

# The Expectation Maximization Algorithm

Gaurav Sharma

University of Rochester

# The Expectation Maximization (EM) Algorithm

- ▶ An iterative algorithm for **Maximum Likelihood (ML)** parameter estimation
- ▶ Formulate Maximum likelihood problems using **incomplete-data** framework
  - ▶ Either in presence of missing data
  - ▶ Or when the model can be simplified by introducing latent variable
- ▶ EM name given by Dempster, who gave the general formulation of the algorithm
- ▶ Wide ranging applications in machine learning, signal processing, statistics, data mining, and many other fields
  - ▶ Learning parameters of finite mixture models (say, mixture of Gaussians)
  - ▶ Model estimation for hidden Markov models (Baum-Welch)

## EM: Toy Example, Mixing Two Coins

- ▶ Random experiment: iid "mixing" of two coins
  - ▶ indexed by  $j = 1, 2$ , Coin  $j$  chosen with probability  $\alpha_j$
  - ▶ Coin characteristics:  $p_j = \Pr(H \mid \text{Coin} = j) = \Pr(1 \mid \text{Coin} = j)$
  - ▶ Parameters  $\theta = [\alpha_1, \alpha_2, p_1, p_2]$ , Constraint  $\sum_j \alpha_j = 1$
- ▶ Observations: Series of outcomes of mixing experiment  
 $x = [x_1, x_2, \dots, x_N]$
- ▶ Want ML estimate of parameters

$$\hat{\theta} = \arg \max_{\theta} p(x \mid \theta)$$

- ▶ What is  $p(x \mid \theta)$ ?
  - ▶ For example, for a specific string  $x = 0010110001?$
- ▶ Recall ML estimation for Bernoulli  $\mathcal{B}$  random variable
  - ▶ Expression for  $p(x \mid \theta)$  was straightforward
  - ▶ Why can't we do the same here?

## EM: Toy Example, Mixing Two Coins

- ▶ In addition to observation  $x$ , need information on which coin used for each outcome to write expression for likelihood
  - ▶ EM Approach: Complete data by introducing latent variables
    - ▶ Latent random variable  $Z^i$  as a  $2 \times 1$  vector for the  $i^{th}$  outcome indicating which coin was used for the  $i^{th}$  toss
    - ▶  $Z_j^i = 1$  if  $j^{th}$  coin produced the  $i^{th}$  outcome
    - ▶  $Z^i$ 's iid vectors having one entry as 1 others 0 ( $\equiv$  Bernoulli RV)
  - ▶ Complete likelihood is simple

$$p(x, z \mid \theta) = \prod_{i=1}^N p(x_i, z^i \mid \theta) = \prod_{i=1}^N \sum_j z_j^i \alpha_j \mathcal{B}(x_i, p_j)$$
$$\mathcal{B}(x, p) = p^x (1 - p)^{1-x}$$

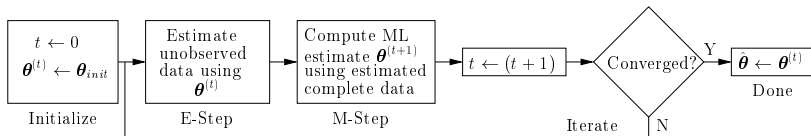
- ▶ How could you use the complete likelihood?

# EM Intuition for Coin Mixing Problem

- ▶ Don't know  $z$  so also estimate it and iteratively update along with parameter estimates
- ▶ An approach:
  - ▶ Say you have initial estimate  $\hat{\theta}^0$  of the parameters  $\theta$
  - ▶ Initialize iteration count  $t \leftarrow 0$
  - ▶ Estimate: Use current estimate of parameters  $\hat{\theta}^t$  to classify outcomes to obtain an estimate  $\hat{z}$  of  $z$
  - ▶ Maximize: Update parameters  $\hat{\theta}^{t+1}$  to maximize  $p(x, \hat{z} \mid \theta)$
  - ▶ Iterate till convergence
- ▶ EM is a refinement of this approach based on the same intuition
- ▶ Question: How well will this work for the coin mixing problem?
  - ▶ Are there parameters for which the approach fails?

