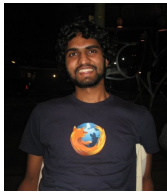


The χ -Divergence for Approximate Inference

Adji Bousso Dieng
Columbia University

Collaborators



Bayesian inference

data - \mathbf{x}

latent variables - \mathbf{z}

probabilistic model - $p(\mathbf{x}, \mathbf{z})$

Goal: Compute $p(\mathbf{z}|\mathbf{x})$



Bayesian inference



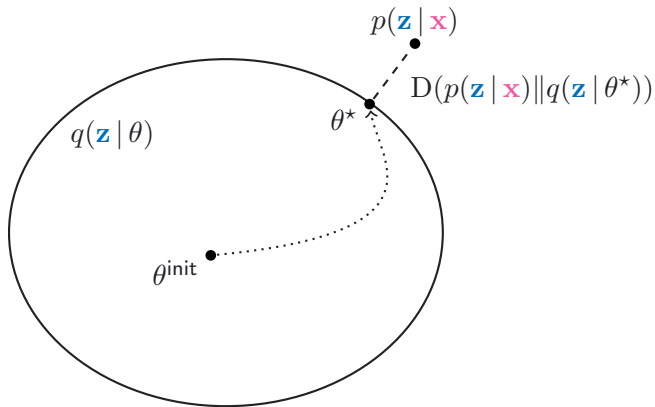
How to compute $p(\mathbf{z}|\mathbf{x})$

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

→ easy for conjugate models

→ intractable for most cases

Variational inference



$$D(p(\mathbf{z} | \mathbf{x}) || q(\mathbf{z} | \theta)) \geq 0$$

$$D(p(\mathbf{z} | \mathbf{x}) || q(\mathbf{z} | \theta)) = 0 \iff p(\mathbf{z} | \mathbf{x}) = q(\mathbf{z} | \theta) \text{ a.e.}$$

Variational inference

Pick a tractable family of approximating distribution $q(\mathbf{z} | \theta)$

Pick a divergence measure $D(p(\mathbf{z} | \mathbf{x}) || q(\mathbf{z} | \theta))$

Optimize this divergence over θ

$$q^*(\mathbf{z} | \theta^*) = \arg \min_{q(\mathbf{z} | \theta)} D(p(\mathbf{z} | \mathbf{x}) || q(\mathbf{z} | \theta))$$

Use $q^*(\mathbf{z} | \theta^*)$ as a surrogate to $p(\mathbf{z} | \mathbf{x})$

Which divergence measure to choose?

Divergence measures

$$D_f(p||q) = \int f\left(\frac{p(x)}{q(x)}\right)q(x)dx$$

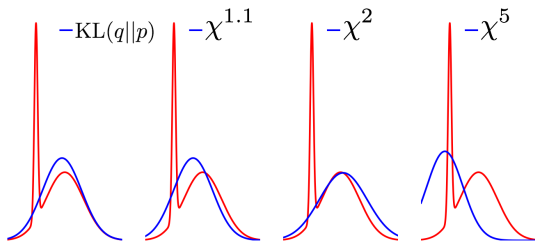
→ f convex and $f(1) = 0$

→ Special cases:

$KL(q||p)$ (**BBVI**) and $KL(p||q)$ (\approx **EP**)

$D_\alpha(q||p)$ (**VR $_\alpha$**) and $D_\alpha(p||q)$ (**CHIVI**)

Divergence measures



Three types of behaviors:

zero-forcing, zero-avoiding, neither

Variational inference with $KL(q||p)$

$$KL(q(\mathbf{z}|\theta)||p(\mathbf{z}|\mathbf{x})) = E_{q(\mathbf{z}|\theta)} \left[\log \left(\frac{q(\mathbf{z}|\theta)}{p(\mathbf{z}|\mathbf{x})} \right) \right]$$

→ Maximizes the evidence **lower bound**

$$ELBO(\theta) = E_{q(\mathbf{z}|\theta)} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\theta)]$$

→ Perform **Stochastic optimization** using Monte Carlo estimates of $\nabla_{\theta} ELBO(\theta)$

Variational inference with $KL(q||p)$

$$ELBO(\theta) = E_{q(\mathbf{z}|\theta)} \left[\log p(\mathbf{x}|\mathbf{z}) \right] - KL(q(\mathbf{z}|\theta)||p(\mathbf{z}))$$

$$\nabla_{\theta} ELBO(\theta) = E_{q(\mathbf{z}|\theta)} \left[\nabla_{\theta} \left(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\theta) \right) \right]$$

- uses **unbiased gradients**
- very efficient and works well in many settings
- **underdispersion** ... too confident
- model overpruning in VAEs

Variational inference with $D_\alpha(q||p)$

$$D_\alpha(q(\mathbf{z}|\theta)||p(\mathbf{z}|\mathbf{x})) = \frac{1}{\alpha-1} \log \int p(\mathbf{z}|\mathbf{x})^{1-\alpha} q(\mathbf{z}|\theta)^\alpha d\mathbf{z}$$

