

- HW4 due Wed
- Quiz due Fri

COMPSCI 688: Probabilistic Graphical Models

Lecture 18: Variational Inference

Dan Sheldon

Manning College of Information and Computer Sciences
University of Massachusetts Amherst

Partially based on materials by Benjamin M. Marlin (marlin@cs.umass.edu) and Justin Domke (domke@cs.umass.edu)

1 / 22

Models

BNs, MRFs

learning + inference

Approx. inference

MCMC

Variational inference $\left\{ \begin{array}{l} \text{learning} \\ \text{models: VAEs} \end{array} \right.$

Additional

- GPs?

- normalizing flows

2 / 22

Variational Inference (VI) Overview

- ▶ Variational inference is an approximate inference approach (alternative to MCMC)
- ▶ Variational inference is at the core of a large family of techniques, **all of which start with the same mathematical idea**
 - ▶ mean-field and structured VI
 - ▶ black-box VI ✓
 - ▶ expectation maximization (EM) ✓
 - ▶ variational EM ✓
 - ▶ variational Bayes
 - ▶ variational auto-encoders ✓
 - ▶ loopy belief propagation and advanced message-passing algorithms

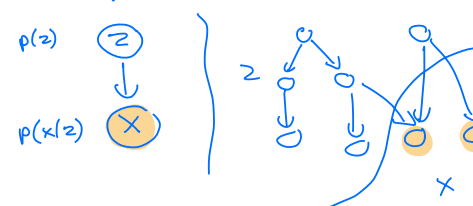
3 / 22

Problem Setting

Assume we have an unnormalized probability model over z . Two examples:

1. Bayesian model $p(z|x)$ for latent z , observed x , unknown $p(x)$
2. Unnormalized model $p(z) = \frac{1}{Z} \tilde{p}(z)$ with unknown Z (e.g., loopy MRF)

1. Bayesian



$\tilde{p}(z)$ - can evaluate at any z

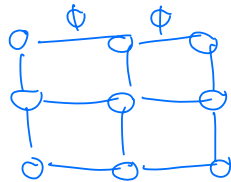
"hard" want: $p(z|x) = \frac{\tilde{p}(z) p(x|z)}{p(x)}$

← "easy" (pointing to $\tilde{p}(z)$)
← "hard" (pointing to $p(x)$)

usually $= \frac{p(z)p(x|z)}{p(x)}$

4 / 22

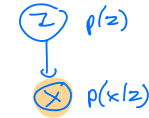
2. Unnormalized model "hard"



$$p(z) = \frac{1}{Z} \tilde{p}(z) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \phi_c(z_c)$$

"hard" "easy"
 ↑ ↑

Problem Setting



For concreteness, henceforth we'll assume the Bayesian model setting:

- ▶ $p(z, x) = p(z)p(x|z)$ easy to compute
- ▶ We observe x , but not z
- ▶ We want to approximate

$$p(z|x) = \frac{p(z, x)}{p(x)}$$

← unnormalized density $\tilde{p}(z)$

but don't know the normalization constant $p(x)$

General Strategy

1. Let $q_\phi(z)$ be a "simple" distribution from some family with parameters ϕ
(Gaussian, diagonal covariance, Gaussian dist over z w/ independent components)
2. Try to optimize

$$\min_{\phi} \text{KL}(q_{\phi}(z) \| p(z|x))$$

"posterior"

("reverse KL")

Then use $q_\phi(z)$ in place of $p(z|x)$

Why use VI?

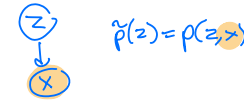


- ▶ Can often get reasonable approximations faster than MCMC
 - get "whole" (approx.) dist instead of samples
- ▶ Gives a bound on $p(x)$ (or " Z "), useful for learning (more later)

Downside: "bias" due to simplifying assumption

The ELBO Decomposition

Big Idea: ELBO Decomposition



This is the math trick that is at the heart of all VI methods:

$$\log p(x) = \underbrace{\sum_z q_\phi(z) \log \frac{p(z, x)}{q_\phi(z)}}_{\text{ELBO}(q_\phi(z) \parallel p(z, x))} + \underbrace{\sum_z q_\phi(z) \log \frac{q_\phi(z)}{p(z|x)}}_{\text{KL}(q_\phi(z) \parallel p(z|x))}$$

↑
unnormalized

- ▶ ELBO: "Evidence Lower Bound" (will explain later)
- ▶ KL: what we want to minimize

Derivation

Claim:

$$\log p(x) = \sum_z q_\phi(z) \log \frac{p(z, x)}{q_\phi(z)} + \sum_z q_\phi(z) \log \frac{q_\phi(z)}{p(z|x)}$$

