# Lecture Stat 461-561 Expectation-Maximization Algorithm

AD

#### January 2008

• Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of the sample  $\mathbf{X} = (X_1, ..., X_n)$ . Then given that  $\mathbf{X} = \mathbf{x}$  is observed, the likelihood function is given by  $f(\mathbf{x}|\theta) = L(\theta|\mathbf{x})$  and  $I(\theta) = \log L(\theta|\mathbf{x})$ .

- Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of the sample  $\mathbf{X} = (X_1, ..., X_n)$ . Then given that  $\mathbf{X} = \mathbf{x}$  is observed, the likelihood function is given by  $f(\mathbf{x}|\theta) = L(\theta|\mathbf{x})$  and  $I(\theta) = \log L(\theta|\mathbf{x})$ .
- The Maximum Likelihood Estimate (MLE) is defined by

$$\widehat{ heta} = \mathop{\mathrm{arg\,max}}_{ heta \in \Theta} I\left( heta
ight).$$

- Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of the sample  $\mathbf{X} = (X_1, ..., X_n)$ . Then given that  $\mathbf{X} = \mathbf{x}$  is observed, the likelihood function is given by  $f(\mathbf{x}|\theta) = L(\theta|\mathbf{x})$  and  $I(\theta) = \log L(\theta|\mathbf{x})$ .
- The Maximum Likelihood Estimate (MLE) is defined by

$$\widehat{ heta} = \mathop{\mathrm{arg\,max}}_{ heta \in \Theta} I\left( heta
ight).$$

• Under regularity assumptions, the MLE is

- Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of the sample  $\mathbf{X} = (X_1, ..., X_n)$ . Then given that  $\mathbf{X} = \mathbf{x}$  is observed, the likelihood function is given by  $f(\mathbf{x}|\theta) = L(\theta|\mathbf{x})$  and  $I(\theta) = \log L(\theta|\mathbf{x})$ .
- The Maximum Likelihood Estimate (MLE) is defined by

$$\widehat{ heta} = \mathop{\mathrm{arg\,max}}_{ heta \in \Theta} I\left( heta
ight).$$

- Under regularity assumptions, the MLE is
  - consistent, i.e.  $\hat{\theta}_n \xrightarrow{\mathsf{P}} \theta_*$  where  $\theta_*$  is the true value

- Let  $f(\mathbf{x}|\theta)$  denote the joint pdf or pmf of the sample  $\mathbf{X} = (X_1, ..., X_n)$ . Then given that  $\mathbf{X} = \mathbf{x}$  is observed, the likelihood function is given by  $f(\mathbf{x}|\theta) = L(\theta|\mathbf{x})$  and  $I(\theta) = \log L(\theta|\mathbf{x})$ .
- The Maximum Likelihood Estimate (MLE) is defined by

$$\widehat{ heta} = egin{argmmatrix} & ext{argmax} \ I\left( heta
ight). \ & heta\in \Theta \ \end{split}$$

- Under regularity assumptions, the MLE is
  - consistent, i.e.  $\widehat{\theta}_n \xrightarrow{\mathsf{P}} \theta_*$  where  $\theta_*$  is the true value
  - asymptotically efficient, i.e. the MLE has the smallest variance, at least for large samples.

#### Standard Numerical Methods

• Assume  $\theta = (\theta_1, ..., \theta_p)^T$ , then under regularity assumptions the MLE are the values that solve

$$rac{\partial I\left( heta
ight)}{\partial heta_{i}} = 0 ext{ for } i = 1, ..., p. ext{ (1)}$$

#### Standard Numerical Methods

 Assume θ = (θ<sub>1</sub>, ..., θ<sub>p</sub>)<sup>T</sup>, then under regularity assumptions the MLE are the values that solve

$$\frac{\partial I\left( heta
ight)}{\partial heta_{i}} = 0 ext{ for } i = 1, ..., p. ext{ (1)}$$

• Example: if  $X_i \overset{\text{i.i.d.}}{\sim} \mathcal{N}\left( heta, 1
ight)$  then

$$L(\theta | \mathbf{x}) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x_i - \theta)^2}{2}\right)$$

and

$$\frac{dl\left(\theta\right)}{d\theta}=0\Leftrightarrow\sum_{i=1}^{n}\left(x_{i}-\theta\right)=0\Rightarrow\widehat{\theta}=\frac{1}{n}\sum_{i=1}^{n}x_{i}.$$

One can check  $\frac{d^2 I(\theta)}{d\theta^2}\Big|_{\widehat{\theta}} < 0$ , hence it is a maximum.

#### Standard Numerical Methods

 Assume θ = (θ<sub>1</sub>, ..., θ<sub>p</sub>)<sup>T</sup>, then under regularity assumptions the MLE are the values that solve

$$\frac{\partial I\left( heta
ight)}{\partial heta_{i}} = 0 ext{ for } i = 1, ..., p. ext{ (1)}$$

• Example: if  $X_i \overset{\text{i.i.d.}}{\sim} \mathcal{N}\left( heta, 1
ight)$  then

$$L(\theta | \mathbf{x}) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x_i - \theta)^2}{2}\right)$$

and

$$\frac{dl\left(\theta\right)}{d\theta}=0\Leftrightarrow\sum_{i=1}^{n}\left(x_{i}-\theta\right)=0\Rightarrow\widehat{\theta}=\frac{1}{n}\sum_{i=1}^{n}x_{i}.$$

One can check  $\left. \frac{d^2 I(\theta)}{d\theta^2} \right|_{\widehat{\theta}} < 0$ , hence it is a maximum.

• In more complex examples, we cannot solve (1) easily and we need to rely on numerical methods.

• Iterative algorithm to find the MLE of  $\theta = (\theta_1, ..., \theta_p)^{\mathsf{T}}$ .

