

Gradient Estimator Summary

May 26, 2018

In general, we want to be able to estimate the gradient of

$$\mathbb{E}_{p(b|\theta)}[f(b)]$$

w.r.t. parameter θ .

Note that

$$\partial_{\theta} \mathbb{E}_{p(b|\theta)}[f(b)] = \partial_{\theta} \int_{-\infty}^{\infty} f(b) p(b|\theta) db \quad (1)$$

$$= \int_{-\infty}^{\infty} f(b) \partial_{\theta} p(b|\theta) db \quad (2)$$

The problem here is that equation (2) doesn't have the form of expectation. So we can't do MCMC here.

One way of dealing with this is that using Score Function

$$\int_{-\infty}^{\infty} f(b) \partial_{\theta} p(b|\theta) db \quad (3)$$

$$= \int_{-\infty}^{\infty} f(b) p(b|\theta) \partial_{\theta} \log(p(b|\theta)) db \quad (4)$$

So equation (4) can be seen as expectation of $f(b) \partial_{\theta} \log(p(b|\theta))$ over distribution $p(b|\theta)$.

$$\partial L = \int_{-\infty}^{\infty} f(b)p(b|\theta)\partial_{\theta}\log(p(b|\theta))db \quad (5)$$

$$= \mathbb{E}_{p(b|\theta)}[f(b)\partial_{\theta}\log(p(b|\theta))] \quad (6)$$

$$\approx \frac{1}{S} \sum_s f(b^{(s)})\partial_{\theta}\log(p(b^{(s)}|\theta)) \quad (7)$$

Where $b^{(b)}$ is sample from $p(b|\theta)$ based on current parameter θ

- Pros of this method is that it doesn't require function f to be differentiable as we never differentiate f and b can be either discrete or continuous as long as its pdf is differentiable
- However, this method to be known has high variance and unable to update parameter θ efficiently

Reparameterization

Alternatively, since our original problem arises from that θ parametrizes the expectation. So instead of differentiate θ , we differentiate a different parameters on which θ depends.

- sample $\epsilon \sim N(0, 1)$
- choose a smooth function h and write $b = h(\epsilon, \theta)$
-

$$\partial_{\theta} L = \partial_{\theta} \mathbf{E}_{p(b|\theta)}[f(b)] \quad (8)$$

$$= \partial_{\theta} \mathbf{E}_{\epsilon \sim p(\epsilon)}[f(g(\theta, \epsilon))] \quad (9)$$

$$= \mathbf{E}_{\epsilon \sim N(0,1)}[\partial_{\theta} f(g(\theta, \epsilon))] \quad (10)$$

$$= \mathbf{E}_{\epsilon \sim N(0,1)}[\partial_{\theta} f(g(\theta, \epsilon)) \partial_{\theta} g(\theta, \epsilon)] \quad (11)$$

$$\approx \frac{1}{S} \sum_s^S \partial_{\theta} f(g(\theta, \epsilon^{(s)})) \partial_{\theta} g(\theta, \epsilon^{(s)}) \quad (12)$$

Reparameterization

- This reparameterization trick gives us lower variances
- However, if we look at the term inside the sum. We notice that we need function f to be differentiable w.r.t. ϕ ; and we need a differentiable function g to represent b .
- But what if f is not differentiable ?
- What if b is discrete ?

So we see that the Pros of Score Function becomes Cons in reparameterization and vice verse.

Therefore, we need something more powerful.

- More flexibility of f : non differentiable
- More flexibility of random variable b : can be either discrete or continuous
- Low variance of gradient estimator
- Gradient Estimator is unbiased

We introduce a control variate function C to alleviate the high variance of score function

$$\mathbb{E}_{p(b|\theta)}[(f(b) - C)\partial_{\theta}\log(p(b|\theta))] \quad (13)$$

Note that if C is a constant, then (13) estimator is unbiased, but if C is function depends on b , it's biased (Shown next slides)

So we need a correction of bias

$$\mathbb{E}_{p(b|\theta)}[(f(b) - C(b))\partial_{\theta}\log(p(b|\theta))] + \partial_{\theta}\mathbb{E}_{p(b|\theta)}[C(b)] \quad (14)$$

Justification

Before we propose the next gradient estimator; let's justify why equation (14) make sense

