

Theoretical Basis for EM

Sargur Srihari

srihari@cedar.buffalo.edu

EM Algorithm

- Given joint distribution $p(X, Z | \theta)$ over observed variables X and latent variables Z governed by parameters θ
goal is to find θ that maximizes likelihood function $p(X | \theta)$

- Step 1:** Choose an initial setting for the parameters θ^0
- Step 2: E Step: Evaluate** $p(Z | X, \theta^t)$ Can be done since joint $p(X, Z | \theta)$ is known
- Step 3: M Step: Evaluate** θ^{t+1} given by

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

where

$$Q(\theta, \theta^t) = \sum_Z p(Z | X, \theta^t) \ln p(X, Z | \theta)$$

Summation due to expectation

- Check for convergence
 - of either log-likelihood or parameter values
- If not satisfied then let $\theta^t \leftarrow \theta^{t+1}$
- Return to **Step 2**

General EM Algorithm defined heuristically

- EM is a general technique
 - Find *maximum likelihood* solutions (parameters) for probabilistic models with latent variables
- EM was defined heuristically
 - It can be proved to maximize the likelihood function
 - Monotonically increases observed data *log-likelihood* until it reaches a local *maximum*
 - Proof involves obtaining lower bound on log-likelihood function
 - Dempster 1977

Likelihood Function with Latent Variables

- Given observed variables $X = \{x_1, \dots, x_N\}$ and hidden variables $Z = \{z_1, \dots, z_N\}$, likelihood function is

$$p(X | \theta) = \sum_Z p(X, Z | \theta) = \sum_{z_i} \prod_{i=1}^N p(x_i, z_i | \theta) = \prod_{i=1}^N \sum_{z_i} p(x_i, z_i | \theta)$$

- Log-likelihood is obtained by taking log of above:

$$\ell(\theta) = \ln p(X | \theta) = \sum_{i=1}^N \ln \left[\sum_{z_i} p(x_i, z_i | \theta) \right]$$

- Goal is to maximize log-likelihood $\ell(\theta)$
 - Direct optimization of $\ell(\theta)$, i.e., $p(X | \theta)$ is difficult, whereas complete-data likelihood $p(X, Z | \theta)$ is easier (since component forms are tractable)
 - Example of GMM: $p(x, z | \theta) = p(x | z, \theta) p(z | \theta)$ and $p(x | z, \theta) = \prod_{k=1}^K N(x | \mu_k, \Sigma_k)^{z_k}$, $p(z | \theta) = \prod_{k=1}^K \pi_k^{z_k}$

Decomposition of $p(X|\theta)$

- For any distribution $q(Z)$, this decomposition holds:

$$\ln p(X | \theta) = L(q, \theta) + KL(q \| p)$$

where we define

$$L(q, \theta) = \sum_z q(Z) \ln \left\{ \frac{p(X, Z | \theta)}{q(Z)} \right\}$$

$$KL(q \| p) = - \sum_z q(Z) \ln \left\{ \frac{p(Z | X, \theta)}{q(Z)} \right\}$$

L is a *functional* that takes a function as input and produces a value as output, like entropy
It contains joint distribution of X and Z

KL is the Kullback-Leibler Divergence
Contains conditional distribution of Z given X

Algebraic Derivation of expression:

$$\begin{aligned} p(X, Z | \theta) &= p(Z | X, \theta) p(X | \theta) \text{ from product rule} \\ \ln p(X, Z | \theta) &= \ln p(Z | X, \theta) + \ln p(X | \theta) \\ \ln p(X, Z | \theta) - \ln q(Z) &= \ln p(Z | X, \theta) - \ln q(Z) + \ln p(X | \theta) \\ \ln \frac{p(X, Z | \theta)}{q(Z)} &= \ln \frac{p(Z | X, \theta)}{q(Z)} + \ln p(X | \theta) \\ q(Z) \ln \frac{p(X, Z | \theta)}{q(Z)} &= q(Z) \ln \frac{p(Z | X, \theta)}{q(Z)} + q(Z) \ln p(X | \theta) \\ \sum_z q(Z) \ln \frac{p(X, Z | \theta)}{q(Z)} &= \sum_z q(Z) \ln \frac{p(Z | X, \theta)}{q(Z)} + \sum_z q(Z) \ln p(X | \theta) \\ \sum_z q(Z) \ln \frac{p(X, Z | \theta)}{q(Z)} &= \sum_z q(Z) \ln \frac{p(Z | X, \theta)}{q(Z)} + \ln p(X | \theta) \\ \ln p(X | \theta) &= \sum_z q(Z) \ln \frac{p(X, Z | \theta)}{q(Z)} - \sum_z q(Z) \ln \frac{p(Z | X, \theta)}{q(Z)} \end{aligned}$$

