P(y=1), \dots, P(Y=k)} \prod_i P(y_i, x_i)$$

But we don't know y_i 's!!!

Maximize marginal likelihood:

$$\begin{aligned} \operatorname{argmax} \prod_j P(x_j) &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i, x_j) \\ &= \operatorname{argmax} \prod_j \sum_i P(y_j=i) p(x_j | y_j=i) \end{aligned}$$



K-means: Generative model

Maximize marginal likelihood:

$$\begin{aligned}\operatorname{argmax} \prod_j P(x_j) &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i, x_j) \\ &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i) p(x_j | y_j=i)\end{aligned}$$

$$P(y_j = i, x_j) \propto P(y_j = i) \exp\left[-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right]$$

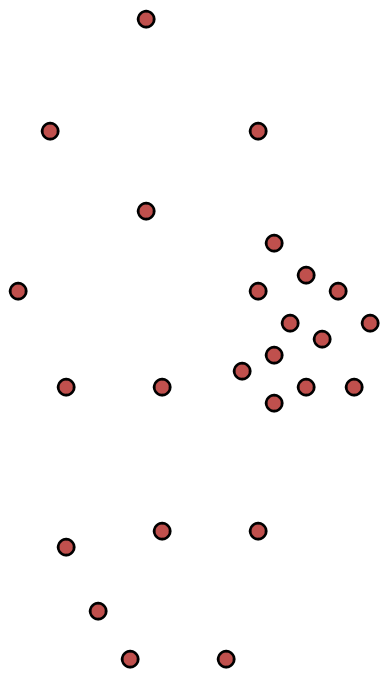
If each x_j belongs to one class $C(j)$ (hard assignment), marginal likelihood:

$$P(y_j=i) = 1 \text{ or } 0 \quad 1 \text{ if } i = C(j)$$

$$\prod_{j=1}^m \sum_{i=1}^k P(y_j = i, x_j) \propto \prod_{j=1}^m \exp\left[-\frac{1}{2\sigma^2} \|x_j - \mu_{C(j)}\|^2\right] = \sum_{j=1}^m -\frac{1}{2\sigma^2} \|x_j - \mu_{C(j)}\|^2$$

Same as K-means!!!

(One) bad case for K-means



- Clusters may not be linearly separable
- Clusters may overlap
- Some clusters may be “wider” than others

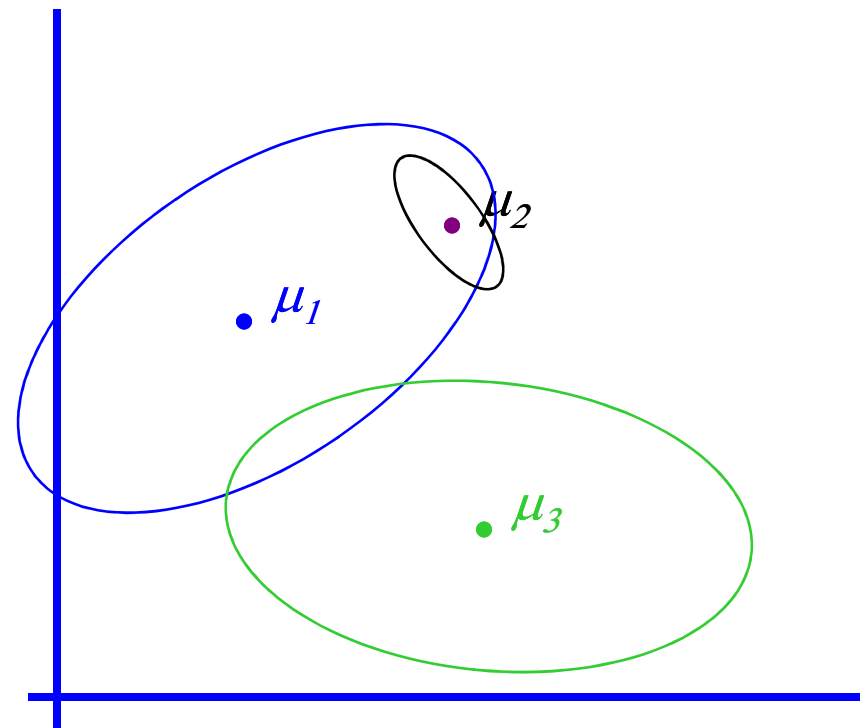
General GMM

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 $P(y=i)$
- 2) Datapoint $x \sim N(\mu_i, \Sigma_i)$



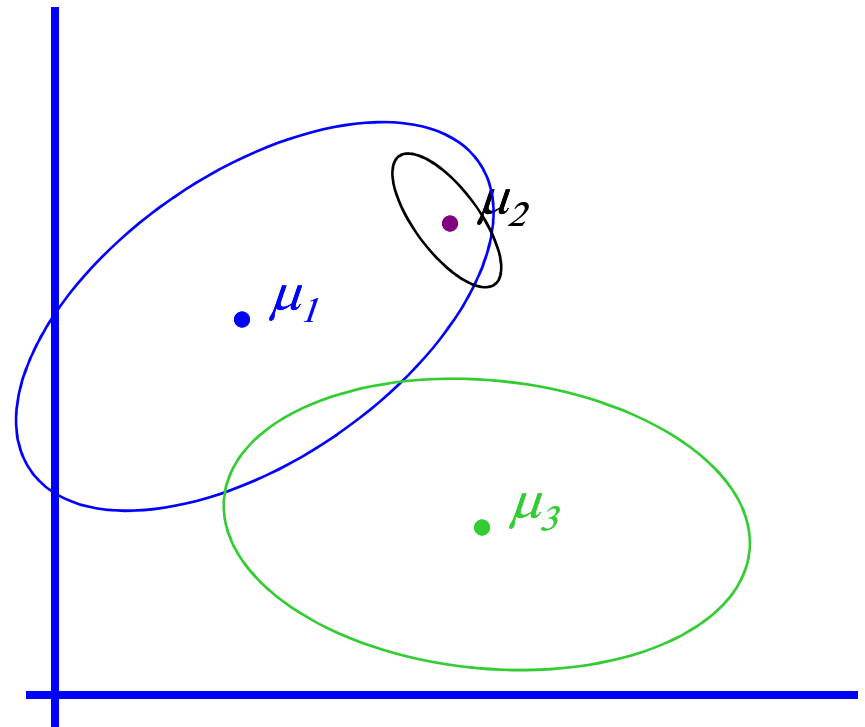
General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

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

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

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



General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

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

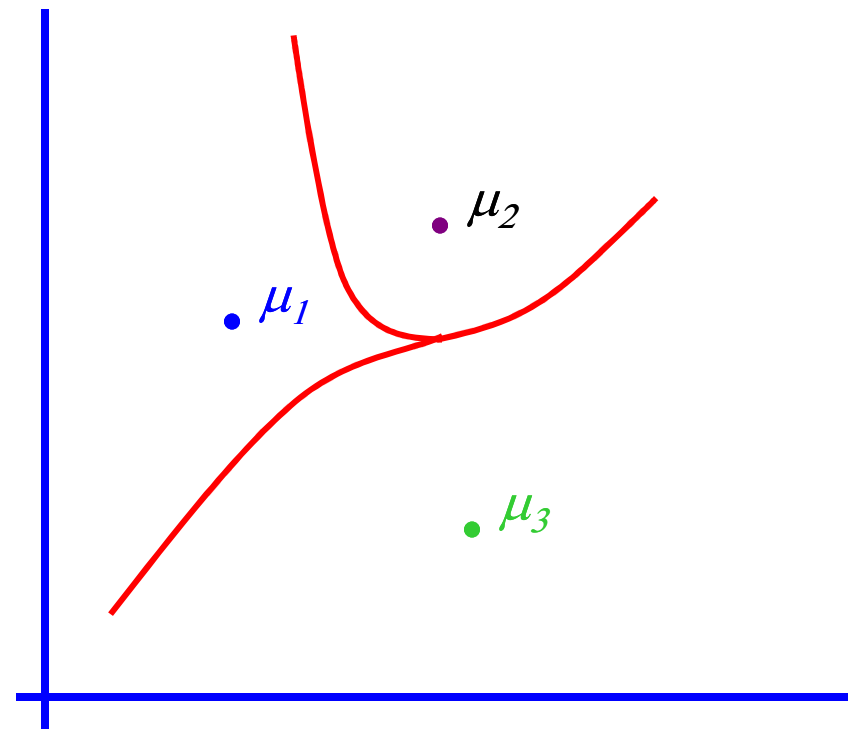
Gaussian Bayes Classifier:

$$\log \frac{P(y = i | x)}{P(y = j | x)}$$

$$= \log \frac{p(x | y = i)P(y = i)}{p(x | y = j)P(y = j)}$$

$$= x^T \mathbf{W} x + \mathbf{w}^T x$$

Depend on $\mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K, P(y=1), \dots, P(Y=k)$



