



Joint Institute of Engineering



Design and Implementation of Speech Recognition Systems

Fall 2014

Ming Li

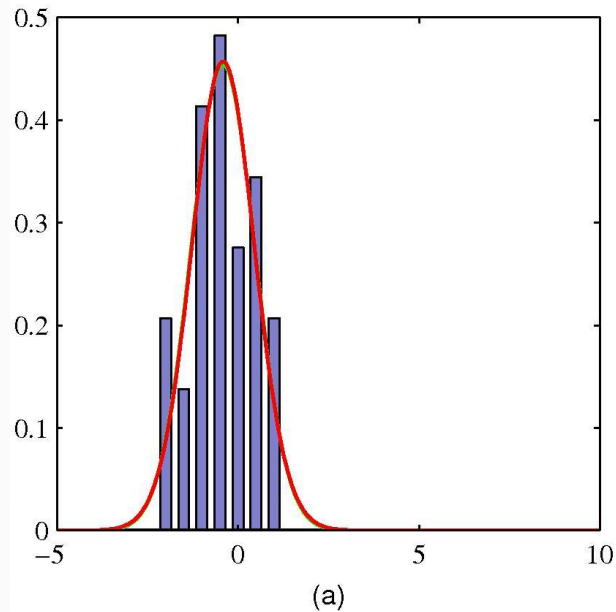
Special topic: the Expectation-Maximization algorithm and GMM

Sep 30 2014

Some graphics are from Pattern Recognition and Machine learning, Bishop, Copyright © Springer

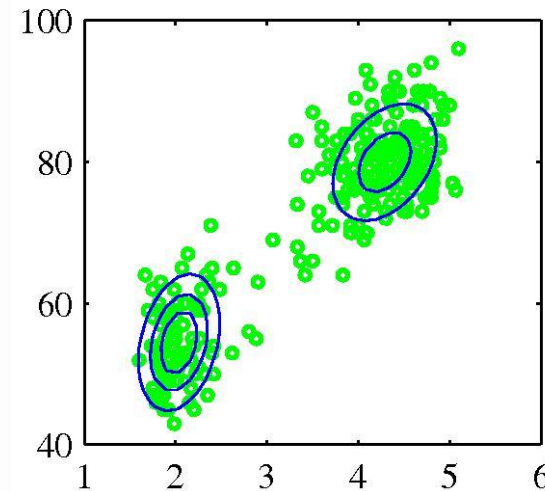
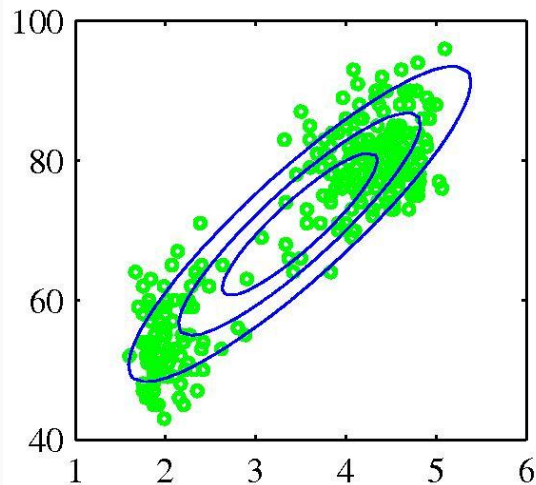
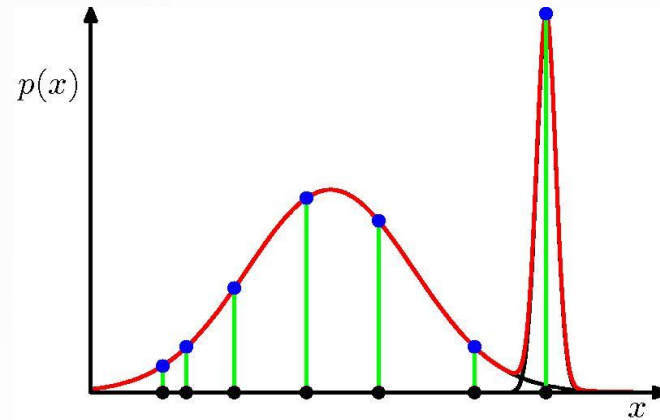
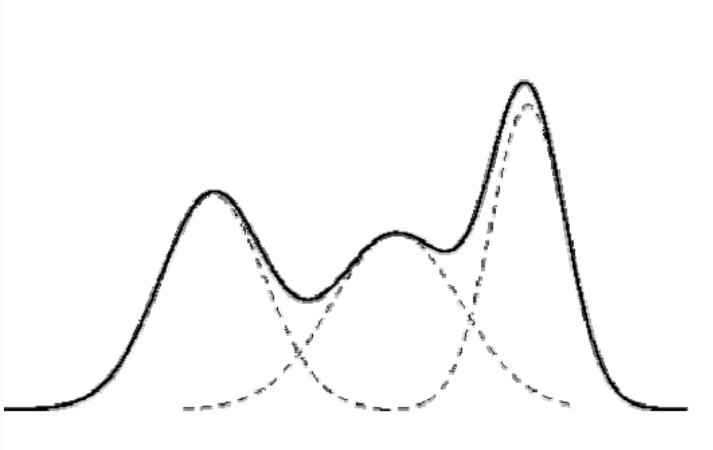
Some equations are from the slides edited by Muneem S.