Iterative algorithm to find the MLE of θ = (θ<sub>1</sub>, ..., θ<sub>p</sub>)<sup>T</sup>.
Assume θ<sup>(k)</sup> is the estimate at iteration k then

$$\frac{\partial I(\theta)}{\partial \theta} \approx \frac{\partial I\left(\theta^{(k)}\right)}{\partial \theta} + \frac{\partial^2 I\left(\theta^{(k)}\right)}{\partial \theta \partial \theta^{\mathsf{T}}} \left(\theta - \theta^{(k)}\right).$$

Iterative algorithm to find the MLE of θ = (θ<sub>1</sub>, ..., θ<sub>p</sub>)<sup>T</sup>.
Assume θ<sup>(k)</sup> is the estimate at iteration k then

$$\frac{\partial I\left(\theta\right)}{\partial \theta} \approx \frac{\partial I\left(\theta^{(k)}\right)}{\partial \theta} + \frac{\partial^{2} I\left(\theta^{(k)}\right)}{\partial \theta \partial \theta^{\mathsf{T}}} \left(\theta - \theta^{(k)}\right).$$

By writing

$$g(\theta) = \left(\frac{\partial I(\theta)}{\partial \theta_1}, \frac{\partial I(\theta)}{\partial \theta_2}, ..., \frac{\partial I(\theta)}{\partial \theta_p}\right)^{\mathsf{T}}, \\ H(\theta) = \frac{\partial^2 I(\theta)}{\partial \theta \partial \theta^{\mathsf{T}}} = \left(\frac{\partial^2 I(\theta)}{\partial \theta_i \partial \theta_i}\right)$$

this means that

$$\mathbf{0} = g\left(\theta\right) \approx g\left(\theta^{(k)}\right) + H\left(\theta^{(k)}\right)\left(\theta - \theta^{(k)}\right)$$

Iterative algorithm to find the MLE of θ = (θ<sub>1</sub>, ..., θ<sub>p</sub>)<sup>T</sup>.
Assume θ<sup>(k)</sup> is the estimate at iteration k then

$$\frac{\partial I\left(\theta\right)}{\partial \theta} \approx \frac{\partial I\left(\theta^{(k)}\right)}{\partial \theta} + \frac{\partial^{2} I\left(\theta^{(k)}\right)}{\partial \theta \partial \theta^{\mathsf{T}}} \left(\theta - \theta^{(k)}\right).$$

By writing

$$g(\theta) = \left(\frac{\partial I(\theta)}{\partial \theta_1}, \frac{\partial I(\theta)}{\partial \theta_2}, ..., \frac{\partial I(\theta)}{\partial \theta_p}\right)^{\mathsf{T}}, \\ H(\theta) = \frac{\partial^2 I(\theta)}{\partial \theta \partial \theta^{\mathsf{T}}} = \left(\frac{\partial^2 I(\theta)}{\partial \theta_i \partial \theta_i}\right)$$

this means that

$$0 = g(\theta) \approx g\left(\theta^{(k)}\right) + H\left(\theta^{(k)}\right)\left(\theta - \theta^{(k)}\right)$$

This suggests

$$\theta^{(k+1)} = \theta^{(k)} - H\left(\theta^{(k)}\right)^{-1} g\left(\theta^{(k)}_{c}\right) \in \mathbb{R} \quad \text{for } k \in \mathbb{R} \quad \text{for } k \in \mathbb{R}$$

January 2008

Iterative algorithm to find the MLE of θ = (θ<sub>1</sub>, ..., θ<sub>p</sub>)<sup>T</sup>.
Assume θ<sup>(k)</sup> is the estimate at iteration k then

$$\frac{\partial I\left(\theta\right)}{\partial \theta} \approx \frac{\partial I\left(\theta^{(k)}\right)}{\partial \theta} + \frac{\partial^{2} I\left(\theta^{(k)}\right)}{\partial \theta \partial \theta^{\mathsf{T}}} \left(\theta - \theta^{(k)}\right).$$

By writing

$$g(\theta) = \left(\frac{\partial I(\theta)}{\partial \theta_1}, \frac{\partial I(\theta)}{\partial \theta_2}, ..., \frac{\partial I(\theta)}{\partial \theta_p}\right)^{\mathsf{T}}, \\ H(\theta) = \frac{\partial^2 I(\theta)}{\partial \theta \partial \theta^{\mathsf{T}}} = \left(\frac{\partial^2 I(\theta)}{\partial \theta_i \partial \theta_i}\right)$$

this means that

$$0 = g(\theta) \approx g\left(\theta^{(k)}\right) + H\left(\theta^{(k)}\right)\left(\theta - \theta^{(k)}\right)$$

This suggests

$$\theta^{(k+1)} = \theta^{(k)} - H\left(\theta^{(k)}\right)^{-1} g\left(\theta^{(k)}_{c}\right) \in \mathbb{R} \quad \text{for } k \in \mathbb{R} \quad \text{for } k \in \mathbb{R}$$

January 2008

• It may prove difficult to calculate the Hessian matrix.

- It may prove difficult to calculate the Hessian matrix.
- At each iteration, it requires computing a new Hessian matrix.

- It may prove difficult to calculate the Hessian matrix.
- At each iteration, it requires computing a new Hessian matrix.
- Depending on the initial value, the method may converge or diverge.

- It may prove difficult to calculate the Hessian matrix.
- At each iteration, it requires computing a new Hessian matrix.
- Depending on the initial value, the method may converge or diverge.
- Quasi-Newton methods have been proposed to mitigate these problems: it avoids calculating the Hessian  $H\left(\theta^{(k)}\right)$  and step widths can be introduced to accelerate convergence/prevent divergence.

• Determine a search direction vector  $d_k = H_k^{-1} g\left( heta^{(k)} 
ight)$ .

- Determine a search direction vector  $d_k = H_k^{-1} g\left( heta^{(k)} 
  ight)$  .
- Determine the optimum step width  $\lambda_k$  maximizing  $I\left( heta^{(k)}+\lambda d_k
  ight)$  .

- Determine a search direction vector  $d_k = H_k^{-1}g\left( heta^{(k)}
  ight)$ .
- Determine the optimum step width  $\lambda_k$  maximizing  $I\left( heta^{(k)}+\lambda d_k
  ight)$  .

• Set 
$$\theta^{(k+1)} = \theta^{(k)} + \lambda_k d_k$$
 and  $y_k = g\left(\theta^{(k+1)}\right) - g\left(\theta^{(k)}\right)$ .