Lower bound on log-likelihood

$$\ell(\theta) = \ln p(X | \theta) = L(q, \theta) + KL(q || p)$$

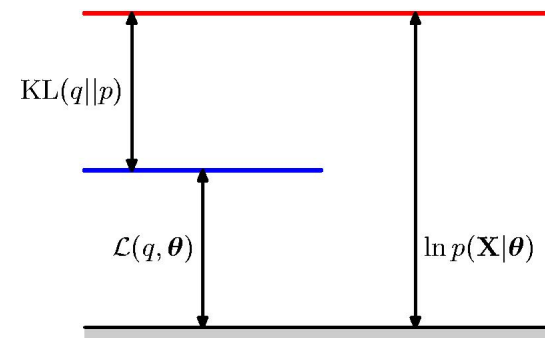
L is the lower bound

- Since $KL(q||p) \geq 0$ It follows that

$$\ln p(X|\theta) \geq L(q, \theta)$$

- Or $L(q, \theta)$ is a lower bound on $\ln p(X|\theta)$

$$\ell(\theta) \geq L(q, \theta)$$



E-step maximizes Lower Bound

- For any distribution $q(z_i)$ the following holds

$$\ell(\theta) = \ln p(X | \theta) = \sum_{i=1}^N \ln \left[\sum_{z_i} p(x_i, z_i | \theta) \right] = \sum_{i=1}^N \ln \left[\sum_{z_i} q(z_i) \frac{p(x_i, z_i | \theta)}{q(z_i)} \right]$$

Multiplying numerator and denominator by $q(z_i)$

- Using Jensen's inequality $f\left(\sum_{i=1}^N \lambda_i x_i\right) \leq \sum_{i=1}^N \lambda_i f(x_i)$ we get a lower bound

$$\ell(\theta) \geq \sum_{i=1}^N \sum_{z_i} q(z_i) \frac{p(x_i, z_i | \theta)}{q(z_i)}$$

- The lower bound is $Q(\theta, q) = \sum_i E_{q_i} [\ln p(x_i, z_i | \theta)] + H(q_i)$ H is the entropy of q_i

- What form of q yields the tightest lower bound?

- Lower bound is also sum of i terms of the form

$$L(\theta, q_i) = \sum_{z_i} q_i(z_i) \ln \frac{p(x_i, z_i | \theta)}{q_i(z_i)} = -KL(q_i(z_i) || p(z_i | x_i, \theta)) + \ln p(x_i | \theta)$$

- We can maximize the lower bound by setting $q_i(z_i) = p(z_i | x_i, \theta_i)$
- Since θ is unknown, instead we use $q_i^t(z_i) = p(z_i | x_i, \theta^t)$
 - Where θ^t is our estimate of parameters at iteration t