Parametric density estimation



- How to estimate those parameters?
 - ML
 - MAP

How about Multimodal cases?



Mixture Models

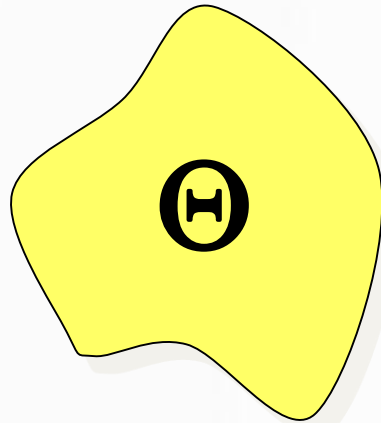
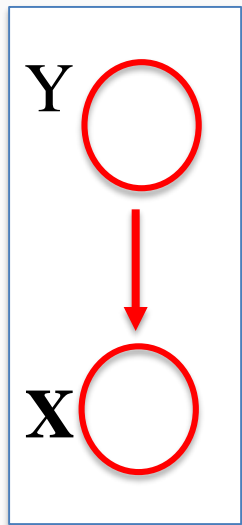
- If you believe that the data set is comprised of **several distinct populations**
- It has the following form:

$$p(\mathbf{x} | \Theta) = \sum_{j=1}^M \alpha_j p_j(\mathbf{x} | \theta_j) \quad \text{with} \quad \sum_{j=1}^M \alpha_j = 1$$

$$\Theta = (\alpha_1, \dots, \alpha_M, \theta_1, \dots, \theta_M)$$

$$p(\mathbf{x} | \Theta) = \sum_{j=1}^M \alpha_j p_j(\mathbf{x} | \theta_j)$$

Mixture Models



$$\mathbf{Y} = \{y_1, y_2, \dots, y_N\}$$

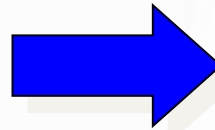
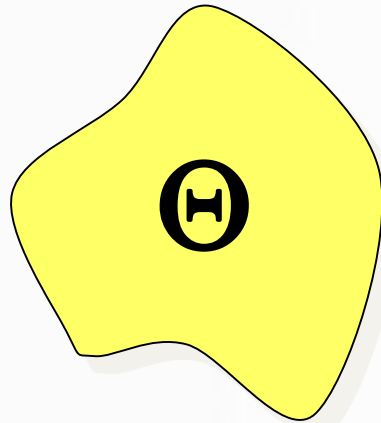
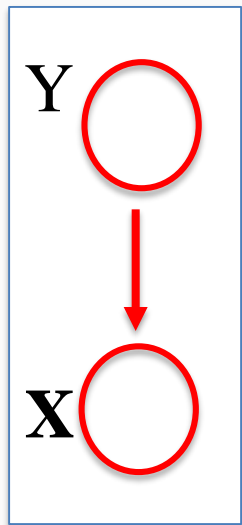


$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

Let $y_i \in \{1, \dots, M\}$ represents the source that generates the data.