- Determine a search direction vector  $d_k = H_k^{-1} g\left( heta^{(k)} 
  ight)$  .
- Determine the optimum step width  $\lambda_k$  maximizing  $I\left( heta^{(k)}+\lambda d_k
  ight)$  .
- Set  $\theta^{(k+1)} = \theta^{(k)} + \lambda_k d_k$  and  $y_k = g\left(\theta^{(k+1)}\right) g\left(\theta^{(k)}\right)$ .
- Obtain an estimate of  $H\left(\theta^{(k+1)}\right)^{-1}$  by using Davidson-Fletcher-Powell (DFP) or Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm

$$\begin{aligned} H_{k+1}^{-1} &= H_{k}^{-1} + \frac{s_{k}s_{k}^{\mathsf{T}}}{s_{k}^{\mathsf{T}}y_{k}} - \frac{H_{k}^{-1}y_{k}y_{k}^{\mathsf{T}}H_{k}^{-1}}{y_{k}^{\mathsf{T}}H_{k}^{-1}y_{k}}, \\ H_{k+1}^{-1} &= H_{k}^{-1} + \frac{s_{k}y_{k}^{\mathsf{T}}H_{k}^{-1}}{s_{k}^{\mathsf{T}}y_{k}} - \frac{H_{k}^{-1}y_{k}s_{k}^{\mathsf{T}}}{s_{k}^{\mathsf{T}}y_{k}} + \left(1 + \frac{y_{k}H_{k}^{-1}y_{k}^{\mathsf{T}}}{s_{k}^{\mathsf{T}}y_{k}}\right) \frac{s_{k}s_{k}^{\mathsf{T}}}{s_{k}^{\mathsf{T}}y_{k}} \end{aligned}$$

where  $s_k = \theta^{(k+1)} - \theta^{(k)}$ .

• The algorithm is initialized with  $\theta^{(0)}$  and  $H_0^{-1}$  where  $H_0^{-1}$  is picked at the identity matrix, an appropriately scale diagonal matrix or an approximate values of  $H\left(\theta^{(0)}\right)^{-1}$ .

- The algorithm is initialized with  $\theta^{(0)}$  and  $H_0^{-1}$  where  $H_0^{-1}$  is picked at the identity matrix, an appropriately scale diagonal matrix or an approximate values of  $H(\theta^{(0)})^{-1}$ .
- In situations where  $g(\theta)$  is also difficult to compute, it can be computed approximately numerically.

- The algorithm is initialized with  $\theta^{(0)}$  and  $H_0^{-1}$  where  $H_0^{-1}$  is picked at the identity matrix, an appropriately scale diagonal matrix or an approximate values of  $H(\theta^{(0)})^{-1}$ .
- In situations where  $g(\theta)$  is also difficult to compute, it can be computed approximately numerically.
- *Example*: Consider  $X_i \stackrel{\text{i.i.d.}}{\sim} f(x | \mu, \tau^2)$  where

$$f(x|\mu, \tau^2) = \frac{1}{\pi} \frac{\tau}{(y-\mu)^2 + \tau^2}$$

so for *n* observations

$$\frac{\partial l(\mu, \tau^2)}{\partial \mu} = 2\sum_{i=1}^{N} \frac{X_i - \mu}{(X_i - \mu)^2 + \tau^2}$$
$$\frac{\partial l(\mu, \tau^2)}{\partial \tau^2} = \frac{n}{2\tau^2} - \sum_{i=1}^{N} \frac{1}{(X_i - \mu)^2 + \tau^2}$$

• Numerical results for  $\theta^{(0)} = \left(\mu^{(0)}, \tau^{2(0)}\right)^{\mathsf{T}} = \left(0, 1\right)^{\mathsf{T}}$ .

k	$\mu^{(k)}$	$\tau^{2(k)}$	$-I\left(\theta^{(k)}\right)$	$\frac{\partial I\left(\theta^{(k)}\right)}{\partial \mu}$	$\frac{\partial l\left(\theta^{(k)}\right)}{\partial \tau^2}$
0	0.000	1.000	48.126	-0.839	-1.098
1	0.231	1.302	47.874	0.188	-0.144
2	0.180	1.357	47.865	-0.046	-0.040
3	0.189	1.379	47.865	0.002	-0.001
4	0.189	1.380	47.865	-0.001	0.000
5	0.189	1.380	47.865	0.000	0.000

• Numerical results for  $\theta^{(0)} = \left(\mu^{(0)}, \tau^{2(0)}\right)^{\mathsf{T}} = \left(0, 1\right)^{\mathsf{T}}$ .

k	$\mu^{(k)}$	$\tau^{2(k)}$	$-I\left(\theta^{(k)}\right)$	$\frac{\partial I\left(\theta^{(k)}\right)}{\partial \mu}$	$\frac{\partial l\left(\theta^{(k)}\right)}{\partial \tau^2}$
0	0.000	1.000	48.126	-0.839	-1.098
1	0.231	1.302	47.874	0.188	-0.144
2	0.180	1.357	47.865	-0.046	-0.040
3	0.189	1.379	47.865	0.002	-0.001
4	0.189	1.380	47.865	-0.001	0.000
5	0.189	1.380	47.865	0.000	0.000

• The quasi-Newton method converges in 5 iterations here.

• Numerical results for  $\theta^{(0)} = \left(\mu^{(0)}, \tau^{2(0)}\right)^{\mathsf{T}} = \left(0, 1\right)^{\mathsf{T}}$ .

k	$\mu^{(k)}$	$\tau^{2(k)}$	$-I\left(\theta^{(k)}\right)$	$\frac{\partial l(\theta^{(k)})}{\partial \mu}$	$\frac{\partial l\left(\theta^{(k)}\right)}{\partial \tau^2}$
0	0.000	1.000	48.126	-0.839	-1.098
1	0.231	1.302	47.874	0.188	-0.144
2	0.180	1.357	47.865	-0.046	-0.040
3	0.189	1.379	47.865	0.002	-0.001
4	0.189	1.380	47.865	-0.001	0.000
5	0.189	1.380	47.865	0.000	0.000

- The quasi-Newton method converges in 5 iterations here.
- Many statisticians are not big fans of such methods and prefer using the Expectation-Maximization algorithm. (Note: statisticians should follow optimization courses).

