

LVMs (Contd), Expectation Maximization (2)

CS771: Introduction to Machine Learning

Piyush Rai

What is EM Doing?

2

Maximizing ILL

Assuming \mathbf{Z} to be discrete, else replace it by an integral

- The MLE problem was $\Theta_{MLE} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}|\Theta) = \operatorname{argmax}_{\Theta} \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$
- What EM (and ALT-OPT in a crude way) did is max of CLL: $\Theta_{MLE} = \operatorname{argmax}_{\Theta} \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$
- But we did not solve the original problem. Is it okay?
- Assume $p_{\mathbf{Z}} = p(\mathbf{Z}|\mathbf{X}, \Theta)$ and $q(\mathbf{Z})$ to be some prob distribution over \mathbf{Z} , then

Function of a distribution q and parameter Θ

$$\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_{\mathbf{Z}})$$

May verify this identity

- In the above $\mathcal{L}(q, \Theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \right\}$ and $KL(q||p_{\mathbf{Z}}) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \Theta)}{q(\mathbf{Z})} \right\}$
- Since KL is always non-negative $\log p(\mathbf{X}|\Theta) \geq \mathcal{L}(q, \Theta)$, so $\mathcal{L}(q, \Theta)$ is a lower-bound on ILL
- Thus if we maximize $\mathcal{L}(q, \Theta)$, it will also improve $\log p(\mathbf{X}|\Theta)$



What is EM Doing?

- As we saw, $\mathcal{L}(q, \Theta)$ depends on q and Θ
- Let's maximize $\mathcal{L}(q, \Theta)$ w.r.t. q with Θ fixed at Θ^{old}

The posterior distribution of \mathbf{Z} given older parameters Θ^{old} (will need this posterior to get the expectation of CLL)

Since $\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_z)$ is constant when Θ is held fixed at Θ^{old}

$$\hat{q} = \operatorname{argmax}_q \mathcal{L}(q, \Theta^{\text{old}}) = \operatorname{argmin}_q KL(q||p_z) = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$$

- Now let's maximize $\mathcal{L}(q, \Theta)$ w.r.t. Θ with q fixed at $\hat{q} = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$

$$\begin{aligned} \Theta^{\text{new}} &= \operatorname{argmax}_{\Theta} \mathcal{L}(\hat{q}, \Theta) = \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} \right\} \\ &= \operatorname{argmax}_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) \\ &= \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)] \\ &= \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{\text{old}}) \end{aligned}$$

Maximization of expected CLL w.r.t. the posterior distribution of \mathbf{Z} given older parameters Θ^{old}



The EM Algorithm in its general form..

- Maximization of $\mathcal{L}(q, \Theta)$ w.r.t. q and Θ gives the EM algorithm (Dempster, Laird, Rubin, 1977)

The EM Algorithm

- 1 Initialize Θ as $\Theta^{(0)}$, set $t = 1$
- 2 Step 1: Compute **posterior** of latent variables given current parameters $\Theta^{(t-1)}$

$$p(\mathbf{z}_n^{(t)} | \mathbf{x}_n, \Theta^{(t-1)}) = \frac{p(\mathbf{z}_n^{(t)} | \Theta^{(t-1)}) p(\mathbf{x}_n | \mathbf{z}_n^{(t)}, \Theta^{(t-1)})}{p(\mathbf{x}_n | \Theta^{(t-1)})} \propto \text{prior} \times \text{likelihood}$$

- 3 Step 2: Now maximize the **expected complete data log-likelihood** w.r.t. Θ

$$\Theta^{(t)} = \arg \max_{\Theta} \mathcal{Q}(\Theta, \Theta^{(t-1)}) = \arg \max_{\Theta} \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n^{(t)} | \mathbf{x}_n, \Theta^{(t-1)})} [\log p(\mathbf{x}_n, \mathbf{z}_n^{(t)} | \Theta)]$$

- 4 If not yet converged, set $t = t + 1$ and go to step 2.