- for any fix c , $\mathbb{E}_{p(b|\theta)}[(f(b) - c)\partial_\theta \log(p(b|\theta))]$ is unbiased

$$\mathbb{E}_{p(b|\theta)}[(f(b) - c)\partial_\theta \log(p(b|\theta))] \quad (15)$$

$$= \int p(b|\theta)(f(b) - c)\partial_\theta \log(p(b|\theta)) \quad (16)$$

$$= \int p(b|\theta)(f(b) - c)\frac{1}{p(b|\theta)}\partial_\theta(p(b|\theta)) \quad (17)$$

$$= \int (f(b) - c)\partial_\theta(p(b|\theta)) \quad (18)$$

$$= \int f(b)\partial_\theta p(b|\theta) - \int c\partial_\theta p(b|\theta) \quad (19)$$

$$= \int f(b)\partial_\theta p(b|\theta) = \partial_\theta \mathbb{E}_{p(b|\theta)}[f(b)] \quad (20)$$

Note that $\int c\partial_\theta p(b|\theta) = c\partial_\theta \int p(b|\theta) = c\partial_\theta 1 = 0$. The integral is over the range b

Notice from above (equation (19)), that if c is not constant but a function of random variable b , above results might not hold (eg. if $c(b) = b$, then second term in equation (19) is $\partial_{\theta} \mathbb{E}_{p(b|\theta)}(b)$). Equality in (20) doesn't hold.

Continues from equation (19) with $C(b)$

$$\int f(b) \partial_{\theta} p(b|\theta) - \int c(b) \partial_{\theta} p(b|\theta) \quad (21)$$

$$= \partial_{\theta} \mathbb{E}_{p(b|\theta)}[f(b)] - \partial_{\theta} \mathbb{E}_{p(b|\theta)}[c(b)] \quad (22)$$

This justify why we need a correction term (which is simply add one more term to cancel the negative bias). So, Equation (14) is an unbiased estimator

- Notice, that when $C = f$; equation(14) becomes exactly as regular reparameterization.
- Hence, C acts as a balance of variances between score function and reparameterization (Variance can be as low as reparameterization trick)
- Also, if $C \neq f$, no need to differentiate through f . We no longer requires f to be differentiable

Objective

- We want C to be differentiable w.r.t. θ
- We want to approximate C such that it makes the equation (14) with lowest variances
- Since one of our goal is to deal with discrete b ; the equation (14) requires us to find a good reparameterization function h

High Level Procedure

- Find a differentiable function h for random variable b
- update θ that optimization the equation (14)
- update the parameters of function C such that minimized the variance of equation (14)

First Approach (Concrete)

Let's solve each of our goals step by step. First, we aim at dealing discrete b by find a good smooth function for re-parameter. Let $b \sim p(b|\phi)$ be a Bernoulli random variable for simple illustration

- sample $u \sim \text{Unif}(0, 1)$
- deterministically define $z = h(u, \theta)$ given u (e.g.
 $z = \log(\frac{\theta}{1-\theta}) + \log(\frac{u}{1-u})$)
- Now; we reparameterize $b = H(z) \approx \sigma_\lambda(z) = (1 + \exp(-\frac{z}{\lambda}))^{-1}$
where $H(z) = 1$ if $z \geq 0$ and $H(z) = 0$ otherwise
Hence,

$$\partial_\theta \mathbb{E}_{p(b|\theta)}[f(b)] \approx \partial_\theta \mathbb{E}_{p(z|\theta)}[f(\sigma_\lambda(z))] \quad (23)$$

$$= \partial_\theta \mathbb{E}_{p(u)}[f(\sigma_\lambda(h(u, \theta)))] \quad (24)$$

$$= \mathbb{E}_{p(u)}[\partial_\theta f(\sigma_\lambda(h(u, \theta)))] \quad (25)$$

- Now, we can estimate equation (25) with reparameterization trick

Weakness of Concrete

- We see that Concrete Relax overcome one major issues: discrete random variable
- However, notice that in this approach, chosen hyperparameter λ is hard
- If λ is large, σ_λ is more horizontally flat(stretched), so lower variance for different values of z but it less accurate as an approximator of Bernoulli b
- As $\lambda \rightarrow 0$, more squashed, it approximate Bernoulli more accurately, but the variance goes to very high

REBAR approach

Rebar Estimator is an alternative that deal with this bias and variance trade-off by combine concrete relax with equation (14)

Notice from equation (23)

$$\partial_{\theta} \mathbb{E}_{p(z|\theta)} f(\sigma_{\lambda}(z)) = \mathbb{E}_{p(z|\theta)} [f(\sigma_{\lambda}(z)) \partial_{\theta} \log(p(z|\theta))] \quad (26)$$

Also, notice that in equation (14)

$$\mathbb{E}_{p(b|\theta)} [f(b) \partial_{\theta} \log(p(b))] = \partial_{\theta} \mathbb{E}_{p(b|\theta)} [f(b)] \quad (27)$$

$$= \partial_{\theta} \mathbb{E}_{p(z|\theta)} [f(H(z))] \quad (28)$$

$$= \mathbb{E}_{p(z|\theta)} [f(H(z)) \partial_{\theta} \log(p(z|\theta))] \quad (29)$$