“Quadratic Decision boundary” – second-order terms don’t cancel out

General GMM

Maximize marginal likelihood:

$$\begin{aligned}\operatorname{argmax} \prod_j P(x_j) &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i, x_j) \\ &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i) p(x_j | y_j=i)\end{aligned}$$

Uncertain about class of each x_j (soft assignment), $P(y_j=i) = P(y=i)$

$$\prod_{j=1}^m \sum_{i=1}^k P(y_j = i, x_j) \propto \prod_{j=1}^m \sum_{i=1}^k P(y = i) \frac{1}{\sqrt{\det(\Sigma_i)}} \exp\left[-\frac{1}{2} (x_j - \mu_i)^T \Sigma_i (x_j - \mu_i)\right]$$

How do we find the μ_i 's which give max. marginal likelihood?

* Set $\frac{\partial}{\partial \mu_i} \log \text{Prob} (\dots) = 0$ and solve for μ_i 's. Non-linear non-analytically solvable

* Use gradient descent: Often slow but doable

Expectation-Maximization (EM)

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

- EM is an optimization strategy for objective functions that can be interpreted as likelihoods in the presence of missing data.
- It is much simpler than gradient methods:
 - No need to choose step size.
 - Enforces constraints automatically.
 - Calls inference and fully observed learning as subroutines.
- EM is an Iterative algorithm with two linked steps:
 - E-step: fill-in hidden values using inference
 - M-step: apply standard MLE/MAP method to completed data
- We will prove that this procedure monotonically improves the likelihood (or leaves it unchanged). Thus it always converges to a local optimum of the likelihood.

Expectation-Maximization (EM)

A simple case:

We have unlabeled data $\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_m$

We know there are k classes

We know $P(y=1), P(y=2) P(y=3) \dots P(y=K)$

We don't know $\mu_1 \mu_2 \dots \mu_k$

We know common variance σ^2

We can write $P(\text{ data } | \mu_1 \dots \mu_k)$

$$= p(x_1 \dots x_m | \mu_1 \dots \mu_k)$$

$$= \prod_{j=1}^m p(x_j | \mu_1 \dots \mu_k)$$

Independent data

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

Marginalize over class

$$\propto \prod_{j=1}^m \sum_{i=1}^k \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y = i)$$

Expectation (E) step

If we know μ_1, \dots, μ_k \rightarrow easily compute prob. point x_j belongs to class $y=i$

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

Simply evaluate gaussian and normalize

Maximization (M) step

If we know prob. point x_j belongs to class $y=i$

→ MLE for μ_i is weighted average

imagine multiple copies of each x_j , each with weight $P(y=i|x_j)$:

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

EM for spherical, same variance GMMs

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

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

EM for axis-aligned GMMs

$$\Sigma_i = \begin{pmatrix} \sigma_{i,1}^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_{i,2}^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_{i,3}^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{i,m-1}^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_{i,m}^2 \end{pmatrix}$$

Iterate. On iteration t let our estimates be

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

E-step

Compute “expected” classes of all datapoints for each class

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

Just evaluate a Gaussian at x_j

M-step

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

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

$m = \#$ data points

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

E-step

Compute "expected" classes of all datapoints for each class

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

Just evaluate a Gaussian at x_j

M-step

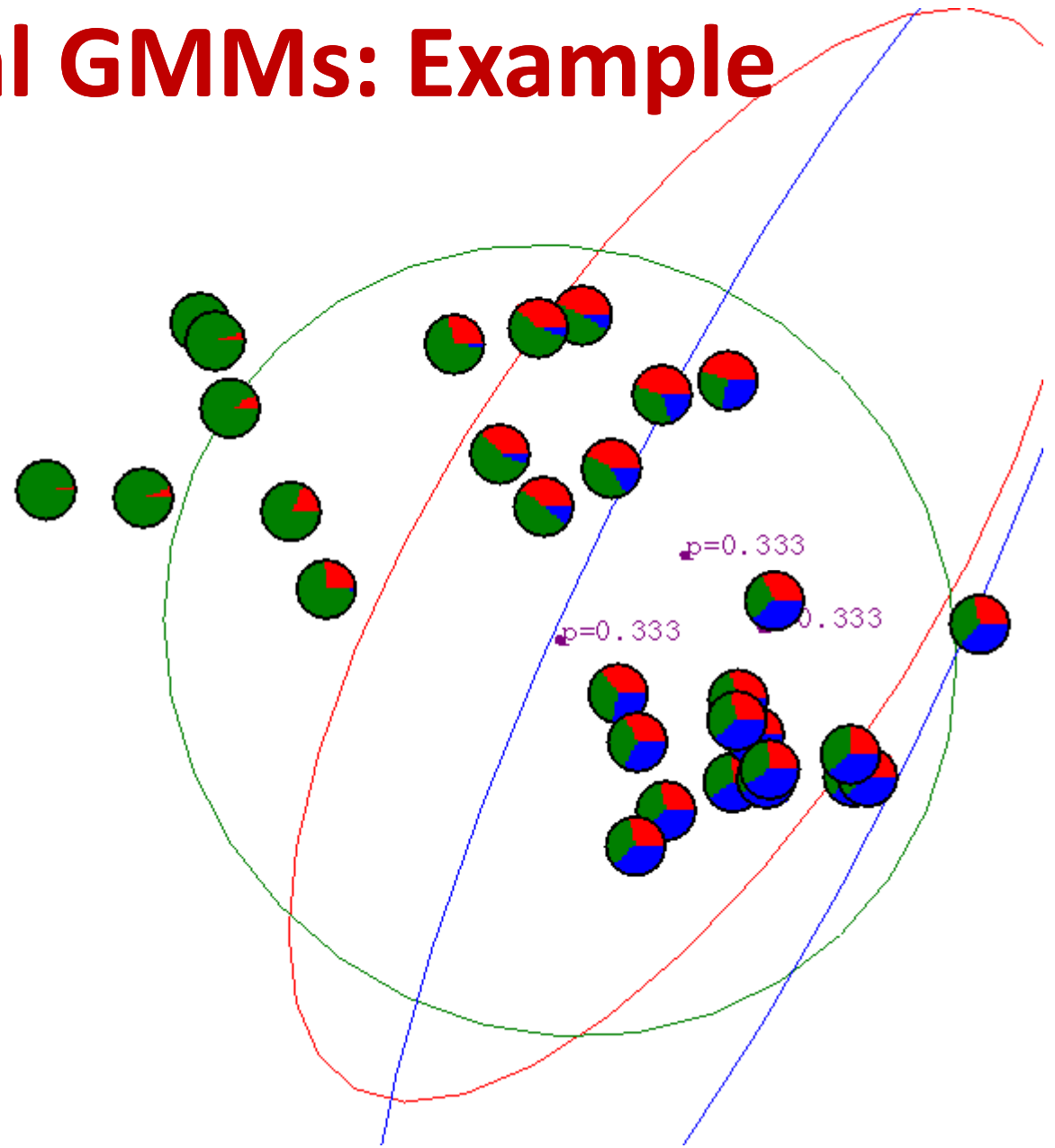
Compute MLEs given our data's class membership distributions (weights)

