

Expectation Maximization (Contd)

CS698X: Topics in Probabilistic Modeling and Inference

Piyush Rai

The Expectation-Maximization (EM) Algorithm

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

The EM Algorithm

Primarily designed for doing point estimation of the parameters Θ but also gives (CP of) latent variables z_n

Usually computing CP + expected CLL is referred to as the **E step**, and max. of exp-CLL w.r.t. Θ as the **M step**



① Initialize Θ as $\Theta^{(0)}$, set $t = 1$

② Step 1: Compute **posterior** of latent variables given current parameters $\Theta^{(t-1)}$

Conditional posterior of each latent variable z_n

Latent variables also assumed indep. a priori

$$p(z_n^{(t)} | x_n, \Theta^{(t-1)}) = \frac{p(z_n^{(t)} | \Theta^{(t-1)}) p(x_n | z_n^{(t)}, \Theta^{(t-1)})}{p(x_n | \Theta^{(t-1)})} \propto \text{prior} \times \text{likelihood}$$

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

Assuming the (expected) CLL $\mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \Theta^{\text{old}})} [\log p(\mathbf{X}, \mathbf{Z} | \Theta)]$ factorizes over all observations

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

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

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



What's Going On?

4

- As we saw, the maximization of lower bound $\mathcal{L}(q, \Theta)$ had two steps
- Step 1 finds the optimal q (call it \hat{q}) by setting it as 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

KL becomes zero and $\mathcal{L}(q, \Theta)$ becomes equal to $\log p(\mathbf{X}|\Theta)$; thus their curves touch at current Θ

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

Local optima found for Θ_{MLE}

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

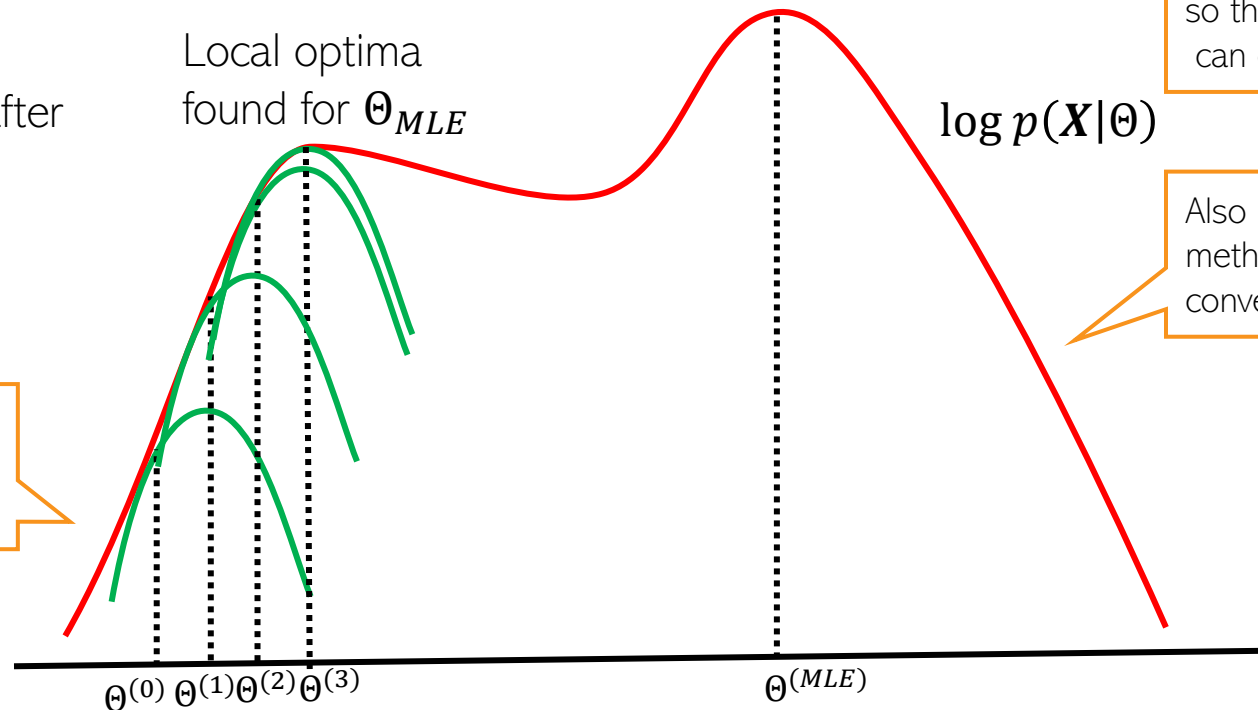


Also kind of similar to Newton's method (and has second order like convergence behavior in some cases)

