# EM Algorithm

Aarti Singh

Slides courtesy:  Eric Xing, Carlos Guetrin

Machine Learning 10-701/15-781
Oct 27, 2010

# K-means Recap ...

- Randomly initialize $k$ centers
  - ☐ $\mu^{(0)} = \mu_1^{(0)}, \ldots, \mu_k^{(0)}$

- **Classify**: Assign each point $j \in \{1, \ldots m\}$ to nearest center:
  - ☐ $C^{(t)}(j) \leftarrow \arg \min_{i=1,\ldots,k} \|\mu_i^{(t)} - x_j\|^2$

- **Recenter**: $\mu_i$ becomes centroid of its points:
  - ☐ $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j:C^{(t)}(j)=i} \|\mu - x_j\|^2 \qquad i \in \{1, \ldots, k\}$
  - ☐ Equivalent to $\mu_i \leftarrow$ average of its points!

# What is K-means optimizing?

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

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

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

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

# K-means algorithm

- Optimize potential function:

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

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

- **K-means algorithm:**

**(1)** Fix $\mu$, optimize C

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

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

# K-means algorithm

- Optimize potential function:

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

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

- **K-means algorithm:**

**(2)** Fix C, optimize $\mu$

$$\min_{\mu_1, \mu_2, \ldots \mu_K} \sum_{i=1}^{K} \sum_{j:C(j)=i} ||\mu_i - x_j||^2 = \sum_{i=1}^{K} \min_{\mu_i} \underbrace{\sum_{j:C(j)=i} ||\mu_i - x_j||^2}$$

**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 descent on F)

**(1)** Fix $\mu$, optimize C          **Expectation step**

**(2)** Fix C, optimize $\mu$          **Maximization step**

Today, we will see a generalization of this approach:
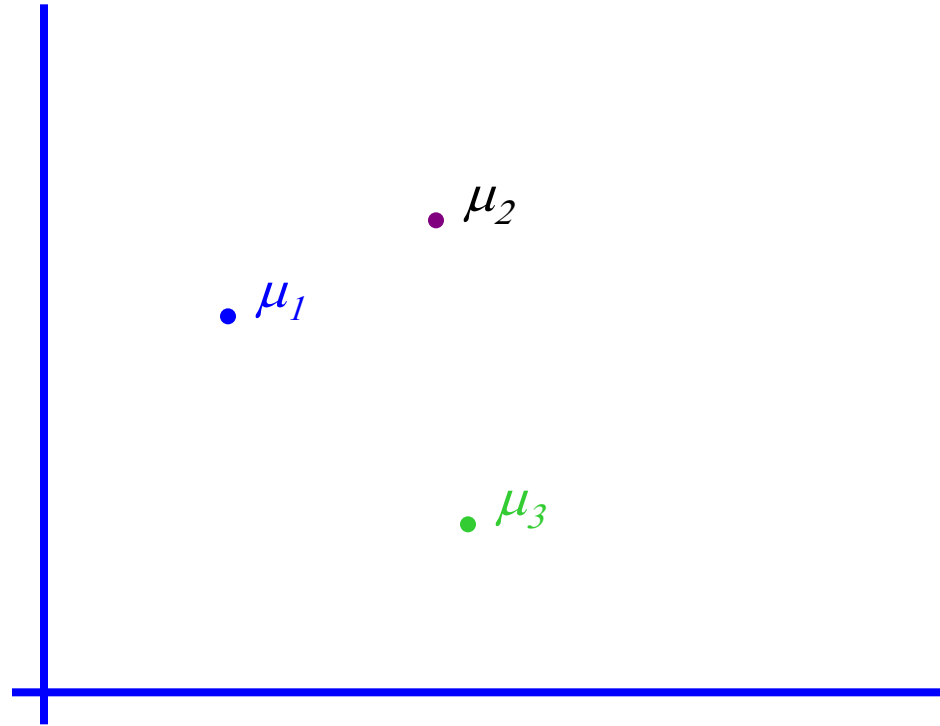
**EM algorithm**

# 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
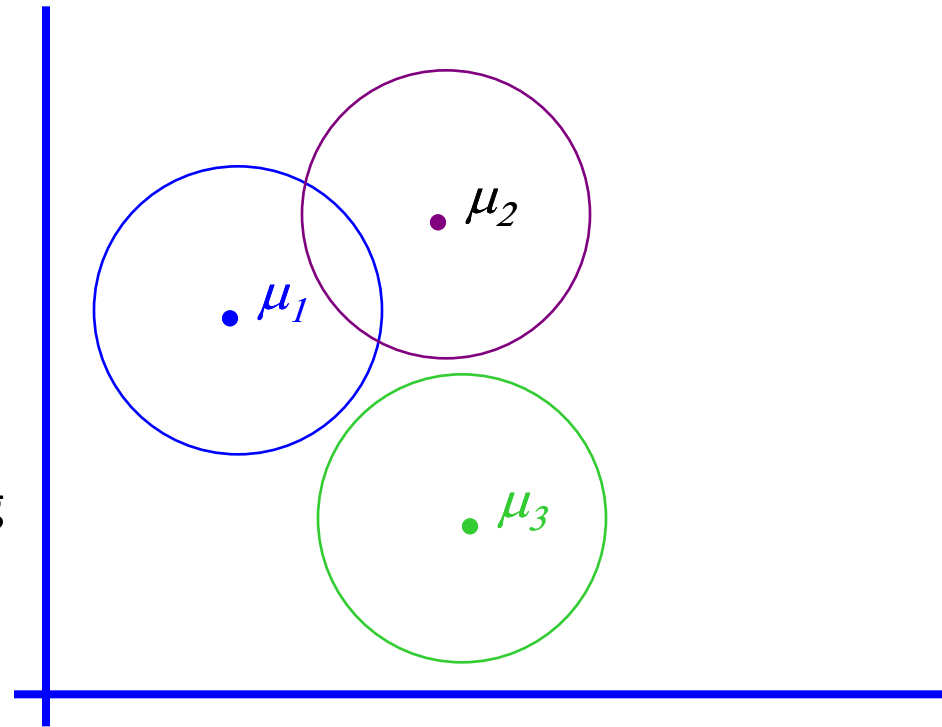
Generative approach

# Gaussian Mixture Model

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

- There are k components

- Component $i$ has an associated mean vector $\mu_i$

$\mu_2$

$\mu_1$

$\mu_3$

# Gaussian Mixture Model

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

- There are k components

- Component *i* has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$

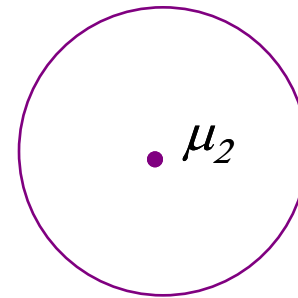Each data point is generated according to the following recipe:

# Gaussian Mixture Model

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

- There are k components

- Component $i$ has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$

Each data point is generated according to the following recipe:

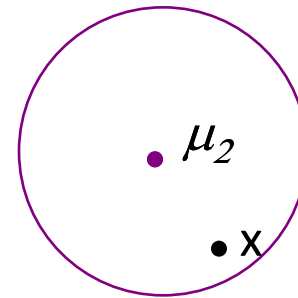1) Pick a component at random: Choose component i with probability $P(y=i)$

$\mu_2$

# Gaussian Mixture Model

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

- There are k components

- Component $i$ has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ 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 $P(y=i)$
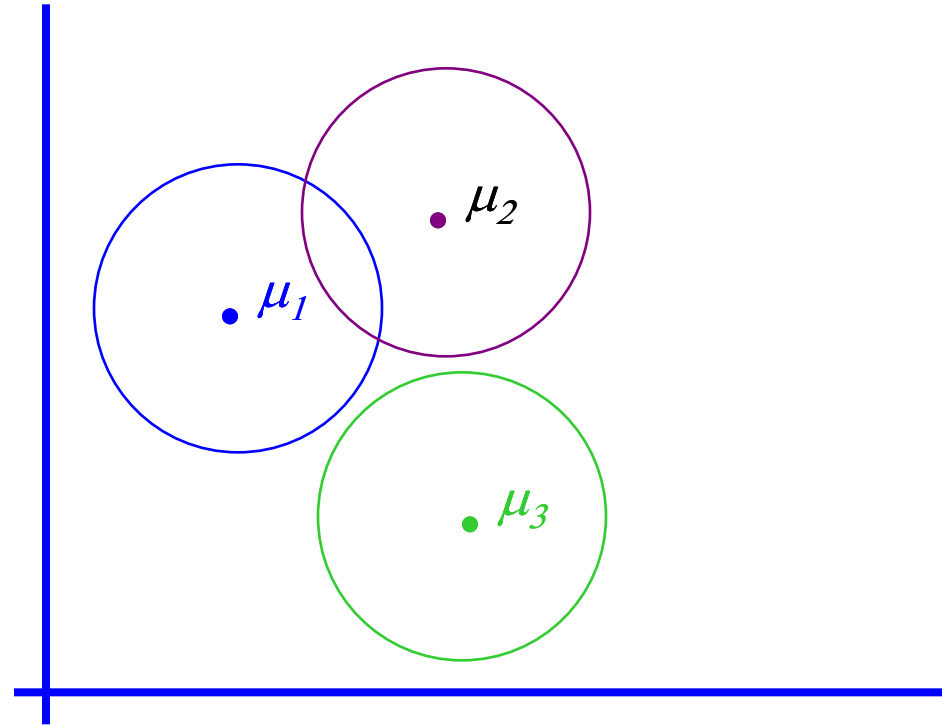
2) Datapoint x ~ $N(\mu_i, \sigma^2 I)$

# Gaussian Mixture Model

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

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

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

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

# Recall: Gaussian Bayes Classifier

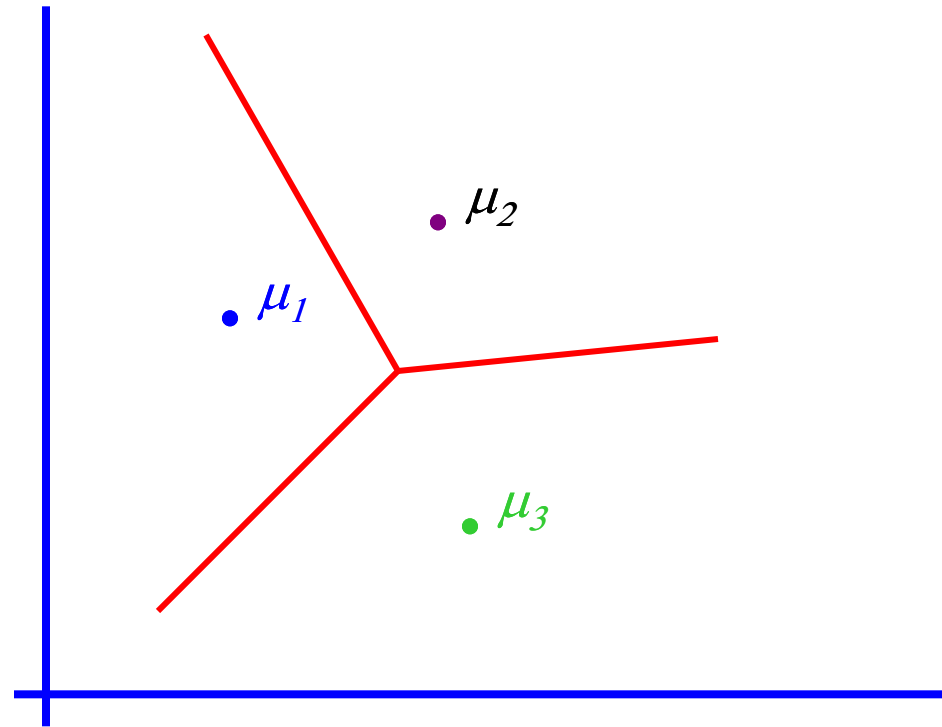Mixture of K Gaussians distributions:  (Multi-modal distribution)

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

Gaussian Bayes Classifier:

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

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

$$= \mathrm{w}^{\mathrm{T}}\mathrm{x}$$

Depends on $\mu_1,\ \mu_2,\ ..\ ,\ \mu_K,\ \sigma^2$ , P(y=1),…, P(Y=k)

**"Linear Decision boundary"** – Recall that second-order terms cancel out
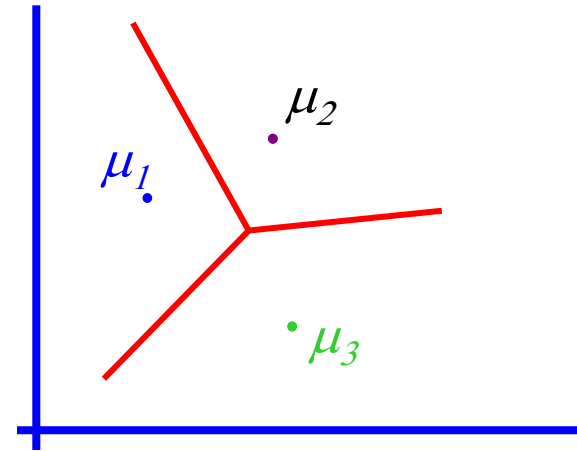
# MLE for GMM

Maximum Likelihood Estimate (MLE)

argmax $\prod_{j=1}^{m}$ P(y$_j$,x$_j$)

$\mu_1, \mu_2, .., \mu_K, \sigma^2,$
P(y=1),..., P(Y=k)

But we don't know y$_j$'s!!!

