

# Inference in Latent Variable Models: The EM Algorithm

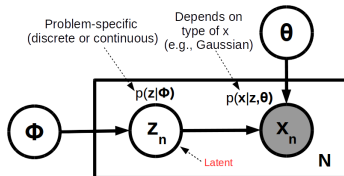
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Feb 8, 2018

# Parameter Estimation in Latent Variable Models

- Assume each observation  $\mathbf{x}_n$  to be associated with a “local” latent variable  $\mathbf{z}_n$



- 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
- The MLE would be  $\hat{\Theta} = \arg \max_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n|\Theta)$  where

$$\text{if } \mathbf{z}_n \text{ is discrete: } \log p(\mathbf{x}_n|\Theta) = \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n|\Theta) = \log \sum_{k=1}^K p(\mathbf{x}_n|\mathbf{z}_n = k, \Theta) p(\mathbf{z}_n = k|\Theta)$$

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

- MLE difficult:**  $p(\mathbf{x}_n|\Theta)$  is usually **not in exp-family**  $\Rightarrow \log p(\mathbf{x}_n|\Theta)$  won't have a simple form

# Parameter Estimation in Latent Variable Models

- Rather than MLE on  $\sum_{n=1}^N \log p(\mathbf{x}_n|\Theta)$ , how about doing **MLE on  $\sum_{n=1}^N \log p(\mathbf{x}_n, \mathbf{z}_n|\Theta)$**  ?
  - One reason:  $\log p(\mathbf{x}_n, \mathbf{z}_n|\Theta) = \log p(\mathbf{x}_n|\mathbf{z}_n, \Theta)p(\mathbf{z}_n|\Theta)$  usually has a **much simpler expression**
  - .. simpler especially when  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  and  $p(\mathbf{z}_n|\Theta)$  are **exponential family distributions**
- $p(\mathbf{x}_n, \mathbf{z}_n|\Theta)$  known as **complete data likelihood**,  $p(\mathbf{x}_n|\Theta)$  known as the **incomplete data likelihood**
- Is MLE on  $\sum_{n=1}^N \log p(\mathbf{x}_n, \mathbf{z}_n|\Theta)$  instead of our original objective  $\sum_{n=1}^N \log p(\mathbf{x}_n|\Theta)$  right thing?
  - **Yes :-)** Will see the justification shortly
- Denoting  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ , consider the new MLE problem

$$\hat{\Theta} = \arg \max_{\Theta} \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- However since  $\mathbf{Z}$  is latent, we will actually maximize the expectation  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$

$$\hat{\Theta} = \arg \max_{\Theta} \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$$

.. where the expectation is w.r.t. to an “optimal” distribution over  $\mathbf{Z}$  (will see shortly what it is)

# An Important Identity

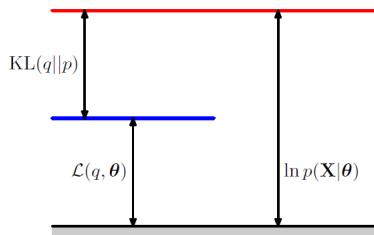
- Define  $p_z = p(\mathbf{Z}|\mathbf{X}, \Theta)$  and let  $q(\mathbf{Z})$  be some distribution over  $\mathbf{Z}$
- Assume discrete  $\mathbf{Z}$ , the identity below holds for any choice of the distribution  $q(\mathbf{Z})$

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

$$\mathcal{L}(q, \Theta) = \sum_{\mathbf{z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \right\}$$

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

(Exercise: Verify the above identity)



- Since  $\text{KL}(q||p_z) \geq 0$ ,  $\mathcal{L}(q, \Theta)$  is a **lower-bound** on  $\log p(\mathbf{X}|\Theta)$
- Maximizing  $\mathcal{L}(q, \Theta)$  won't decrease  $\log p(\mathbf{X}|\Theta)$ . Can maximize  $\log p(\mathbf{X}|\Theta)$  by maximizing  $\mathcal{L}(q, \Theta)$
- Note:  $\mathcal{L}(q, \Theta)$  depends on  $q$  and  $\Theta$ . Consider **alternating maximization** w.r.t each fixing the other

# An Alternating Optimization Scheme for $\mathcal{L}(q, \Theta)$

- The identity we had was  $\log p(\mathbf{X}|\Theta) = \mathcal{L}(q, \Theta) + \text{KL}(q||p_z)$  with

$$\mathcal{L}(q, \Theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \right\} \quad \text{and} \quad \text{KL}(q||p_z) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \Theta)}{q(\mathbf{Z})} \right\}$$