- Note: If we can take the MAP estimate $\hat{\mathbf{z}}_n$ of \mathbf{z}_n (not full posterior) in Step 1 and maximize the CLL in Step 2 using that, i.e., do $\arg \max_{\Theta} \sum_{n=1}^N [\log p(\mathbf{x}_n, \hat{\mathbf{z}}_n^{(t)} | \Theta)]$ this will be ALT-OPT



The Expected CLL

- Expected CLL in EM is given by (assume observations are i.i.d.)

$$\begin{aligned} Q(\Theta, \Theta^{old}) &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \Theta^{old})} [\log p(\mathbf{x}_n, \mathbf{z}_n|\Theta)] \\ &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \Theta^{old})} [\log p(\mathbf{x}_n|\mathbf{z}_n, \Theta) + \log p(\mathbf{z}_n|\Theta)] \end{aligned}$$

Was indeed the case of GMM: $p(\mathbf{z}_n|\Theta)$ was multinoulli, $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$ was Gaussian

- If $p(\mathbf{z}_n|\Theta)$ and $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$ are exp-family distributions, $Q(\Theta, \Theta^{old})$ has a very simple form
- In resulting expressions, replace terms containing \mathbf{z}_n 's by their respective expectations, e.g.,
 - \mathbf{z}_n replaced by $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})}[\mathbf{z}_n]$
 - $\mathbf{z}_n \mathbf{z}_n^T$ replaced by $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})}[\mathbf{z}_n \mathbf{z}_n^T]$
- However, in some LVMs, these expectations are intractable to compute and need to be approximated (beyond the scope of CS771)



EM: An Illustration

- As we saw, EM maximizes the lower bound $\mathcal{L}(q, \Theta)$ in two steps
- Step 1 finds the optimal q (call it \hat{q}) by setting it the posterior of \mathbf{Z} given current Θ
- Step 2 maximizes $\mathcal{L}(\hat{q}, \Theta)$ w.r.t. Θ which gives a new Θ .

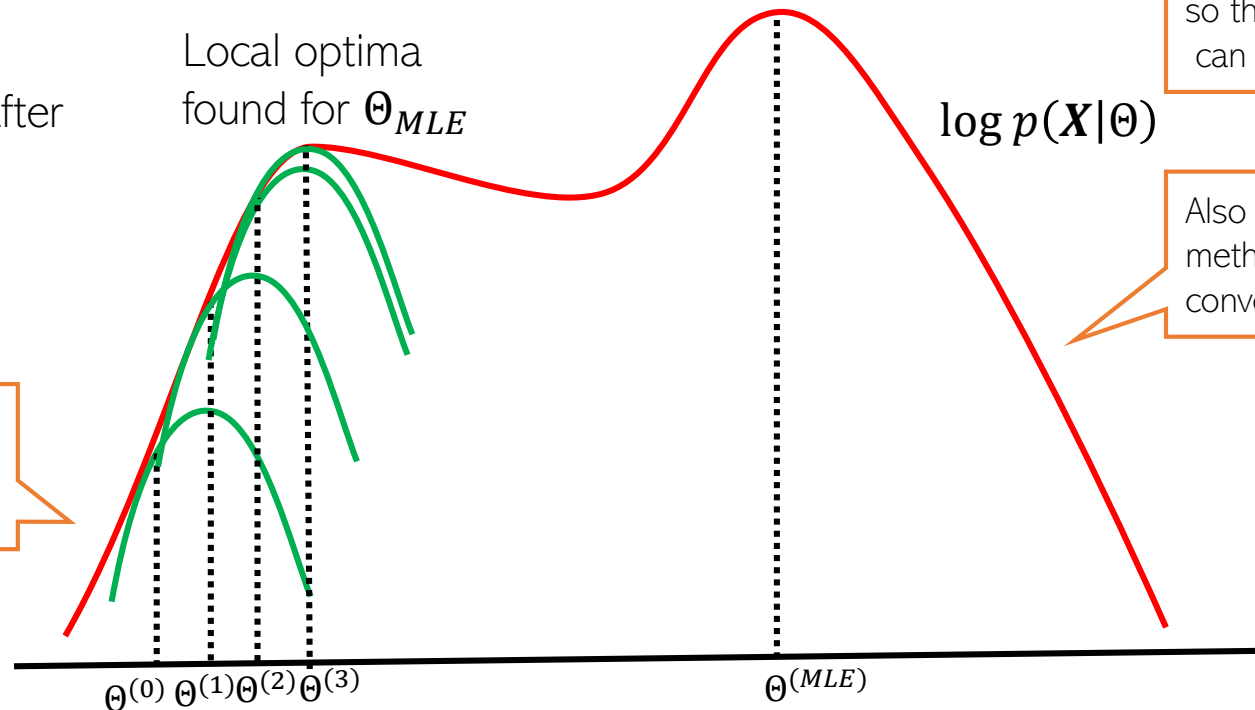
Alternating between them until convergence to some local optima