Maximize marginal likelihood:

argmax $\prod_j$ P(x$_j$) = argmax $\prod_j \sum_{i=1}^{K}$ P(y$_j$=i,x$_j$)

$\qquad$ = argmax $\prod_j \sum_i^{K}$ P(y$_j$=i)p(x$_j$|y$_j$=i)

$\qquad$ = argmax $\prod_j \sum_{i=1}^{K}$ P(y$_j$=i) $\exp\left[ -\frac{1}{2\sigma^2} \left\| x_j - \mu_i \right\|^2 \right]$

# K-means and GMM

**"Linear" Decision Boundaries**



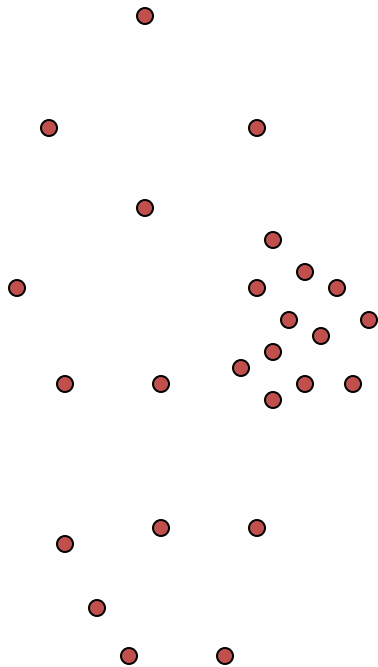Assume data comes from a mixture of K Gaussians distributions with same variance

Hard assignment:

$$P(y_j = i) = 1 \quad \text{if} \quad i = C(n)$$
$$\qquad\quad = 0 \quad \text{otherwise}$$

Maximize marginal likelihood:

$$\text{argmax} \prod_j P(x_j) \quad \equiv \arg\min_{\mu_1,\,..\,,\,\mu_K \atop C(1),\,...,\,C(m)} \sum_{j=1}^{m} \left\| x_j - \mu_{C(j)} \right\|^2$$

$\mu_1, \mu_2, .. , \mu_K, \sigma^2,$
$P(y=1), ..., P(y=k)$

**Same as K-means!!!**

# (One) bad case for K-means

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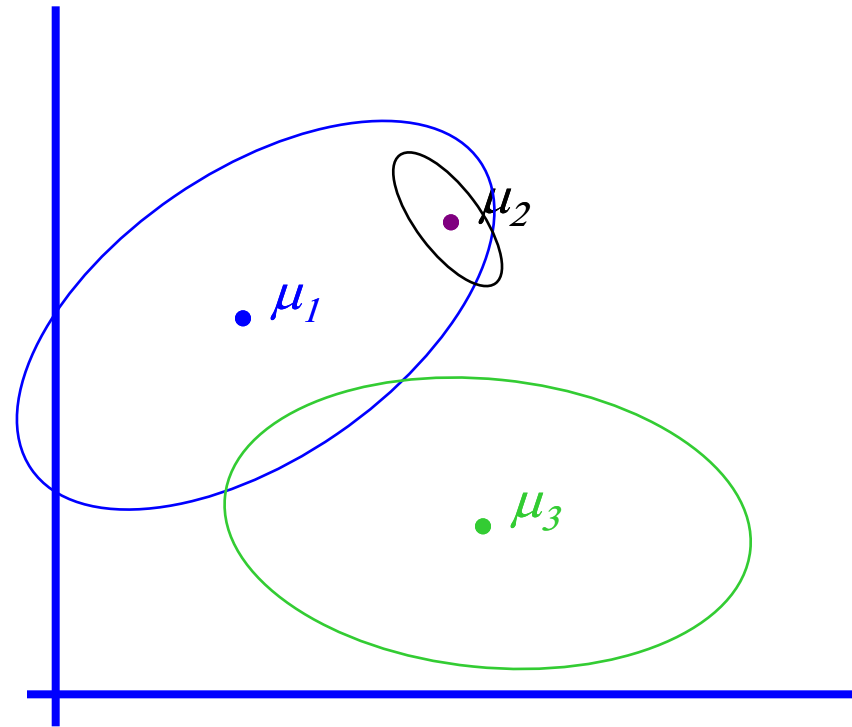- Some clusters may be "wider" than others

# General GMM

GMM – Gaussian Mixture Model  (Multi-modal distribution)

- There are k components

- Component *i* has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\Sigma_i$

Each data point is generated according to the following recipe:

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

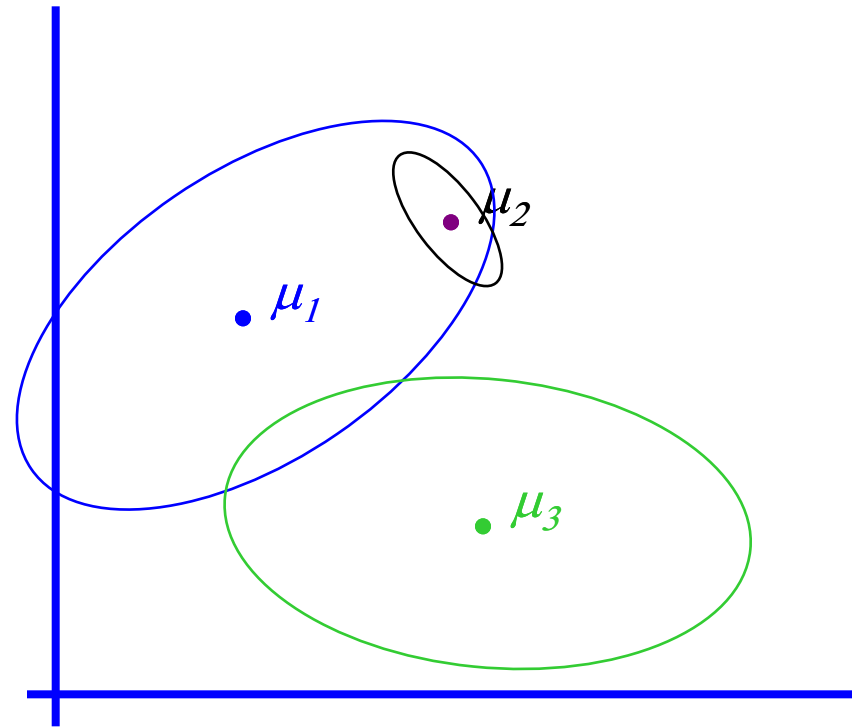2) Datapoint x ~ $\mathrm{N}(\mu_i,\ \Sigma_i)$

# General GMM

GMM – Gaussian Mixture Model  (Multi-modal distribution)

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

$p(x) = \sum\limits_i p(x|y=i) \, P(y=i)$

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

# General GMM

GMM – Gaussian Mixture Model  (Multi-modal distribution)