- With  $\Theta$  fixed at  $\Theta^{old}$ ,  $\log p(\mathbf{X}|\Theta)$  won't change. The  $q$  that maximizes  $\mathcal{L}(q, \Theta)$  will be

$$\hat{q} = \arg \max_q \mathcal{L}(q, \Theta^{old})$$

.. which is attained when  $\text{KL}(q||p_z) = 0$ . Therefore  $\hat{q} = p_z = p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$

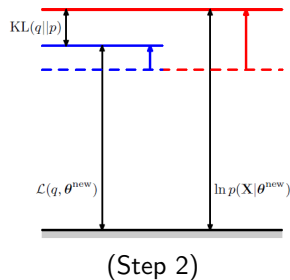
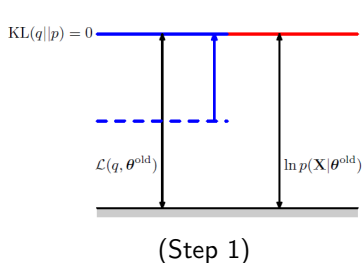
- With  $q$  fixed at  $\hat{q} = p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ , the  $\Theta$  that maximizes  $\mathcal{L}(q, \Theta)$  will be

$$\Theta^{new} = \arg \max_{\Theta} \mathcal{L}(\hat{q}, \Theta) = \arg \max_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})} = \arg \max_{\Theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

.. therefore,  $\Theta^{new} = \arg \max_{\Theta} \mathcal{Q}(\Theta, \Theta^{old})$  where  $\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$

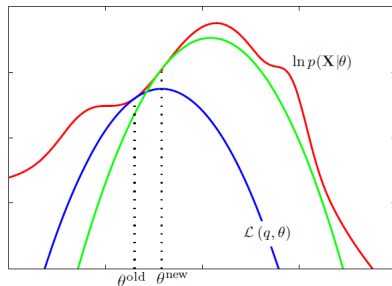
# Why Alternating Optimization Works Here?

- The two-step alternating optimization scheme we saw never decreases  $p(\mathbf{X}|\Theta)$
- To see this consider both steps: (1) Optimize  $q$  with  $\Theta = \Theta^{old}$ ; (2) Optimize  $\Theta$  using this  $q$



- Step 1 keeps  $\Theta$  fixed, so  $p(\mathbf{X}|\Theta)$  obviously can't decrease (stays unchanged in this step)
- Step 2 maximizes the lower bound  $\mathcal{L}(q, \Theta)$  w.r.t  $\Theta$ . Thus  $p(\mathbf{X}|\Theta)$  can't decrease!

# Convergence: A View in the Parameter ( $\Theta$ ) Space



- E-step: Update of  $q$  makes the  $\mathcal{L}(q, \Theta)$  curve touch the  $\log p(\mathbf{X}|\Theta)$  curve at  $\Theta^{old}$
- M-step gives the maxima  $\Theta^{new}$  of  $\mathcal{L}(q, \Theta^{old})$
- Next E-step readjusts  $\mathcal{L}(q, \Theta^{old})$  curve (green) to meet  $\log p(\mathbf{X}|\Theta)$  curve again, now at  $\Theta^{new}$
- This continues until a local maxima of  $\log p(\mathbf{X}|\Theta)$  is reached

# The Expectation Maximization (EM) Algorithm

Initialize the parameters:  $\Theta^{old}$ . Then alternate between these steps:

- **E (Expectation) step:**

- Compute the posterior distribution  $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$  over latent variables  $\mathbf{Z}$  using  $\Theta^{old}$
- Compute the **expected complete data log-likelihood** w.r.t. *this* posterior distribution

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

- **M (Maximization) step:**

- **Maximize** the expected complete data log-likelihood w.r.t.  $\Theta$

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

- If the incomplete log-lik  $p(\mathbf{X}|\Theta)$  not yet converged then set  $\Theta^{old} = \Theta^{new}$  and go to the E step.



# The Expected CLL

- Deriving the EM algorithm requires finding the expression of the expected CLL

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

- If  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  and  $p(\mathbf{z}_n|\Theta)$  are exp-family distributions, expected CLL will have a simple form
- Finding the expression for the expected CLL in such cases is fairly straightforward
  - First write down the expressions for  $p(\mathbf{x}_n|\mathbf{z}_n, \Theta)$  and  $p(\mathbf{z}_n|\Theta)$  and simplify as much as possible
  - In the resulting expressions, replace all terms containing  $\mathbf{z}_n$ 's by their respective expectations, e.g.,
    - $\mathbf{z}_n$  replaced by  $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \Theta^{old})}[\mathbf{z}_n]$ , i.e., the posterior mean of  $\mathbf{z}_n$
    - $\mathbf{z}_n \mathbf{z}_n^\top$  replaced by  $\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \Theta^{old})}[\mathbf{z}_n \mathbf{z}_n^\top]$
    - .. and so on..

