



# Brief Intro to Variational Inference

*WHY, WHAT and HOW*

*Only basic statistics background needed*

**(Nov 02, 2016)**

**YANG Jiancheng**



# Outline

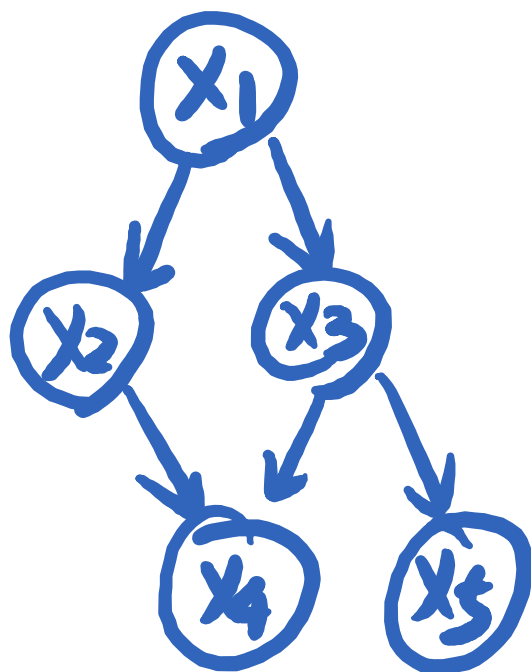
- **I. WHY**
  - **Probability Review**
  - **Graphic Models**
  - **The Point of Variational Inference**
- **II. WHAT**
  - **Ways to infer**
  - **Short Words**
- **III. HOW**
  - **Information and Entropy**
  - **KL divergence**
  - **Derivation**
  - **Example**



# • I. WHY

## • Probability Review

pdf / pmf = probability density / mass function



Bayes Rules

joint pdf

$$p(x|y)$$

conditional pdf

$$= \frac{p(x, y)}{p(y)}$$

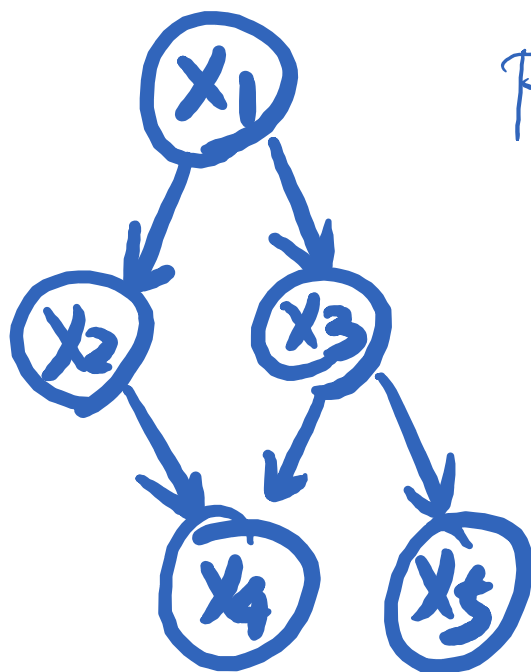
prior

$$\frac{p(x, y)}{\int_x p(x, y) dx}$$



# • I. WHY

## • Graphic Models



Joint pdf **EASY**

$$P(X_1, X_2, X_3, X_4, X_5) = P(X_5 | X_3) \cdot P(X_4 | X_2, X_3) \\ \cdot P(X_3 | X_1) \cdot P(X_2 | X_1) \cdot P(X_1)$$

Conditional pdf

$$P(X_3, X_4 | X_1, X_2, X_5) = \frac{P(X_1, X_2, X_3, X_4, X_5)}{\text{unknown} \leftarrow P(X_1, X_2, X_5)}$$

**HARD**

hard  $\leftarrow$

$$= \int \int_{X_3, X_4} P(X_1, X_2, X_3, X_4, X_5) dX_3 dX_4$$



## • I. WHY

- The Point of Variation Inference

JOINT <sup>?</sup> → CONDITIONAL  
easy hard.

Answer is YES.



## • II. WHAT

- Ways to infer

### a) MCMC

- Gibbs Sampling
- Metropolis Hastings

Exact  
Sampling

Easy  
Slow

### b) Variational Inference

Good Approximate Hard  
Deterministic Fast

### c) Laplace Inference

Poor Approximate Easy  
Deterministic Fast



- II. WHAT

- Short Words

We try to use a simple distribution

$$q(z)$$

to approximate a complex conditional

$$p(z|x)$$

When we don't/hardly know  $p(x)$



### • III. HOW

- Information and Entropy

Information  $I = -\log P(x)$ ,  $x$  is event.

