

Approximate Inference Using Expectation Maximization

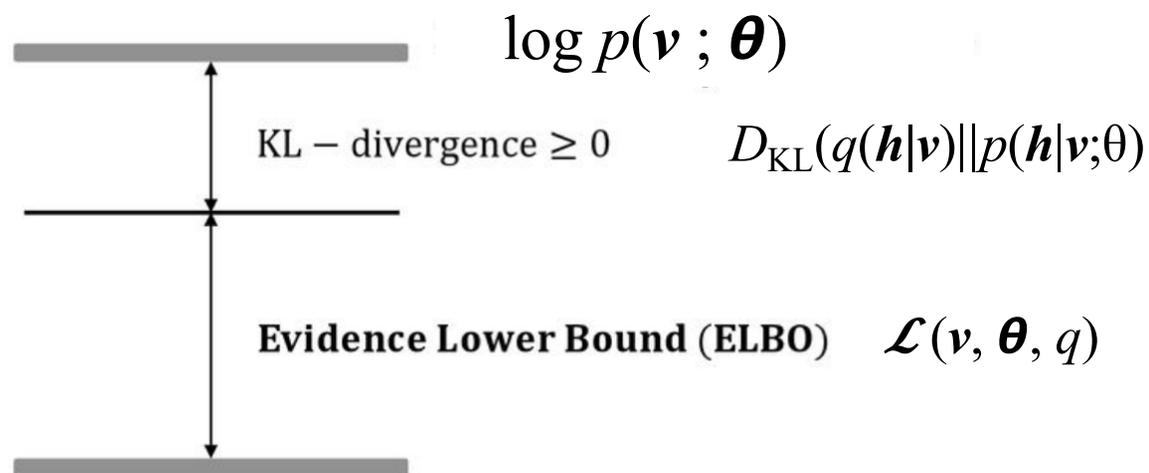
Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Approximate Inference

- Task of Inference
- Intractability in Inference
 1. Inference as Optimization
 2. Expectation Maximization
 3. MAP Inference and Sparse Coding
 4. Variational Inference and Learning
 5. Learned Approximate Inference

Approximate Inference

- Inference problem $p(\mathbf{v}; \boldsymbol{\theta})$ converted into optimization problem of finding a tractable q



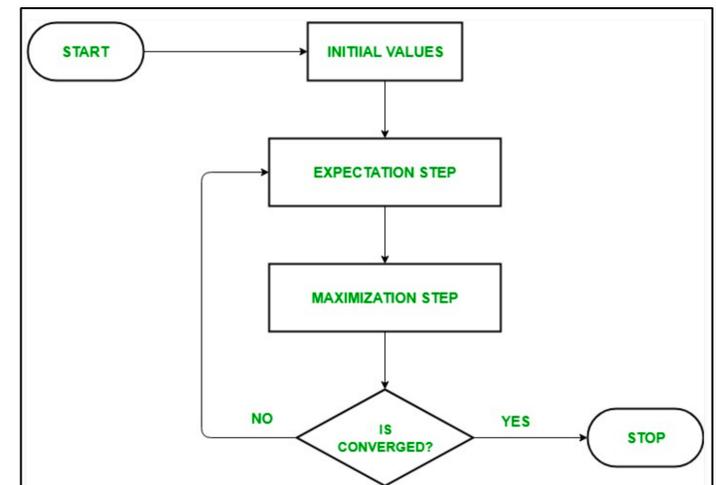
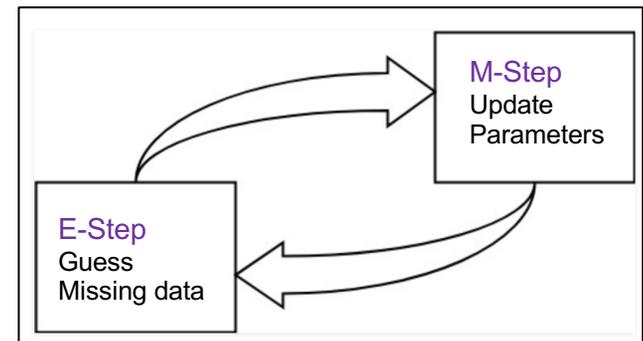
- Expectation Maximization Approach

E-step: maximize \mathcal{L} wrt q

M-step: maximize \mathcal{L} wrt $\boldsymbol{\theta}$

Expectation Maximization

- A training algorithm for models with latent variables
- Unlike other approximate inference algorithms, EM is not approximate inference, but an approach to learning with an approximate posterior
- Alternates two steps until convergence
 - E-Step (Expectation Step)
 - M-Step (Maximization Step)



The E Step (Expectation)

- Let $\theta^{(0)}$ denote value of parameters at the beginning of the step
- Set $q(h^{(i)} | \mathbf{v}) = p(h^{(i)} | \mathbf{v}^{(i)}; \theta^{(0)})$ for all indices i of the training examples $\mathbf{v}^{(i)}$ we want to train on
 - (both batch and minibatch variants are valid)
 - By this we mean q is defined in terms of the *current* parameter value of $\theta^{(0)}$;
 - if we vary θ then $p(\mathbf{h} | \mathbf{v}; \theta)$ will change but $q(\mathbf{h} | \mathbf{v})$ will remain equal to $p(\mathbf{h} | \mathbf{v}; \theta^{(0)})$

Note: $\theta^{(0)}$ is the current parameter value.

It is updated in the M-Step.

At the next E-step we update q based on the new value of p as determined by the new $\theta^{(0)}$

The M-Step (Maximization)

- Completely or partially maximize

$$\sum_i \mathcal{L}(v^{(i)}, \theta, q)$$

training examples i

– wrt using θ your optimization algorithm of choice

- This can be viewed as a coordinate ascent algorithm to maximize \mathcal{L}
- On one step, we maximize \mathcal{L} with respect to q , and on the other, we maximize \mathcal{L} wrt θ

Exact and Approximate Inference

- Although E-step involves exact inference, EM uses approximate inference in some sense
- M-step assumes same value of q for all θ
 - Introduces gap between \mathcal{L} and true $\log p(\mathbf{v})$ as the M-step moves further away from $\theta^{(0)}$ used in E-step
- Fortunately, E-step reduces gap to zero again as we enter loop next time

SGD is a special case of EM

- Stochastic gradient ascent on latent variable models is a special case of EM where the M step consists of taking a single gradient step
- Other variants of the EM algorithm can make much larger steps
- For some model families, the M step can even be performed analytically, jumping all the way to the optimal solution for θ given the current q

First Insight due to EM: Max. Likelihood

- In learning model parameters are updated to improve likelihood
 - Missing variables provided by estimated posterior distribution
- This insight is not unique to EM
 - E.g., SGD to maximize log-likelihood has same property
 - Log-likelihood gradient computations require taking expectations wrt the posterior distribution over hidden units

Second insight of EM

- We can continue to use one value of q even after we have moved to a different value of θ
- This insight used throughout ML to derive large M-step updates
- Deep learning models are too complex to admit a tractable solution for an optimal large M-step update, so this second insight which is more unique to EM is rarely used