$p(z|x) = \frac{p(z, x)}{p(x)}$

Proof. Start with RHS and simplify:

$$\begin{aligned} \text{RHS} &= \sum_z q_\phi(z) \left[\log p(z, x) - \log q_\phi(z) + \log q_\phi(z) - \log p(z|x) \right] \\ &= \sum_z q_\phi(z) \left[\log p(z, x) - \log p(z|x) \right] \\ &= \sum_z q_\phi(z) \log p(x) \\ &= \log p(x) \sum_z q_\phi(z) \\ &= \log p(x) \end{aligned}$$

ELBO Significance

$$\log p(x) = \underbrace{\sum_z q_\phi(z) \log \frac{p(z, x)}{q_\phi(z)}}_{\text{ELBO}(q_\phi(z) \parallel p(z, x))} + \underbrace{\sum_z q_\phi(z) \log \frac{q_\phi(z)}{p(z|x)}}_{\text{KL}(q_\phi(z) \parallel p(z|x))}$$

↑
unnormalized

↑
target "hard" (normalized)

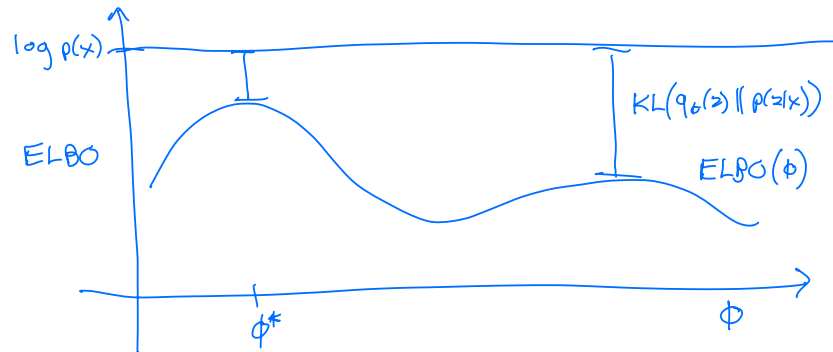
VI: max ELBO($q_\phi \parallel p$) wrt ϕ

1. KL is "hard": can't evaluate the *normalized* distribution $p(z|x)$
2. ELBO is "easy" (ish). Uses *unnormalized* distribution $p(z, x)$. Can often evaluate or approximate it, e.g., by Monte Carlo:

$$\text{sample } z^{(1)}, \dots, z^{(N)} \sim q_\phi(z), \text{ then compute } \frac{1}{N} \sum_{i=1}^N \log \frac{p(z^{(i)}, x)}{q_\phi(z^{(i)})}$$

3. KL is non-negative
4. Therefore $\log p(x) \geq \text{ELBO}$ ("Evidence lower bound")
5. Therefore, choosing ϕ to maximize the ELBO is the same as choosing ϕ to minimize the KL (since $\log p(x)$ is constant with respect to ϕ)

ELBO Interpretation: Picture

 $\phi \mapsto q_\phi$ Fixed $p(z, x)$ 

13 / 22

Variational Inference

14 / 22

Uses of VI

There are two different uses of VI

target simple

1. Approximate a posterior distribution: $p(z|x) \approx q_\phi(z)$
2. Bound the log-likelihood: $\log p_\theta(x) \geq \text{ELBO}(q_\phi(z) || p_\theta(z, x))$, usually in a learning procedure for $p_\theta(x)$ (details to come)

MLE: find θ to max $\log p_\theta(x)$

15 / 22

Basic VI Algorithm

1. Input: $p(z, x)$ and fixed x Want: $p(z|x)$
2. Choose some approximating family $q_\phi(z)$
3. Maximize $\text{ELBO}(q_\phi(z) || p(x, z))$ wrt ϕ
4. Use $q_\phi(z)$ as a proxy for $p(z|x)$

- Many choices for
- ▶ Target $p(z|x)$
 - ▶ Approximating family q_ϕ
 - ▶ How to estimate ELBO
 - ▶ How to do optimization