# The Expectation Maximization (EM) Algorithm (Precursor)

- ▶ Iterate between
  - ▶ E-Step: Estimating “unobserved” data
  - ▶ M-Step: Maximum likelihood estimation of  $\theta$  from “completed” data



## EM: Complete vs Incomplete Likelihood

- ▶ If we cannot even evaluate  $p(x | \theta)$  readily, ML estimation seems to be hard
- ▶ EM: Complete the data and iterate between "estimating the missing data" and maximizing the complete likelihood with the "estimated missing data"
  - ▶ Coin example
- ▶ Data is often "incomplete" in many practical situations
  - ▶ Missing data, partial observations, indirect observation
- ▶ EM addresses this class of problems
  - ▶ Provides an indirect approach for ML estimation

# The Expectation Maximization (EM) Algorithm

- ▶ ML estimates of parameters

$$\hat{\theta} = \arg \max_{\theta} p(x | \theta)$$

- ▶ The EM algorithm addresses scenarios where it is hard to evaluate/formulate

$$p(x | \theta)$$

but we can complete the observations to obtain  $y$  such that  $p(y | \theta)$  is easy to evaluate and

$$x = f(y)$$

for some (typically many-to-one) mapping  $f()$ .



# Likelihood for Incomplete Data from Complete Data

- Likelihood for observed “incomplete” data from “complete” data

$$p(x | \theta) = \int_{y: f(y)=x} p(y | \theta) dy$$

- Log likelihood

$$l_x(\theta) = \ln(p(x | \theta))$$

# EM Algorithm Intuition

- ▶ Idea: If complete data was available, would like to maximize  $p(y | \theta) \equiv \text{maximize } \ln(p(y | \theta))$
- ▶ Since  $y$  is unavailable, maximize expectation of  $\ln(p(y | \theta))$  given data  $x$  and current estimate of parameters  $\theta^{(t)}$
- ▶ Two step procedure
  - ▶ E-Step: Compute the expectation

$$Q(\theta, \theta^t) = E [\ln(p(y | \theta)) | x, \theta^t]$$

- ▶ M-Step: Update the parameters to maximize the expectation

$$\theta^{t+1} = \arg \max_{\theta} Q(\theta, \theta^t)$$

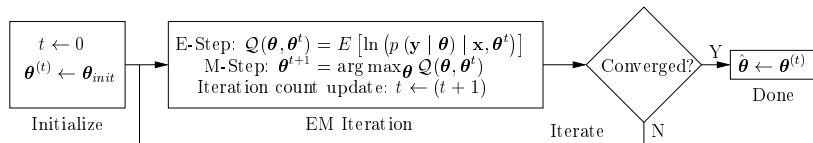
# The Expectation Maximization (EM) Algorithm: Formal Statement

- ▶ Observed data  $\mathbf{x}$ , Full data  $\mathbf{y}$ 
  - ▶ E-Step: Compute the expectation

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = E [\ln (p(\mathbf{y} \mid \boldsymbol{\theta})) \mid \mathbf{x}, \boldsymbol{\theta}^t]$$

- ▶ M-Step: Update parameter estimate to maximize the expectation

$$\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$$



## Formal EM algorithm for two coin mixing problem

- Recall, complete likelihood

$$p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) = \prod_{i=1}^N p(x_i, z^i \mid \boldsymbol{\theta}) = \prod_{i=1}^N \sum_j z_j^i \alpha_j \mathcal{B}(x_i, p_j)$$

- E Step:

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) &= E[\ln(p(\mathbf{y} \mid \boldsymbol{\theta})) \mid \mathbf{x}, \boldsymbol{\theta}^t] \\ &= \sum_{i=1}^N E \left[ \ln \left( \sum_j z_j^i \alpha_j \mathcal{B}(x_i, p_j) \right) \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)} \right] \\ &= \sum_{i=1}^N \sum_j \gamma_j^i \ln(\alpha_j \mathcal{B}(x_i, p_j)) \quad \text{Why?} \\ \gamma_j^i &= E[z_j^i \mid \mathbf{x}, \boldsymbol{\theta}^t] = Pr[z_j^{(i)} = 1 \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)}] \end{aligned}$$

$\gamma_j^i$  is posterior prob. that  $i^{th}$  outcome came from  $j^{th}$  coin ( $\mid \mathbf{x}, \boldsymbol{\theta}^t$ )

