

# BACKPROPAGATING THROUGH THE VOID: OPTIMIZING DIFFERENTIABLE CONTROL VARIATES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Gradient-based optimization is the foundation of deep learning and reinforcement learning. Even when the mechanism being optimized is unknown or not differentiable, optimization using high-variance or biased gradient estimates is still often the best strategy. We introduce a general framework for learning low-variance, unbiased gradient estimators for black-box functions of random variables. These estimators can be jointly trained with model parameters or policies, and are applicable in both discrete and continuous settings. We give unbiased, adaptive analogs of state-of-the-art reinforcement learning methods such as deep deterministic policy gradients and advantage actor-critic. We also demonstrate this framework for training discrete latent-variable models.

## 1 INTRODUCTION

Gradient-based optimization has been key to most recent advances in machine learning and deep reinforcement learning. The back-propagation algorithm (BPTT), also known as reverse-mode automatic differentiation (AD) computes exact gradients of deterministic, differentiable objective functions. The reparameterization trick (RPT) allows backpropagation to give unbiased, low-variance estimates of gradients of expectations of continuous random variables. This has allowed effective stochastic optimization of large probabilistic latent-variable models.

Unfortunately, backpropagation cannot be easily applied to problems involving discrete random variables, or when the function being optimized is a black box (BB). This is the case in most reinforcement learning settings, or when fitting probabilistic models with discrete latent variables. Much recent work has been devoted to constructing gradient estimators for these situations. In reinforcement learning, advantage actor-critic methods (AAC) give unbiased gradient estimates with reduced variance obtained by jointly optimizing the policy parameters with an estimate of the value function. In discrete latent-variable models, low-variance but biased gradient estimates can be given by continuous relaxations of discrete variables (CR).

A recent advance by (Schmitt et al., 2017) used a continuous relaxation to construct a control variate (CV) for functions of discrete random variables. Low-variance estimates of the expectation of the control variate can be computed using the reparameterization trick to produce an unbiased estimator with lower variance than previous methods.

Furthermore, (Schmitt et al., 2017) showed how to tune the free parameters of these relaxations to minimize their variance during training.

We generalize the method of (Schmitt et al., 2017) to learn a free-form control variate parameterized by a neural network, giving a lower-variance, unbiased gradient estimator which can be applied to a wider variety of problems with greater flexibility.

Most notably, our method is applicable even when no continuous relaxation is available, as in reinforcement learning or black box function optimizing.

We derive improved variants of popular reinforcement learning methods with unbiased gradients and more stable training dynamics.

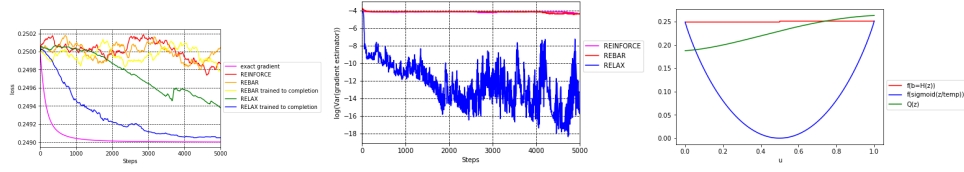


Figure 1: *Left*: Estimator losses. *Centre*: Estimator variance. *Right*: Learned relaxation function.

## 2 MOTIVATION

(WE SHOULD BREAK THIS UP INTO MOTIVATION AND BACKGROUND) IE WHY DO WE WANT THIS, RL, DISCRETE RANDOM VARIABLE MODELS, YADDA, YADDA, YADDA

THEN HOW HAS IT BEEN DONE IN THE PAST

THEN WHATS WRONG WITH THAT, LIMITATIONS, SHIT LIKE THAT

### 2.1 OPTIMIZING CONTINUOUS BLACK-BOX FUNCTIONS

Deep deterministic policy gradients (?)

Also: (?)

## 3 BACKGROUND: GRADIENT ESTIMATORS

How can we choose the parameters of a distribution to maximize an expectation? This problem comes up in reinforcement learning, where we must choose a policy  $\pi$  on actions  $a$  to maximize the expected reward  $\mathbb{E}_{p(a|\pi)}[r(a)]$ . It also comes up in fitting latent-variable models, when we wish to maximize the marginal probability  $p(x|\theta) = \sum p(x|z)p(z|\theta) = \mathbb{E}_{p(z|\theta)}[p(x|z)]$ . In this paper, we'll consider the general problem of optimizing

$$\mathcal{L}(\theta) = \mathbb{E}_{p_\theta(b)}[f(b)]. \quad (1)$$

Later, we will discuss the case when  $f(b)$  depends directly on  $\theta$ .

When the parameters  $\theta$  are high-dimensional, gradient-based optimization is appealing because it provides information about how to adjust each parameter individually. Stochastic optimization is essential for large problems, but is only guaranteed to converge to a fixed point of the original objective (?) when the stochastic gradients  $\hat{g}$  are unbiased, i.e.  $\mathbb{E}[\hat{g}] = \nabla_\theta \mathbb{E}_{p(b|\theta)}[f(b)]$ .

