Latent Variable Models and EM Algorithm

CS772A: Probabilistic Machine Learning Piyush Rai

Hybrid Inference (posterior infer. + point est.)

- In many models, we infer posterior on some unknowns and do point est. for others
- We have already seen MLE-II for lin reg. which alternates between
 - Inferring CP over the main parameter given the point estimates of hyperparams
 - Maximizing the marginal lik. to do point estimation for hyperparams
- The Expectation-Maximization algorithm (will see today) also does something similar
 - In E step, the CP of latent variables is inferred, given <u>current</u> point-est of params
 - M step maximizes expected complete data log-lik. to get point estimates of params
- If we can't (due to computational or other reasons) infer posterior over all unknowns, how to decide which variables to infer posterior on, and for which to do point-est?
- Usual approach: Infer posterior over local vars and point estimates for global vars
 - Reason: We typically have plenty of data to reliably estimate the global variables so it is okay even if we just do point estimation for those

CP of $w: p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}, \hat{\lambda}, \hat{\beta})$

 $\{\hat{\lambda}, \hat{\beta}\} = \operatorname{argmax}_{\lambda,\beta} p(\boldsymbol{y}|\boldsymbol{X}, \lambda, \beta)$

Nomenclature/Notation Alert

- Why call some unknowns as parameters and others as latent variables?
- Well, no specific reason. Sort of a convention adopted by some algorithms
 - EM: Unknowns estimated in E step referred to as latent vars; those in M step as params
 - Usually: Latent vars (Conditional) posterior computed; parameters point estimation
- Some algos won't make such distinction and will infer posterior over all unknowns
- Sometimes the "global" or "local" unknown distinction makes it clear
 - Local variables = latent variables, global variables = parameters
- But remember that this nomenclature isn't really cast in stone, no need to be confused so long as you are clear as to what the role of each unknown is, and how we want to estimate it (posterior or point estimate) and using what type of inference algorithm

CS772A: PML

Inference/Parameter Estimation in Latent Variable Models using Expectation-Maximization (EM)



Parameter Estimation in Latent Variable Models

• Assume each observation x_n to be associated with a "local" latent variable z_n



- Although we can do fully Bayesian inference for all the unknowns, suppose we only want a point estimate of the "global" parameters $\Theta = (\theta, \phi)$ via MLE/MAP
- Such MLE/MAP problems in LVMs are difficult to solve in a "clean" way
 - Would typically require gradient based methods with no closed form updates for Θ
 - However, EM gives a clean way to obtain closed form updates for Θ



Why MLE/MAP of Params is Hard for LVMs?

- Suppose we want to estimate $\Theta = (\theta, \phi)$ via MLE. If we knew \mathbf{z}_n , we could solve
- Easy to solve $\Theta_{MLE} = \arg \max_{\Theta} \sum_{n=1}^{n} \log p(\mathbf{x}_n, \mathbf{z}_n | \Theta) = \arg \max_{\Theta} \sum_{n=1}^{n} [\log p(\mathbf{z}_n | \phi) + \log p(\mathbf{x}_n | \mathbf{z}_n, \theta)]$ = Easy. Usually closed form if $p(\mathbf{z}_n | \phi)$ and $p(\mathbf{x}_n | \mathbf{z}_n, \theta)$ have simple forms experiment. In particular, if they are exp-fam distributions
- However, since in LVMs, \mathbf{z}_n is hidden, the MLE problem for Θ will be the following Basically, the marginal (elihood after tegrating out z_n $\Theta_{MLE} = \arg \max_{\Theta} \sum_{n=1}^n \log p(\mathbf{x}_n | \Theta) = \arg \max_{\Theta} \log p(\mathbf{X} | \Theta)$ $\log p(\mathbf{x}_n | \Theta)$ will not have a simple expression since $p(\mathbf{x}_n | \Theta)$ requires sum/integral likelihood after integrating out z_n

$$p(\mathbf{x}_n|\Theta) = \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n|\Theta)$$
 ... or if \mathbf{z}_n is continuous: $p(\mathbf{x}_n|\Theta) = \int p(\mathbf{x}_n, \mathbf{z}_n|\Theta) d\mathbf{z}_n$

• MLE now becomes difficult (basically MLE-II now), no closed form expression for Θ .

• Can we maximize some other quantity instead of $\log p(x_n | \Theta)$ for this MLE?

CS772A: PML

An Important Identity

• Assume $p_z = p(Z|X, \Theta)$ and q(Z) to be some prob distribution over Z, then

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

• In the above
$$\mathcal{L}(q, \Theta) = \sum_{Z} q(Z) \log \left\{ \frac{p(X, Z | \Theta)}{q(Z)} \right\}$$

•
$$KL(q||p_z) = -\sum_Z q(\mathbf{Z}) \log\left\{\frac{p(\mathbf{Z}|\mathbf{X},\Theta)}{q(\mathbf{Z})}\right\}$$

Assume \mathbf{Z} discrete

- KL is always non-negative, so $\log p(X|\Theta) \ge \mathcal{L}(q,\Theta)$
- Thus $\mathcal{L}(q, \Theta)$ is a lower-bound on $\log p(X|\Theta)$
- Thus if we maximize $\mathcal{L}(q, \Theta)$, it will also improve $\log p(X|\Theta)$
- Also, as we'll see, it's easier to maximize $\mathcal{L}(q, \Theta)$



Verify the identity





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 Expected CLL

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

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

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

CS772A: PML

What's Going On?

Alternating between them until convergence to some local optima

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

- 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 Z given current Θ
- Step 2 maximizes $\mathcal{L}(\hat{q}, \Theta)$ w.r.t. Θ which gives a new Θ .



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 $oldsymbol{Z}$, e.g.,

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

- 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
 - ${\ensuremath{\,^\circ}}$ 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

CS772A: PML

EM also provides the conditional posterior over the latent variables Z (from E step)

Some Applications of EM

- Mixture Models and Dimensionality Reduction/Representation Learning
 - Mixture Models: Mixture of Gaussians, Mixture of Experts, etc
 - Dim. Reduction/Representation Learning: Probabilistic PCA, Variational Autoencoders
- Problems with missing features or missing labels (which are treated as latent variables)
 - $\widehat{\Theta} = \operatorname{argmax}_{\Theta} \log p(\mathbf{x}^{obs} | \Theta) = \operatorname{argmax}_{\Theta} \log \int p([\mathbf{x}^{obs}, \mathbf{x}^{miss}] | \Theta) d\mathbf{x}^{miss}$
 - $\widehat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^{N} \log p(x_n, y_n | \Theta) + \sum_{n=N+1}^{N+M} \log \sum_{c=1}^{K} p(x_n, y_n = c | \Theta)$

Hyperparameter estimation in probabilistic models (an alternative to MLE-II)

MLE-II estimates hyperparams by maximizing the marginal likelihood, e.g.,

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

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

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



For a Bayesian linear

regression model

Expectations w.r.t.

the CP of **w**