$P \downarrow, I \uparrow$

Entropy (Average Information)  $\rightarrow$  The measure of uncertainty.

discrete:  $H = -\sum_i P_i \log P_i$

continuous:  $H = -\int_{\mathcal{X}} P(x) \log P(x) dx$





### • III. HOW

- KL divergence

Kullback-Leibler divergence

$$KL(P \parallel q) = - \sum_i P_i \log \frac{q_i}{P_i}$$

$p, q$  is distribution

$$= - \sum_i P_i \log q_i - (- \sum_i P_i \log P_i)$$

$$\underline{H(p, q)} - H(p)$$



Cross-entropy



### • III. HOW

- KL divergence

but not symmetrical!

Measure the "distance" of two distributions  
Information inequality

$$KL(P, Q) = H(P, Q) - H(P) \geq 0$$
$$KL = 0 \text{ iff } P = Q$$

use Jensen's inequality to prove

Cross-entropy as loss function

$$H(P, Q) = KL(P, Q) + H(P)$$

$H(P, Q) \downarrow$

$KL \downarrow$

$P, Q$  similarity  $\uparrow$



### • III. HOW

- KL divergence

Maximum Entropy Principle

(the principle of insufficient reasoning)

The discrete distribution with max entropy is  
Uniform distribution  $u(x)$

$$0 \leq KL(p||u) = \sum_i p_i \log \frac{p_i}{u_i} = \sum_i p_i \log p_i - \sum_i p_i \log u_i$$

$$\Rightarrow H(p) \leq H(u) = \log |X|$$

$H(p)$  get max iff  $p = u$ .

$\rightarrow |X| =$  the number  
of state



### • III. HOW

- Derivation

We want to use  $q(z)$  to estimate  $P(z|x)$   
known  $P(z,x)$

$$\min_q KL(q(z) || P(z|x))$$



### • III. HOW

#### • Derivation

$$KL = - \sum_z q(z) \log \frac{P(z|x)}{q(z)}$$

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

$$\Rightarrow KL = - \sum_z q(z) \log \frac{P(x, z)}{q(z)} + \underbrace{\sum_z q(z) \cdot \log P(x)}_{\substack{\text{sum to 1} \\ \text{given } x, \text{ fixed.}}}$$

$$\therefore \underbrace{KL}_{\geq 0} + \boxed{\sum_z q(z) \cdot \log \frac{P(x, z)}{q(z)}} = \log P(x) \leq 0$$

Lower Bound ( $L$ )  $\leq 0$



### • III. HOW

#### • Derivation

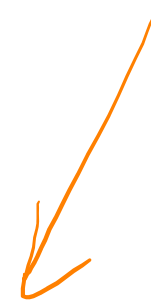
$$KL + \mathcal{L} = \log P(x) \text{ (fixed)}$$

known



$$\boxed{P(x, z)}$$

$$\therefore \min_q KL \Leftrightarrow \max_q \mathcal{L} = \sum_z q(z) \log \frac{P(x, z)}{q(z)}$$



$P(x|z)$  disappear!

Though, estimate  $q(z)$  is still not easy.



### • III. HOW

#### • Derivation

$q$  is a "constructed" distribution, so we ASSUME:

$z = (z_1, z_2, z_3)$ ,  $z_1, z_2, z_3$  is independent.

thus,

$$q(z) = q(z_1, z_2, z_3) = q(z_1) \cdot q(z_2) \cdot q(z_3)$$

and,

$$\mathcal{L} = \sum_{z_1} \sum_{z_2} \sum_{z_3} q(z_1) \cdot q(z_2) \cdot q(z_3) (\log p(x, z) - \log q(z_1) - \log q(z_2) - \log q(z_3))$$



### • III. HOW

#### • Derivation

Still difficult, how about imagining:

We know  $z_2, z_3$ , want to solve  $z_1$ .

$$\sum_{z_1} \sum_{z_2} \sum_{z_3} q(z_1) q(z_2) q(z_3) \cdot \textcircled{1} \log P(x, z)$$

$$\textcircled{2} - \log q(z_1)$$

$$\textcircled{3} - \log q(z_2) - \log q(z_3)$$

$$\textcircled{1} = \sum_{z_1} q(z_1) \cdot \left[ \sum_{z_2} \sum_{z_3} q(z_2) \cdot q(z_3) \log P(x, z) \right] = \sum_{z_1} q(z_1) \cdot E_{z_2, z_3} \log P(x, z)$$

$$\textcircled{2} = - \sum_{z_1} q(z_1) \cdot \log q(z_1) \cdot \left[ \sum_{z_2} \sum_{z_3} q(z_2) q(z_3) \right] = 1$$

$$\textcircled{3} = - \sum_{z_1} q(z_1) \cdot \sum_{z_2} \sum_{z_3} q(z_2) q(z_3) (\log q(z_2) + \log q(z_3)) = - \sum_{z_1} q(z_1) \cdot K = K?$$





### • III. HOW

#### • Derivation

$$\mathcal{L} = \sum_{z_1} q(z_1) \left[ \mathbb{E}_{z_2, z_3} \log P(x, z) - K \right] - \sum_{z_1} q(z_1) \cdot \log q(z_1)$$

Let  $f(x, z) = C_1 e^{\mathbb{E}_{z_2, z_3} \log P(x, z)}$  is a distribution.