$p(x|y=i) \sim \mathrm{N}(\mu_i, \Sigma_i)$
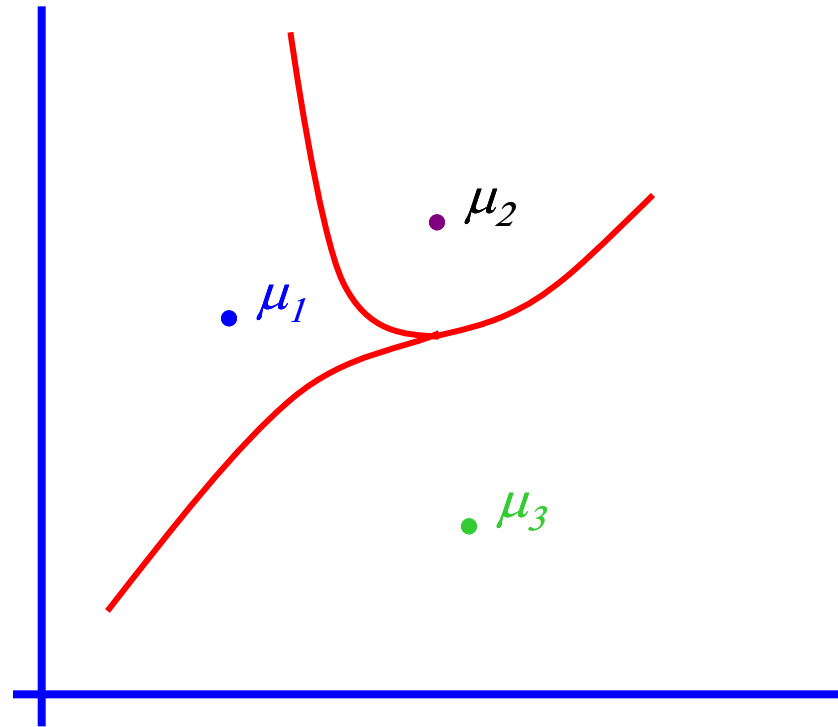
Gaussian Bayes Classifier:

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

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

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

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

**"Quadratic Decision boundary"** – second-order terms don't cancel out

# General GMM

Maximize marginal likelihood:

$$\text{argmax} \prod_j P(x_j) = \text{argmax} \prod_j \sum_{i=1}^{K} P(y_j=i, x_j)$$

$$= \text{argmax} \prod_j \sum_{i=1}^{K} P(y_j=i) p(x_j | y_j=i)$$

Soft assignment: $P(y_j=i) = P(y=i)$

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

How do we find the $\mu_i$'s and P(y=i)s which give max. marginal likelihood?

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

 * Use gradient descent:    Doable, but often slow

# 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.

- 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 $x_1\ x_2\ \ldots\ x_m$

We know there are k classes

We know P(y=1), P(y=2) P(y=3) … P(y=K)

We <u>don't</u> know $\mu_1\ \mu_2\ ..\ \mu_k$

We know common variance $\sigma^2$

We can write P( data | $\mu_1$…. $\mu_k$) $= \mathrm{p}\left(x_1...x_m\big|\mu_1...\mu_k\right)$

$$= \prod_{j=1}^{m} \mathrm{p}\left(x_j\big|\mu_1...\mu_k\right) \qquad \text{\color{red}Independent data}$$

$$= \prod_{j=1}^{m}\sum_{i=1}^{k} \mathrm{p}\left(x_j\big|y=i,\mu_i\right)\mathrm{P}\left(y=i\right) \qquad \text{\color{red}Marginalize over class}$$

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

# Expectation (E) step

If we know $\mu_1, ..., \mu_k \rightarrow$ easily compute prob. point $x_j$ belongs to class y=i

For each point $x_j$, j = 1, …, m

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

simply evaluate gaussian and normalize

Equivalent to assigning clusters to each data point in K-means

# Maximization (M) step

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

$\rightarrow$ MLE for $\mu_i$ is weighted average

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

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

Equivalent to updating cluster centers in K-means

# EM for spherical, same variance GMMs

**E-step**

Compute "expected" classes of all datapoints for each class

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

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)

$$\mu_i = \frac{\sum\limits_{j=1}^{m} P\left(y = i \middle| x_j\right)x_j}{\sum\limits_{j=1}^{m} P\left(y = i \middle| x_j\right)}$$

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

**E-step**

Compute "expected" classes of all datapoints for each class

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

Just evaluate a Gaussian at $x_j$

**M-step**

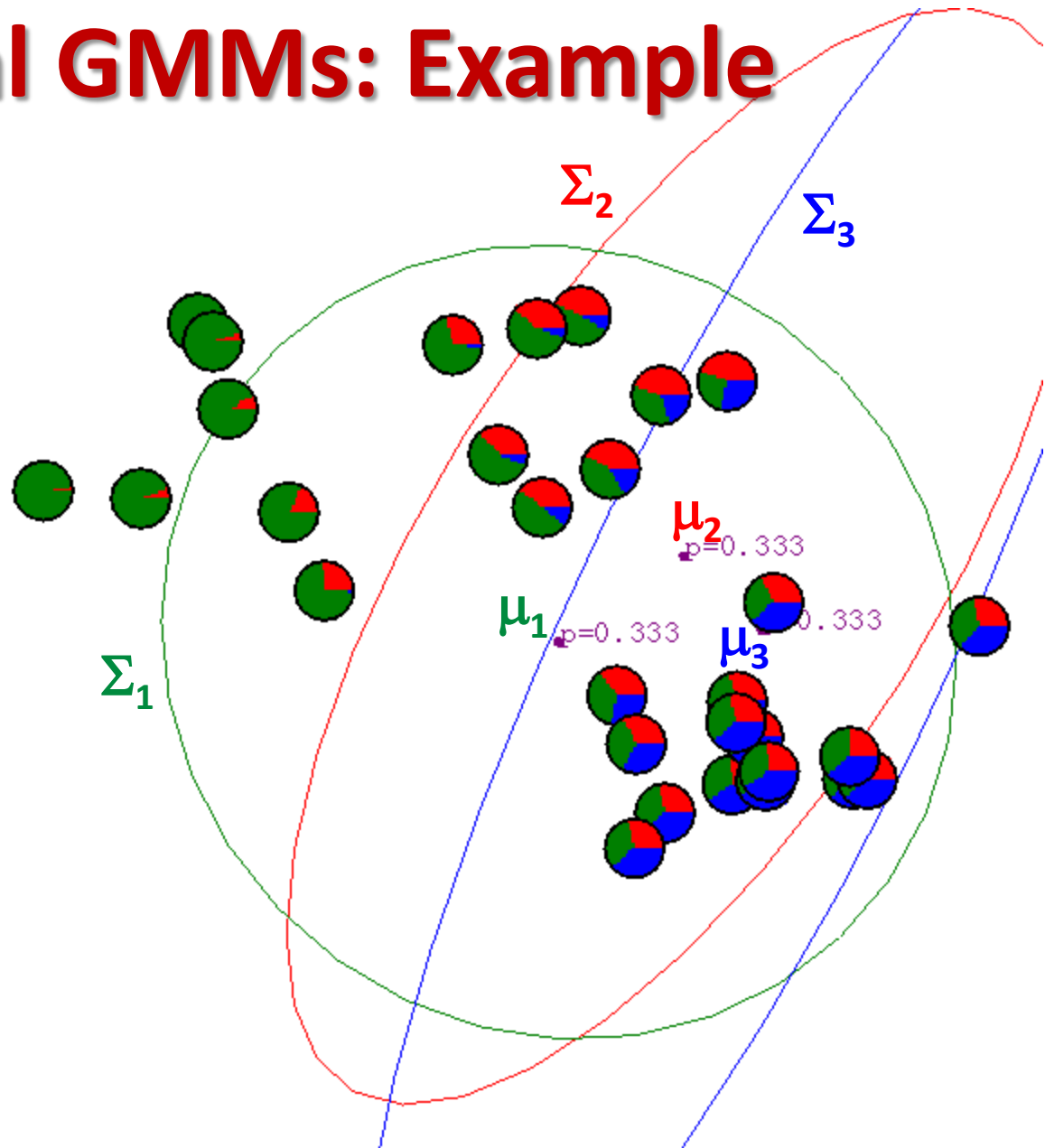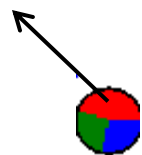Compute MLEs given our data's class membership distributions (weights)