 Although the EM algorithm does not apply to all models, it is powerful and elegant: one of the most popular algorithms in statistics.

- Although the EM algorithm does not apply to all models, it is powerful and elegant: one of the most popular algorithms in statistics.
- It is well-suited to so-called missing data problems where you are interested in maximizing with respect to θ the likelihood function

$$L\left(\left.\theta\right|\mathbf{y}\right) = g\left(\left.\mathbf{y}\right|\theta\right)$$

where

$$g(\mathbf{y}|\theta) = \int f(\mathbf{y}, \mathbf{x}|\theta) d\mathbf{x}.$$

- Although the EM algorithm does not apply to all models, it is powerful and elegant: one of the most popular algorithms in statistics.
- It is well-suited to so-called missing data problems where you are interested in maximizing with respect to θ the likelihood function

$$L\left(\left.\theta\right|\mathbf{y}\right) = g\left(\left.\mathbf{y}\right|\theta\right)$$

where

$$g\left( \left. \mathbf{y} \right| \mathbf{ heta} 
ight) = \int f\left( \left. \mathbf{y}, \mathbf{x} \right| \mathbf{ heta} 
ight) d\mathbf{x}.$$

We call

- Although the EM algorithm does not apply to all models, it is powerful and elegant: one of the most popular algorithms in statistics.
- It is well-suited to so-called missing data problems where you are interested in maximizing with respect to θ the likelihood function

$$L\left(\left.\theta\right|\mathbf{y}\right) = g\left(\left.\mathbf{y}\right|\theta\right)$$

where

$$g(\mathbf{y}|\theta) = \int f(\mathbf{y}, \mathbf{x}|\theta) d\mathbf{x}.$$

We call

• Y the incomplete data (i.e. the observed data) and X the missing data or the augmented data

- Although the EM algorithm does not apply to all models, it is powerful and elegant: one of the most popular algorithms in statistics.
- It is well-suited to so-called missing data problems where you are interested in maximizing with respect to θ the likelihood function

$$L\left(\left.\theta\right|\mathbf{y}\right) = g\left(\left.\mathbf{y}\right|\theta\right)$$

where

$$g(\mathbf{y}|\theta) = \int f(\mathbf{y}, \mathbf{x}|\theta) d\mathbf{x}.$$

We call

- Y the incomplete data (i.e. the observed data) and X the missing data or the augmented data
- (X, Y) the complete data (although remember that X is not observed!).

- Although the EM algorithm does not apply to all models, it is powerful and elegant: one of the most popular algorithms in statistics.
- It is well-suited to so-called missing data problems where you are interested in maximizing with respect to θ the likelihood function

$$L\left(\left.\theta\right|\mathbf{y}\right) = g\left(\left.\mathbf{y}\right|\theta\right)$$

where

$$g(\mathbf{y}|\theta) = \int f(\mathbf{y}, \mathbf{x}|\theta) d\mathbf{x}.$$

We call

- Y the incomplete data (i.e. the observed data) and X the missing data or the augmented data
- (X, Y) the complete data (although remember that X is not observed!).
- More generally, we can have  $\mathbf{Y} = h(\mathbf{X})$ .

• We introduce the function  $Q: \Theta imes \Theta 
ightarrow \mathbb{R}$ 

$$Q\left(\theta, \theta'\right) = \mathbb{E}\left[\left|\log f\left(\left.\mathbf{y}, \mathbf{x} \right| \theta\right)\right| \mathbf{y}, \theta'\right] = \int \log f\left(\left.\mathbf{y}, \mathbf{x} \right| \theta\right) \ . \ f\left(\left.\mathbf{x} \right| \mathbf{y}, \theta'\right) d\mathbf{x}.$$

< ロ > < 同 > < 三 > < 三

• We introduce the function  $Q: \Theta imes \Theta 
ightarrow \mathbb{R}$ 

$$Q\left(\theta,\theta'\right) = \mathbb{E}\left[\left|\log f\left(\left|\mathbf{y},\mathbf{x}\right|\theta\right)\right|\mathbf{y},\theta'\right] = \int \log f\left(\left|\mathbf{y},\mathbf{x}\right|\theta\right) \ . \ f\left(\left|\mathbf{x}\right|\mathbf{y},\theta'\right)d\mathbf{x}.$$

• The EM algorithm is an iterative algorithm defined at iteration k + 1 by

$$\theta^{(k+1)} = \underset{\theta \in \Theta}{\operatorname{arg\,max}} Q\left(\theta, \theta^{(k)}\right).$$
(2)

• We introduce the function  $Q: \Theta imes \Theta o \mathbb{R}$ 

$$Q\left(\theta,\theta'\right) = \mathbb{E}\left[\log f\left(\left|\mathbf{y},\mathbf{x}\right|\theta\right)\right|\mathbf{y},\theta'\right] = \int \log f\left(\left|\mathbf{y},\mathbf{x}\right|\theta\right) \ . \ f\left(\left|\mathbf{x}\right|\mathbf{y},\theta'\right)d\mathbf{x}.$$

 The EM algorithm is an iterative algorithm defined at iteration k + 1 by

$$\theta^{(k+1)} = \underset{\theta \in \Theta}{\operatorname{arg\,max}} Q\left(\theta, \theta^{(k)}\right).$$
(2)

• **Theorem**: The sequence  $\left\{ \theta^{(k)} \right\}$  defined by (2) satisfies

$$L\left(\left.\theta^{(k+1)}\right|\mathbf{y}\right) \geq L\left(\left.\theta^{(k)}\right|\mathbf{y}\right).$$

• Example: Let

$$Y_i \stackrel{\text{i.i.d.}}{\sim} \theta m(y) + (1 - \theta) h(y)$$

where  $\theta \in [0, 1]$  is unknown whereas m(y) and h(y) are known pdfs.

