

Review: VI has us pick a density family $q(w|v)$ w/ parameters $v \in V$
 We then optimize an objective $d(v)$ 01

BBVI

BbB

~~"score function trick"~~

~~"log derivative trick"~~

how to compute $d(v)$?
 (what identity/technique)

monte carlo

monte carlo

how to compute $\nabla_v d(v)$

score function trick

~~log derivative~~
 reparameterization trick

write formula for it

$$\nabla_v d(v) \approx \frac{1}{S} \sum_{s=1}^S \nabla_v [\log q(v)]$$

$(\log p(w^s, y|x) - \log q(w^s))$

do we need log prior to be differentiable

X

✓

log lik

X

✓

log q

~~X~~ ✓

✓

$$v = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix} \quad \text{w/ mean } \mu, \text{ variance } \sigma^2$$

02

Idea: if $g(w/v)$ is Gaussian, consider two ways to sample from it

(1) draw $w \sim \text{Normal}(\mu, \sigma^2)$

(2) draw $\varepsilon \sim \mathcal{N}(0, 1)$ standardized r.v.
 $w \leftarrow \sigma \varepsilon + \mu$ plus a transform

Can prove these two procedures produce same distribution.

Key idea: define $w \hat{=} t(v, \varepsilon)$

deterministic mapping of parameters v
 and a std. r.v. ε
w/ distr. $p(\varepsilon)$

Now, for any function $f(w, v)$ that depends on w and v ,

we can rewrite expectations

$$\mathbb{E}_{g(w/v)} [f(w, v)] = \mathbb{E}_{p(\varepsilon)} [f(t(v, \varepsilon), v)]$$

why? now gradient w.r.t. v can go inside

$$\nabla_v \mathbb{E}_{g(w/v)} [f(w, v)] = \nabla_v \mathbb{E}_{p(\varepsilon)} [f(t(v, \varepsilon), v)]$$

~~$\mathbb{E}_{p(\varepsilon)} [\nabla_v f(t(v, \varepsilon), v)]$~~
 chain rule \rightarrow

chain rule \downarrow

$$\nabla_v \mathbb{E}_{q(w|v)} [f(w,v)] = \mathbb{E}_{p(\varepsilon)} \left[\frac{\partial f(w,v)}{\partial w} \frac{\partial w}{\partial v} + \frac{\partial f(w,v)}{\partial v} \right]$$

Given S samples, can compute as

$$= \frac{1}{S} \sum_s \left. \frac{\partial f(w,v)}{\partial w} \right|_{w=w^s} \left. \frac{\partial w}{\partial v} \right|_{w=w^s} + \frac{\partial f(w,v)}{\partial v}$$

Illustration:

What is gradient wrt μ of expected value of ?
 Very easy w/ reparameterization ^{1D} normal?

First way:

$$\begin{aligned} \nabla_{\mu} \mathbb{E}_{w \sim N(\mu, \sigma^2)} [w] \\ &= \nabla_{\mu} \int p(w|\mu, \sigma^2) w \, dw \\ &= \nabla_{\mu} \int \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} e^{-\frac{1}{2\sigma^2}(w-\mu)^2} w \, dw \end{aligned}$$

HARD!

Second way:

w can be parameterized as $w \sim \mu + \sigma \varepsilon$ $\varepsilon \sim N(0, 1)$

$$\begin{aligned} \nabla_{\mu} \mathbb{E}_{\varepsilon \sim N(0,1)} [\mu + \sigma \varepsilon] \\ &= \nabla_{\mu} \mu + \nabla_{\mu} \mathbb{E}[\sigma \varepsilon] \xrightarrow{\rightarrow 0} = 1 \end{aligned}$$

const wrt μ

Algorithm: Applying BB to BNNs

What we want to optimize:

$$\max_{\nu \in V} \mathcal{L}(\nu), \quad \mathcal{L}(\nu) = \frac{1}{S} \sum_S \log p(y, w^S) - \log q(w^S | \nu)$$

$$q(w | \nu) = \prod_j N(w_j | \mu_j, \tau_j)$$

$$\nu_j = \begin{bmatrix} \mu_j \\ \tau_j \end{bmatrix}$$

mean can be any real
stddev must be positive

wait! will this work for SGD?
want unconstrained parameter could do $s \hat{=} \log(\tau)$
 $p \hat{=} \log(1 + e^\tau)$

Using S samples, can compute gradient as

$$\nabla_{\mu_j} \mathbb{E} \left[\log p(y, w) - \log q(w | \mu, \tau) \right]$$

$$= \frac{1}{S} \sum_{s=1}^S \mathbb{E}_{p(\epsilon)} \left[\frac{\partial}{\partial \mu_j} \left(\log p(y, w(\epsilon) | x) - \log q(w(\epsilon) | \mu, \tau) \right) \right]$$

$$= \frac{1}{S} \sum_{s=1}^S \frac{\partial f}{\partial \mu} + \frac{\partial f}{\partial w} \frac{\partial w}{\partial \mu}$$

$$\nabla_p \mathbb{E} \left[\log p(y, w) - \log q(w | \mu, p) \right]$$

can be computed in same way

Alg: Init: means μ , log stddevs S
for each iter:

$$\mu \leftarrow \mu - \alpha \nabla_{\mu} \mathcal{L}$$

$$p \leftarrow p - \alpha \nabla_p \mathcal{L}$$

need many iters
need to tune step size

All these can produce samples from approx. posterior which we can use for prediction

Methods:

$$P(y_* | x_*) = \int P(y_* | w, x_*) P(w | x_*) dw$$

HMC

BBVI

Bayes By Backprop

requires	HMC	BBVI	Bayes By Backprop
diff'able log prior PDF	✓		✓
diff'able log lik PDF	✓		✓
a chosen "q" approx posterior density	N/A	✓	✓
diff'able log PDF for q		✓	✓
q must be Normal (at least, for purposes of this class)			✓
estimated gradients		BIG variance	smaller variance

defined by param v ,
For any $q(w/v)$, we draw S samples w^1, w^2, \dots, w^S

BBVI gradient estimator

$$\nabla_v d(v) = \frac{1}{S} \sum_{s=1}^S \left(\nabla_v \log q(w^s/v) \right) \left(\log P(y, w^s/x) - \log q(w^s/v) \right)$$

BbB gradient estimator

Assume $q(w/v) \sim \text{Normal}(\mu_v, \Sigma_v)$. Draw S samples $\epsilon^1, \dots, \epsilon^S$ from $N(0, 1)$.

$$\nabla_v d(v) = \frac{1}{S} \sum_{s=1}^S \nabla_v \left[\log P(y, w(\epsilon^s, \mu_v, \Sigma_v)) - \log q(w(\epsilon^s, \mu_v, \Sigma_v)/v) \right]$$