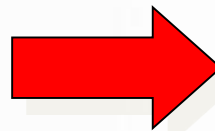
$$p(\mathbf{x}_i | \Theta) = \sum_{j=1}^M \alpha_j p_j(\mathbf{x}_i | \theta_j)$$

Mixture Models



$$y_i = j$$

$$p(y_i = j | \Theta) = \alpha_j$$



$$\mathbf{X}_i$$

$$p(\mathbf{x}_i | y = j, \Theta) = p_j(\mathbf{x}_i | \theta_j)$$

Let $y_i \in \{1, \dots, M\}$ represents the source that generates the data.

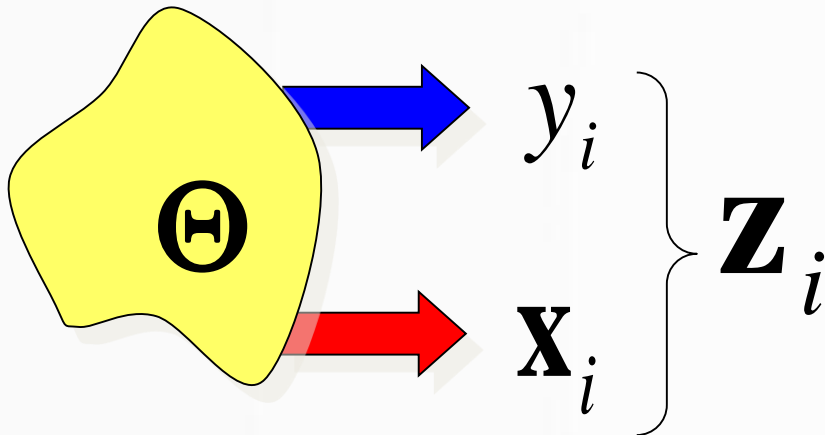
Mixture Models

$$p(\mathbf{x}_i | y = j, \Theta) = p_j(\mathbf{x}_i | \theta_j)$$

$$p(y_i = j | \Theta) = \alpha_j$$



$$p(\mathbf{x}_i | \Theta) = \sum_{y_i=1}^M p(\mathbf{x}_i, y_i | \Theta) = \sum_{y_i=1}^M p(y_i = j | \Theta) p(\mathbf{x}_i | y = j, \Theta) = \sum_{j=1}^M \alpha_j p_j(\mathbf{x}_i | \theta_j)$$



$$p(\mathbf{z}_i | \Theta) = p(\mathbf{x}_i, y_i | \Theta) = p(y_i | \mathbf{x}_i, \Theta) p(\mathbf{x}_i | \Theta)$$

Mixture Models

$$p(\mathbf{x}_i | y = j, \Theta) = p_j(\mathbf{x}_i | \theta_j)$$

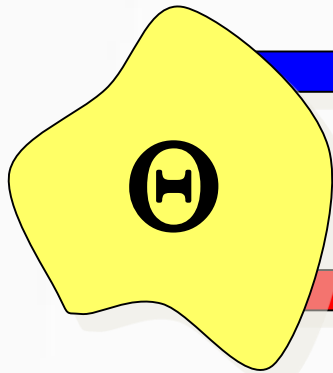
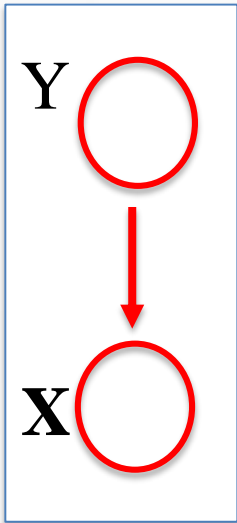
$$p(y_i = j | \Theta) = \alpha_j$$

$$p(\mathbf{z}_i | \Theta) = p(\mathbf{x}_i, y_i | \Theta) = p(y_i | \mathbf{x}_i, \Theta) p(\mathbf{x}_i | \Theta)$$

$$\begin{aligned} p(y_i | \mathbf{x}_i, \Theta) &= \frac{p(\mathbf{x}_i, y_i, \Theta)}{p(\mathbf{x}_i, \Theta)} = \frac{p(\mathbf{x}_i | y_i, \Theta) p(y_i, \Theta)}{p(\mathbf{x}_i, \Theta)} \\ &= \frac{p(\mathbf{x}_i | y_i, \Theta) p(y_i | \Theta) p(\Theta)}{p(\mathbf{x}_i | \Theta) p(\Theta)} = \frac{p(\mathbf{x}_i | y_i, \Theta) p(y_i | \Theta)}{p(\mathbf{x}_i | \Theta)} \\ &= \frac{p_{y_i}(\mathbf{x}_i | \theta_{y_i}) \alpha_{y_i}}{\sum_{j=1}^M \alpha_j p_j(\mathbf{x} | \theta_j)} \end{aligned}$$

Given \mathbf{x} and Θ , the conditional density of y can be computed.

