

Variational Inference and Learning

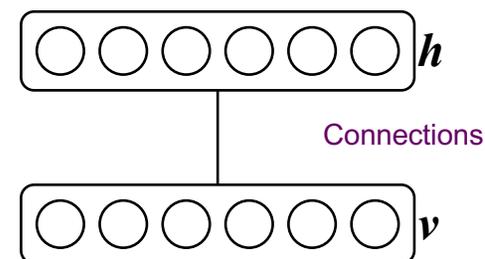
Sargur N. Srihari
srihari@cedar.buffalo.edu

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

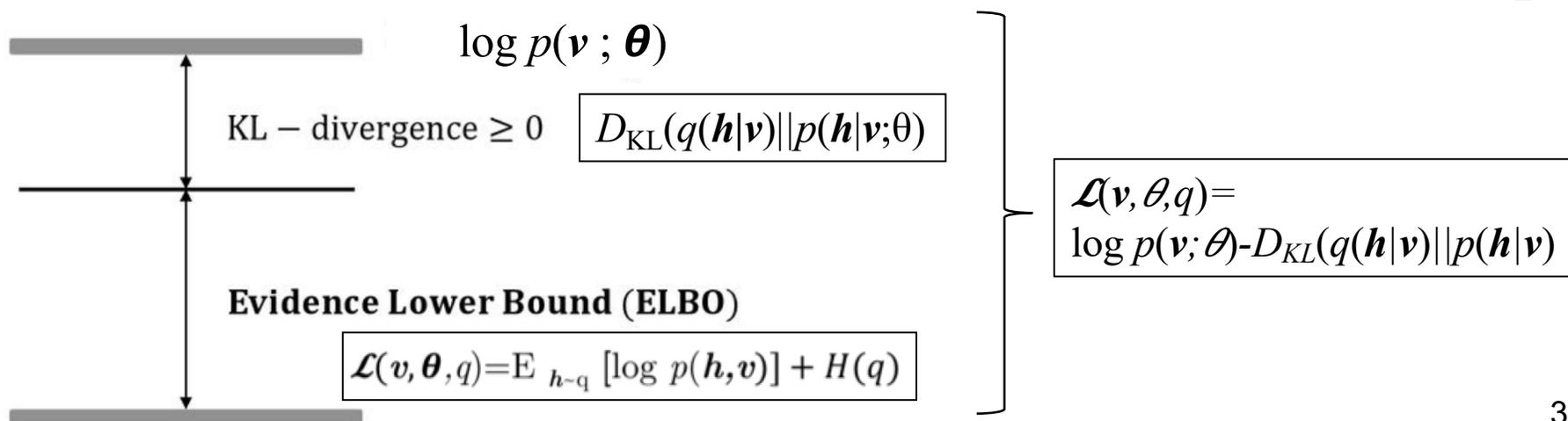
Setup for Variational Inference

- Observed variables \mathbf{v} and latent \mathbf{h}
- Known joint distribution $p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})$
- We would like to compute $\log p(\mathbf{v}; \boldsymbol{\theta})$ so as to choose $\boldsymbol{\theta}$ to increase it



– It is too costly to marginalize \mathbf{h} using $p(\mathbf{v}; \boldsymbol{\theta}) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})$

- Instead compute a lower bound $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$ on $p(\mathbf{v}; \boldsymbol{\theta})$
- Approx inference is approximate optimization to find q



Summary of Approximate Inference

- ELBO $\mathcal{L}(\mathbf{v}, \theta, q)$ is a lower bound on $\log p(\mathbf{v}; \theta)$
 - Inference viewed as maximizing \mathcal{L} wrt q
 - Learning viewed as maximizing \mathcal{L} wrt θ
- 1. EM allows large learning steps with a fixed q
- 2. Learning based on MAP inference
 - Learn using a point estimate of $p(\mathbf{h} | \mathbf{v})$ rather than inferring entire distribution
- 3. Next: more general approach to variational learning

Common Variational Approaches

1. *Mean field approach*: q is a factorial distribution

$$q(h|\mathbf{v}) = \prod_i q(h_i|\mathbf{v})$$

2. *Structured variational approach*:

– Impose any PGM on q

Variational: Discrete vs Continuous

- Specifying how to factorize means:
 - Discrete case:
 - Optimize finite no. of variables describing q distribution
 - Continuous case:
 - Use calculus of variations over a space of functions
 - Determine which function should be used to represent q
- Although calculus of variations is not used in discrete case, it is still called as variational
- Calculus of variations is powerful
 - Only need specify factorization, not design of q

Maximizing \mathcal{L} means fitting q to p

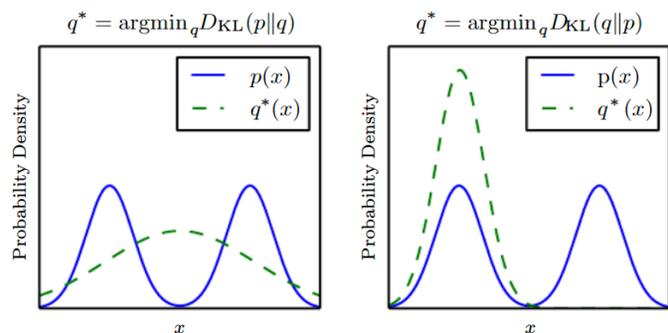
- Maximizing \mathcal{L} wrt q is equivalent to minimizing $D_{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$
- In this sense, we are fitting q to p
- However, we are doing so with the opposite direction of the KL divergence than we are used to using for fitting an approximation
- When we use maximum likelihood learning to fit a model to data, we minimize $D_{KL}(p_{data}||p_{model})$

Optimization matches low probabilities

- Maximum likelihood encourages model to have high probability wherever data has high probability
- Optimization-based inference encourages q to have low probability wherever true posterior has low probability
 - Illustrated in figure next

Asymmetry of KL divergence

- We wish to approximate p with q
 - We have a choice of using $D_{KL}(p||q)$ or $D_{KL}(q||p)$
 - Bimodal Gaussian for p and Gaussian for q



- (a) Effect of minimizing $D_{KL}(p||q)$
 - select q that has high probability where p has high probability
 - q chooses to blur the two modes together
- (b) Effect of minimizing $D_{KL}(q||p)$
 - select q that has low probability where p has low probability
 - Avoids putting probability mass in low probability areas between modes

Choice of KL Direction

- Directions of KL have different properties
- Choice depends on application
 - For inference optimization use $D_{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$ for computational reasons.
 - Computing $D_{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$ involves evaluating expectations wrt q , so by designing q to be simple, we can simplify the required expectations
 - Opposite direction of KL requires computing expectations wrt true posterior
 - Because the form of the true posterior is determined by the choice of model, we cannot design a reduced-cost approach to computing $D_{KL}(p(\mathbf{h}|\mathbf{v})||q(\mathbf{h}|\mathbf{v}))$ exactly

Topics in Variational Inference

1. Discrete Latent Variables
 2. Calculus of Variations
 3. Continuous Latent Variables
 4. Interactions between Learning and Inference
- We discuss a subset of the above topics here

Calculus of Variations

- A function of a function f is known as a functional $J[f]$
- We can take *functional derivatives* wrt to individual values of of the function $f(\mathbf{x})$ at any specific value of \mathbf{x}
- Denoted by $\frac{\delta}{\delta f(\mathbf{x})} J$

Functional Derivative Identity

- For differentiable functions $f(\mathbf{x})$ and differentiable functions $g(y, \mathbf{x})$ with continuous derivatives

$$\frac{\delta}{\delta f(\mathbf{x})} \int g(f(\mathbf{x}), \mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial y} g(f(\mathbf{x}), \mathbf{x})$$

- Intuition: $f(\mathbf{x})$ is an infinite vector indexed by \mathbf{x}
 - Identity is same as for a vector $\theta \in R^n$ indexed by positive integers

$$\frac{\partial}{\partial \theta_i} \sum_j g(\theta_j, j) = \frac{\partial}{\partial \theta_i} g(\theta_i, i)$$

- More general is the Euler-Lagrange equation
 - Allows g to depend on derivatives of f as well as value of f , but not needed for deep learning