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

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



Online/Incremental EM

Have to do it in each iteration of EM

- Computing CP of latent variable \mathbf{z}_n of each observation \mathbf{x}_n can be expensive
 - Before we do the M step to update Θ , we must wait for all CPs to be computed ☹
- Recall that, for i.i.d. case, the expected CLL is often a sum over all data points

$$Q(\Theta, \Theta^{old}) = \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \Theta)] = \sum_{n=1}^N \mathbb{E}[\log p(\mathbf{x}_n | \mathbf{z}_n, \theta)] + \mathbb{E}[\log p(\mathbf{z}_n | \phi)]$$

- Can compute this quantity recursively using small minibatches of data

Expected CLL from iteration $t-1$

$$Q_t = (1 - \gamma_t) Q_{t-1} + \gamma_t \left[\sum_{n=1}^{N_t} \mathbb{E}[\log p(\mathbf{x}_n | \mathbf{z}_n, \theta)] + \mathbb{E}[\log p(\mathbf{z}_n | \phi)] \right]$$

Expected CLL from this minibatch

Only requires CP for the latent variables from this minibatch of observations

where $\gamma_t = (1 + t)^{-\kappa}$, $0.5 < \kappa \leq 1$ is a decaying learning rate

- MLE on above Q_t can be shown to be equivalent to a simple recursive updates for Θ

$$\Theta^{(t)} = (1 - \gamma_t) \times \Theta^{(t-1)} + \gamma_t \times \arg \max_{\Theta} \underbrace{Q(\Theta, \Theta^{t-1})}_{\substack{\text{computed using only} \\ \text{the } N_t \text{ examples} \\ \text{from this minibatch}}}$$



How M Step uses Sufficient Statistics

- Recall the batch EM algorithm for a K component Gaussian mixture model
 - Cluster id z_n s.t. $z_{nk} = 1$ if x_n belongs to cluster k , and 0 otherwise
 - The conditional posterior of z_{nk} is $p(z_{nk} = 1 | x_n, \Theta) \propto \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$
- Denoting current iteration by t , and expectation computed in E step: $\mathbb{E} [z_{nk}^{(t)}] = \gamma_{nk}^{(t)}$
- The M step updates for params $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ are

$$\mu_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} x_n$$

$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} (x_n - \mu_k^{(t)})(x_n - \mu_k^{(t)})^\top$$

$$\pi_k^{(t)} = \frac{\sum_{n=1}^N \gamma_{nk}^{(t)}}{N}$$

- Each update depends on sum of **expected** sufficient statistics (ESS). For each x_n, z_n

ESS for μ_k is $\gamma_{nk}^{(t)} x_n$; ESS for Σ_k is $\gamma_{nk}^{(t)} (x_n - \mu_k^{(t)})(x_n - \mu_k^{(t)})^\top$; ESS for π_k is $\gamma_{nk}^{(t)}$



Batch EM in terms of Sufficient Statistics

- Denote the **sum of ESS** as $\mathbf{S} = \sum_{n=1}^N \mathbf{s}_n$ where each ESS

$$\mathbf{s}_n = \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \Theta, \mathbf{x}_n) \phi(\mathbf{x}_n, \mathbf{z}_n)$$

SS for \mathbf{x}_n and \mathbf{z}_n

- M step updates of Θ are like computing a function of \mathbf{S} , i.e., $\Theta = f(\mathbf{S})$

Batch EM in terms of ESS

- Initialize \mathbf{S} and compute parameters $\Theta = f(\mathbf{S})$
- For $t = 1 : T$ (or until convergence)
 - $\mathbf{S}^{new} = 0$ (fresh sum of ESS; will be computed in this iteration)
 - For $n = 1 : N$

$$\begin{aligned} \mathbf{s}_n &= \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \Theta) \phi(\mathbf{x}_n, \mathbf{z}_n) = \mathbb{E}[\phi(\mathbf{x}_n, \mathbf{z}_n)] \\ \mathbf{S}^{new} &= \mathbf{S}^{new} + \mathbf{s}_n \end{aligned}$$

- $\mathbf{S} = \mathbf{S}^{new}$
- Recompute parameters $\Theta = f(\mathbf{S})$

- Note: In general, there may be more than one sum of ESS (one for each param)
 - E.g., for GMM, one for π_k , one for μ_k , one for Σ_k



Online EM in terms of Sufficient Statistics

- Works in a similar way as batch EM except we perform online updates for \mathbf{S}
- Can be done in one of the two manners (Liang and Klein, 2009)
 - Stepwise EM (based on recursively updating the sum of ESS)
 - Incremental EM (based on deleting old and adding new ESS of each data point)