## Formal EM algorithm for two coin mixing problem

- ▶ E Step: Equivalent to estimating  $\gamma_j^i$ 
  - ▶ From Bayes rule:

$$\begin{aligned}\gamma_j^i &= Pr \left[ z_j^{(i)} = 1 | x^{(i)}, \theta^{(t)} \right] = \frac{Pr \left[ x^{(i)}, z_j^{(i)} = 1 | \theta^{(t)} \right]}{\sum_{z_i} Pr \left[ x^{(i)}, z_i | \theta^{(t)} \right]} \\ &= \frac{\alpha_j^{(t)} \mathcal{B}(x_i, p_j^{(t)})}{\sum_{l=1}^2 \alpha_l^{(t)} \mathcal{B}(x_i, p_l^{(t)})}\end{aligned}$$

- ▶ Soft as opposed to "hard" categorization of outcomes to coins

$$\begin{aligned}\mathcal{Q}(\theta, \theta^t) &= \sum_{i=1}^N \sum_j \gamma_j^i \ln (\alpha_j \mathcal{B}(x_i, p_j)) \\ &= \sum_{i=1}^N \sum_j \gamma_j^i (\ln \alpha_j + x_i \ln p_j + (1 - x_i) \ln(1 - p_j))\end{aligned}$$

## Formal EM algorithm for two coin mixing problem

- M Step: Maximize

$$\begin{aligned}\boldsymbol{\theta}^{t+1} &\stackrel{\text{def}}{=} [\alpha_1^{t+1}, \alpha_2^{t+1}, p_1^{t+1}, p_2^{t+1}] = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \sum_j \gamma_j^i (\ln \alpha_j + x_i \ln p_j + (1 - x_i) \ln(1 - p_j))\end{aligned}$$

- Calculus of variations (constraints important)

$$\begin{aligned}\frac{\nabla \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^t)}{\nabla \alpha_1} = 0 &\equiv \sum_{i=1}^N \left( \gamma_1^i \frac{1}{\alpha_1} - (1 - \gamma_1^i) \frac{1}{1 - \alpha_1} \right) = 0 \\ \alpha_j^{t+1} &= \frac{1}{N} \sum_{i=1}^N \gamma_j^i & p_j^{t+1} &= \frac{\sum_{i=1}^N \gamma_j^i x_i}{\sum_{i=1}^N \gamma_j^i}\end{aligned}$$

## Formal EM algorithm for two coin mixing problem

- E Step: Estimate posterior probabilities  $\gamma_j^i$

$$\begin{aligned}\gamma_j^i &= \frac{\alpha_j^{(t)} \mathcal{B}(x_i, p_j^{(t)})}{\sum_{l=1}^2 \alpha_l^{(t)} \mathcal{B}(x_i, p_l^{(t)})} \\ &= \frac{\alpha_j^{(t)} \left(p_j^{(t)}\right)^{x_i} \left(1 - p_j^{(t)}\right)^{1-x_i}}{\sum_{l=1}^2 \alpha_l^{(t)} \left(p_l^{(t)}\right)^{x_i} \left(1 - p_l^{(t)}\right)^{1-x_i}}\end{aligned}$$

- M Step:

$$\begin{aligned}\alpha_j^{t+1} &= \frac{1}{N} \sum_{i=1}^N \gamma_j^i \\ p_j^{t+1} &= \frac{\sum_{i=1}^N \gamma_j^i x_i}{\sum_{i=1}^N \gamma_j^i}\end{aligned}$$

## EM algorithm for two coin mixing problem

- ▶ Final EM algorithm relatively simple and intuitive
  - ▶ E Step: Compute posterior probability that outcome  $i$  came from coin  $j$
  - ▶ M Step: Update prob to posterior probability weighted mean fraction of heads in outcomes weight of  $i^{th}$  outcome is posterior probability that outcome  $i$  came from coin  $j$
- ▶ Contrast with abstraction
  - ▶ E-Step: Compute the expectation

$$Q(\theta, \theta^t) = E [\ln (p(y | \theta)) | x, \theta^t]$$

- ▶ M-Step: Update the parameters to maximize the expectation

$$\theta^{t+1} = \arg \max_{\theta} Q(\theta, \theta^t)$$

- ▶ Simple intuitive structure of EM algorithm extends to many problems



# Actual EM Application Examples

- ▶ WAMI to roadmap alignment
  - ▶ Formulating maximum likelihood directly was challenging
  - ▶ Different behavior of on-road vs spurious vehicles
  - ▶ Introduction of latent variable simplified things
    - ▶ Complete likelihood involving latent variables
    - ▶ Actual likelihood = marginal of complete likelihood over latent variables
- ▶ Color Barcodes
  - ▶ Similar story
  - ▶ Additional approximation/constraints from Physics
- ▶ Common theme in applications of EM
  - ▶ Incomplete Data/Latent Variable Problems