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

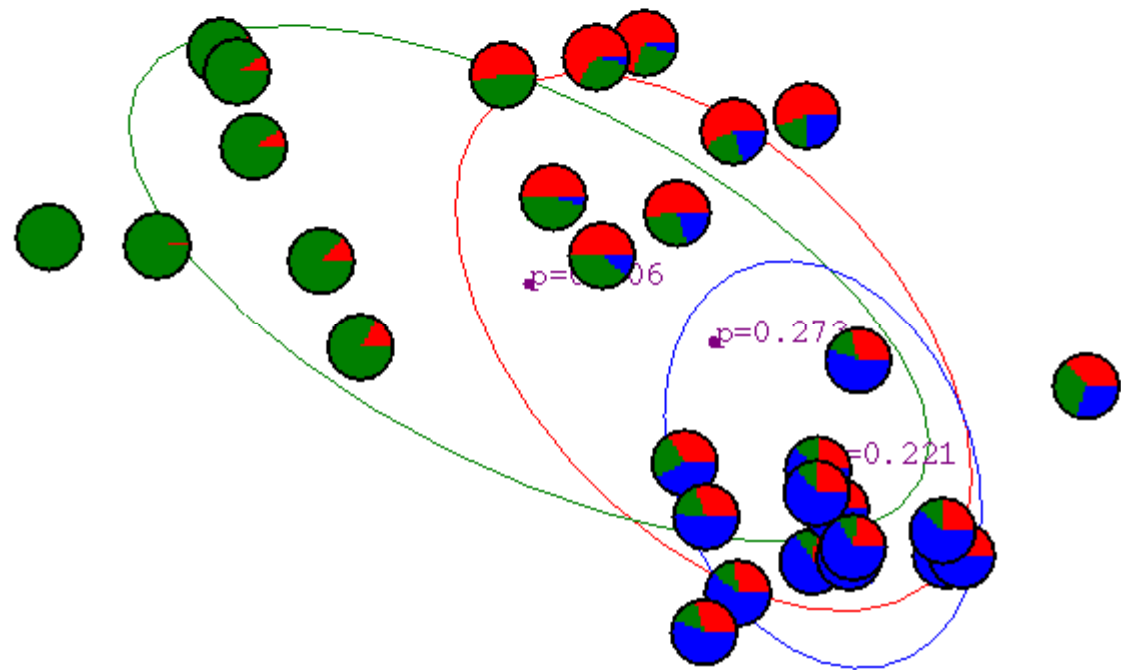
$$p_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t)}{m}$$