Online EM as Stepwise EM

- Initialize the sum of ESS \mathbf{S} and compute $\Theta = f(\mathbf{S})$
- For $t = 1 : T$ (or until convergence)
 - Set “learning rate” γ_t , pick a random example n and compute its sufficient statistics

$$\mathbf{s}_n = \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \Theta) \phi(\mathbf{x}_n, \mathbf{z}_n)$$

$$\mathbf{S} = (1 - \gamma_t) \mathbf{S} + \gamma_t \mathbf{s}_n$$

- Recompute $\Theta = f(\mathbf{S})$



Online EM in terms of Sufficient Statistics

- The other Online EM approach “Incremental EM” needs no learning rate

Online EM as Incremental EM

- Initialize each ESS \mathbf{s}_n , $n = 1, \dots, N$, $\mathbf{S} = \sum_{n=1}^N \mathbf{s}_n$, and compute $\Theta = f(\mathbf{S})$
- For $t = 1 : T$ (or until convergence)
 - Pick a random example n and update its exp. sufficient statistics

$$\mathbf{s}_n^{\text{new}} = \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \Theta) \phi(\mathbf{x}_n, \mathbf{z}_n)$$

$$\mathbf{S} = \mathbf{S} + \mathbf{s}_n^{\text{new}} - \mathbf{s}_n$$

$$\mathbf{s}_n = \mathbf{s}_n^{\text{new}}$$

- Recompute $\Theta = f(\mathbf{S})$

- However, incremental EM requires keeping a track of sum of ESS \mathbf{S} as well as each \mathbf{s}_n
- In practice, stepwise EM outperforms batch EM as well as incremental EM on many problems (can refer to Liang and Klein, 2009 for some examples of models where these algos were tried)



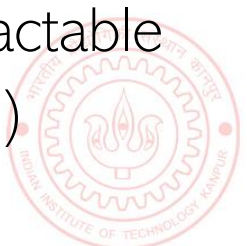
EM vs Gradient-based Methods

- Can also estimate params using gradient-based optimization instead of EM
 - We can usually explicitly sum over or integrate out the latent variables \mathbf{Z} , e.g.,

$$\mathcal{L}(\Theta) = \log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- Now we can optimize $\mathcal{L}(\Theta)$ using first/second order optimization to find the optimal Θ
- EM is usually preferred over this approach because
 - The M step has often simple closed-form updates for the parameters Θ
 - Often constraints (e.g., PSD matrices) are automatically satisfied due to form of updates
 - In some cases[†], EM usually converges faster (and often like second-order methods)
 - E.g., Example: Mixture of Gaussians with when the data is reasonably well-clustered
 - EM applies even when the explicit summing over/integrating out is expensive/intractable
 - EM also provides the conditional posterior over the latent variables \mathbf{Z} (from E step)

[†]Optimization with EM and Expectation-Conjugate-Gradient (Salakhutdinov et al, 2003), On Convergence Properties of the EM Algorithm for Gaussian Mixtures (Xu and Jordan, 1996), Statistical guarantees for the EM algorithm: From population to sample-based analysis (Balakrishnan et al, 2017)



Some Applications of EM

- Mixture Models (each data-point comes from one of K mixture components)
 - Examples: Mixture of Gaussians, Mixture of Experts, etc
- Latent variable models for dimensionality reduction or representation learning
 - Probabilistic PCA, Factor Analysis, Variational Autoencoders, etc
- Problems models with missing features/labels (treated as latent variables)
 - An example of problem with missing labels: [Semi-supervised learning](#)
- [Hyperparameter estimation](#) in probabilistic models (an alternative to MLE-II)
 - MLE-II estimates hyperparams by maximizing the marginal likelihood, e.g.,

$$\{\hat{\lambda}, \hat{\beta}\} = \operatorname{argmax}_{\lambda, \beta} p(\mathbf{y}|\mathbf{X}, \lambda, \beta) = \operatorname{argmax}_{\lambda, \beta} \int p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

For a Bayesian linear regression model

- With EM, can treat \mathbf{w} as latent var, and λ, β as “parameters”
 - E step will estimate the CP of \mathbf{w} given current estimates of λ, β
 - M step will re-estimate λ, β by maximizing the expected CLL

$$\mathbb{E}[\log p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda)] = \mathbb{E}[\log p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \beta) + \log p(\mathbf{w}|\lambda)]$$

Expectations w.r.t. the CP of \mathbf{w}



EM: Some Comments

- The E and M steps may not always be possible to perform exactly. Some reasons

- The conditional posterior of latent variables $p(\mathbf{Z}|\mathbf{X}, \Theta)$ may not be easy to compute
 - Will need to approximate $p(\mathbf{Z}|\mathbf{X}, \Theta)$ using methods such as MCMC or variational inference

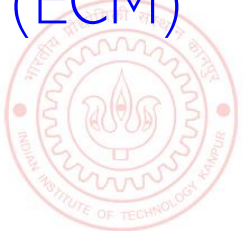
- Even if $p(\mathbf{Z}|\mathbf{X}, \Theta)$ is easy, the expected CLL may not be easy to compute

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

Can often be approximated by Monte-Carlo using sample from the CP of \mathbf{Z}

Results in Monte-Carlo EM

- 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)
- Other advanced probabilistic inference algos are based on ideas similar to EM
 - E.g., Variational Bayes (VB) inference, a.k.a. Variational Inference (VI)



Coming Up Next

- Variational Inference