# EM Algorithm for Channel-wise Color Barcodes

- Recall model relation

$$\mathbf{D}_{3 \times 3} \mathbf{l}_{3 \times N} \approx \mathbf{d}_{3 \times N}$$

where  $\mathbf{l}_{3 \times N}$  is indicator variable indicating printing in  $\{R, G, B\}$  channels for the  $N$  -pixels,  $\mathbf{D}_{3 \times 3}$  is the unknown channel cross-interference matrix,  $\mathbf{d}_{3 \times N}$  are the observed densities corresponding to the pixels.

# EM Algorithm for Channel-wise Color Barcodes

- ▶ Probabilistic model formulation
  - ▶ Model noise in each pixel as iid zero mean Gaussian random variable

$$\mathbf{d}_{3 \times N} = \mathbf{D}_{3 \times 3} \mathbf{l}_{3 \times N} + \boldsymbol{\eta}_{3 \times N}$$

where  $\boldsymbol{\eta}_{3 \times N}$  is the noise in the  $N$  pixels.

- ▶ Incomplete data  $\mathbf{x} \equiv \mathbf{d}$ , Complete data  $\mathbf{y} \equiv (\mathbf{d}, \mathbf{l})$ , Parameters  $\boldsymbol{\theta} \equiv \mathbf{D}$
- ▶ Complete data log likelihood

$$l_d(\mathbf{D}) = -C \|\mathbf{d} - \mathbf{D}\mathbf{l}\|^2$$

where  $C$  is a positive constant independent of  $\mathbf{D}$ .

# EM Algorithm for Channel-wise Color Barcodes

## ► E-Step:

$$\begin{aligned} Q(D, D^{(t)}) &= E \left[ l_d(D) \mid d, D^{(t)} \right] = -C \left\| d - DE \left[ I \mid d, D^{(t)} \right] \right\|^2 \\ &= -C \left\| d - D\tilde{I} \right\|^2 \\ \tilde{I} &= E \left[ I \mid d, D^{(t)} \right] = \left( D^{(t)} \right)^{-1} d \end{aligned} \tag{1}$$

used fact that noise is assumed to be zero mean

## ► M-Step:

$$\begin{aligned} D^{(t+1)} &= \arg \max_D Q(D, D^{(t)}) \\ &= \arg \min_D \left\| d - D\tilde{I} \right\|^2 \end{aligned} \tag{2}$$

- Note: imposing non-negativity of  $D$  and  $\tilde{I}$  helps convergence to correct local minima

## EM Algorithm: For Exponential Family PDFs

- ▶ Exponential family complete likelihood

$$p(y | \theta) = \frac{1}{a(\theta)} b(y) \exp \left( c^T(\theta) s(y) \right)$$

- ▶  $a(\theta)$  and  $c(\theta)$  are vectors that are functions of the parameters
- ▶  $s(y)$  = sufficient statistic
- ▶ Exponential family of distributions includes many distributions of common interest
  - ▶ Gaussian, Poisson, binomial, uniform, Rayleigh, etc

## Example: Multivariate Gaussian as an Exponential Family PDF

- ▶ Recall: Multivariate Gaussian, Parameters  $\theta = (\mu, \Sigma)$

$$p(y | \theta) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left( -\frac{(y - \mu)^T \Sigma^{-1} (y - \mu)}{2} \right)$$

- ▶ Multivariate Gaussian as an exponential family PDF
  - ▶ Observe that the exponent is a polynomial with terms involving  $y_i$ ,  $y_i y_j$ , and constants
  - ▶ Can be expressed in the form

$$\begin{aligned} p(y | \theta) &= \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left( \sum_i \alpha_i y_i + \sum_{i,j} \beta_{ij} y_i y_j + \gamma \right) \\ &= \frac{1}{a(\theta)} \exp \left( c^T(\theta) s(y) \right) \end{aligned}$$