Complete-Data Likelihood Function



$$\mathbf{y} = \{y_1, \dots, y_N\}$$

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

$$\mathbf{z}_i = (\mathbf{x}_i, y_i)$$

$$\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$$

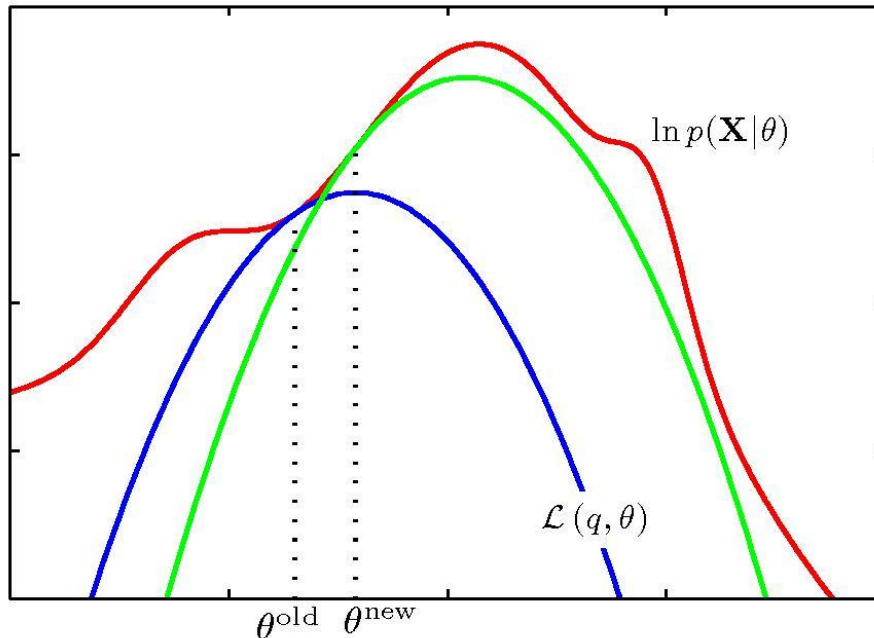
$$p(\mathbf{Z} / \Theta) = p(\mathbf{X}, \mathbf{y} / \Theta)$$

But we don't know
the latent variable y

$$= \prod_{i=1}^N \underbrace{p(\mathbf{x}_i | y_i, \Theta)}_{\theta_{y_i}} \underbrace{p(y_i)}_{\alpha_{y_i}} = \prod_{i=1}^N p_{y_i}(\mathbf{x}_i | \theta_{y_i})$$

Expectation-Maximization (1)

- If we want to find the local maxima of $f(x)$
 - Define an auxiliary function $A(x, x^t)$
 - Satisfy $f(x) \geq A(x, x^t), \forall x$ and $f(x^t) = A(x^t, x^t)$



$$x^{t+1} = \arg \max A(x, x^t)$$



$$A(x^{t+1}, x^t) \geq A(x^t, x^t) = f(x^t)$$

$$f(x^{t+1}) \geq A(x^{t+1}, x^t)$$

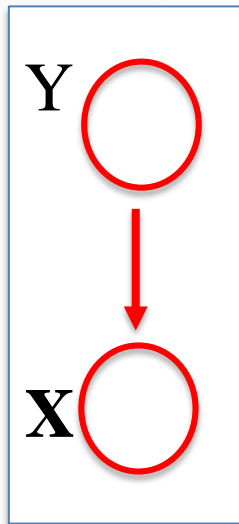


$$f(x^{t+1}) \geq f(x^t)$$

$$f(x^{t+2}) \geq f(x^{t+1})$$

Expectation-Maximization (2)