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

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

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

$m$ = #data points

# EM for general GMMs: Example



$\Sigma_2$

$\Sigma_3$

$\mu_2$
p=0.333
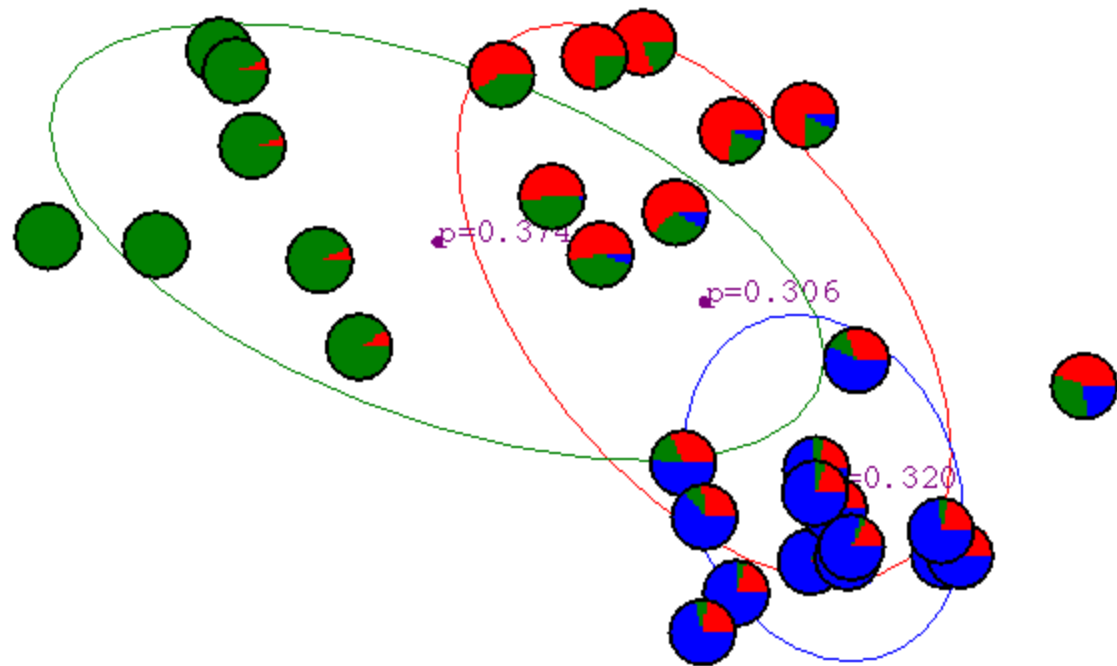
$\mu_1$
p=0.333   0.333

$\mu_3$

$\Sigma_1$

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

# After 1st iteration

# After 2nd iteration



p=0.374

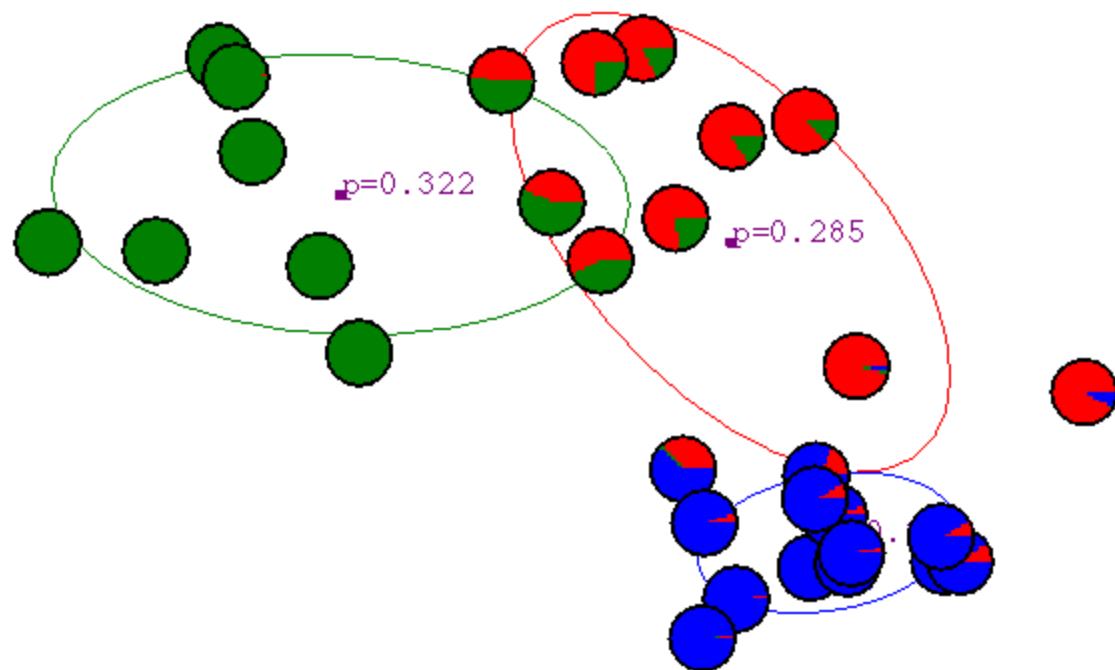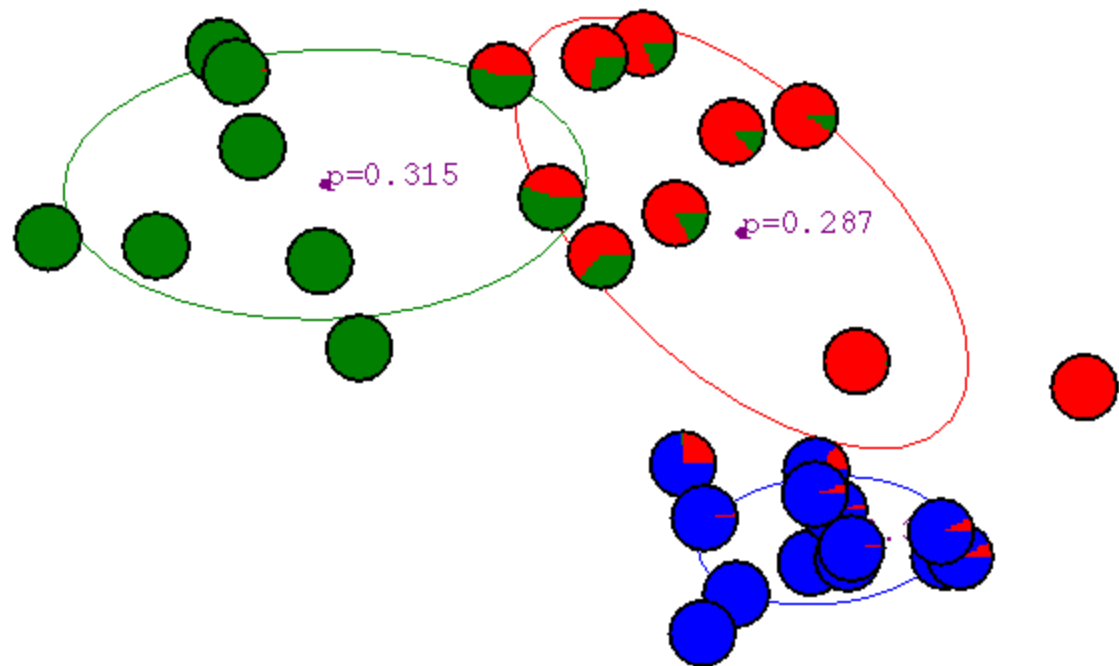p=0.306