$m = \#$ data points

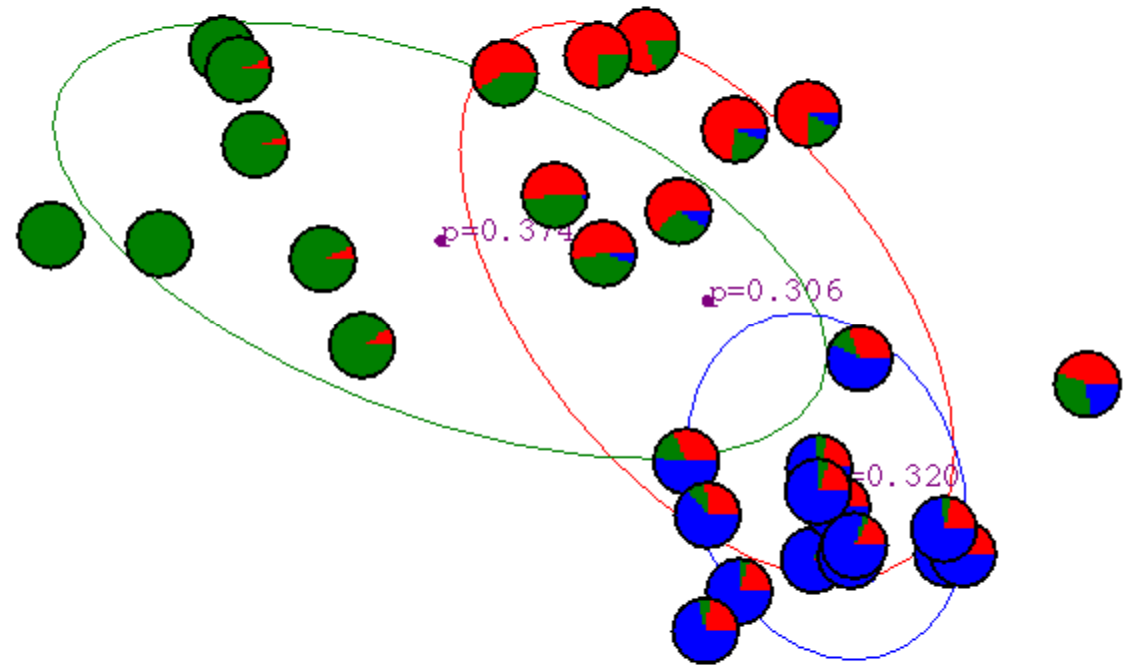
EM for general GMMs: Example



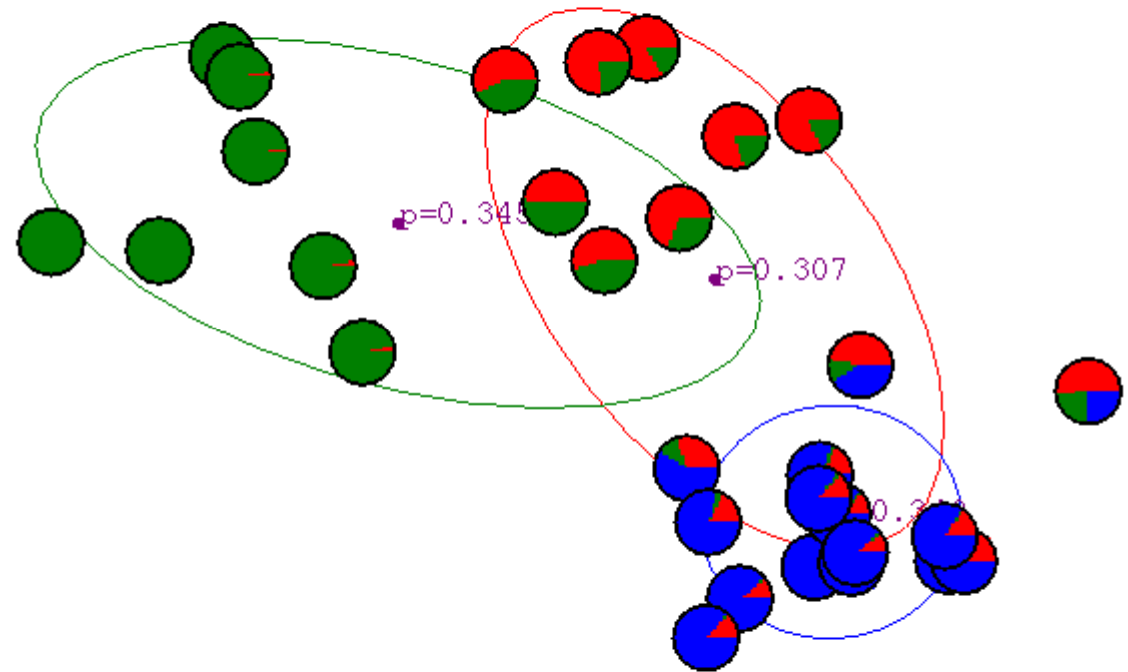
After 1st iteration



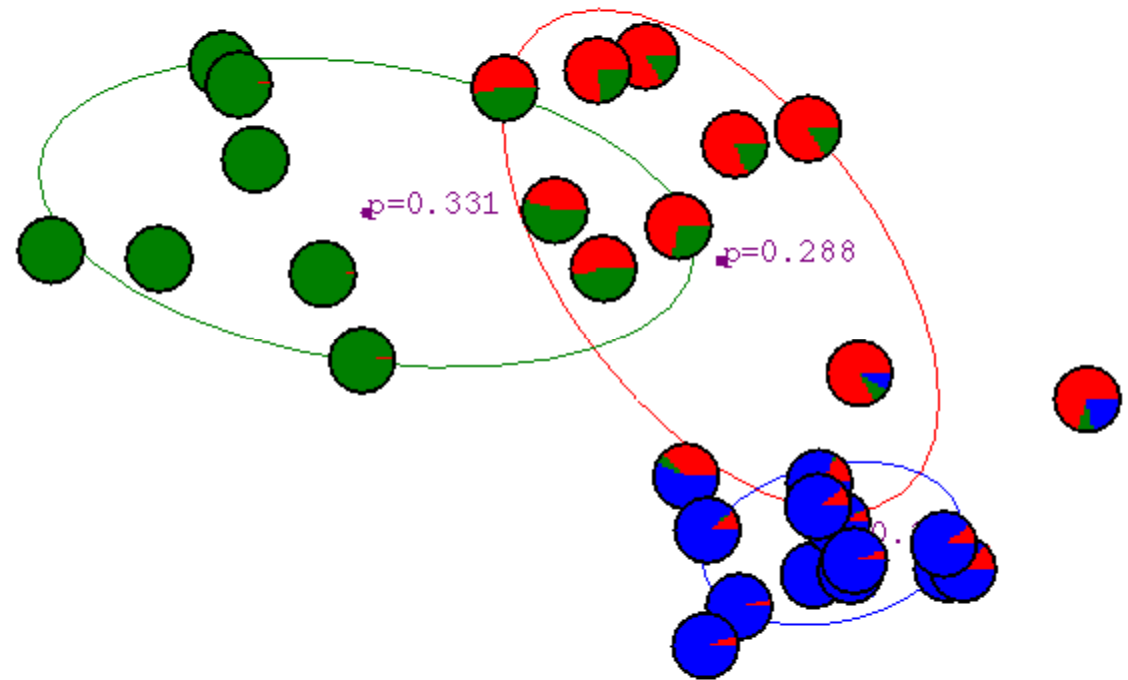
After 2nd iteration



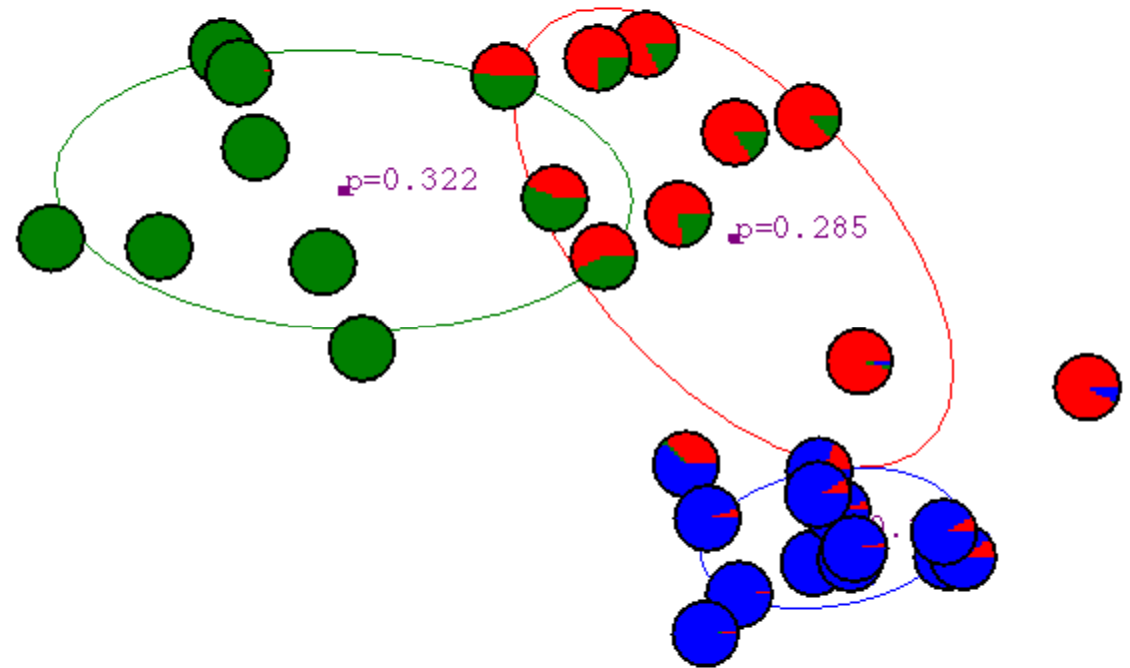
After 3rd iteration



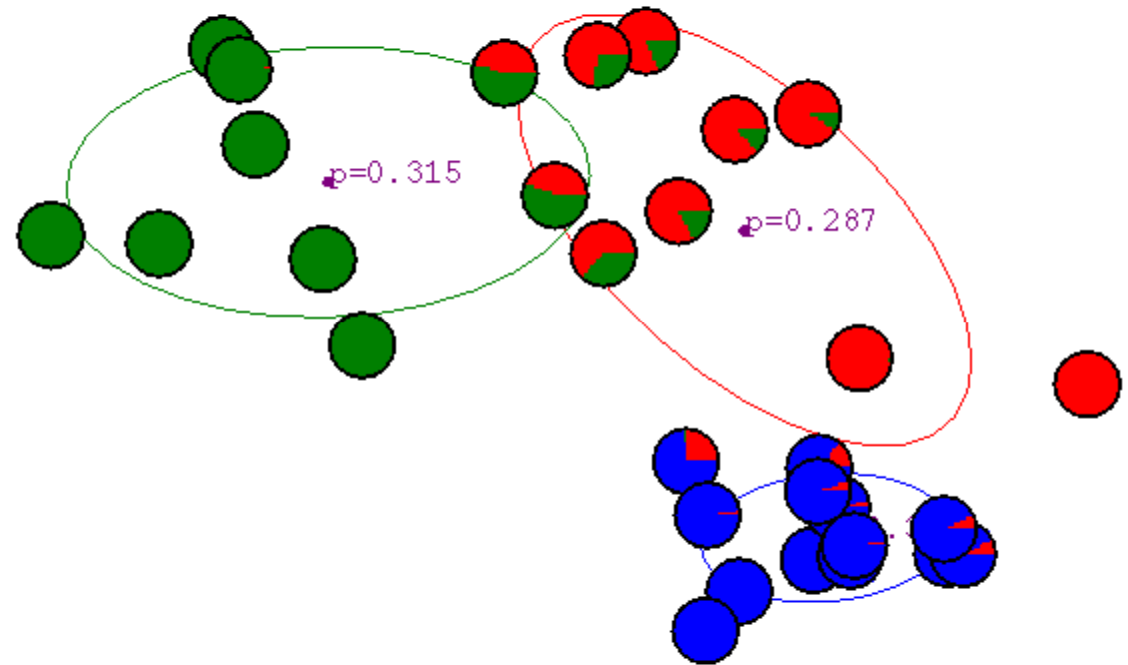
After 4th iteration



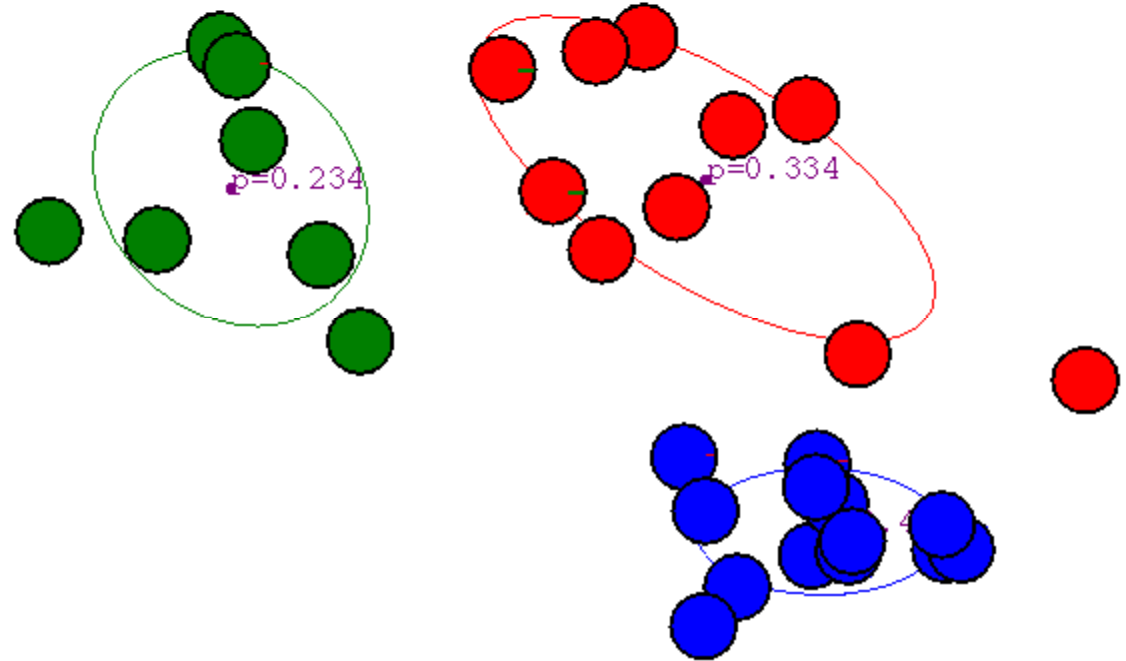
After 5th iteration



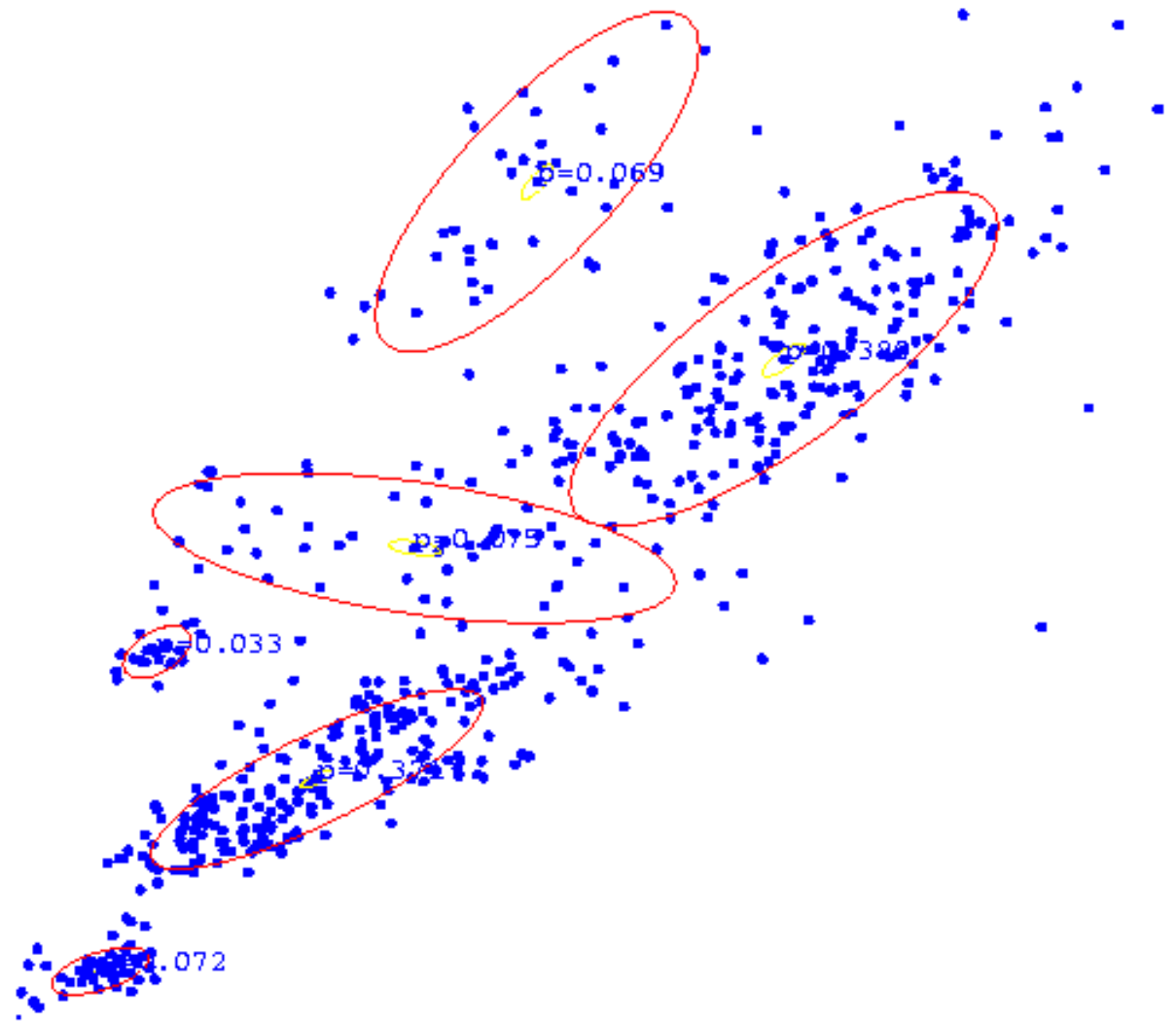
After 6th iteration



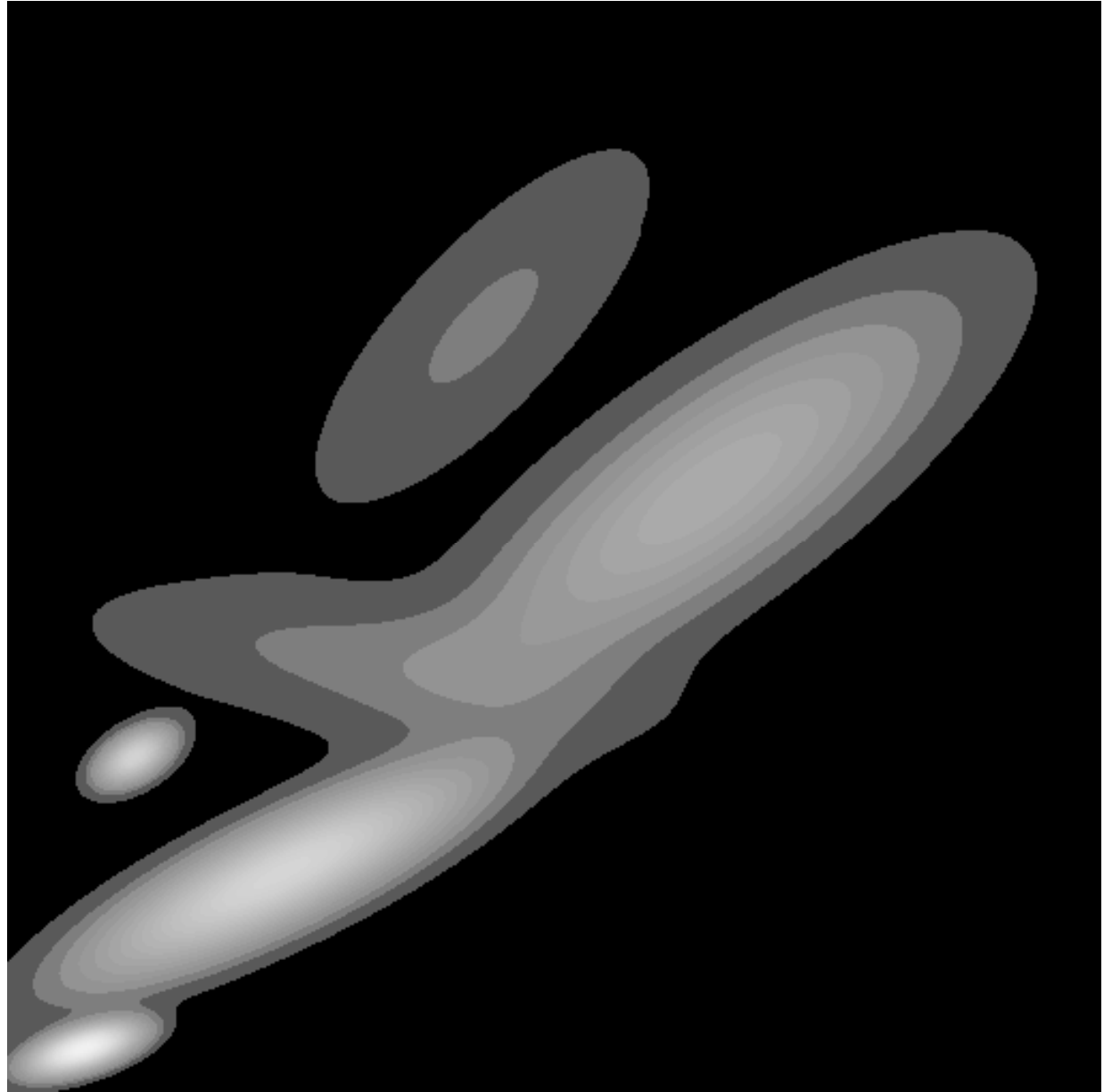
After 20th iteration



GMM clustering of the assay data

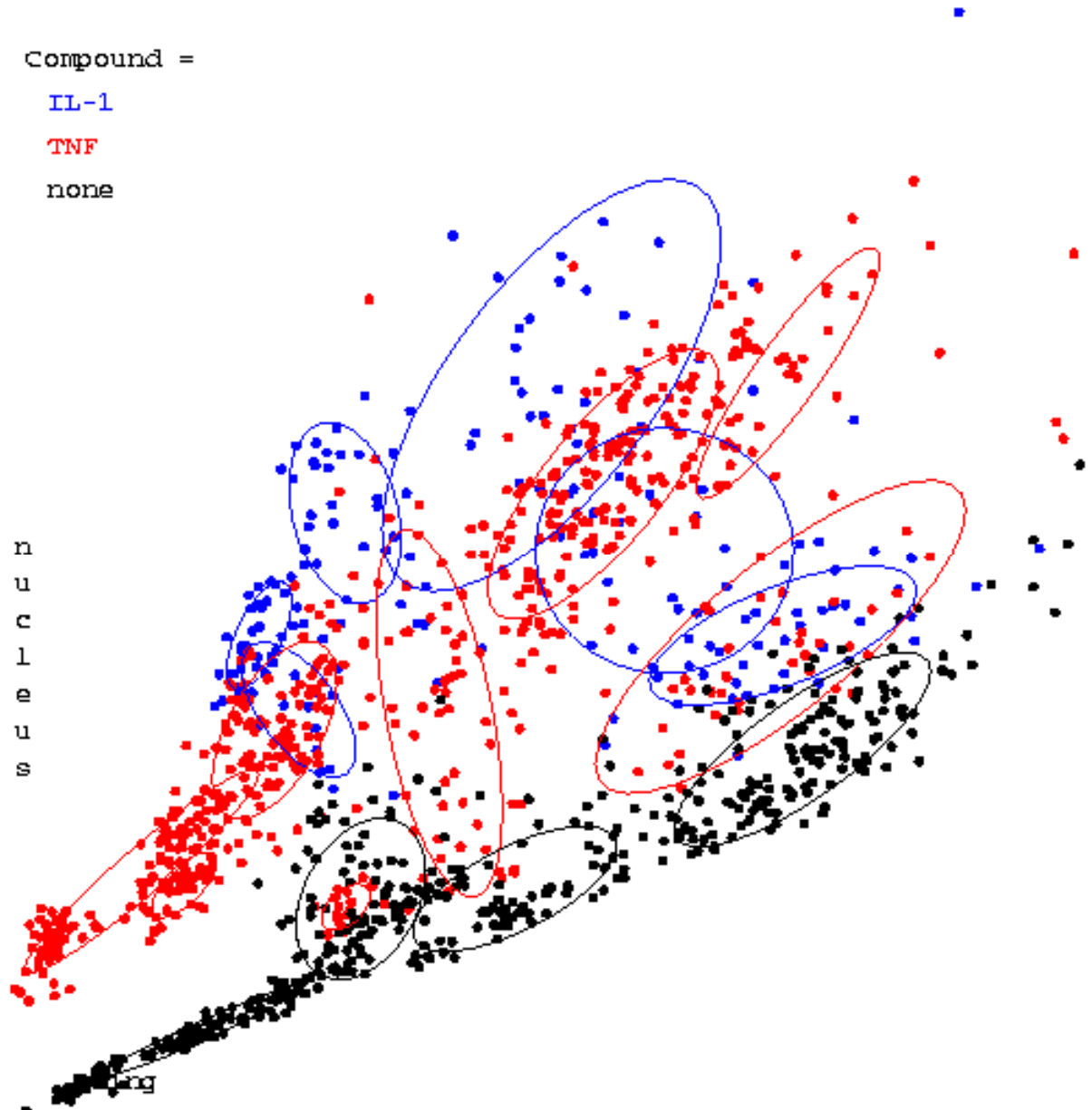


Resulting Density Estimator

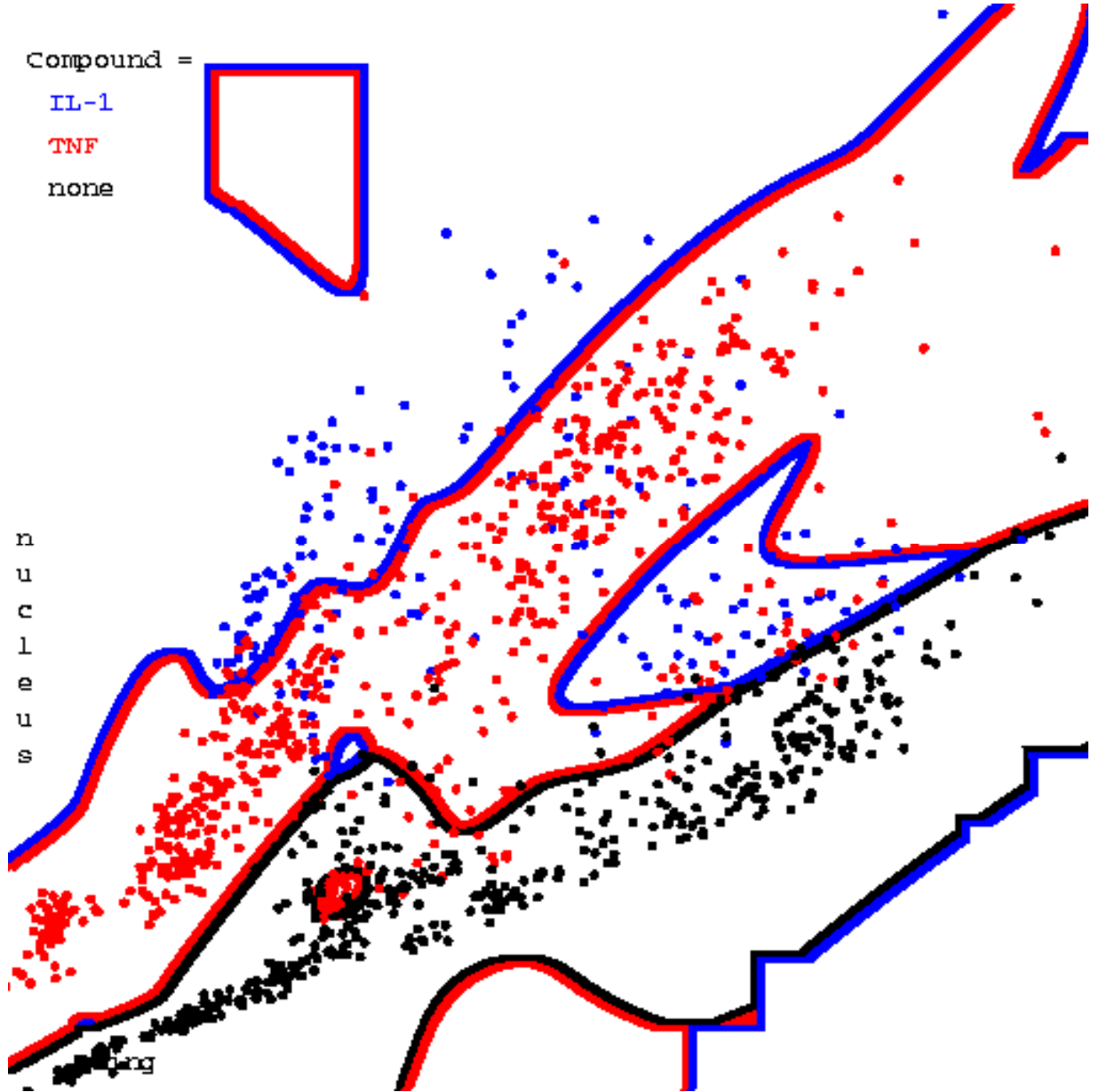


Three classes of assay

(each learned with its own mixture model)



Resulting Bayes Classifier



General EM algorithm

Marginal likelihood – \mathbf{x} is observed, \mathbf{z} is missing:

$$\begin{aligned} P(\mathbf{D}; \theta) &= \log \prod_{j=1}^m P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta) \end{aligned}$$

E step

\mathbf{x} is observed, \mathbf{z} is missing

