

Introduction oooooo	KL Minimization and ELBO oooooo	Variational Inference oooo	Variational Learning oooo
------------------------	------------------------------------	-------------------------------	------------------------------

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 / 24

Introduction ●ooooo	KL Minimization and ELBO oooooo	Variational Inference oooo	Variational Learning oooo
------------------------	------------------------------------	-------------------------------	------------------------------

Introduction

2 / 24

Introduction ●ooooo	KL Minimization and ELBO oooooo	Variational Inference oooo	Variational Learning oooo
------------------------	------------------------------------	-------------------------------	------------------------------

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 / 24

Introduction ●ooooo	KL Minimization and ELBO oooooo	Variational Inference oooo	Variational Learning oooo
------------------------	------------------------------------	-------------------------------	------------------------------

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)

4 / 24

Introduction  
oooo●ooo

KL Minimization and ELBO  
ooooooo

Variational Inference  
ooooo

Variational Learning  
oooo

5 / 24

## General Strategy

1. Let  $q_\phi(z)$  be a “simple” distribution from some family with parameters  $\phi$
2. Try to optimize

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

where  $D$  is some “distance”. Then use  $q_\phi(z)$  in place of  $p(z|x)$

7 / 24

Introduction  
oooo●ooo

KL Minimization and ELBO  
ooooooo

Variational Inference  
ooooo

Variational Learning  
oooo

## 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

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

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

6 / 24

## Why use VI?

- Can often get reasonable approximations faster than MCMC
- Gives a bound on  $p(x)$  (or “ $Z$ ”), useful for learning (more later)

8 / 24

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

9 / 24

## KL Minimization and ELBO

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

10 / 24

**Idea: “Distance” Minimization**

We want  $q_\phi(z) \approx p(z|x)$ .

**Idea:** define a “distance”  $D(q_\phi(z) \| p(z|x))$  and choose  $\phi$  to minimize it.

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

11 / 24

## KL Divergence

A widely used “distance” between distributions is the Kullback-Leibler divergence:

$$\text{KL}(q||p) = \int q(\mathbf{z}) \log \left( \frac{q(\mathbf{z})}{p(\mathbf{z})} \right) d\mathbf{z}$$

It is a *divergence* because it only satisfies some properties of a distance metric. It satisfies:

- ▶  $\text{KL}(q||p) \geq 0$  for all  $q$  and  $p$
- ▶  $\text{KL}(q||p) = 0$  if and only if  $q = p$

It does **not** satisfy:

- ▶  $\text{KL}(q||p) = \text{KL}(p||q)$  for all  $q, p$
- ▶  $\text{KL}(q||p) \leq \text{KL}(q||s) + \text{KL}(s||p)$  for all  $q, p, s$

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

12 / 24

## Big Idea: ELBO Decomposition

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

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

- ▶ ELBO: “Evidence Lower BOund” (will explain later)
- ▶ KL: what we want to minimize

Introduction oooooooo	KL Minimization and ELBO oooooo●	Variational Inference ooooo	Variational Learning oooo
<b>Derivation</b>			
<b>Claim:</b> $\log p(x) = \sum_z q_\phi(z) \log \frac{p(z, x)}{q_\phi(z)} + \sum_z q_\phi(z) \log \frac{q_\phi(z)}{p(z x)}$			
<b>Proof.</b> Start with RHS and simplify: $\begin{aligned} \text{RHS} &= \sum_z q_\phi(z) [\log p(z, x) - \log q_\phi(z) + \log q_\phi(z) - \log p(z x)] \\ &= \sum_z q_\phi(z) [\log p(z, x) - \log p(z, x) + \log p(x)] \\ &= \sum_z q_\phi(z) \log p(x) \\ &= \log p(x) \sum_z q_\phi(z) \\ &= \log p(x) \end{aligned}$			
13 / 24			

Introduction oooooooo	KL Minimization and ELBO ooooo●○	Variational Inference ooooo	Variational Learning oooo
<b>ELBO Significance</b>			
$\log p(x) = \underbrace{\sum_z q_\phi(z) \log \frac{p(z, x)}{q_\phi(z)}}_{\text{ELBO}(q_\phi(z) \parallel p(z, x))} + \underbrace{\sum_z q_\phi(z) \log \frac{q_\phi(z)}{p(z x)}}_{\text{KL}(q_\phi(z) \parallel p(z x))}$			
<ol style="list-style-type: none"> <li>1. KL is "hard": can't evaluate the <i>normalized</i> distribution <math>p(z x)</math></li> <li>2. ELBO is "easy"(ish). Uses <i>unnormalized</i> distribution <math>p(z, x)</math>. Can often evaluate or approximate it, e.g., by Monte Carlo:            sample <math>z^{(1)}, \dots, z^{(N)} \sim q_\phi(z)</math>, then compute <math>\frac{1}{N} \sum_{i=1}^N \log \frac{p(z^{(i)}, x)}{d_\phi(z^{(i)})}</math></li> <li>3. KL is non-negative</li> <li>4. Therefore <math>\log p(x) \geq \text{ELBO}</math> ("Evidence lower bound")</li> <li>5. Therefore, choosing <math>\phi</math> to maximize the ELBO <b>is the same</b> as choosing <math>\phi</math> to minimize the KL (since <math>\log p(x)</math> is constant with respect to <math>\phi</math>)</li> </ol>			
14 / 24			

Introduction oooooooo	KL Minimization and ELBO oooooo●	Variational Inference ooooo	Variational Learning oooo
<b>ELBO Interpretation: Picture</b>			
15 / 24			

Introduction oooooooo	KL Minimization and ELBO oooooo	Variational Inference ●oooo	Variational Learning oooo
<b>Variational Inference</b>			
16 / 24			

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

## Uses of VI

There are two different uses of VI

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

17 / 24

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

## Basic VI Algorithm

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

Many choices for

- ▶ Approximating family  $q_\phi$
- ▶ How to estimate ELBO
- ▶ How to do optimization

18 / 24

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

## ELBO Intuition

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

- ▶ 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

19 / 24

Introduction

KL Minimization and ELBO

Variational Inference

Variational Learning

## ELBO Intuition

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

20 / 24

Introduction ooooooo  
KL Minimization and ELBO ooooooo  
Variational Inference ooooo  
Variational Learning ●ooo

## Variational Learning

21 / 24

Introduction ooooooo  
KL Minimization and ELBO ooooooo  
Variational Inference ooooo  
Variational Learning o●oo

## Expectation Maximization (EM): VI + Learning

EM is a classical algorithm for maximum-likelihood learning with latent variables  
**Goal:** choose  $\theta$  to maximize  $\log p_\theta(x) = \log \sum_z p_\theta(z, x)$  given observed  $x$

**Usual lower-bound derivation**

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

**EM Algorithm**

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

**Gives local maximum of  $\log p_\theta(x)$  wrt  $\theta$**   
(Jensen's inequality)

22 / 24

Introduction ooooooo  
KL Minimization and ELBO ooooooo  
Variational Inference ooooo  
Variational Learning ooo●

23 / 24

Introduction ooooooo  
KL Minimization and ELBO ooooooo  
Variational Inference ooooo  
Variational Learning ooo●

## Variational EM

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  $\phi$  of the approximating distribution and parameters  $\theta$  of the model ("simultaneous inference and learning")

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

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

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

- ▶ (Stochastic) gradient ascent
- ▶ Alternating (partial) optimization steps

24 / 24