Makes $\mathcal{L}(q, \Theta)$ equal to $\log p(\mathbf{X}|\Theta)$; thus the curves touch at current Θ

Green curve: $\mathcal{L}(\hat{q}, \Theta)$ after setting q to \hat{q}

Local optima found for Θ_{MLE}

Good initialization matters; otherwise would converge to a poor local optima



Note that Θ only changes in Step 2 so the objective $\log p(\mathbf{X}|\Theta)$ can only change in Step 2



Also kind of similar to Newton's method (and has second order like convergence behavior in some cases)

Unlike Newton's method, we don't construct and optimize a quadratic approximation, but a lower bound

Even though original MLE problem $\text{argmax}_{\Theta} \log p(\mathbf{X}|\Theta)$ could be solved using gradient methods, EM often works faster and has cleaner updates

Recap: ALT-OPT vs EM

- ALT-OPT does the following

- 1 Initialize $\Theta = \hat{\Theta}$
- 2 Estimate \mathbf{Z} as $\hat{\mathbf{Z}} = \arg \max_{\mathbf{Z}} \log p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$
- 3 Estimate Θ as $\hat{\Theta} = \arg \max_{\Theta} \log p(\mathbf{X}, \hat{\mathbf{Z}}|\Theta)$
- 4 Go to step 2 if not converged

This step could potentially throw away a lot of information about the latent variable \mathbf{Z}

- EM addresses it using “soft” version of ALT-OPT

- 1 Initialize $\Theta = \hat{\Theta}$
- 2 Compute the posterior distribution of \mathbf{Z} , i.e., $p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$
- 3 Estimate Θ by maximizing the expected CLL $\hat{\Theta} = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$
- 4 Go to step 2 if not converged

ALT-OPT can be seen as an approximation of EM – the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta)$ is replaced by a point mass at its mode



EM: Some Comments

- The E and M steps may not always be possible to perform exactly. Some reasons
 - Posterior of latent variables $p(\mathbf{Z}|\mathbf{X}, \Theta)$ may not be easy to find and may require approx.
 - Even if $p(\mathbf{Z}|\mathbf{X}, \Theta)$ is easy, expected CLL, i.e., $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$ may still not be tractable

$$\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \int \log p(\mathbf{X}, \mathbf{Z}|\Theta) p(\mathbf{Z}|\mathbf{X}, \Theta) d\mathbf{Z}$$

Monte-Carlo EM

..and may need to be approximated, e.g., using [Monte-Carlo expectation](#)

Gradient methods may still be needed for this step

- Maximization of the expected CLL may not be possible in closed form
- EM works even if the M step is only solved approximately ([Generalized EM](#))
- If M step has multiple parameters whose updates depend on each other, they are updated in an alternating fashion - called [Expectation Conditional Maximization \(ECM\)](#) algorithm
- Other advanced probabilistic inference algorithms are based on ideas similar to EM
 - E.g., Variational Bayesian inference a.k.a. [Variational Inference \(VI\)](#)
- EM is also related to non-convex optimization algorithms [Majorization-Maximization \(MM\)](#)