Compute probability of missing data given current choice of θ

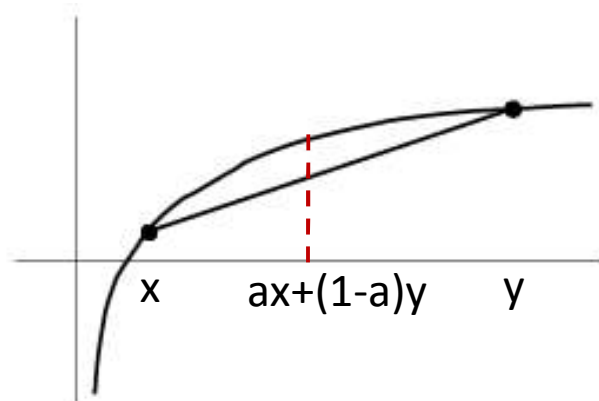
$$Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) = P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$$

$$\text{E.g., } P(y = i | \mathbf{x}_j, \lambda_t)$$

Lower-bound on marginal likelihood

$$\begin{aligned} P(\mathbf{D}; \theta) &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} \underbrace{Q(\mathbf{z} | \mathbf{x}_j)}_{P(\mathbf{z})} \underbrace{\frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}}_{f(\mathbf{z})} \end{aligned}$$

Jensen's inequality: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$



log: concave function

$$\log(ax+(1-a)y) \geq a \log(x) + (1-a) \log(y)$$

Lower-bound on marginal likelihood

$$\begin{aligned} P(\mathbf{D}; \theta) &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} \underbrace{Q(\mathbf{z} | \mathbf{x}_j)}_{P(\mathbf{z})} \underbrace{\frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}}_{f(\mathbf{z})} \end{aligned}$$

Jensen's inequality: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\begin{aligned} &\geq \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m.H(Q) \end{aligned}$$