## EM Algorithm: For Exponential Family PDFs

- Expectation (E) step for exponential family

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) &= E \left[ \ln (p(y | \boldsymbol{\theta})) \mid \mathbf{x}, \boldsymbol{\theta}^t \right] \\ &= E \left[ \ln b(y) + \left( \mathbf{c}^T(\boldsymbol{\theta}) s(y) \right) - \ln a(\boldsymbol{\theta}) \mid \mathbf{x}, \boldsymbol{\theta}^t \right] \\ &= E \left[ \ln b(y) \mid \mathbf{x}, \boldsymbol{\theta}^t \right] + \mathbf{c}^T(\boldsymbol{\theta}) E \left[ s(y) \mid \mathbf{x}, \boldsymbol{\theta}^t \right] - \ln a(\boldsymbol{\theta}) \end{aligned}$$

- Note  $E \left[ s(y) \mid \mathbf{x}, \boldsymbol{\theta}^t \right]$  is a conditional estimate of the sufficient statistic given the observed data and current estimated parameters
- EM sometimes called estimation maximization algorithm

## EM Algorithm: For Exponential Family PDFs

- Recall: Expectation

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = E [\ln b(y) \mid \mathbf{x}, \boldsymbol{\theta}^t] + \mathbf{c}^T(\boldsymbol{\theta}) E [s(y) \mid \mathbf{x}, \boldsymbol{\theta}^t] - \ln a(\boldsymbol{\theta})$$

- Observe that  $E [\ln b(y) \mid \mathbf{x}, \boldsymbol{\theta}^t]$  does not depend on  $\boldsymbol{\theta}$
- EM simplifies to
  - E-Step (equivalent):

$$\mathbf{s}^{t+1} = E [s(y) \mid \mathbf{x}, \boldsymbol{\theta}^t]$$

- M-Step:

$$\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} \mathbf{c}^T(\boldsymbol{\theta}) \mathbf{s}^{t+1} - \ln a(\boldsymbol{\theta})$$



# Convergence of EM Algorithm

- ▶ Locally convergent
  - ▶ A contraction map
- ▶ Not necessarily to right point!

# Gaussian Mixture Models

- ▶  $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$  an  $N$  i.i.d random vectors that follows  $K$  -component Gaussian mixture distribution
- ▶ Probability density function:

$$p(\mathbf{x}^{(i)}|\boldsymbol{\theta}) = \sum_{j=1}^K \alpha_j \mathcal{N}(\mathbf{x}^{(i)}|\mu_j, \Sigma_j) \quad (3)$$

- ▶  $\mathcal{N}(\cdot)$  denotes the normal distribution
- ▶  $\alpha_j$  denotes the mixing coefficient of the  $j$  -th Gaussian
- ▶  $\mu_j$  and  $\Sigma_j$  are the mean and covariance matrix of  $j$  -th Gaussian

# Gaussian Mixture Model: Example

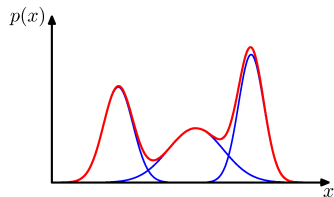


Figure: Mixture of 3 1D Gaussians (pdf)

# GMM: ML Parameter Estimate

- ▶ Likelihood function:

$$L(\mathbf{X}, \boldsymbol{\theta}) = \sum_{i=1}^N \log \left( \sum_{j=1}^K \alpha_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j) \right) \quad (4)$$

- ▶ Goal: Compute the ML estimate of parameters  $\boldsymbol{\theta} = \{\alpha_j, \mu_j, \Sigma_j\}_{j=1}^K$  that maximizes  $L(\mathbf{X}, \boldsymbol{\theta})$

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} L(\mathbf{X}, \boldsymbol{\theta}) \quad (5)$$

- ▶ Analytical solution is not possible
  - ▶ Use iterative EM algorithm
  - ▶  $\mathbf{X}$  is treated as incomplete data
  - ▶ Introduce latent variables to simplify the optimization problem