=0.320

# After 3rd iteration

# After 4th iteration

# After 5<sup>th</sup> iteration
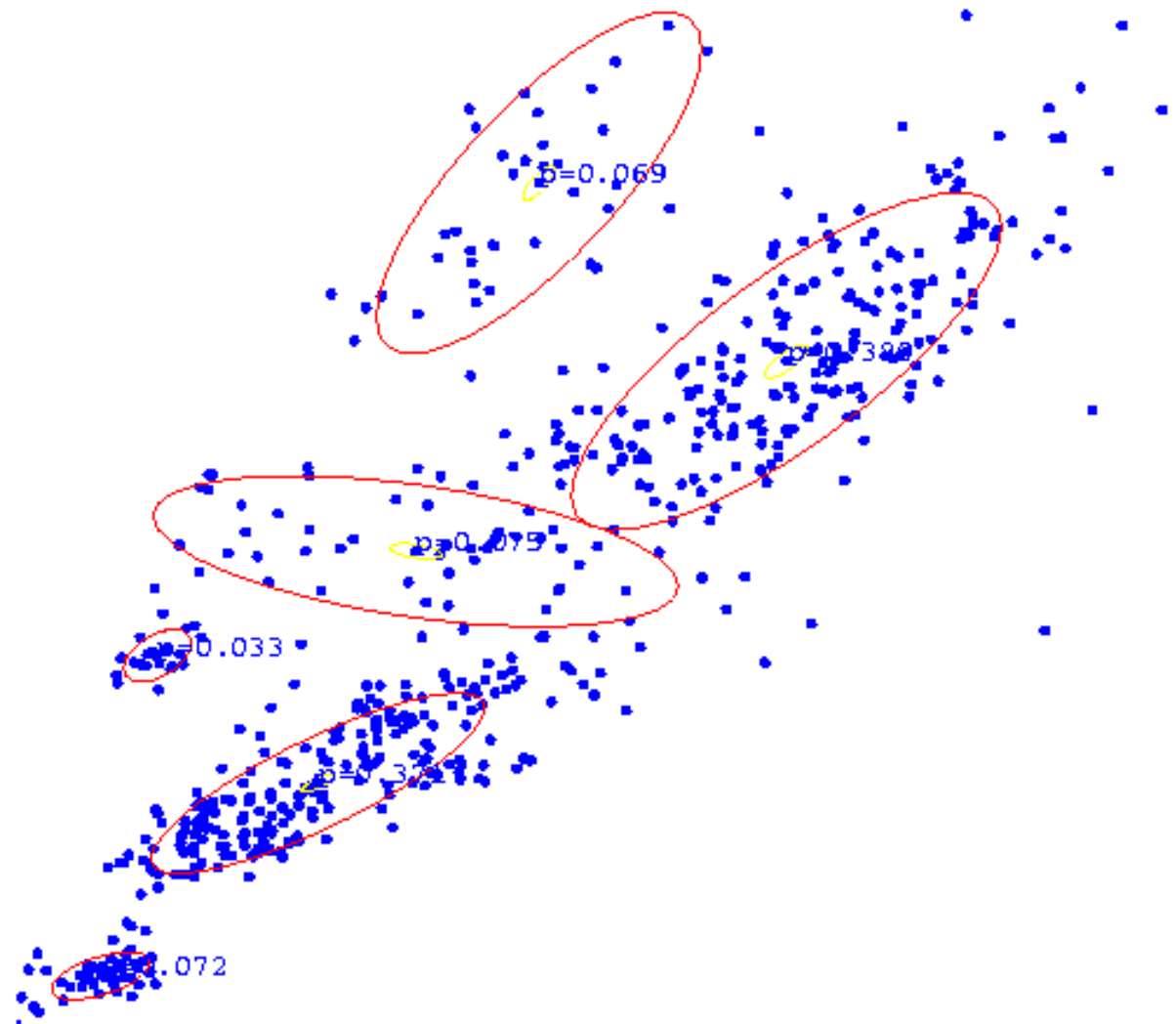
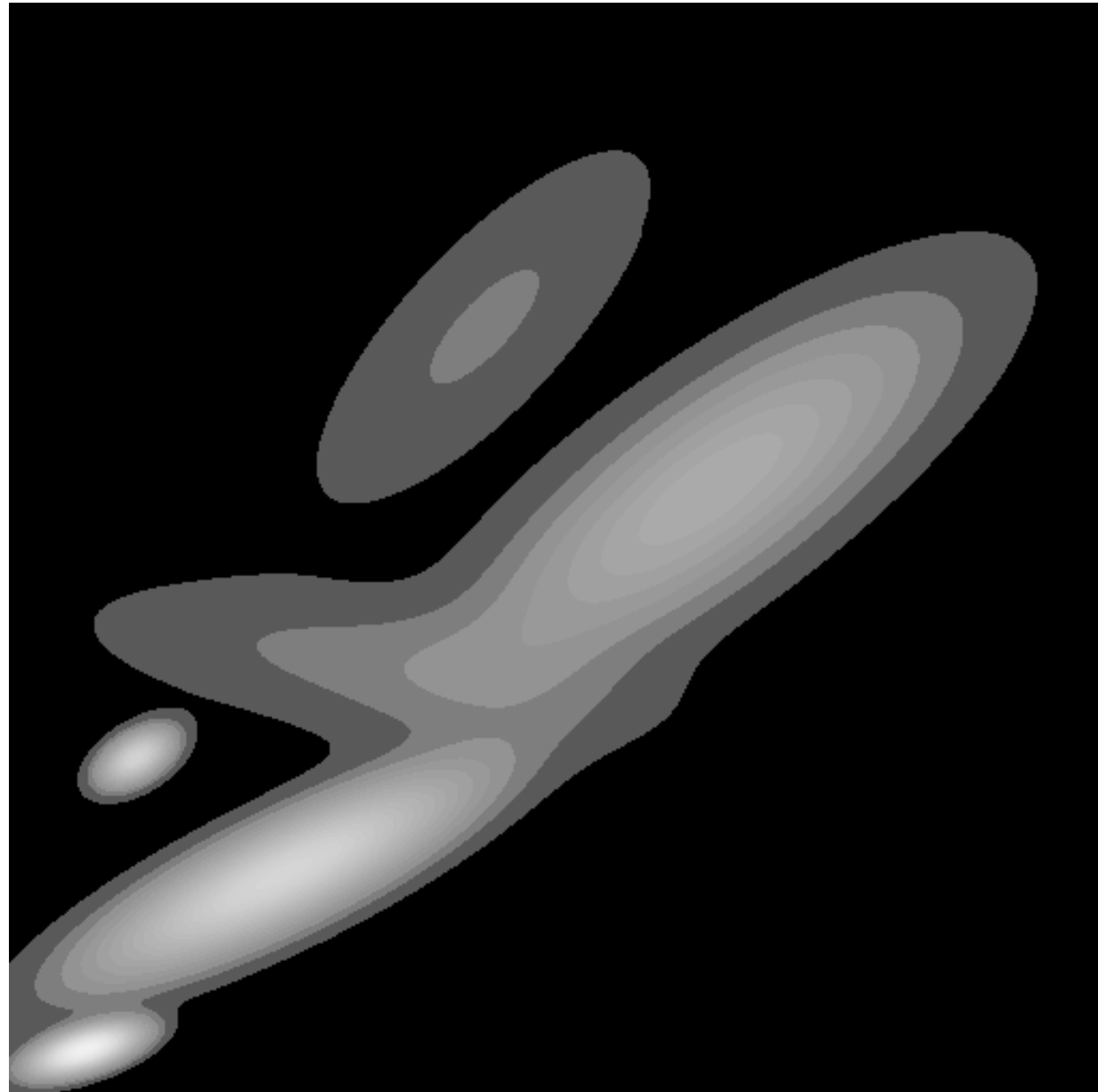# After 6<sup>th</sup> iteration