- The goal: finding ML solutions for probabilistic models with latent variables



$$\log p(\mathbf{x} | \Theta) = \log \sum_y p(\mathbf{x}, y | \Theta)$$

Difficult, auxiliary function

$$\log \sum_y q(y) \frac{p(\mathbf{x}, y | \Theta)}{q(y)} \geq \sum_y q(y) \log \left(\frac{p(\mathbf{x}, y | \Theta)}{q(y)} \right)$$

$$\sum_y q(y) = 1$$

Auxiliary function

$$\log p(\mathbf{x} | \Theta) = \log \sum_y p(\mathbf{x}, y | \Theta) \geq \sum_y q(y) \log \left(\frac{p(\mathbf{x}, y | \Theta)}{q(y)} \right)$$

Expectation-Maximization (2)

- What is the gap?

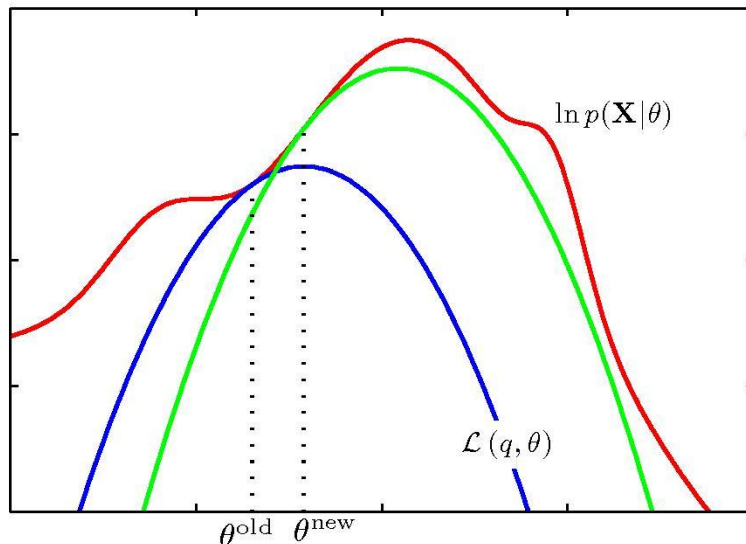
$$\begin{aligned}
 & \log p(\mathbf{x} | \Theta) - \sum_y q(y) \log \left(\frac{p(\mathbf{x}, y | \Theta)}{q(y)} \right) = \\
 & \sum_y q(y) \left[\log p(\mathbf{x} | \Theta) - \log \left(\frac{p(\mathbf{x}, y | \Theta)}{q(y)} \right) \right] = \\
 & = \sum_y q(y) \log \left(\frac{q(y) p(\mathbf{x} | \Theta)}{p(\mathbf{x}, y | \Theta)} \right) = \sum_y q(y) \log \left(\frac{q(y)}{p(y | x, \Theta)} \right) \\
 & = KL(q(y) \| p(y | x, \Theta)) \geq 0 \qquad D_{KL}(P \| Q) = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right) \\
 & KL(q(y) \| p(y | x, \Theta)) = 0 \quad \text{iff } q(y) = p(y | x, \Theta)
 \end{aligned}$$

Expectation-Maximization (2)

- We have shown that

$$\log p(\mathbf{x} | \Theta) - \underbrace{\sum_y q(y) \log\left(\frac{p(\mathbf{x}, y | \Theta)}{q(y)}\right)}_{A(\Theta, \Theta^{old})} = KL(q(y) \| p(y | x, \Theta)) \geq 0$$

$$q(y) = p(y | x, \Theta^{old})$$



$$\begin{aligned} \sum_y q(y) \log\left(\frac{p(\mathbf{x}, y | \Theta^{old})}{q(y)}\right) &= \\ \sum_y p(y | x, \Theta^{old}) \log\left(\frac{p(\mathbf{x}, y | \Theta^{old})}{p(y | x, \Theta^{old})}\right) &= \\ \sum_y p(y | x, \Theta^{old}) \log(p(\mathbf{x} | \Theta^{old})) &= \log(p(\mathbf{x} | \Theta^{old})) \end{aligned}$$

Expectation-Maximization (3)

- Any physical meaning of the **auxiliary function**?