# EM for Gaussian Mixture Models

- ▶ Latent variable  $\mathbf{z}^{(i)} = [z_1^{(i)}, \dots, z_K^{(i)}]$  for each data vector  $\mathbf{x}^{(i)}$ 
  - ▶  $z^{(i)}$  indicates which component produced  $\mathbf{x}^{(i)}$
  - ▶ If  $\mathbf{x}^{(i)}$  is produced by the  $m$ -th mixture component, then  $z_m^{(i)} = 1$  and  $z_p^{(i)} = 0, \forall p \neq m$
- ▶ Complete data,  $\mathbf{Y} = \{\mathbf{X}, \mathbf{Z}\}$  and complete log-likelihood:

$$\begin{aligned} L_c(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z}) &= \sum_{i=1}^n \log \left[ \sum_{j=1}^K z_j^{(i)} \alpha_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j) \right] \\ &= \sum_{i=1}^n \sum_{j=1}^K z_j^{(i)} \log \left[ \alpha_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j) \right] \end{aligned} \quad (6)$$

- ▶ Why is the last step in the above equation valid?
- ▶ Actual value of  $z_j^{(i)}$  is unknown
- ▶ Use EM framework

## EM for Gaussian Mixture Models: E-Step

- ▶ Complete data log-likelihood:

$$L_c(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z}) = \sum_{i=1}^n \sum_{j=1}^K z_j^{(i)} \log \left[ \alpha_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j) \right] \quad (7)$$

- ▶ E-Step:

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(t)}) &= E \left[ L_c(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z}) | \mathbf{X}, \boldsymbol{\theta}^{(t)} \right] \\ &= \sum_{i=1}^N \sum_{j=1}^K \gamma_j^{(i)} \log \left[ \alpha_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j) \right] \end{aligned} \quad (8)$$

- ▶  $\gamma_j^{(i)} = E \left[ z_j^{(i)} | \mathbf{X}, \boldsymbol{\theta}^{(t)} \right] = Pr \left[ z_j^{(i)} = 1 | \mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)} \right]$ . Posterior probability that  $i^{\text{th}}$  observation came from  $j^{\text{th}}$  mixture component
- ▶  $\boldsymbol{\theta}^{(t)}$  is the current estimate of parameters

## EM for Gaussian Mixture Models: M-Step

- ▶ Maximize conditional expectation  $Q(\theta|\theta^{(t)})$

$$Q(\theta|\theta^{(t)}) = \sum_{i=1}^N \sum_{j=1}^K \gamma_j^{(i)} \log \left[ \alpha_j \mathcal{N}(\mathbf{x}^{(i)}|\mu_j, \Sigma_j) \right] \quad (9)$$

- ▶ M-Step: Re-estimate parameters  $\theta^{(t+1)} = \arg \max_{\theta} Q(\theta|\theta^{(t)})$ 
  - ▶ Simple calculus (with constraints)

$$\alpha_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_j^{(i)}}{N} \quad \mu_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_j^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^N \gamma_j^{(i)}} \quad (10)$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_j^{(i)} (\mathbf{x}^{(i)} - \mu_j^{(t+1)})(\mathbf{x}^{(i)} - \mu_j^{(t+1)})^T}{\sum_{i=1}^N \gamma_j^{(i)}} \quad (11)$$

# EM for Gaussian Mixture Models: Final Algo. Summary

- ▶ Initial parameter estimate  $\theta^{(0)}$
- ▶ Iteratively apply alternate E and M-steps
  - ▶ E-step: Estimate posterior probabilities  $\gamma_j^{(i)}$  and the expected value of complete data likelihood,

$$\gamma_j^{(i)} = \frac{\alpha_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j)}{\sum_{l=1}^K \alpha_l \mathcal{N}(\mathbf{x}^{(i)} | \mu_l, \Sigma_l)} \quad (12)$$

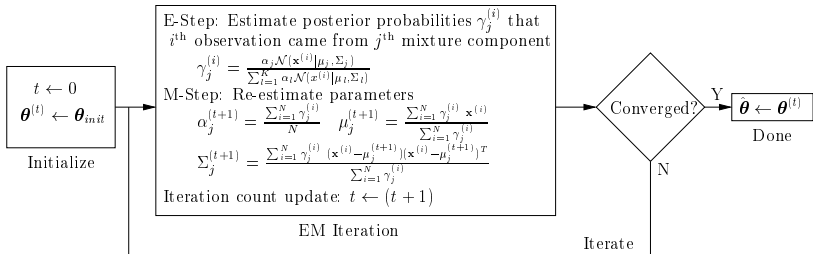
- ▶ M-step: Re-estimate parameters  $\theta^{(t+1)}$

$$\alpha_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_j^{(i)}}{N} \quad \mu_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_j^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^N \gamma_j^{(i)}} \quad (13)$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_j^{(i)} (\mathbf{x}^{(i)} - \mu_j^{(t+1)})(\mathbf{x}^{(i)} - \mu_j^{(t+1)})^T}{\sum_{i=1}^N \gamma_j^{(i)}} \quad (14)$$



