

Introduction ooooooo	KL Minimization and ELBO ooooooo	Variational Inference ooooo	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 ooooooo	Variational Inference ooooo	Variational Learning oooo
------------------------	-------------------------------------	--------------------------------	------------------------------

**Introduction**

2 / 24

Introduction ○●ooooo	KL Minimization and ELBO ooooooo	Variational Inference ooooo	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  $\leftarrow$  HW5
  - ▶ loopy belief propagation and advanced message-passing algorithms

3 / 24

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

**Problem Setting**

$\tilde{p}(z)$

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

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

Have:  $p(z|x) := \tilde{p}(z)$

easy

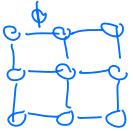
hard

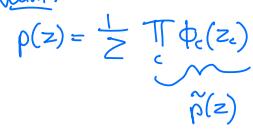
4 / 24

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

Setting 2. unnormalized  $p(z) = \frac{1}{Z} \tilde{p}(z)$

e.g. MRF    want:  $p(z) = \sum_c \prod_c \phi_c(z_c)$





5 / 24

Introduction ooooo●    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

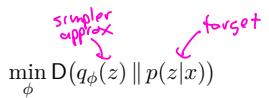
Introduction ooooo●    KL Minimization and ELBO ooooooo    Variational Inference ooooo    Variational Learning oooo

## General Strategy

Want:  $p(z|x) \approx q_\phi(z)$



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 ooooo●    KL Minimization and ELBO ooooooo    Variational Inference ooooo    Variational Learning oooo

## 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 ooooooo

KL Minimization and ELBO ●oooooo

Variational Inference ooooo

Variational Learning oooo

## KL Minimization and ELBO

9 / 24

Introduction ooooooo

KL Minimization and ELBO ●oooooo

Variational Inference ooooo

Variational Learning oooo

### 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.

space of all dists over Z

10 / 24

Introduction ooooooo

KL Minimization and ELBO ●oooooo

Variational Inference ooooo

Variational Learning oooo

## KL Divergence

$$\text{dive} \sum_z q(z) \log \frac{q(z)}{p(z)}$$

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

$$\text{KL}(q||p) = \int q(z) \log \left( \frac{q(z)}{p(z)} \right) dz = \mathbb{E}_{q(z)} \left[ \log \frac{q(z)}{p(z)} \right]$$

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

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

It does **not** satisfy:

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

11 / 24

Introduction ooooooo

KL Minimization and ELBO ●oooooo

Variational Inference ooooo

Variational Learning oooo

### Big Idea: ELBO Decomposition

$$\tilde{p}(z) = \frac{p(z, x)}{p(x)} \Leftarrow 'z'$$

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

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

↑  
approx unnormalized

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

12 / 24

Introduction ooooooo    KL Minimization and ELBO ooooooo    Variational Inference ooooo    Variational Learning oooo

### Derivation

**Claim:**

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

**Proof.** Start with RHS and simplify:

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

13 / 24

Introduction ooooooo    KL Minimization and ELBO ooooooo    Variational Inference ooooo    Variational Learning oooo

### ELBO Significance

$$= \mathbb{E}_{q_\phi(z)} \left[ \log \frac{p(z, x)}{q_\phi(z)} \right]$$

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

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

$\frac{q_\phi(z)}{p(z|x)}$   $\leftarrow \text{normalized}$

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:  
 sample  $z^{(1)}, \dots, z^{(N)} \sim q_\phi(z)$ , then compute  $\frac{1}{N} \sum_{i=1}^N \log \frac{p(z^{(i)}, x)}{d_\phi(z^{(i)})}$   
 3. KL is non-negative  
 4. Therefore  $\log p(x) \geq \text{ELBO}$  ("Evidence lower bound")  
 5. Therefore, choosing  $\phi$  to maximize the ELBO **is the same** as choosing  $\phi$  to minimize the KL (since  $\log p(x)$  is constant with respect to  $\phi$ )

14 / 24

Introduction ooooooo    KL Minimization and ELBO ooooooo    Variational Inference ooooo    Variational Learning oooo

### ELBO Interpretation: Picture

$q_\phi(z) \approx p(z|x)$

Model  $p(z|x)$

15 / 24

Introduction ooooooo    KL Minimization and ELBO ooooooo    Variational Inference ooooo    Variational Learning oooo

### Variational Inference

16 / 24

Introduction ooooooo    KL Minimization and ELBO ooooooo    Variational Inference o●ooo    Variational Learning oooo

## Uses of VI

There are two different uses of VI

1. Approximate a posterior distribution:  $p(z|x) \approx q_\phi(z)$  target simple
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 ooooooo    KL Minimization and ELBO ooooooo    Variational Inference o●ooo    Variational Learning oooo

## Basic VI Algorithm

target

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

- Model  $p(z, x)$
- Approximating family  $q_\phi$
- How to estimate ELBO
- How to do optimization

18 / 24

Introduction ooooooo    KL Minimization and ELBO ooooooo    Variational Inference o●ooo    Variational Learning oooo

## ELBO Intuition

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

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

- 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 ooooooo    KL Minimization and ELBO ooooooo    Variational Inference o●ooo    Variational Learning oooo

## ELBO Intuition

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

$p(z, x)$      $q_\phi(z)$      $p(z|x)$

dist  $q_\phi$  that maximizes energy

dist that maximizes entropy + energy

$\tilde{p}(z) = p(z|x)$

z

20 / 24

Introduction oooooooo

KL Minimization and ELBO oooooooo

Variational Inference ooooo

Variational Learning ●○○

**Variational Learning**

21 / 24

Introduction oooooooo

KL Minimization and ELBO oooooooo

Variational Inference ooooo

Variational Learning ○○○

**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**

$$\begin{aligned} \log p_\theta(x) &= \log \sum_z p_\theta(z, x) \\ &= \log \sum_z q_\phi(z) \cdot \frac{p_\theta(z, x)}{q_\phi(z)} \\ &\stackrel{\text{log concave}}{\geq} \sum_z q_\phi(z) \log \frac{p_\theta(z, x)}{q_\phi(z)} \end{aligned}$$

(Jensen's inequality)

**ELBO**

**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  $\theta$**

22 / 24

Introduction oooooooo

KL Minimization and ELBO oooooooo

Variational Inference ooooo

Variational Learning ○○○

23 / 24

Introduction oooooooo

KL Minimization and ELBO oooooooo

Variational Inference ooooo

Variational Learning ○○○●

**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