$$A(\Theta, \Theta^{old}) = \sum_y p(y | x, \Theta^{old}) \log\left(\frac{p(\mathbf{x}, y | \Theta)}{p(y | x, \Theta^{old})}\right) =$$

$$\sum_y p(y | x, \Theta^{old}) \log(p(\mathbf{x}, y | \Theta)) - \underbrace{\sum_y p(y | x, \Theta^{old}) p(y | x, \Theta^{old})}_{\text{Independent of } \Theta}$$

$$= \sum_y p(y | x, \Theta^{old}) \log(p(\mathbf{x}, y | \Theta)) - \text{const}$$

Conditional expectation of the
complete data log likelihood

Expectation-Maximization (4)

1. Choose an initial setting for the parameters Θ^{old}
2. E step: evaluate $p(y | x, \Theta^{old})$

$$Q(\Theta, \Theta^{old}) = \sum_y p(y | x, \Theta^{old}) \log(p(\mathbf{x}, y | \Theta))$$

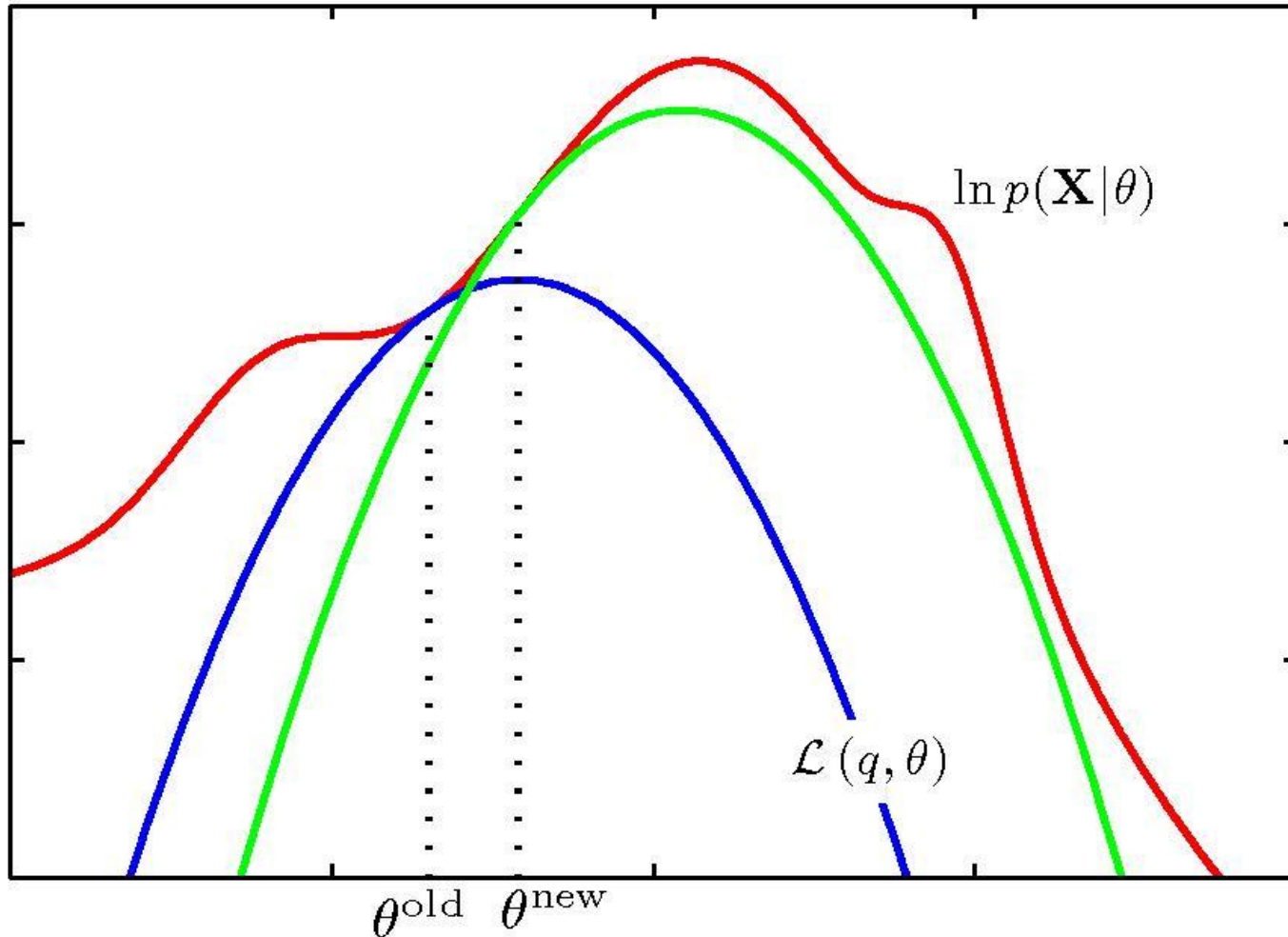
3. M step: evaluate Θ^{new} given by

$$\Theta^{new} = \arg \max_{\Theta} Q(\Theta, \Theta^{old})$$

4. Check for convergence

if not, $\Theta^{old} \leftarrow \Theta^{new}$ and return to step 2

Expectation-Maximization (4)



Guassain Mixture Model (GMM)

Guassain model of a d -dimensional source, say j :

$$p_j(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_j|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right]$$

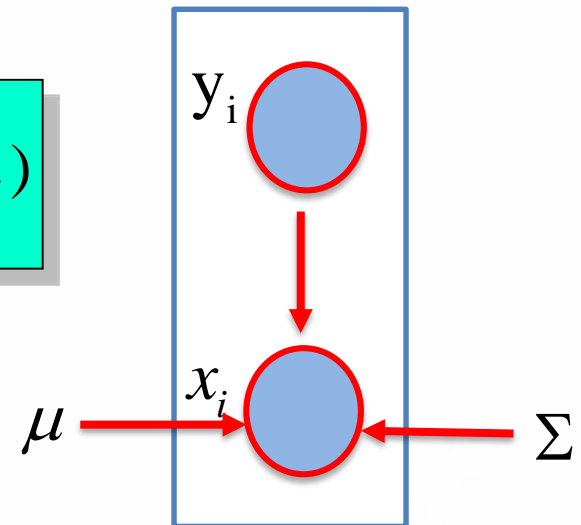
$$\theta_j = (\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

GMM with M sources:

$$p(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M) = \sum_{j=1}^M \alpha_j p_j(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

$$\alpha_j \geq 0$$

$$\sum \alpha_j = 1$$



Mixture Model