→ Maximizes the variational Renyi **lower bound**

$$\mathcal{L}(\theta) = \frac{1}{1-\alpha} \log E_{q(\mathbf{z}|\theta)} \left[\left(\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\theta)} \right)^{1-\alpha} \right]$$
$$\hat{\mathcal{L}}(\theta) = \frac{1}{1-\alpha} \log \frac{1}{K} \sum_{k=1}^K \left[\left(\frac{p(\mathbf{x}, \mathbf{z}_k)}{q(\mathbf{z}_k|\theta)} \right)^{1-\alpha} \right] \text{ where } \mathbf{z}_k \sim q(\mathbf{z}|\theta)$$

→ Perform **Stochastic optimization** using noisy estimates of $\nabla_\theta \mathcal{L}(\theta)$

Variational inference with $D_\alpha(q||p)$

$$\nabla_\theta \mathcal{L}(\theta) = E_{q(\mathbf{z}|\theta)} \left[\frac{w^{1-\alpha}}{E_{q(\mathbf{z}|\theta)}(w^{1-\alpha})} \nabla_\theta \log w \right]$$

$$\log w = \log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\theta)$$

- uses **biased gradients**
- works well in many settings
- **underdispersion** by the nature of the divergence
- cannot handle upper bounds

Variational inference with $D_\alpha(p||q)$

→ Equivalent to minimizing:

$$D_\chi(p(\mathbf{z}|\mathbf{x})||q(\mathbf{z}|\theta)) = \frac{1}{n} E_{q(\mathbf{z}|\theta)} \left[\left(\frac{p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\theta)} \right)^n \right], \quad n = 1 + \alpha$$

→ Equivalent to minimizing the χ upper bound:

$$\text{CUBO}_n(\theta) = \frac{1}{n} \log E_{q(\mathbf{z}|\theta)} \left[\left(\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\theta)} \right)^n \right]$$

→ Equivalent to minimizing:

$$\mathcal{L}(\theta) = \exp(n * \text{CUBO}(\theta))$$

→ Perform **Stochastic optimization** using noisy estimates of $\nabla_\theta \mathcal{L}(\theta)$

Variational inference with $D_\alpha(p||q)$

$$\nabla_\theta \mathcal{L}(\theta) = E_{q(\mathbf{z}|\theta)} \left[w^n \nabla_\theta \log w \right]$$

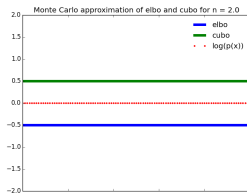
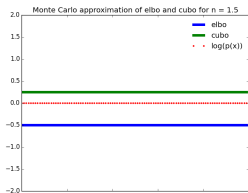
- uses **unbiased gradients**
- **overdispersed** posterior approximations
- performs upper bound minimization
- enables sandwich estimation of the evidence
- black box alternative to EP

Sandwich estimation

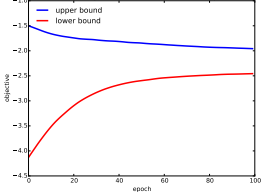
$$\text{ELBO}(\theta) \leq \log p(\mathbf{x}) \leq \text{CUBO}_n(\theta)$$

$$\lim_{n \rightarrow \infty} \text{CUBO}_n(\theta) = \text{ELBO}(\theta)$$

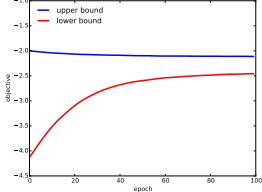
Sandwich estimation



Sandwich Plot Using CHIVI and BBVI On Ionsphere Dataset



Sandwich Plot Using CHIVI and BBVI On Ionsphere Dataset



CHIVI for probit regression

Table: Test error for Bayesian probit regression.

Dataset	BBVI	EP	CHIVI
Pima	0.235 ± 0.006	0.234 ± 0.006	0.222 ± 0.048
Ionos	0.123 ± 0.008	0.124 ± 0.008	0.116 ± 0.05
Madelon	0.457 ± 0.005	0.445 ± 0.005	0.453 ± 0.029
Covertypes	0.157 ± 0.01	0.155 ± 0.018	0.154 ± 0.014

CHIVI for Cox processes

Table: Average L_1 error for posterior uncertainty estimates (ground truth from HMC).

-	Curry	Demarcus	Lebron	Duncan
CHIVI	0.060	0.073	0.0825	0.0849
BBVI	0.066	0.082	0.0812	0.0871

Take-away message

CHIVI

- uses **unbiased gradients**
- favors **overdispersed** posterior approximations
- performs upper bound minimization
- enables sandwich estimation of the evidence
- is a black box alternative to EP
- needs variance reduction techniques