

# Gradient Estimator Summary

Zhexin(Ellery) Lai

University of Toronto

June 10, 2018

# Motivation

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

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

w.r.t. parameter  $\theta$ .

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 use simple Monte Carlo to estimate the expectation

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  $\theta$

- 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  $\theta$  efficiently

# Reparameterization

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

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

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

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

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

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

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

# 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.  $\phi$ ; 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.  $\theta$
- 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  $\theta$  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 equation (25) is biased if  $b$  is discrete, because  $\sigma_\lambda(z) \rightarrow H(z)$  only as  $\lambda \rightarrow 0$
- Consider Bernoulli case:  $b = 1$  or  $0$
- If  $\lambda$  is large,  $\sigma_\lambda$  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.
- This gives us a situation where we have to tune the  $\lambda$  to deal with bias-variance trade-off

# 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 solve it with score function (no need to differentiate nor reparameterize it); 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.

How can we so sure that equation (26) will be an effective variance control?

- Intuition:  $\mathbb{E}_{p(b|\theta)} f(b) = \mathbb{E}_{p(z|\theta)} f(H(z)) \approx \mathbb{E}_{p(z|\theta)} f(\sigma_\lambda(z))$
- Implies that
$$\mathbb{E}_{p(z|\theta)} f(H(z) \partial_\theta \log(p(z|\theta))) \approx \mathbb{E}_{p(z|\theta)} f(\sigma_\lambda(z) \partial_\theta \log(p(z|\theta)))$$
- Then, the score function (left hand side) is **STRONGLY** correlated with the right hand side (That's why we choose right hand side as control variate)

We can directly noise Monte Carlo estimate equation (26)

But, we can do better than that. Let's recall the Law of total Variance

- $\mathbb{E}[\text{Var}(z|b)] = \text{Var}(z) - \text{Var}(\mathbb{E}(z|b))$
- It means that if we can make random variable  $z$  depends on some random variable  $b$ , on average, the variance of  $z|b$  will be even smaller (namely, equation (26) will be more stable)
- In the following slides, we show how to make  $z$  to depends on  $b$  in the equation (26)

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))

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

# REBAR Estimator

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

Varaince Control Objective Function

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

Note that due to its unbiased property

$$\mathbb{E}[\hat{g}_\theta] = f(H(z)) \implies \partial_\lambda f(H(z)) = 0$$

Our gradient estimator for varaince is

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

$$= \partial_\lambda \mathbb{E}[\hat{g}_\theta^2] \quad (50)$$

Update rule for  $\lambda$

$$\lambda \leftarrow \lambda - \alpha * \hat{g}_\lambda$$



- We see how Concrete approach dealing with discrete random variable
- REBAR, based on Concrete, reduce the variances (adding control variate)

But, we still face some challenges

- Notice that equation(47) requires  $f$  to be differentiable (recall one of our original goals)
- There is only a single parameter for variance control, that is  $\lambda$  (too simple, under fit)

# 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  $C_\phi$  where  $\phi$  is the weight and bias of neural net. (We can have as many parameters as we want to control the variance now)
- Instead of update the parameter  $\lambda$  in REBAR, we update the parameter  $\phi$  ( $\lambda$  is incorporated into neural network as weights )
- Even if  $f$  is undifferentiable, we backpropagation on neural network w.r.t. its weight  $\phi$  to control the variance. So more flexible of function  $f$
- So in equation (47), we can replace  $f$  with  $C_\phi$  to get our final version of graident 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  $C_\phi$ , step size  $\alpha_1, \alpha_2$ , 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  $\partial_{\theta} L$

# Experiment

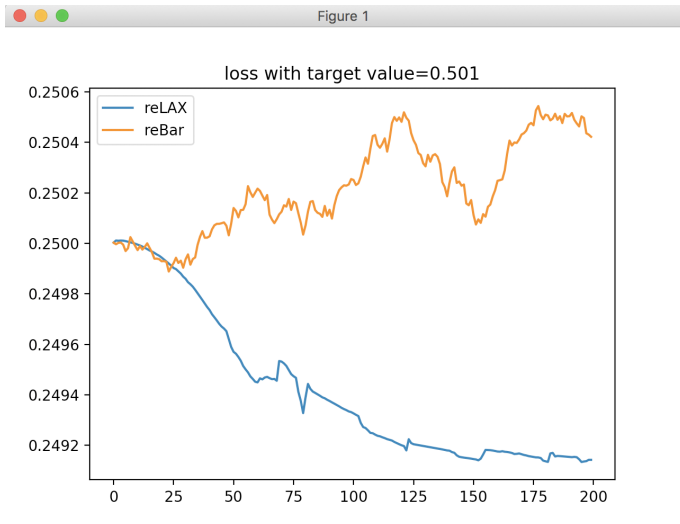
- Consider  $b \sim \text{Bernoulli}(\theta)$ .
- objective to minimize  $\mathbb{E}_b[(0.501 - b)^2]$
- Challenges: (1) Discrete Random Variable. (2) Very small loss value even for completely wrong move (we could have very high variance for each update) (3) We make the task even more challenge, each iteration of gradient estimator, we only sample one point

Let's compare reLAX estimator vs reBAR estimator

# Experiment

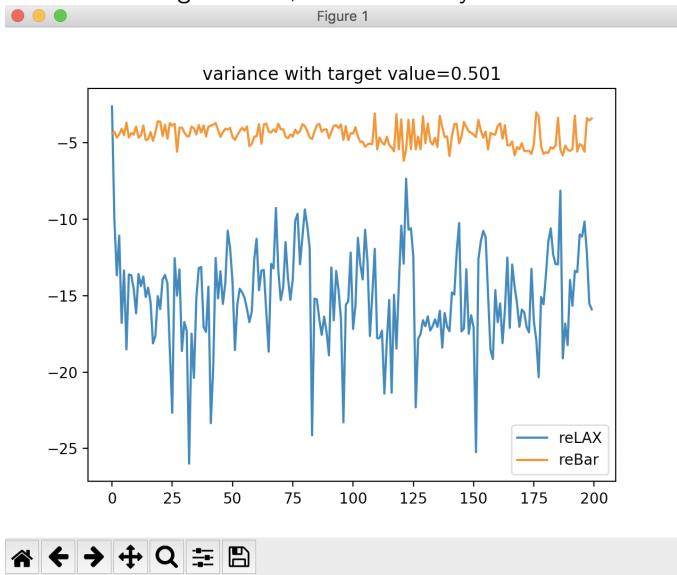
We see that under challenge case, reBar fails for the task over 10000 iterations.

While reLAX successfully complete the task



# Experiment

For the log variance, reLAX is always below reBar





- 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.