• Example: Let

$$Y_i \stackrel{\text{i.i.d.}}{\sim} \theta m(y) + (1 - \theta) h(y)$$

where  $\theta \in [0,1]$  is unknown whereas m(y) and h(y) are known pdfs. • We have

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta) = \prod_{i=1}^{n} (\theta m(y_i) + (1 - \theta) h(y_i)).$$

• Example: Let

$$Y_i \stackrel{\text{i.i.d.}}{\sim} \theta m(y) + (1 - \theta) h(y)$$

where  $\theta \in [0, 1]$  is unknown whereas m(y) and h(y) are known pdfs. • We have

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta) = \prod_{i=1}^{n} (\theta m(y_i) + (1 - \theta) h(y_i)).$$

 Associate to each data Y<sub>i</sub> the latent variable X<sub>i</sub> ∈ {1, 2} which indicates from which distribution Y<sub>i</sub> has been drawn, i.e.

$$Y_i | X_i = 1 \sim m\left(y
ight)$$
 and  $Y_i | X_i = 2 \sim h\left(y
ight)$ 

and  $\Pr(X_i = 1) = 1 - \Pr(X_i = 2) = \theta$ . Thus we have

$$f(\mathbf{y}, \mathbf{x} | \theta) = \prod_{\substack{\{i:x_i=1\}\\ q_i = 1\}}} \theta m(y_i) \prod_{\substack{\{i:x_i=2\}\\ \{i:x_i=1\}}} (1-\theta) h(y_i)$$

where  $n_j = \sum_{i=1}^n \delta_j(x_i)$  with  $\delta_b(a) = 1$  if a = b and zero otherwise.

• From the expression of  $f(\mathbf{y}, \mathbf{x} | \theta)$ , it follows that

$$f\left(\mathbf{x}|\mathbf{y},\theta'\right) = \prod_{i=1}^{n} f\left(x_{i}|y_{i},\theta'\right)$$

where

$$\begin{split} f\left(x_{i}=1|\,y_{i},\theta'\right) &= 1-f\left(x_{i}=2|\,y_{i},\theta'\right) \\ &= \frac{f\left(x_{i}=1,\,y_{i},\theta'\right)}{f\left(x_{i}=1,\,y_{i},\theta'\right)+f\left(x_{i}=2,\,y_{i},\theta'\right)} \\ &= \frac{\theta'm\left(y_{i}\right)}{\theta'm\left(y_{i}\right)+\left(1-\theta'\right)h\left(y_{i}\right)}. \end{split}$$

$$\begin{aligned} Q\left(\theta,\theta'\right) &= \sum_{\mathbf{x}\in\{1,2\}^n} \log f\left(\mathbf{x},\mathbf{y}|\theta\right) \cdot f\left(\mathbf{x}|\mathbf{y},\theta'\right) \\ &= \sum_{\mathbf{x}\in\{1,2\}^n} \left[ \left( \sum_{i=1}^n \delta_1\left(x_i\right) \left(\log\theta + \log m\left(y_i\right)\right) \right) \\ &+ \left( \sum_{i=1}^n \delta_2\left(x_i\right) \right) \left(\log\left(1-\theta\right) + \log h\left(y_i\right)\right) \right] \prod_{i=1}^n f\left(x_i|y_i,\theta'\right) \\ &= \sum_{i=1}^n f\left(x_i = 1|y_i,\theta'\right) \left(\log\theta + \log m\left(y_i\right)\right) \\ &+ \sum_{i=1}^n f\left(x_i = 2|y_i,\theta'\right) \left(\log\left(1-\theta\right) + \log h\left(y_i\right)\right). \end{aligned}$$

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ 三重 - のへの

• We have

$$heta^{(k+1)} = rg\max_{ heta \in \Theta} Q\left( heta, heta^{(k)}
ight).$$

メロト メポト メヨト メヨト



• We have

$$heta^{(k+1)} = rg\max_{ heta \in \Theta} Q\left( heta, heta^{(k)}
ight).$$

• We have

$$\frac{dQ\left(\theta,\theta^{(k)}\right)}{d\theta} = \frac{\sum_{i=1}^{n} f\left(x_{i}=1|y_{i},\theta^{(k)}\right)}{\theta} - \frac{\sum_{i=1}^{n} f\left(x_{i}=2|y_{i},\theta^{(k)}\right)}{1-\theta}.$$
  
By solving  $\frac{dQ\left(\theta,\theta^{(k)}\right)}{d\theta} = 0$  we obtain

$$\theta^{(k+1)} = \frac{\sum_{i=1}^{n} f\left(x_{i} = 1 | y_{i}, \theta^{(k)}\right)}{\sum_{i=1}^{n} f\left(x_{i} = 1 | y_{i}, \theta^{(k)}\right) + \sum_{i=1}^{n} f\left(x_{i} = 2 | y_{i}, \theta^{(k)}\right)}$$

$$= \frac{\sum_{i=1}^{n} \theta^{(k)} m(y_{i})}{\sum_{i=1}^{n} \theta^{(k)} m(y_{i}) + \sum_{i=1}^{n} \left(1 - \theta^{(k)}\right) h(y_{i})}.$$

メロト メポト メヨト メヨト

We have

$$heta^{(k+1)} = rg\max_{ heta \in \Theta} Q\left( heta, heta^{(k)}
ight).$$

• We have

\_

$$\frac{dQ\left(\theta,\theta^{(k)}\right)}{d\theta} = \frac{\sum_{i=1}^{n} f\left(x_{i}=1|y_{i},\theta^{(k)}\right)}{\theta} - \frac{\sum_{i=1}^{n} f\left(x_{i}=2|y_{i},\theta^{(k)}\right)}{1-\theta}.$$

By solving  $\frac{dQ(0,0)}{d\theta} = 0$  we obtain

$$\theta^{(k+1)} = \frac{\sum_{i=1}^{n} f\left(x_{i} = 1 | y_{i}, \theta^{(k)}\right)}{\sum_{i=1}^{n} f\left(x_{i} = 1 | y_{i}, \theta^{(k)}\right) + \sum_{i=1}^{n} f\left(x_{i} = 2 | y_{i}, \theta^{(k)}\right)} \\ = \frac{\sum_{i=1}^{n} \theta^{(k)} m(y_{i})}{\sum_{i=1}^{n} \theta^{(k)} m(y_{i}) + \sum_{i=1}^{n} \left(1 - \theta^{(k)}\right) h(y_{i})}.$$

• This updating equation is simple to implement and it is guaranteed that  $L\left(\theta^{(k+1)} \middle| \mathbf{y}\right) \geq L\left(\theta^{(k)} \middle| \mathbf{y}\right)$ .