- This is the output of the E step, i.e., we determine $q_i^t(z_i)$

M step computes $\ell(\theta^t)$

- After the E step the lower bound is

$$Q(\theta, q^t) = \sum_i E_{q_i^t} [\ln p(x_i, z_i | \theta)] + H(q_i^t)$$

- Second term is constant wrt θ
- So the M step becomes

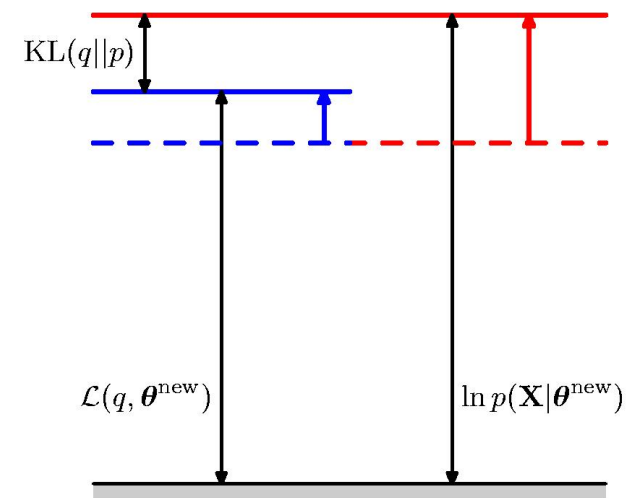
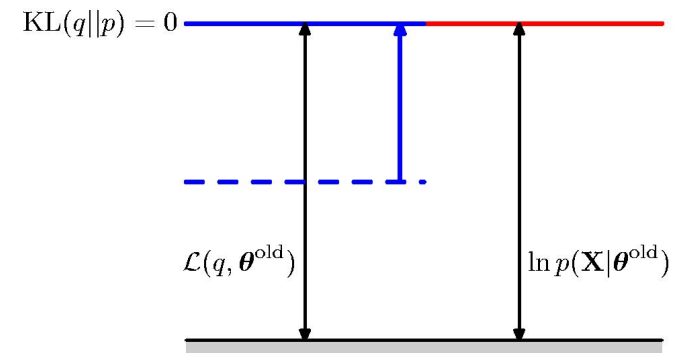
$$\theta^{t+1} = \arg \max_{\theta} Q(\theta, q^t) = \arg \max_{\theta} \sum_i E_{q_i^t} [\ln p(x_i, z_i | \theta)]$$

- Since we choose $q_i(z_i) = p(z_i | x_i, \theta_i)$ KL divergence is zero
 $L(\theta^t, q_i) = \ln p(x_i | \theta^t)$ and hence lower bound reaches $\ell(\theta^t)$

$$Q(\theta^t, q^t) = \sum_i \ln p(x_i | \theta^t) = \ell(\theta^t)$$

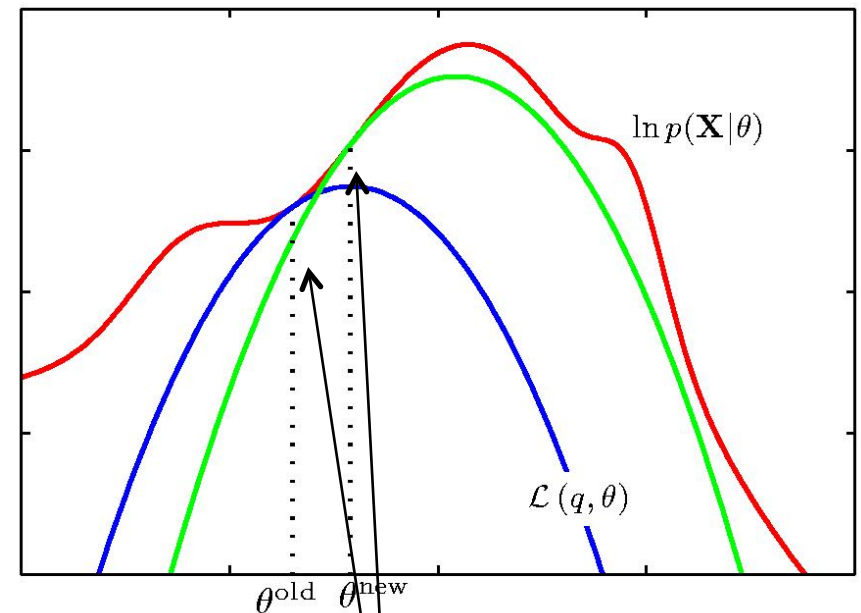
Illustration of E and M steps

- **E step** (Lower bound maximized keeping θ^{old} fixed)
 - θ distribution is set to posterior distribution for current parameter values θ^{old}
 - Causing lower bound to move to same value as log-likelihood with KL vanishing
- **M step**
 - Lower bound is maximized wrt q to give revised value θ^{new}
 - Causes log-likelihood to increase by as much as lower bound



View of EM in parameter space

- EM involves alternately computing lower bound on log-likelihood for the current parameter values
- And maximizing this bound to obtain new parameter values
- Note that the lower bound is a convex function with a unique maximum



Lower bound on log-likelihood function tangential at θ^{old} and θ^{new}

Generalized EM (GEM)

- EM breaks own potentially difficult problem of maximizing the likelihood function into two stages, the E step and the M step
- One or both may remain intractable
- GEM addresses the problem of the intractable M step
 - Instead of maximizing $L(q, \theta)$ wrt θ it changes parameters so as to increase its value
- Similar generalization of the E step can be made