# The Expected CLL via an example

- Let's consider a **latent factor model for dimensionality reduction**

$$p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}) = \mathcal{N}(\mathbf{W}\mathbf{z}_n, \sigma^2 \mathbf{I}_D) \quad p(\mathbf{z}_n) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$$

- A **linear Gaussian model**: Low-dim  $\mathbf{z}_n \in \mathbb{R}^K$  mapped to high-dim  $\mathbf{x}_n \in \mathbb{R}^D$  via  $\mathbf{W} \in \mathbb{R}^{D \times K}$
- The complete data log-likelihood for this model will be

$$\log p(\mathbf{X}, \mathbf{Z} | \mathbf{W}, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \mathbf{W}, \sigma^2) = \log \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) p(\mathbf{z}_n) = \sum_{n=1}^N \{ \log p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2) + \log p(\mathbf{z}_n) \}$$

- Plugging in the expressions for  $p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma^2)$  and  $p(\mathbf{z}_n)$  and simplifying

$$CLL = - \sum_{n=1}^N \left\{ \frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|\mathbf{x}_n\|^2 - \frac{1}{\sigma^2} \mathbf{z}_n^\top \mathbf{W}^\top \mathbf{x}_n + \frac{1}{2\sigma^2} \text{tr}(\mathbf{z}_n \mathbf{z}_n^\top \mathbf{W}^\top \mathbf{W}) + \frac{1}{2} \text{tr}(\mathbf{z}_n \mathbf{z}_n^\top) \right\}$$

- Expected CLL will require replacing  $\mathbf{z}_n$  by  $\mathbb{E}[\mathbf{z}_n]$  and  $\mathbf{z}_n \mathbf{z}_n^\top$  by  $\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top]$ 
  - These expectations can be easily obtained from the posterior  $p(\mathbf{z}_n | \mathbf{x}_n)$  (computed in E step)
- The M step maximizes the expected CLL w.r.t. the parameters ( $\mathbf{W}, \sigma^2$  in this case)

# Online or Incremental EM

- Needn't compute  $p(\mathbf{z}_n|\mathbf{x}_n)$  for every  $\mathbf{x}_n$  in each EM iteration (computational/storage efficiency)
  - Recall that the expected CLL is often a sum over all data points

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

- Can compute this quantity **recursively** using **small minibatches** of data

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

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

- Requires computing  $p(\mathbf{z}_n|\mathbf{x}_n)$  only for data in current mini-batch (computational/storage efficiency)
- MLE on above  $\mathcal{Q}_t$  can be shown to be equivalent to a simple **recursive updates for  $\Theta$**

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

# Some Applications of EM

- Mixture of (multivariate) Gaussians/Bernoullis, multinoullis, Mixture of experts models
- Problems with missing labels/features (treat these as latent variables)
- Note that EM not only gives estimates of the parameters  $\Theta$  but also infers latent variables  $\mathbf{Z}$
- **Hyperparameter estimation** in probabilistic models (an alternative to MLE-II)
  - We've already seen MLE-II where we did MLE on marginal likelihood, e.g., for linear regression

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

- As an alternative, can treat  $\mathbf{w}$  as a latent variable and  $\beta, \lambda$  as parameters and **use EM to learn these**
- Note: In this case, the latent variable  $\mathbf{w}$  is not “local” (but EM still applies)
- E step computes posterior  $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$  assuming  $\beta, \lambda$  fixed from the previous M step
- M step maximizes  $\mathbb{E}[\log p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda)] = \mathbb{E}[\log p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \beta) + \log p(\mathbf{w}|\lambda)]$  w.r.t.  $\lambda, \beta$ 
  - This requires using expectations of quantities like  $\mathbf{w}$  and  $\mathbf{w}\mathbf{w}^\top$  which can be obtained easily from the posterior  $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$  which we compute in the E step

# The EM Algorithm: Some Comments

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

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

.. which can be approximated, e.g., using Monte-Carlo expectation (called 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 \(ECM\)](#) algorithm
- Other advanced probabilistic inference algorithms are based on ideas similar to EM
  - E.g., [Variational Bayesian \(VB\)](#) inference