• Finite mixture of Gaussians

$$Y_i \overset{\text{i.i.d.}}{\sim} \sum_{k=1}^p \pi_k \mathcal{N}\left(\mu_k, \sigma_k^2\right)$$

where  $\pi_k \geq 0$  and  $\sum_{k=1}^{p} \pi_k = 1$ .

イロト イヨト イヨト

Finite mixture of Gaussians

$$Y_{i} \stackrel{\text{i.i.d.}}{\sim} \sum_{k=1}^{p} \pi_{k} \mathcal{N}\left(\mu_{k}, \sigma_{k}^{2}\right)$$

where  $\pi_k \geq 0$  and  $\sum_{k=1}^{p} \pi_k = 1$ .

• This simple mixture model is widely used in practice for density estimation and clustering.

Finite mixture of Gaussians

$$Y_{i} \overset{\text{i.i.d.}}{\sim} \sum_{k=1}^{p} \pi_{k} \mathcal{N}\left(\mu_{k}, \sigma_{k}^{2}\right)$$

where  $\pi_k \geq 0$  and  $\sum_{k=1}^{p} \pi_k = 1$ .

- This simple mixture model is widely used in practice for density estimation and clustering.
- Given *n* observations  $\mathbf{Y} = \mathbf{y} = (y_1, ..., y_n)$ , we are interested in the MLE estimate of  $\theta = (\pi_1, ..., \pi_p, \mu_1, ..., \mu_p, \sigma_1^2, ..., \sigma_p^2)$  where

$$g\left(\mathbf{y}|\theta\right) = \prod_{i=1}^{n} \left( \sum_{k=1}^{p} \frac{1}{\sqrt{2\pi\sigma_{k}^{2}}} \exp\left(-\frac{\left(y_{i}-\mu_{k}\right)^{2}}{2\sigma_{k}^{2}}\right) \right)$$

• Similarly to the previous lecture, we associate to each  $Y_i$  a latent variable  $X_i \in \{1, ..., p\}$  such that

$$\Pr\left(X_i=j\right) = \pi_j \text{ and } |Y_i| \left(X_i=j\right) \sim \mathcal{N}\left(\mu_j, \sigma_j^2\right).$$

It follows that

$$f(\mathbf{x}, \mathbf{y} | \theta) = \prod_{i=1}^{n} \frac{\pi_{x_i}}{\sqrt{2\pi\sigma_{x_i}^2}} \exp\left(-\frac{\left(y_i - \mu_{x_i}\right)^2}{2\sigma_{x_i}^2}\right)$$

and

$$f(\mathbf{x}|\theta,\mathbf{y}) = \prod_{i=1}^{n} f(x_i|\theta, y_i)$$

where

$$f(x_i = j | \theta, y_i) = \frac{\pi_j / \sigma_j \exp\left(-\left(y_i - \mu_j\right)^2 / \left(2\sigma_j^2\right)\right)}{\sum_{m=1}^p \pi_m / \sigma_m \exp\left(-\left(y_i - \mu_m\right)^2 / \left(2\sigma_m^2\right)\right)}.$$

#### • We have

$$Q\left(\theta, \theta^{(k)}\right) = \sum_{\mathbf{x}} \log f\left(\mathbf{x}, \mathbf{y} | \theta\right) \cdot f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right)$$
  
=  $cst + \sum_{\mathbf{x}} \sum_{i=1}^{n} \left(\log \pi_{x_{i}} - \frac{1}{2} \log \sigma_{x_{i}}^{2} - \frac{\left(y_{i} - \mu_{x_{i}}\right)^{2}}{2\sigma_{x_{i}}^{2}}\right) \prod_{i=1}^{n} f\left(x_{i} | \theta^{(k)}, y_{i}\right)$   
=  $cst + \sum_{m=1}^{p} \left[ \left(\log \pi_{m} - \frac{1}{2} \log \sigma_{m}^{2}\right) \left(\sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)\right) - \sum_{i=1}^{n} \frac{\left(y_{i} - \mu_{m}\right)^{2}}{2\sigma_{m}^{2}} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right) \right].$ 

• We have to maximize  $Q\left(\theta, \theta^{(k)}\right)$  over a constrained set, i.e.  $\sum_{k=1}^{p} \pi_{k} = 1$  . We introduce a Lagrange multiplier  $\lambda$  and maxim

 $\sum_{m=1}^{p}\pi_{m}=1.$  We introduce a Lagrange multiplier  $\lambda$  and maximize instead

$$Q\left( heta, heta^{(k)}
ight) + \lambda\left(\sum_{m=1}^{p}\pi_m - 1
ight).$$

• We have to maximize  $Q\left(\theta, \theta^{(k)}\right)$  over a constrained set, i.e.  $\sum_{m=1}^{p} \pi_m = 1$ . We introduce a Lagrange multiplier  $\lambda$  and maximize instead

$$Q\left(\theta,\theta^{(k)}\right) + \lambda\left(\sum_{m=1}^{p}\pi_m - 1\right).$$

