

Intro to Variational Inference & BBVI

Overview

Observed data
 $\{x_n, y_n\}$

Assume Model to Generate Data

(1) Prior $p(\mathbf{z})$

(2) Likelihood $p(y|\mathbf{z}, x)$

Goal: Estimate posterior $p(\mathbf{z}|X, y)$ given observed data

When can I use posterior inference methods?

	General MCMC	HMC (specific)	General VI	BBVI
both prior + lik need to be Gaussian	✓	✓	✓	✓
prior only is a Gaussian	✓	✓	✓	✓
\mathbf{z} Continuous prior + lik have differentiable log pdf \mathbf{z} is cont.	✓	✓	✓	✓
\mathbf{z} discrete prior + lik have pdfs we cannot differentiate \mathbf{z} has discrete	✓	✗	✓	✓
	✓	✗	✓	✓

What tunable settings does HMC require?

$$PDF_{Cat}(\mathbf{z}) = \prod_{k=1}^K p_k(\mathbf{z}_k)$$

What tunable setting does BBVI require?

#leaping steps, step size

step size, #monte carlo samples

Variational Inference

Goal: Estimate a posterior distribution $q(z|x,y)$
(z or w here)

Any distr. can be represented two possible ways

- set of samples $\{w^s\}_{s=1}^S$

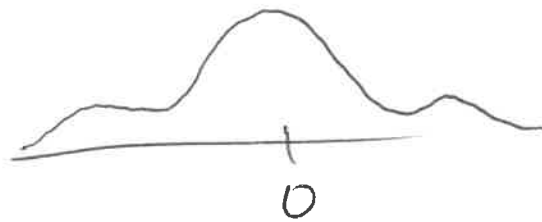
- as a known density function
w/ known parameters

Normal
w/ mean
& std dev

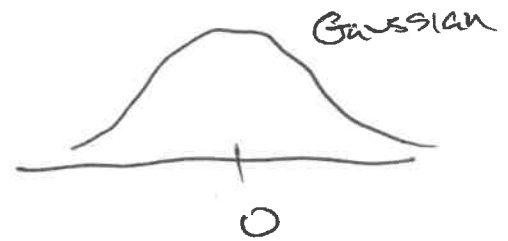
Poisson
w/ mean

etc.

Not all distributions belong to known family,
but can often make approximations



\approx



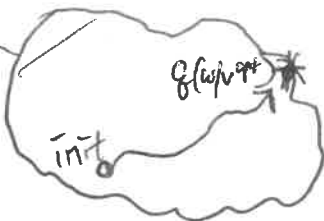
Idea of VI:

Choose an approximating family of density functions
call this $q(w)$ with possible parameters $v \in \mathcal{V}$

Goal: Make $q(w)$ as "close" as possible to $p(w|y,x)$
by finding specific value of v within \mathcal{V}

All distr. over w

distr. which
belong to
 $q(\cdot)$



- $p(w|y,x)$ $\min KL(q||p)$

"Kullback Leibler"

→ textbook says →

$$KL(q \parallel p) \triangleq \mathbb{E}_{q(w)} \left[\log \frac{q(w)}{p(w|x,y)} \right]$$

what is KL if $q=p$ everywhere? exactly 0

what if $q \neq p$ everywhere? $KL(q \parallel p) > 0$

Properties:

not symmetric

Useful to measure "distance" between q & p

Motivated by information theory

Measures expected num. "bits" (if \log_2)

required to code samples from q using best possible code for p

Look at expanding KL

$$KL \triangleq \mathbb{E}_{q(w)} \left[\log \frac{q(w)}{p(w|x,y)} \right]$$

Bayes thm
 $p(w|x,y) = \frac{p(w)p(y|w,x)}{p(y|x)}$

$$= \mathbb{E}_{q(w)} \left[\log q(w) - \log \frac{p(w)p(y|w,x)}{p(y|x)} \right]$$

$$= \mathbb{E}_{q(w)} \left[\log q(w) - \log p(w)p(y|w,x) \right] + \log p(y|x)$$

marginal dist
constant

want to do $\min_q KL(q \parallel p(w,y,x))$

can instead $\min_q \mathbb{E}_q \left[\log q(w) - \log p(w) - \log p(y|w,x) \right]$