p=0.315

p=0.287

# After 20<sup>th</sup> iteration

# GMM clustering of the assay data

# Resulting Density Estimator

# Three classes of assay
**(each learned with it's own mixture model)**



Compound =

IL-1

TNF

nucleus

**Resulting Bayes Classifier**

# General EM algorithm

Marginal likelihood – **x** is observed, **z** is missing:

$$\log \mathrm{P(D;\theta)} = \log \prod_{j=1}^{m} P(\mathbf{x}_j \mid \theta)$$

$$= \sum_{j=1}^{m} \log P(\mathbf{x}_j \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)$$

$$\mathrm{D} = \left\{ \mathrm{x}_j \right\}_{j=1}^{m}$$

$\theta$ - model parameter(s)

# E step

**x** is observed, **z** is missing

Compute probability of missing data given current choice of $\theta$

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

$$\text{E.g., } P\left(y = i \middle| x_j, \lambda_t\right)$$

**M step** – Compute estimate of $\theta$ by maximizing marginal likelihood using $Q^{(t+1)}(\mathbf{z}|\mathbf{x}_j)$

# Lower-bound on marginal likelihood

$$P(D; \theta) = \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} \underbrace{Q(\mathbf{z} \mid \mathbf{x}_j)}_{\text{P}(\mathbf{z})} \underbrace{\frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}}_{\text{f}(\mathbf{z})}$$

**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

$$P(D; \theta) = \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} \underbrace{Q(\mathbf{z} \mid \mathbf{x}_j)}_{P(\mathbf{z})} \underbrace{\frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}}_{f(\mathbf{z})}$$

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

$$\geq \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

$$= \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + \underbrace{m.H(Q)}_{\textbf{Independent of } \theta}$$

# Lower-bound on marginal likelihood

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

**Independent of** $\theta$

$$\|$$

$$\underbrace{\sum_{\mathbf{z}} \sum_{j=1}^{m} \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) Q(\mathbf{z} \mid \mathbf{x}_j)}_{}$$

Expected log likelihood wrt Q

Since **z** is missing, instead take expectation over it (recall: probability of missing data **z** computed in E-step)

# M step

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

Maximize lower bound on marginal likelihood

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

Expected log likelihood wrt $Q^{(t+1)}$

Use expected counts instead of counts when computing MLE:
   If learning requires Count(**x**,**z**), Use $E_{Q(t+1)}$[Count(**x**,**z**)]

# EM as Coordinate Ascent

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

**M-step**: Fix Q, maximize F over $\theta$

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

M-step maximizes lower bound F on marginal likelihood => doesn't decrease the marginal likelihood

**E-step**: Fix $\theta$, maximize F over Q

We'll show this next. Thus,

E-step also maximizes lower bound F on marginal likelihood => doesn't decrease the marginal likelihood

Since marginal likelihood is bounded, Convergence follows!

# Convergence of EM

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

**E-step**: Fix $\theta$, maximize F over Q

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

$$= \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)}) P(\mathbf{x}_j \mid \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

$$= \underbrace{\sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)}}_{-KL(Q(\mathbf{z}|\mathbf{x}_j), P(\mathbf{z}|\mathbf{x}_j, \theta^{(t)}))} + \underbrace{\sum_{j=1}^{m} \sum_{\mathbf{z}} \overset{1}{Q(\mathbf{z}|\mathbf{x}_j)} \log P(\mathbf{x}_j|\theta^{(t)})}_{P(D; \theta^{(t)})}$$

**KL divergence between two distributions**

# Convergence of EM

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

**E-step**: Fix $\theta$, maximize F over Q

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

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

**KL>=0, Maximized if KL divergence = 0        KL(Q,P) = 0 iif Q = P**

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

Thus, E-step also maximizes lower bound F on marginal likelihood => doesn't decrease the marginal likelihood