# EM Algorithm for Gaussian Mixture Models



# Convergence of EM Algorithm

- ▶ Proof of convergence of EM algorithm (Dempster [1])
- ▶ Further extended by (Wu et al [5])
  - ▶ Minor corrections to Dempster's proof
  - ▶ Conditions under which EM converges to stationary points and local maxima
- ▶ EM has linear rate of convergence (Redner & Walker [4])
- ▶ EM shows asymptotic superlinear rate of convergence ([6, 3, 2])
  - ▶ Connection between first order Gradient ascent and EM algorithm
    - ▶ EM has faster convergence than gradient ascent, for well separated Gaussians
    - ▶ Superlinear rate of convergence, as overlap among Gaussian components goes to zero

# Generalized EM and Convergence

- ▶ Generalized EM (GEM): Any algorithm that chooses  $\theta^{(t+1)}$  to increase  $Q(\theta^{(t)}|\theta^{(t)})$  instead of maximizing.
- ▶ Specifically a GEM chooses any  $\theta^{(t+1)} \in \Omega$  (parameter space) such that,

$$Q(\theta^{(t+1)}|\theta^{(t)}) \geq Q(\theta^{(t)}|\theta^{(t)}) \quad (15)$$

- ▶ GEM algorithm  $\theta^{(t+1)} = M(\theta^{(t)})$ .
  - ▶ an implicit mapping:  $\theta \rightarrow M(\theta)$  from parameter space  $\Omega$  to itself
  - ▶ A fixed point  $\theta^*$  satisfies  $M(\theta^*) = \theta^*$

# Convergence Rate

- ▶ Linear Convergence Rate

- ▶ A sequence  $\{x^{(t)}\}$  is said to be linearly convergent if, there exists a number  $r \in (0, 1)$  such that,

$$\lim_{t \rightarrow \infty} \frac{\|x^{(t+1)} - x^*\|}{\|x^{(t)} - x^*\|} = r \quad (16)$$

- ▶ Rate of convergence  $r$  and speed of convergence  $(1 - r)$
- ▶ Superlinear if  $r = 0$

# Expectation Maximization: Summary

- ▶ EM provides a framework for addressing "missing data" problems
  - ▶ Iteratively optimize likelihood of observed "incomplete" data by "completing" it
    - ▶ Alternation between computation of expectation of unobserved variables given current parameter values and determination of maximizing parameters for "full data" likelihood
- ▶ EM algorithm is guaranteed to converge monotonically to a stationary point
- ▶ EM for GMM:
  - ▶ Fast convergence for well-separated Gaussians with similar sizes
  - ▶ Convergence slows down due to:
    - ▶ strong Overlap among Gaussian components
    - ▶ large dynamic range of the mixing coefficient ( $\alpha_j < 0.1$ )

# References I

- [1] A.P. Dempster, N.M. Laird, and D.B. Rubin. "Maximum likelihood from incomplete data via the EM algorithm". In: *Journal of the Royal Statistical Society. Series B (Methodological)* 39.1 (1977), pp. 1–38.
- [2] J. Ma and S. Fu. "On the correct convergence of the EM algorithm for Gaussian mixtures". In: *Pattern Recognition* 38.12 (2005), pp. 2602–2611.
- [3] J. Ma, L. Xu, and M.I. Jordan. "Asymptotic convergence rate of the EM algorithm for Gaussian mixtures". In: *Neural Computation* 12.12 (2000), pp. 2881–2907.
- [4] R. Redner and H. Walker. "Mixture Densities, maximum likelihood, and the EM Algorithm". In: *SIAM Review* 26.2 (1984), pp. 195–239.
- [5] C. F. Jeff Wu. "On the convergence properties of the EM algorithm". In: *Annals of Statistics* 11 (1983), pp. 95–103.

## References II

- [6] L. Xu and M.I. Jordan. “On convergence properties of the EM algorithm for Gaussian mixtures”. In: *Neural computation* 8.1 (1996), pp. 129–151.