∴

$$\mathcal{L} = \sum_{z_1} q(z_1) \log f(x, z) - \sum_{z_1} q(z_1) \log q(z_1) + \text{const.}$$

$$= \sum_{z_1} q(z_1) \cdot \log \frac{f(x, z)}{q(z_1)} + \text{const}$$

$$= -KL(q(z_1) \parallel f(x, z)) + \text{const.}$$



### • III. HOW

#### • Derivation

Got  $\mathcal{L} + KL(q(z_1) || f(x, z)) = \text{const.}$

$\max \mathcal{L} \Leftrightarrow \min KL$

$$\therefore q(z_1) = f(x, z) = C_1 e^{\sum_{z_2, z_3} \log P(x, z)}$$

So as  $q(z_2)$  and  $q(z_3)$

Then,

$$q(z) = q(z_1) \cdot q(z_2) \cdot q(z_3)$$

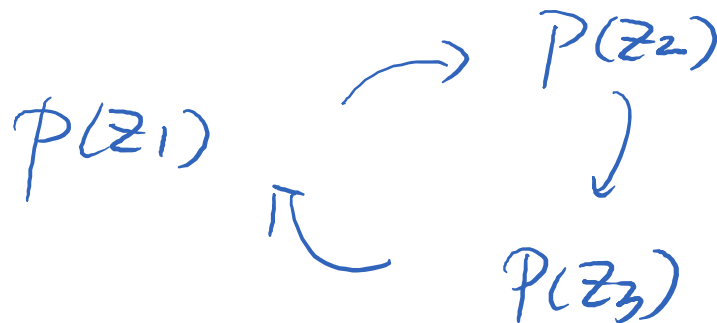


### • III. HOW

- Derivation

In many cases, given  $q$  is distribution,  
we can infer  $q$  directly (use  $\int q dz = 1$ )

But, if we cannot, just iterate





### • III. HOW

#### • Example

Given

$$p(x, y, z) = \lambda_1 \lambda_2 \lambda_3 e^{-\lambda_1 x - \lambda_2 y - \lambda_3 z}$$

calculate  $P(x, y | z)$ .

Simple Bayes Rule Solution =

$$P(x, y | z) = \frac{P(x, y, z)}{P(z)} = \frac{P(x, y, z)}{\int_0^{+\infty} \int_0^{+\infty} P(x, y, z) dx dy} = \lambda_1 \lambda_2 e^{-\lambda_1 x - \lambda_2 y}$$



### • III. HOW

#### • Example

Variational Inference Solution:

Use  $q(x, y) = q(x) \cdot q(y)$  to estimate  $P(x, y | z)$

Use the equation:

$$\begin{aligned}\ln q(x) &= \mathbb{E}_y \ln p(x, y, z) + K \\ &= \mathbb{E}_y (\ln \lambda_1 \lambda_2 \lambda_3 - \lambda_1 x - \lambda_2 y - \lambda_3 z + K) \\ &= -\lambda_1 x - \underbrace{\lambda_2 \mathbb{E}_y y}_{\downarrow} - \underbrace{\lambda_3 z}_{\downarrow} + \underbrace{\ln \lambda_1 \lambda_2 \lambda_3}_{\downarrow} + \underbrace{K}_{\downarrow} \\ &\quad \quad \quad \searrow \quad \downarrow \quad \swarrow \quad \swarrow \\ &\quad \quad \quad \text{Const}\end{aligned}$$



### • III. HOW

- Example

$$\therefore \ln g(x) = -\lambda_1 x + K_1$$

$$\Rightarrow g(x) = C_1 e^{-\lambda_1 x}$$

$$\int g(x) = 1 \Rightarrow C_1 = \lambda_1$$

So as  $g(y)$ .

$$\therefore g(x, y) = \lambda_1 \lambda_2 e^{-\lambda_1 x - \lambda_2 y}$$

The same as Bayes Rule Solution!

A "perfect" approximate!



# Bibliography

- [Variational Inference Tutorial Series by Chieu from NEU](#)
- **Machine Learning: A Probabilistic Perspective (Kevin P. Murphy) Chapter 2**



**Thanks for listening!**