16 / 22

ELBO Intuition

$$= \sum_z q_\phi(z) \log \frac{p(z,x)}{q_\phi(z)}$$

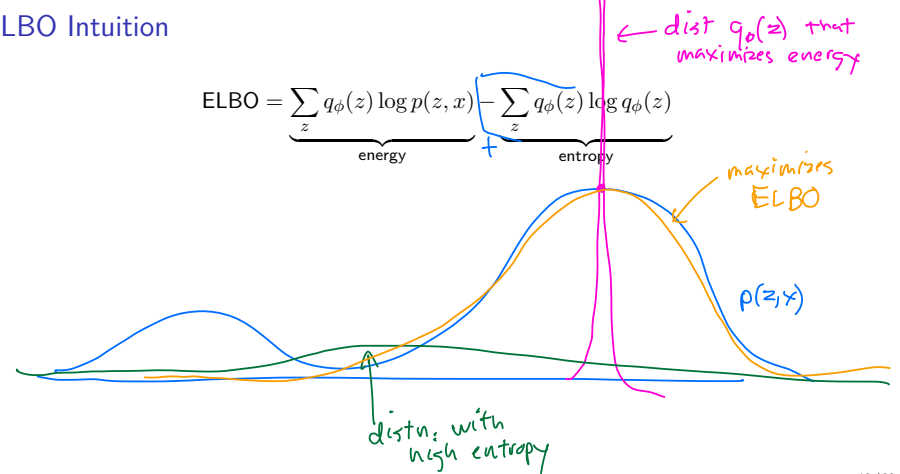
$$\text{ELBO} = \underbrace{\sum_z q_\phi(z) \log p(z,x)}_{\text{energy}} - \underbrace{\sum_z q_\phi(z) \log q_\phi(z)}_{\text{entropy}}$$

energy + entropy

- ▶ energy term encourages $q_\phi(z)$ to be high where $p(z|x)$ is high
- ▶ entropy term encourages $q_\phi(z)$ to be spread out

ELBO Intuition

$$\text{ELBO} = \underbrace{\sum_z q_\phi(z) \log p(z,x)}_{\text{energy}} + \underbrace{-\sum_z q_\phi(z) \log q_\phi(z)}_{\text{entropy}}$$



Variational Learning

MLE: find θ to max $\log p_\theta(\text{data})$



$$\max_{\theta} \log p_{\theta}(x)$$

$$\max_{\theta} \log \sum_z p_{\theta}(z,x)$$

instead

$$\max_{\theta} \text{ELBO}(\theta) \leq \log p_{\theta}(x)$$

Expectation Maximization (EM): VI + Learning

EM is a classical algorithm for maximum-likelihood learning with latent variables

Goal: choose θ to maximize $\log p_{\theta}(x) = \log \sum_z p_{\theta}(z,x)$ given observed x

Usual lower-bound derivation

$$\begin{aligned} \log p_{\theta}(x) &= \log \sum_z p_{\theta}(z,x) \\ &= \log \sum_z \underbrace{q_{\theta}(z)}_{\text{q}} \cdot \underbrace{\frac{p_{\theta}(z,x)}{q_{\theta}(z)}}_{\text{E}} \\ &\geq \sum_z q_{\theta}(z) \log \frac{p_{\theta}(z,x)}{q_{\theta}(z)} \\ &= \text{ELBO} \end{aligned}$$

(Jensen's inequality)

Learn θ

EM Algorithm. Init $\theta \mapsto p_{\theta}(z,x)$

- ▶ Set $q(z) = p_{\theta}(z|x)$ (maximize ELBO wrt q)
- ▶ Maximize $\sum_z q(z) \log \frac{p_{\theta}(z,x)}{q(z)}$ wrt θ
- ▶ Repeat $\underbrace{\sum_z q(z) \log \frac{p_{\theta}(z,x)}{q(z)}}_{\text{ELBO}}$

Gives local maximum of $\log p_{\theta}(x)$ wrt θ

Variational EM

$$\max_q \text{ELBO}(q(z) \parallel p(z, x)) \Leftrightarrow \min_q \text{KL}(q(z) \parallel p(z|x))$$

It is not always possible or practical to compute $p_\theta(z|x)$ exactly in EM.

Variational EM is an extension where the ELBO is maximized jointly with respect to the parameters ϕ of the approximating distribution and parameters θ of the model ("simultaneous inference and learning")

Goal: choose θ to maximize $\log p_\theta(x) = \log \sum_z p_\theta(z, x)$ given observed x .

Define

$$\mathcal{L}(\phi, \theta) = \text{ELBO}(q_\phi(z) \parallel p_\theta(z, x)) = \sum_z q_\phi(z) \log \frac{p_\theta(z, x)}{q_\phi(z)} \leq \log p_\theta(x)$$

then jointly optimize $\mathcal{L}(\phi, \theta)$ with respect to ϕ and θ , e.g.:

- ▶ (Stochastic) gradient ascent
- ▶ Alternating (partial) optimization steps \rightarrow Traditional EM
 $\begin{array}{l} 1. \text{ max fully wrt } q \Rightarrow q(z) = p(z|x) \\ 2. \text{ max wrt } \theta \end{array}$