Optimization wrt a vector

- Procedure to optimize a function wrt a vector
 - Take the gradient of the function wrt the vector
 - Solve for the point where every element of the gradient is equal to zero
- Procedure to optimize a functional
 - Solve for the function where the functional derivative at every point is equal to zero

Example of Optimizing a Functional

- Find maximum entropy distribution over $x \in R$
- Entropy of $p(x)$ is defined as $H[p] = -E_x \log p(x)$
- For continuous values expectation is an integral

$$H[p] = \int p(x) \log p(x) dx$$

Lagrangian constraints for maxent

1. To ensure a distribution, $p(x)$ integrates to 1
 2. Since entropy increases unbounded with variance, seek distribution with variance σ^2
 3. Since problem is undetermined
 - as distribution can be arbitrarily shifted without changing entropy we fix the mean to be μ
- Lagrangian functional is

$$L[p] = \lambda_1 \left(\int p(x) dx - 1 \right) + \lambda_2 (E[x] - \mu) + \lambda_3 (E[(x - \mu)^2] - \sigma^2) + H[p]$$

$$L[p] = \int (\lambda_1 p(x) dx + \lambda_2 p(x)x + \lambda_3 p(x)(x - \mu)^2 - p(x) \log p(x)) dx - \lambda_1 - \mu \lambda_2 - \sigma^2 \lambda_3$$

Functional derivative of Lagrangian

- To minimize Lagrangian wrt p , we set the functional derivatives equal to 0:

$$L[p] = \int (\lambda_1 p(x) dx + \lambda_2 p(x)x + \lambda_3 p(x)(x - \mu)^2 - p(x) \log p(x)) dx - \lambda_1 - \mu \lambda_2 - \sigma^2 \lambda_3$$

$$\boxed{\frac{\delta}{\delta f(\mathbf{x})} \int g(f(\mathbf{x}), \mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial y} g(f(\mathbf{x}), \mathbf{x})}$$

$$y=p(x)$$

$$g_1(y,x)=y$$

$$g_2(y,x)=yx$$

$$g_3(y,x)=y(x-\mu)^2$$

$$g_4(y,x)=y \log y$$

$$\frac{d}{dy}(y \log y) = 1 + \log y$$

$$\forall x, \frac{\delta}{\delta p(x)} L = \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1 - \log p(x) = 0$$

- This condition tells us functional form of $p(x)$
 - By algebraic rearrangement

$$p(x) = \exp(\lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1)$$

- We never assumed directly that $p(x)$ has this form

Choosing λ values

- We must choose λ values so that all constraints are satisfied
- We may set $\lambda_1=1-\log \sigma\sqrt{2\pi}$, $\lambda_2=0$, $\lambda_3=-1/2\sigma^2$ to obtain $p(x)=N(x:\mu,\sigma^2)$
- Because the normal distribution has maximum entropy, we impose the least possible structure by using maximum entropy

Continuous Latent Variables

- When our graphical model contains continuous latent variables, we can perform variational inference and learning by maximizing \mathcal{L}
- We must now use calculus of variations when maximizing \mathcal{L} with respect to $q(\mathbf{h}|\mathbf{v})$
- Not necessary for practitioners to solve calculus of variations problems
- Instead there is a general equation for mean-field fixed point updates

Optimal $q(h_i|\mathbf{v})$ for Mean Field

- If we make the mean field approximation

$$q(\mathbf{h}|\mathbf{v}) = \prod_i q(h_i|\mathbf{v})$$

- The optimal $q(h_i|\mathbf{v})$, if we fix $q(h_j|\mathbf{v})$ for all $j \neq i$, is obtained by normalizing the

equation
$$q(h_i|\mathbf{v}) = \exp(E_{h_{-i} \sim q(h_{-i}|\mathbf{v})} \log p(\mathbf{v}, \mathbf{h}))$$

So long as p does not assign 0 probability to any joint configuration of variables.

Carrying out expectation yields correct functional form of $q(h_i|\mathbf{v})$

Designed to be iteratively applied for each value of i until convergence

Ex: Latent $h \in \mathbb{R}^2$ and one visible variable v