M step

$$P(\mathbf{D}; \theta) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m.H(Q)$$

Maximize lower bound on marginal likelihood

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \underbrace{\sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)}_{\text{Expected log likelihood}}$$

Use expected counts instead of counts:

If learning requires $\text{Count}(\mathbf{x}, \mathbf{z})$, Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{z})]$

Convergence of EM

$$P(\mathbf{D}; \theta) \geq F(\theta, Q)$$

M-step: Fix Q , maximize F over θ

$$P(\mathbf{D}; \theta) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m \cdot H(Q^{(t)})$$

Maximizes lower bound F on marginal likelihood

E-step: Fix θ , maximize F over Q

Convergence of EM

$$P(\mathbf{D}; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$\begin{aligned} P(\mathbf{D}; \theta^{(t)}) &\geq F(\theta^{(t)}, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}) P(\mathbf{x}_j | \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \underbrace{\sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)}}_{-KL(Q(\mathbf{z}|\mathbf{x}_j), P(\mathbf{z}|\mathbf{x}_j, \theta^{(t)}))} + \underbrace{\sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{x}_j | \theta^{(t)})}_{P(\mathbf{D}; \theta^{(t)})} \end{aligned}$$

KL divergence between two distributions

Convergence of EM

$$P(\mathbf{D}; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$\begin{aligned} P(\mathbf{D}; \theta^{(t)}) &\geq F(\theta^{(t)}, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \sum_{j=1}^m -KL(Q(\mathbf{z} | \mathbf{x}_j), P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})) + P(\mathbf{D}; \theta^{(t)}) \end{aligned}$$

KL ≥ 0 , Maximized if KL divergence = 0 **KL(Q,P) = 0** iff **Q = P**

Recall E-step: $Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) = P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$

Convergence of EM

$$P(\mathbf{D}; \theta) \geq F(\theta, Q)$$

M-step: Fix Q , maximize F over θ

$$P(\mathbf{D}; \theta) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m \cdot H(Q^{(t)})$$