Evidence lower bound objective $\max_{v \in V} \mathcal{L}(v)$

$$\mathcal{L}(v) = \mathbb{E}_{q(w|v)} [\log p(w) + \log p(y|w, x) - \log q(w)]$$

Two interps

$$\mathbb{E}_{q(w)} [\log p(y, w|x)] + \mathbb{E}_{q(w)} [-\log q(w)]$$

joint prob. entropy

$$\mathbb{E}_{q(w)} [\log p(y|x, w)] + \mathbb{E}_{q(w)} [\log \frac{p(w)}{q(w)}] \left. \vphantom{\mathbb{E}_{q(w)}} \right\} \text{KL to prior}$$

reconstruction likelihood

Two choices:

- (1) Density family for q w/ V space
- (2) Specific parameters $v \in V$

if $\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \in \mathbb{R}^2$

could be

$$N(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \pm)$$

$$N(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix})$$

$$N(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \Sigma)$$

simple

flexible

want to balance simplicity & flexibility in making approximations

How can we compute $\kappa(\nu)$?

Monte Carlo approximation!

$$\mathbb{E}_{q(w)} [\log p(y|w, x) + \log p(w) - \log q(w)]$$

$$\approx \frac{1}{S} \sum_{s=1}^S \log p(y|x, w^s) + \log p(w^s) - \log q(w^s)$$

when $w^s \stackrel{iid}{\sim} q(w^s | \nu)$

works for any model (lik & prior pdf choice)
as long as we can

(1) sample from q

(2) evaluate logpdf of q at w

How can we find the variational parameter ν that maximizes?

if we could estimate gradients, just go uphill



how to estimate gradients?

can't use samples w^s from MC approx.

Once we've drawn, we've lost dependence on ν

$$\frac{1}{S} \sum_s \nabla_{\nu} [\log p(y|x, w^s) + \dots]$$

doesn't work!

Idea: use small identity that's true for any distrib.

Given $q(\omega|v)$ w/ continuous param. v

$$\nabla_v \log q(\omega|v) = \frac{1}{q(\omega|v)} \nabla_v [q(\omega|v)]$$

by chain rule of log function

How does this help?

$$\begin{aligned} \nabla_v \mathbb{E}_{q(\omega|v)} [\log q(\omega|v)] &= \mathbb{E}_{q(\omega|v)} [\nabla_v \log q(\omega|v)] \\ &= \mathbb{E}_{q(\omega|v)} \left[\frac{1}{q(\omega|v)} \nabla_v q(\omega|v) \right] \\ &= \int \frac{q(\omega|v)}{q(\omega|v)} \nabla_v q(\omega|v) d\omega \\ &= \nabla_v \int q(\omega|v) d\omega \\ &= \nabla_v [1] = 0 \end{aligned}$$

this "trick" lets us simplify expectations of log probs.

Now for ELBO objective $\mathcal{L}(v)$

$$\nabla_v \mathbb{E}_q [\log p(y, w|x) - \log q(w|v)]$$

$$\nabla_v \int q(w) [\log p(y, w|x) - \log q(w|v)] dw$$

$$\int \nabla_v \left[\underbrace{q(w)}_A [\log p(y, w|x) - \log q(w|v)] \right] dw$$

product rule says

$$\int \underbrace{[\nabla_v q(w)]}_{\text{grad A} * B} (\log p(y, w|x) - \log q(w|v)) dw$$

$$+ \int \nabla_v [\log p(y, w|x) - \log q(w|v)] \underbrace{q(w)}_{\text{we know gradient is 0}}$$

const w.r.t v
so grad = 0

$$\underbrace{\text{grad B} * A}$$

whole term equals zero!

so keep first term & apply trick again.

$$\nabla_v \mathcal{L}(v) = \int \underbrace{[\nabla_v q(w)]}_{\text{ID trick}} (\log p(y, w|x) - \log q(w|v)) dw$$

$$= \int \underbrace{[\nabla_v \log q(w|v)]}_{\text{ID trick}} \underbrace{q(w)}_{\text{ID trick}} (\log p(y, w|x) - \log q(w|v)) dw$$

expectation!

$$= \text{something we can estimate via Monte Carlo!}$$

$$\frac{1}{S} \sum_{s=1}^S [\nabla_v \log q(w^s|v)] (\log p(y, w^s|x) - \log q(w^s|v))$$