Recall from our original gradient estimator proposal

$$\mathbb{E}_{p(b|\theta)}[(f(b) * \partial_{\theta} \log(p(b|\theta)) - C(b) * \partial_{\theta} \log(p(b|\theta)))] + \partial_{\theta} \mathbb{E}_{p(b|\theta)}[C(b)]$$

- We use equation (29) to approximate $f(b)$ ($H(z)$ is discontinuous here but it's Ok as we don't reparameterize it; we solve it with score function; then alleviate its variance)
- Now, we need to design a $C(b)$ Since we require $C(b)$ to be differentiable even with discrete input b . We use the trick from Concrete Estimator. Basically, use equation(26) as our control variate.

Note that equation (26) can be written as

$$\mathbb{E}_{p(z|\theta)}[f(\sigma_\lambda(z))\partial_\theta \log(p(z|\theta))] \quad (30)$$

$$= \sum_z p(z)[f(\sigma_\lambda(z))\partial_\theta \log(p(z|\theta))] \quad (31)$$

$$= \sum_z \sum_b [p(b, z)f(\sigma_\lambda(z))\partial_\theta \log(p(z|\theta))] \quad (32)$$

$$= \sum_z \sum_b [p(z|b)p(b)f(\sigma_\lambda(z))\partial_\theta \log(p(z|\theta))] \quad (33)$$

$$= \sum_b p(b) \sum_z [p(z|b)f(\sigma_\lambda(z))\partial_\theta \log(p(z|\theta))] \quad (34)$$

$$= \sum_b p(b) \sum_z \left[p(z|b)f(\sigma_\lambda(z))\partial_\theta \left[\log(p(z|b, \theta)) + \log(p(b|\theta)) \right] \right] \quad (35)$$

$$= \mathbb{E}_{p(b)} \left[\partial_{\theta} \mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))] \right] + \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z)) \partial_{\theta} \log(p(b))] \right] \quad (36)$$

equation(35) to equation(36) are derived from the definition of expectation and log trick

For the gradient inside the first term; we can estimate it with reparameterization

$$\mathbb{E}_{p(b)} \left[\partial_{\theta} \mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))] \right] = \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(v)} [\partial_{\theta} f(\sigma_{\lambda}(\tilde{z}))] \right] \quad (37)$$

where $v \sim \text{Unif}(0, 1)$, $\tilde{z} = \tilde{h}(v, b, \theta)$;

while the second term can be estimated as follows

$$\mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} [f(\sigma_\lambda(z)) \partial_\theta \log(p(b))] \right] \quad (38)$$

$$= \sum_b p(b) \sum_z p(z|b) [f(\sigma_\lambda(z)) \partial_\theta \log(p(b))] \quad (39)$$

$$= \sum_b p(b) \sum_z p(z|b) \left[f(\sigma_\lambda(z)) \frac{1}{p(b)} \partial_\theta p(b) \right] \quad (40)$$

$$= \sum_b p(b) \frac{1}{p(b)} \sum_z p(z|b) [f(\sigma_\lambda(z)) \partial_\theta p(b)] \quad (41)$$

$$= \sum_b \sum_z p(z|b) [f(\sigma_\lambda(z)) \partial_\theta p(b)] \quad (42)$$

$$= \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} [f(\sigma_\lambda(z)) \partial_\theta \log(p(b|\theta))] \right] \quad (43)$$

$$= \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(v)} [f(\sigma_\lambda(\tilde{z})) \partial_\theta \log(p(b|\theta))] \right] \quad (44)$$

Now, we can combined things together, plug equation (44),(37) into (36); then plug (36) into (14) (to replace $C(b)$); finally replace f with (28)
Note that our reparameterization make b deterministically depends on u and v from Uniform distribution.

$$\mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)}(F) \right] = \sum_b \sum_z p(z|b)p(b)(F) \quad (45)$$

$$= \sum_b \sum_z p(z, b)(F) = \mathbb{E}_{p(u,v)}[F] \quad (46)$$

Our gradient estimator (originally equation (21))

$$\hat{g}_\theta = \mathbb{E}_{p(u,v)} \left[\left(f(H(z)) - f(\sigma_\lambda(\tilde{z})) \right) \partial_\theta \log(p(b)) + \partial_\theta f(\sigma_\lambda(z)) - \partial_\theta f(\sigma_\lambda(\tilde{z})) \right] \quad (47)$$

- Instead of having λ to be hyperparameters, REBAR update λ with objection that reduce the variances

Varaince Control Objective Function

$$\text{Var}(\hat{g}_\theta) = \mathbb{E}[\hat{g}_\theta^2] - (\mathbb{E}[\hat{g}_\theta])^2 \quad (48)$$