$$p(\mathbf{x} | \Theta) = \sum_{l=1}^M \alpha_l p_l(\mathbf{x} | \theta_l)$$

$$\Theta = (\alpha_1, \dots, \alpha_M, \theta_1, \dots, \theta_M)$$

Correlated
with α_l only.

subject to

Correlated
with θ_l only.

$$Q(\Theta, \Theta^{old}) = \sum_{l=1}^M p(l | \mathbf{x}, \Theta^{old}) \log(p(\mathbf{x}, y | \Theta))$$

$$Q(\Theta, \Theta^g) = \sum_{l=1}^M p(l | \mathbf{x}_i, \Theta^g) \sum_{i=1}^N \log(\alpha_l p_l(\mathbf{x}_i | \theta_l))$$

$$p(\mathbf{X}, \mathbf{y} | \Theta) = \prod_{i=1}^N \alpha_{y_i} p_{y_i}(\mathbf{x}_i | \theta_{y_i})$$

$$Q(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \sum_{l=1}^M \sum_{i=1}^N \log[p_l(\mathbf{x}_i | \theta_l)] p(l | \mathbf{x}_i, \Theta^g)$$

M step: $\Theta^{new} = \arg \max_{\Theta} Q(\Theta, \Theta^{old})$

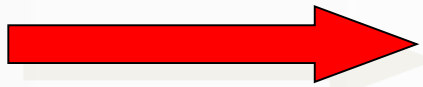


$$Q(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \sum_{l=1}^M \sum_{i=1}^N \log[p_l(\mathbf{x}_i | \theta_l)] p(l | \mathbf{x}_i, \Theta^g)$$

Finding α_l

Due to the constraint on α_l 's, we introduce *Lagrange Multiplier* λ , and solve the following equation.

$$\frac{\partial}{\partial \alpha_l} \left[\sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \lambda \left(\sum_{l=1}^M \alpha_l - 1 \right) \right] = 0, \quad l = 1, \dots, M$$



$$\sum_{i=1}^N \frac{1}{\alpha_l} p(l | \mathbf{x}_i, \Theta^g) + \lambda = 0, \quad l = 1, \dots, M$$



$$\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) + \alpha_l \lambda = 0, \quad l = 1, \dots, M$$