We have

$$\frac{\frac{\partial Q\left(\theta,\theta^{(k)}\right)}{\partial \sigma_{m}^{2}} = -\frac{\left(\sum_{i=1}^{n} f\left(x_{i}=m|\theta^{(k)},y_{i}\right)\right)}{2\sigma_{m}^{2}} + \sum_{i=1}^{n} \frac{\left(y_{i}-\mu_{m}\right)^{2}}{2\sigma_{m}^{4}} f\left(x_{i}=m|\theta^{(k)},y_{i}\right)$$

$$2\sigma_{m}^{2} \frac{\partial Q\left(\theta,\theta^{(k)}\right)}{\partial \mu_{m}} = \mu_{m} \sum_{i=1}^{n} f\left(x_{i}=m|\theta^{(k)},y_{i}\right)$$

$$-\sum_{i=1}^{n} y_{i} f\left(x_{i}=m|\theta^{(k)},y_{i}\right)$$

• It follows that

$$\mu_{m}^{(k+1)} = \frac{\sum_{i=1}^{n} y_{i} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)}{\sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)},$$
  
$$\sigma_{m}^{2(k+1)} = \frac{\sum_{i=1}^{n} \left(y_{i} - \mu_{m}^{(k+1)}\right)^{2} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)}{\sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)}.$$

・ロト ・聞 ト ・ ヨト ・ ヨト

• It follows that

$$\mu_{m}^{(k+1)} = \frac{\sum_{i=1}^{n} y_{i} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)}{\sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)},$$
  
$$\sigma_{m}^{2(k+1)} = \frac{\sum_{i=1}^{n} \left(y_{i} - \mu_{m}^{(k+1)}\right)^{2} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)}{\sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)}.$$

• We have

$$\frac{\partial Q\left(\theta, \theta^{(k)}\right)}{\partial \pi_{m}} = \frac{\left(\sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)\right)}{\pi_{m}} + \lambda,$$
  
$$\frac{\partial Q\left(\theta, \theta^{(k)}\right)}{\partial \lambda} = \sum_{m=1}^{p} \pi_{m} - 1 \Rightarrow \pi_{m}^{(k+1)} = \frac{1}{N} \sum_{i=1}^{n} f\left(x_{i} = m | \theta^{(k)}, y_{i}\right)$$

▲ □ ▶ < □ ▶ < □</p>

• Although the EM is widely used to compute the MLE for finite mixture of Gaussians, we have

$$g(\mathbf{y}|\theta) \to \infty$$
 when  $\exists \sigma_i \to 0$ .

 Although the EM is widely used to compute the MLE for finite mixture of Gaussians, we have

$$g(\mathbf{y}|\theta) \to \infty$$
 when  $\exists \sigma_i \to 0$ .

Hence, the true MLE estimate does not give sensible results. The EM is used to find a sensible local maximum of g (y | θ).

• Although the EM is widely used to compute the MLE for finite mixture of Gaussians, we have

$$g(\mathbf{y}|\theta) \to \infty$$
 when  $\exists \sigma_i \to \mathbf{0}$ .

- Hence, the true MLE estimate does not give sensible results. The EM is used to find a sensible local maximum of g (y | θ).
- A way to circumvent this problem consists of introducing a prior distribution  $p(\theta)$  on  $\theta$  and to use the EM to maximize  $g(\mathbf{y}|\theta) p(\theta)$ . We will discuss this Bayesian approach later.

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta)$$
 where  $g(\mathbf{y} | \theta)$  follows from  $f(\mathbf{z} | \theta)$ .

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta)$$
 where  $g(\mathbf{y} | \theta)$  follows from  $f(\mathbf{z} | \theta)$ .

We call

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta)$$
 where  $g(\mathbf{y} | \theta)$  follows from  $f(\mathbf{z} | \theta)$ .

- We call
  - Y the incomplete data (i.e. the observed data) and Z are the complete data.

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta)$$
 where  $g(\mathbf{y} | \theta)$  follows from  $f(\mathbf{z} | \theta)$ .

- We call
  - Y the incomplete data (i.e. the observed data) and Z are the complete data.
  - We can have Y as a subset of Z, i.e. Z = (X, Y), or Y = h(Z) where h is a many-to-one mapping.

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta)$$
 where  $g(\mathbf{y} | \theta)$  follows from  $f(\mathbf{z} | \theta)$ .

- We call
  - Y the incomplete data (i.e. the observed data) and Z are the complete data.
  - We can have Y as a subset of Z, i.e. Z = (X, Y), or Y = h(Z) where h is a many-to-one mapping.
- The EM proceeds as follows

$$\theta^{(k+1)} = \operatorname*{arg\,max}_{\theta \in \Theta} \int \log f\left(\left. \mathbf{z} \right| \theta\right) \ . \ f\left(\left. \mathbf{z} \right| \mathbf{y}, \theta^{(k)} \right) d\mathbf{z}$$

• Example: Consider the following genetic linkage model where observations

$$(Y_1, Y_2, Y_3, Y_4) \sim \mathcal{M}\left(n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4}\right).$$

Image: Image:

3

• Example: Consider the following genetic linkage model where observations

$$(Y_1, Y_2, Y_3, Y_4) \sim \mathcal{M}\left(n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4}\right).$$

• The observed likelihood function is given by

$$g\left(\mathbf{y}\right|\theta) \propto (2+\theta)^{y_1} \left(1-\theta\right)^{y_2+y_3} \theta^{y_4}$$

• Example: Consider the following genetic linkage model where observations

$$(Y_1, Y_2, Y_3, Y_4) \sim \mathcal{M}\left(n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4}\right).$$

The observed likelihood function is given by

$$g(\mathbf{y}|\theta) \propto (2+\theta)^{y_1} (1-\theta)^{y_2+y_3} \theta^{y_4}$$

• Introduce the artificial missing data  $(X_1, X_2)$  such that  $Y_1 = X_1 + X_2$ and define

$$Z = (Z_1, ..., Z_5) = (X_1, X_2, Y_2, Y_3, Y_4)$$
  
 
$$\sim \mathcal{M}\left(n; \frac{1}{2}, \frac{\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4}\right)$$

In this case we have  $\mathbf{Y} = h(\mathbf{Z})$ ; i.e.  $(Y_1, Y_2, Y_3, Y_4) = (Z_1 + Z_2, Z_3, Z_4, Z_5)$ . This equation defines  $f(\mathbf{z}|\theta) \propto (1-\theta)^{y_2+y_3} \theta^{x_2+y_4}$ 

and clearly we have

$$\sum_{z_1,z_2} f(\mathbf{z}|\theta) \,\delta_{z_1+z_2}(y_1) \,\delta_{z_3}(y_2) \,\delta_{z_4}(y_3) \,\delta_{z_5}(y_4) = g(\mathbf{y}|\theta), \quad \mathbb{R}$$