Take the derivative w.r.t. λ

$$\hat{g}_\lambda = \partial_\lambda (\mathbb{E}[\hat{g}_\theta^2] - (\mathbb{E}[\hat{g}_\theta])^2) \quad (49)$$

$$\lambda^{(t)} \leftarrow \lambda^{(t-1)} - \alpha * \hat{g}_\lambda \quad (50)$$

- We see how Concrete approach dealing with discrete random variable
- REBAR, based on Concrete, reduce the variances (adding control variate)
- However, notice that equation(38) requires f to be differentiable
- The only parameter for variance control is λ

Justification of reLAX

- Note that with correction of bias term, the estimator is unbiased and holds for any choice of C ; and by Universal approximation theorem; we can represent C with a neural network C_ϕ where ϕ is the weight and bias of neural net.
- Instead of update the parameter λ in REBAR, we update the parameter ϕ with the same loss function (the variance of gradient estimator, λ can be incorporated into neural network with lots of other parameters)
- Even if f is undifferentiable, we are backpropagation on neural network w.r.t. its weight ϕ to control the variance. So more flexible of function f
- So in equation (39), we can replace f with C_ϕ to get our final version of gradient estimator

More specific,

$$\hat{g}_{\text{reLAX}} = \mathbb{E}_{p(u,v)} \left[\left[f(b) - [C_\phi(\tilde{z})] \right] \partial_\theta \log p(b|\theta) - \partial_\theta [C_\phi(\tilde{z})] \right] + \partial_\theta [C_\phi(z)] \quad (51)$$

Algorithm

Let's summary all things together

Algorithm 1: reLAX

input : differentiable pdf $p(b|\theta)$, a function $f(\cdot)$, neural network C_ϕ , step size α_1, α_2 , reparameterization distribution q , hard threshold function $H(\cdot)$, smooth function $h(u, \theta)$, $\tilde{h}(v, H(z), \theta)$

```
1 while not converge do
2    $u, v \sim q$  ;                               // sample reparam random variable
3    $z_i \leftarrow h(u, \theta)$  ;                 // Input for Concrete Relax
4    $\tilde{z}_i \leftarrow \tilde{h}(v, H(z_i), \theta)$  ;    // Input for conditional reparam on  $H(z)$ 
5    $\hat{g}_\theta \leftarrow \text{equation (50)}$ 
6    $\hat{g}_\phi \leftarrow \partial_\phi \text{Var}(\hat{g}_\theta)$ 
7    $\theta \leftarrow \theta - \alpha_1 * \hat{g}_\theta$  ;      // update the distribution parameter
8    $\phi \leftarrow \phi - \alpha_2 * \hat{g}_\phi$  ;    // update control variate neural network weight
9 end
```

In Stochastic Variational Inference, we are optimizing our objective function (minimized the negative ELBO)

$$L = -\mathbb{E}_{q(z|x, \theta)}[\log(p(x|z) + \log(p(z)) - \log(q(z|x, \theta)))]$$

by performing gradient descent to update the variational parameters

$$\theta^t \leftarrow \theta^{(t-1)} - \alpha * \partial_{\theta} L$$

But compute $\partial_{\theta} L$ can be hard as most of case, it unlikely for it to be computed analytically.

Therefore, we can use the gradient estimator reLAX proposed to estimate $\partial_{\theta} L$

Further Consideration

- In above equation, we assume f doesn't take the parameter θ . In general, such as SVI, we see that function inside the expectation is $f(b, \theta)$
- By chain rule:

$$\partial_{\theta} \mathbb{E}_{p(b|\theta)} [f(b, \theta)] \quad (52)$$

$$= \int \partial_{\theta} (p(b|\theta) f(b, \theta)) \quad (53)$$

$$= \int (f(b, \theta) \partial_{\theta} p(b|\theta)) + \int (p(b|\theta) \partial_{\theta} f(b, \theta)) \quad (54)$$

$$= \mathbb{E}_{p(b|\theta)} [f(b, \theta) \partial_{\theta} \log(p(b|\theta))] + \mathbb{E}_{p(b|\theta)} [\partial_{\theta} f(b, \theta)] \quad (55)$$

- Note that first term in equation (52) is what we have already solved in previous. The second term can be estimated with simple MCMC estimator.

- George Tucker, Andriy Mnih, Chris J Maddison, John Lawson, and Jascha Sohl-Dickstein. REBAR: Low-variance, unbiased gradient estimates for discrete latent variable models. In Advances in Neural Information Processing Systems, pages 2624-2633, 2017.
- Will Grathwohl, Dami Choi, Yuhuai Wu, Geoff Roeder, and David Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. arXiv:1711.00123, 2017.