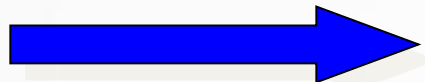


$$Q(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \sum_{l=1}^M \sum_{i=1}^N \log[p_l(\mathbf{x}_i | \theta_l)] p(l | \mathbf{x}_i, \Theta^g)$$

Finding $\alpha_l \rightarrow \lambda = -N$



$$\sum_{l=1}^M \sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) + \lambda \sum_{l=1}^M \alpha_l = 0$$



$$\underbrace{\sum_{i=1}^N \sum_{l=1}^M p(l | \mathbf{x}_i, \Theta^g)}_N + \lambda \underbrace{\sum_{l=1}^M \alpha_l}_1 = 0$$

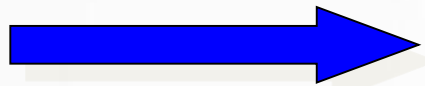


$$\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) + \alpha_l \lambda = 0, \quad l = 1, \dots, M$$



$$Q(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \sum_{l=1}^M \sum_{i=1}^N \log[p_l(\mathbf{x}_i | \theta_l)] p(l | \mathbf{x}_i, \Theta^g)$$

Finding $\alpha_l \rightarrow \lambda = -N$



$$\alpha_l = \frac{1}{N} \sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)$$

$$p(l | \mathbf{x}_i, \Theta^g) = \frac{\alpha_l^g p_l(\mathbf{x}_i | \theta_l^g)}{\sum_{j=1}^M \alpha_j^g p_j(\mathbf{x} | \theta_j^g)}$$



$$\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) + \alpha_l \lambda = 0, \quad l = 1, \dots, M$$



$$Q(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \sum_{l=1}^M \sum_{i=1}^N \log[p_l(\mathbf{x}_i | \theta_l)] p(l | \mathbf{x}_i, \Theta^g)$$

Finding θ_l Only need to maximize this term

Consider GMM

$$p_l(\mathbf{x} | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_l|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_l)^T \boldsymbol{\Sigma}_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l) \right]$$

$$\theta_l = (\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)$$

$$\log[p_l(\mathbf{x} | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)] = \underbrace{-\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_l|^{1/2}}_{\text{unrelated}} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_l)^T \boldsymbol{\Sigma}_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l)$$



$$Q(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \log(\alpha_l) p(l | \mathbf{x}_i, \Theta^g) + \sum_{l=1}^M \sum_{i=1}^N \log[p_l(\mathbf{x}_i | \theta_l)] p(l | \mathbf{x}_i, \Theta^g)$$

Finding θ_l Only need to maximize this term

Therefore, we want to maximize:

$$Q'(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \left(-\frac{1}{2} \log |\Sigma_l|^{1/2} - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_l) \right) p(l | \mathbf{x}_i, \Theta^g)$$

How?

$$p(l | \mathbf{x}_i, \Theta^g) = \frac{\alpha_l^g p_l(\mathbf{x}_i | \theta_l^g)}{\sum_{j=1}^M \alpha_j^g p_j(\mathbf{x} | \theta_j^g)}$$

Finding μ_l

Therefore, we want to maximize:

$$Q'(\Theta, \Theta^g) = \sum_{l=1}^M \sum_{i=1}^N \left(-\frac{1}{2} \log |\Sigma_l|^{1/2} - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right) p(l | \mathbf{x}_i, \Theta^g)$$

$$\mu_l = \frac{\sum_{i=1}^N \mathbf{x}_i p(l | \mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

$$\Sigma_l = \frac{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

Summary

EM algorithm for GMM

Given an initial guess Θ^g , find Θ^{new} as follows

$$p(l | \mathbf{x}_i, \Theta^g) = \frac{\alpha_l^g p_l(\mathbf{x}_i | \theta_l^g)}{\sum_{j=1}^M \alpha_j^g p_j(\mathbf{x}_i | \theta_j^g)}$$

$$\alpha_l^{new} = \frac{1}{N} \sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)$$

$$\mu_l^{new} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l | \mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

$$\Sigma_l^{new} = \frac{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \mu_l^{new})(\mathbf{x}_i - \mu_l^{new})^T}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

Not converge

$$\Theta^g \leftarrow \Theta^{new}$$