Maximizes lower bound F on marginal likelihood

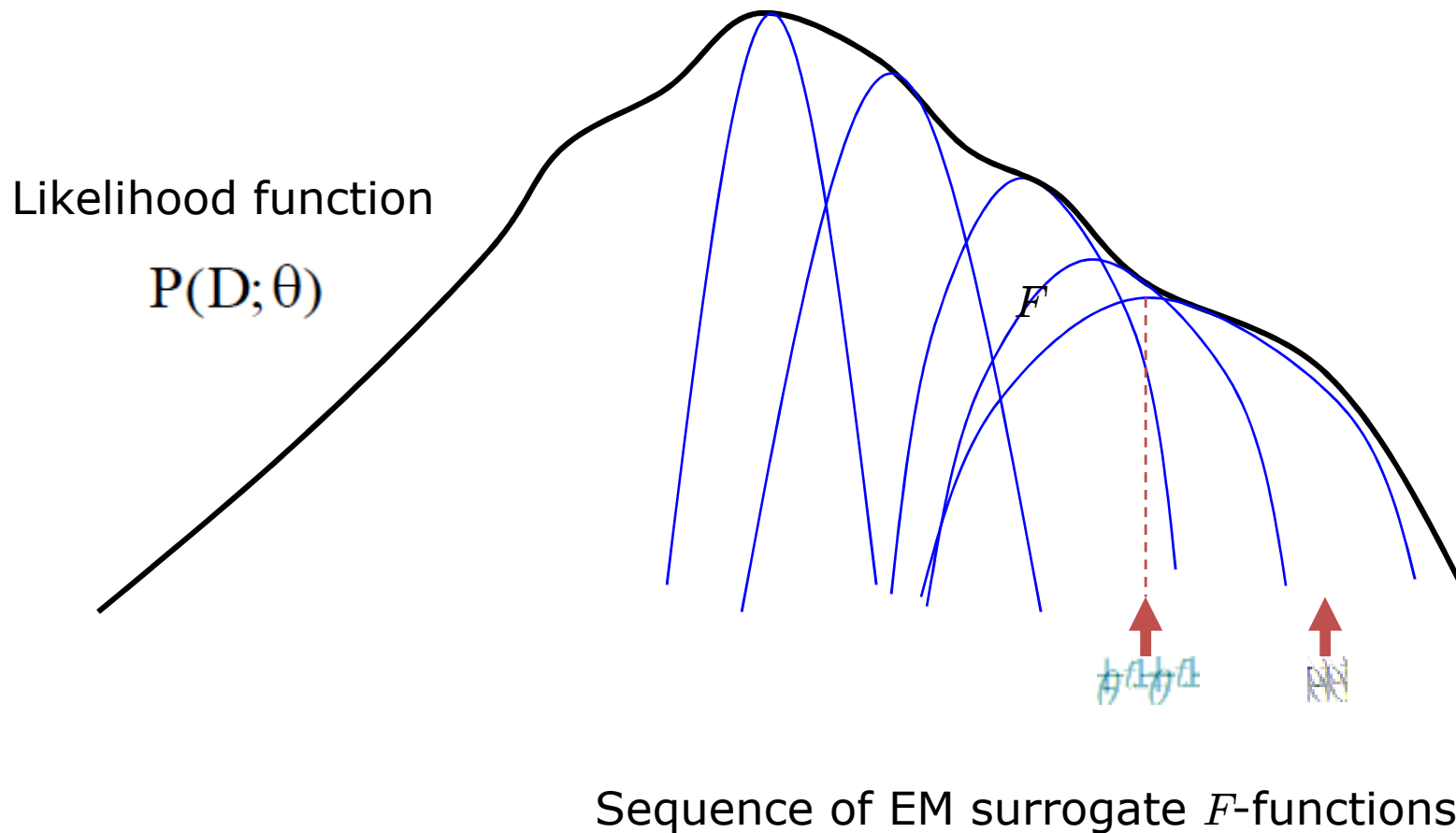
E-step: Fix θ , maximize F over Q

$$P(\mathbf{D}; \theta^{(t)}) \geq F(\theta^{(t)}, Q) = P(\mathbf{D}; \theta^{(t)}) - \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))$$

Re-aligns F with marginal likelihood

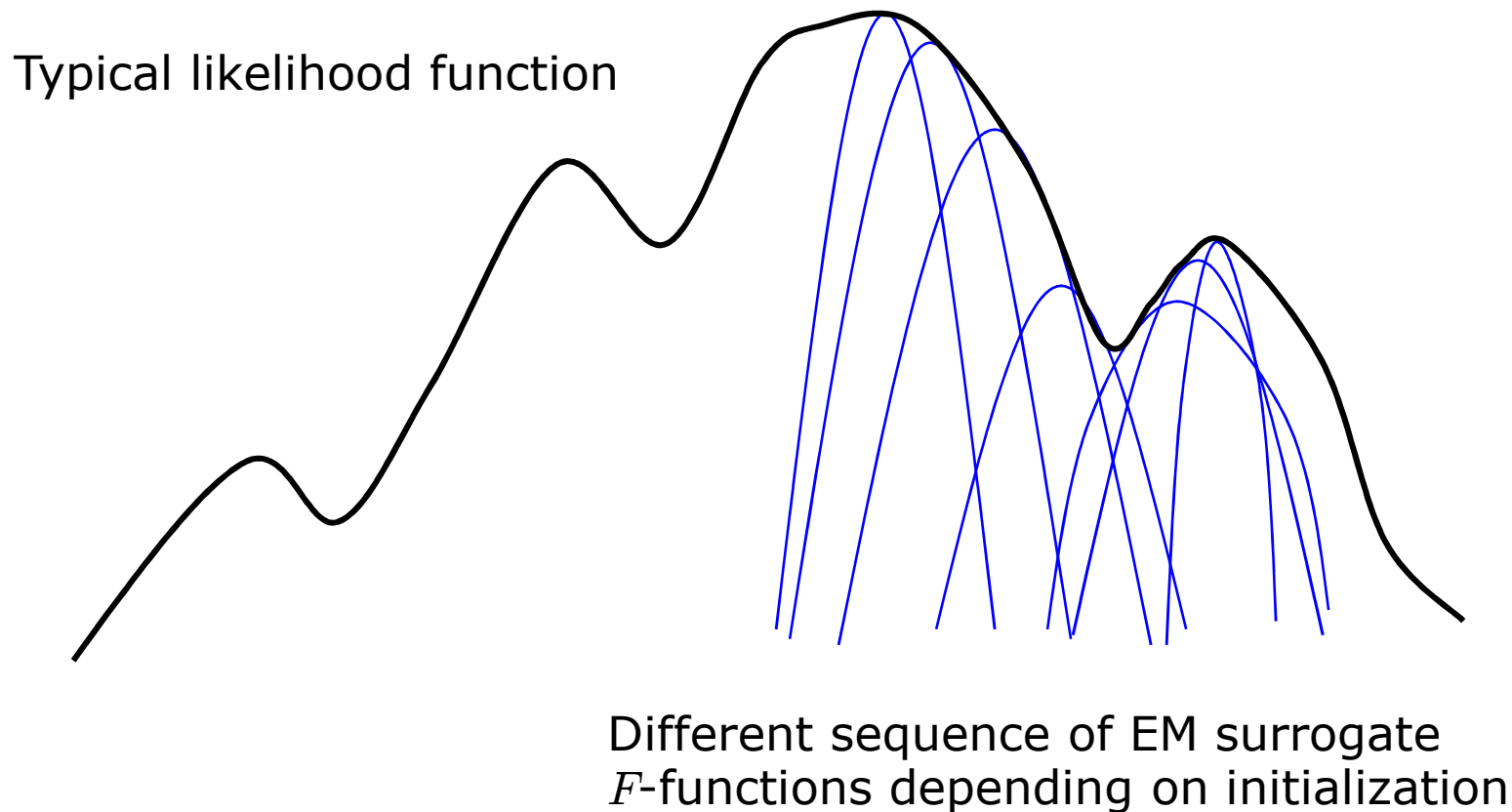
$$F(\theta^{(t)}, Q^{(t+1)}) = P(\mathbf{D}; \theta^{(t)})$$

Monotonic convergence of EM



EM monotonically converges to a local maximum of likelihood !

Monotonic convergence of EM



Use multiple, randomized initializations in practice

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: $Q^{t+1} = \arg \max_Q F(\theta^t, Q)$
 2. M-step: $\theta^{t+1} = \arg \max_{\theta} F(\theta, Q^{t+1})$
- In the M-step we optimize a lower bound on the likelihood. In the E-step we close the gap, making bound=likelihood.
- EM performs coordinate ascent on F, can get stuck in local minima.
- BUT Extremely popular in practice.