# Convergence of EM

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

**M-step**: Fix Q, maximize F over θ

$$\mathrm{P}(\mathrm{D};\theta) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t)})$$
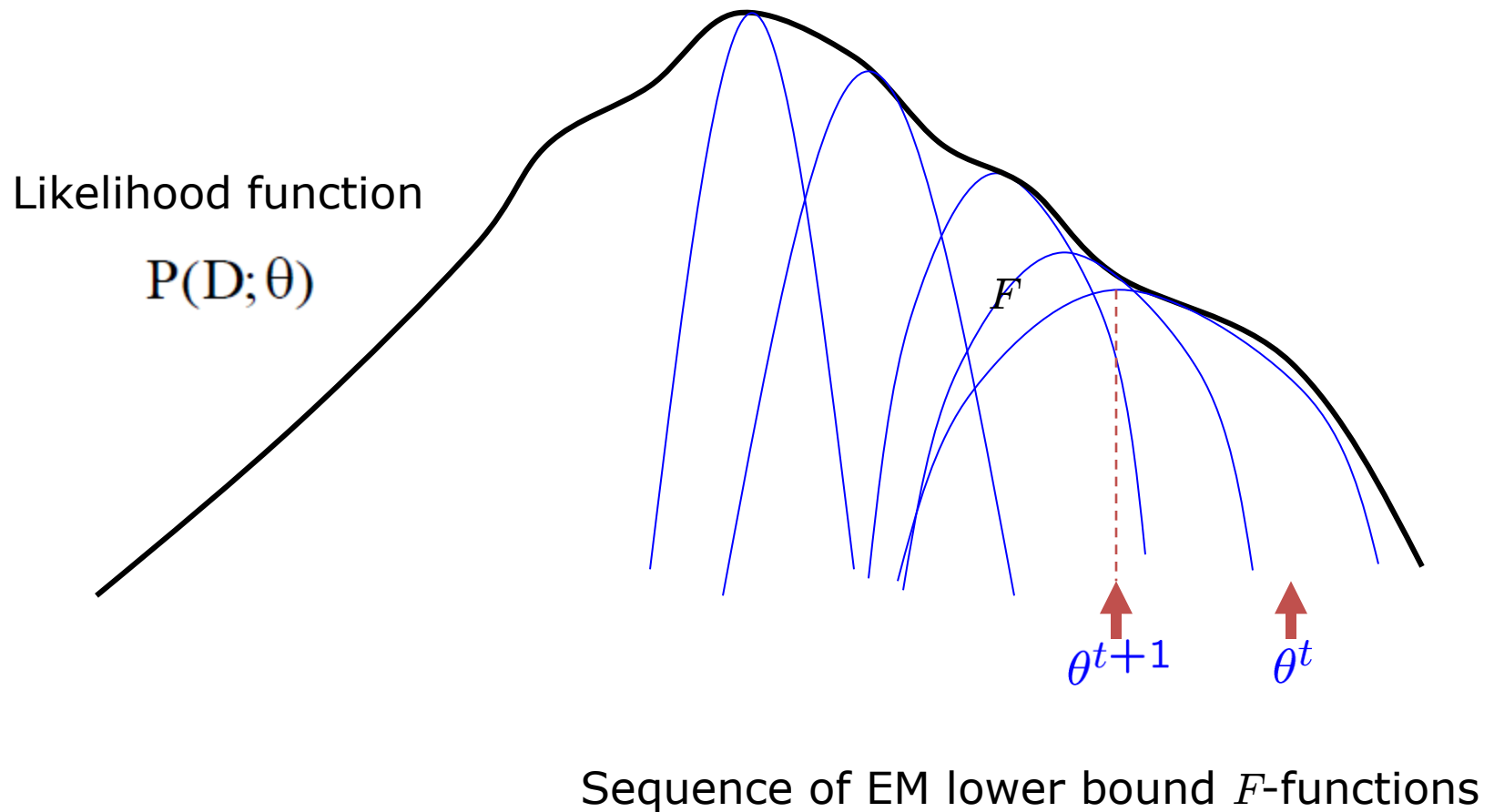
Maximizes lower bound F on marginal likelihood

**E-step**: Fix θ, maximize F over Q

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

Re-aligns F with marginal likelihood

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

# Convergence of EM



Likelihood function

$$P(D; \theta)$$

$F$

$\theta^{t+1}$      $\theta^t$

Sequence of EM lower bound $F$-functions
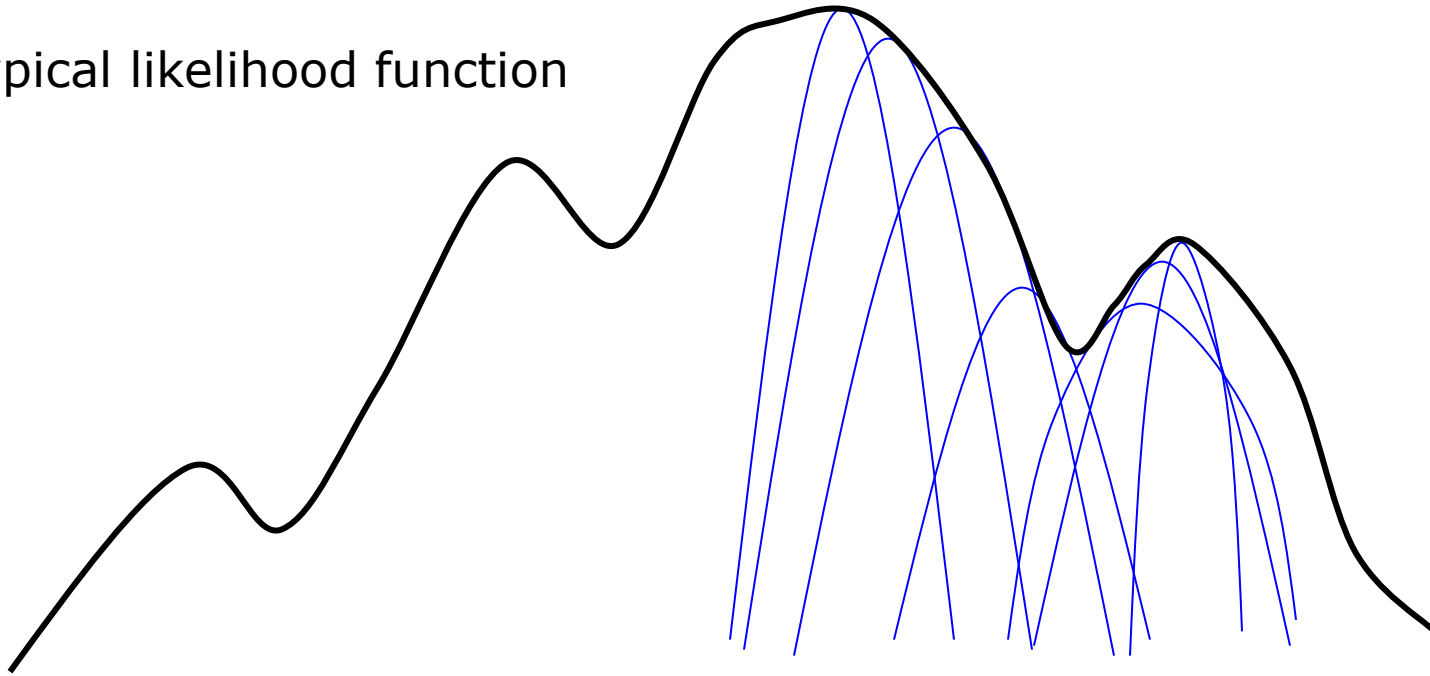
**EM monotonically converges to a local maximum of likelihood !**

# Convergence of EM

Typical likelihood function

Different sequence of EM lower bound
$F$-functions depending on initialization

**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.