January 2008

• We have 
$$f(\mathbf{z}|\mathbf{y}, \theta') = f(x_1, x_2|\mathbf{y}, \theta') \delta_{y_2}(z_3) \delta_{y_3}(z_4) \delta_{y_4}(z_5)$$
 where  

$$\begin{aligned} f(x_1, x_2|\mathbf{y}, \theta') &= f(x_1, x_2|y_1, \theta') \\ &= \mathcal{M}\left(y_1; \frac{\frac{1}{2}}{\frac{1}{2} + \frac{\theta'}{4}}, \frac{\frac{\theta'}{4}}{\frac{1}{2} + \frac{\theta'}{4}}\right).\end{aligned}$$

Now we have

$$\begin{aligned} Q\left(\theta,\theta^{(k)}\right) &= \sum_{x_1,x_2} \left[ cst + (y_2 + y_3) \log\left(1 - \theta\right) + (x_2 + y_4) \log\theta \right] \\ &\times f\left(x_1, x_2 | y_1, \theta^{(k)}\right) \\ &= cst + (y_2 + y_3) \log\left(1 - \theta\right) + \left(\mathbb{E}\left(X_2 | y_1, \theta^{(k)}\right) + y_4\right) \log\theta \end{aligned}$$
where  $\mathbb{E}\left(X_2 | y_1, \theta^{(k)}\right) &= y_1 \frac{\theta^{(k)}}{2 + \theta^{(k)}}$ , thus
$$\frac{dQ(\theta, \theta^{(k)})}{d\theta} &= 0 \Leftrightarrow -\frac{(y_2 + y_3)}{(1 - \theta)} + \frac{y_1 \frac{\theta^{(k)}}{2 + \theta^{(k)}} + y_4}{\theta} = 0 \\ &\Rightarrow \theta^{(k+1)} = \left(y_1 \frac{\theta^{(k)}}{2 + \theta^{(k)}} + y_2 + y_3 + y_4\right)^{-1} \left(y_1 \frac{\theta^{(k)}}{2 + \theta^{(k)}} + y_4\right). \end{aligned}$$

### Proof of Theorem for Expectation-Maximization Algorithm

• We want to show that  $L\left(\left.\theta^{(k+1)}\right|\mathbf{y}\right) \geq L\left(\left.\theta^{(k)}\right|\mathbf{y}\right)$  for  $\theta^{(k+1)} = \underset{\theta \in \Theta}{\arg \max} Q\left(\theta, \theta^{(k)}\right)$ .

#### Proof of Theorem for Expectation-Maximization Algorithm

• We want to show that  $L\left(\left.\theta^{(k+1)}\right|\mathbf{y}\right) \ge L\left(\left.\theta^{(k)}\right|\mathbf{y}\right)$  for  $\theta^{(k+1)} = \underset{\substack{\theta \in \Theta\\ \theta \in \Theta}}{\operatorname{arg\,max}} Q\left(\theta, \theta^{(k)}\right)$ .

$$f(\mathbf{x}|\theta, \mathbf{y}) = \frac{f(\mathbf{x}, \mathbf{y}|\theta)}{g(\mathbf{y}|\theta)} \Leftrightarrow g(\mathbf{y}|\theta) = L(\theta|\mathbf{y}) = \frac{f(\mathbf{x}, \mathbf{y}|\theta)}{f(\mathbf{x}|\theta, \mathbf{y})}$$

thus

$$\log L\left(\left.\theta\right|\mathbf{y}\right) = \log f\left(\left.\mathbf{x},\mathbf{y}\right|\theta\right) - \log f\left(\left.\mathbf{x}\right|\theta,\mathbf{y}\right)$$
  
and for any value  $\theta^{(k)}$ 

$$\log L(\theta | \mathbf{y}) = \underbrace{\int \log f(\mathbf{x}, \mathbf{y} | \theta) . f(\mathbf{x} | \theta^{(k)}, \mathbf{y}) d\mathbf{x}}_{=Q(\theta, \theta^{(k)})} - \int \log f(\mathbf{x} | \theta, \mathbf{y}) . f(\mathbf{x} | \theta^{(k)}, \mathbf{y}) d\mathbf{x}.$$

• Now the EM ensures by construction that  

$$Q\left(\theta^{(k+1)}, \theta^{(k)}\right) \ge Q\left(\theta^{(k)}, \theta^{(k)}\right).$$
So if we can prove that  

$$\int \log f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right) . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$$

$$\le \int \log f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$$

then the Theorem is proved.

• Now the EM ensures by construction that  

$$Q\left(\theta^{(k+1)}, \theta^{(k)}\right) \ge Q\left(\theta^{(k)}, \theta^{(k)}\right).$$
So if we can prove that  

$$\int \log f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right) . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$$

$$\le \int \log f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$$

then the Theorem is proved.

• We have thanks to Jensen's inequality

$$\int \log \frac{f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right)}{f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right)} . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$$

$$\leq \log \int \frac{f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right)}{f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right)} . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x} = 0$$

• Now the EM ensures by construction that  $Q\left(\theta^{(k+1)}, \theta^{(k)}\right) \ge Q\left(\theta^{(k)}, \theta^{(k)}\right).$ So if we can prove that  $\int \log f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right) . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$   $\le \int \log f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) . f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$ 

then the Theorem is proved.

• We have thanks to Jensen's inequality

$$\int \log \frac{f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right)}{f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right)} \cdot f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x}$$

$$\leq \log \int \frac{f\left(\mathbf{x} | \theta^{(k+1)}, \mathbf{y}\right)}{f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right)} \cdot f\left(\mathbf{x} | \theta^{(k)}, \mathbf{y}\right) d\mathbf{x} = 0$$

 There are numerous variations of the EM in the literature. We will discuss some later.