**The score-function gradient estimator** One of the most generally-applicable gradient estimators is known as the score-function estimator, or REINFORCE (?):

$$\hat{g}_{\text{reinforce}} = f(b) \nabla_\theta \log p(b|\theta), \quad b \sim p(b|\theta) \quad (2)$$

This estimator is unbiased, but in general has high variance. Intuitively, this estimator is limited by the fact that it doesn't use any information about how  $f$  depends on  $b$ , only on the final outcome of evaluating  $f(b)$ .

**The reparameterization trick** When  $f$  is continuous and differentiable, and the latent variables  $b$  can be written as a deterministic, differentiable function of a random draw from a fixed distribution, the reparameterization trick (???) creates a low-variance, unbiased gradient estimator by making the dependence of  $b$  on  $\theta$  explicit:

$$\hat{g}_{\text{reparam}} = \nabla_\theta f(b(\theta, \epsilon)), \quad \epsilon \sim p(\epsilon) \quad (3)$$

This gradient estimator is used when training high-dimensional, continuous latent-variable models, such as variational autoencoders or GANs. One intuition for why this gradient estimator is preferable to REINFORCE is that it depends on  $\partial f / \partial b$ , which directly exposes the dependence of  $f$  on  $b$ .

**Concrete relaxation** When  $b$  is discrete and  $f$  is known, one general approach is to differentiate a continuous relaxation of the discrete random variables. ? and ? developed a differentiable relaxation of the categorical distribution, called the concrete distribution:

$$\hat{g}_{\text{concrete}} = \nabla_{\theta} f(\sigma_{\lambda}(\log \theta - \log(-\log \mathbf{u}))), \quad \mathbf{u} \sim \text{uniform}[0, 1] \quad (4)$$

where  $\sigma_{\lambda}$  is the softmax function with temperature  $\lambda$  [TODO: Reference this in appendix]. This gradient estimator is fairly effective in practice, but produces biased gradients. Additionally, it is not clear how to set the temperature  $\lambda$ .

**Control variates** Control variates are a general trick for reducing the variance of a Monte Carlo estimator. Given an estimator  $f(b)$ , a control variate is a function  $\hat{f}(b)$  with a known mean  $\mathbb{E}_{p(b)}[\hat{f}]$ . Subtracting the control variate from our estimator and adding its mean gives us a new estimator:

$$\hat{f}_{\text{new}}(b) = f(b) - \hat{f}(b) + \mathbb{E}_{p(b)}[\hat{f}(b)] \quad (5)$$

This new estimator has the same expectation as the old one:

$$\mathbb{E}_{p(b)}[\hat{f}_{\text{new}}] = \mathbb{E}_{p(b)}[f(b) - \hat{f}(b) + \mathbb{E}_{p(b)}[\hat{f}(b)]] = \mathbb{E}_{p(b)}[f(b)] \quad (6)$$

Importantly, the new estimator has lower variance than  $f(b)$  if  $\hat{f}(b)$  is positively correlated with  $f(b)$ .

**REBAR** The recently-developed REBAR method (?) reduces the variance of the REINFORCE gradient estimator for functions of discrete random variables by adding a control variate based on the concrete relaxation (??). ALL THIS STUFF I HAVE COMMENTED OUT SHOULD BE CUT DOWN AND PLACED HERE

## 4 SCOPE AND LIMITATIONS

One major limitation of the REBAR estimator and the concrete relaxation is that they requires the function being optimized, whose input is only defined at discrete inputs, to also accept continuous inputs, and to be differentiable w.r.t. those inputs. This makes REBAR and the concrete relaxation inapplicable for optimizing black-box functions, as in reinforcement learning settings where the reward is an unknown function of the environment.

Following ?, the following overview focuses on a single discrete Bernoulli random variable. However, the generalization to categorical variables is straightforward. (LETS ADD THE REST OF OUR CRITICISMS GRIEVANCES HERE)

## 5 A GENERAL FAMILY OF GRADIENT ESTIMATORS

In this work, we seek to synthesize and generalize the approaches introduced above. Mainly, we desire a simple and generic gradient estimator that can be applied to expectations of known or black-box functions of discrete or random variables. We would like these approaches to be simple to implement and add minimal computational overhead.

The foundation of our approach is the score function gradient estimator with a control variate who's expectation can be estimated with low variance using the reparameterization trick. This control variate is neural network which is trained directly to minimize the variance of the estimated gradients. The central result of this paper is that learning the function in the control variate leads to even better convergence properties and lower variance gradient estimates.

We first outline the algorithm when we are optimizing  $E_{p_{\theta}(b)}[f(b)]$  where  $b$  is a continuous random variable and then when  $b$  is discrete.

### 5.1 CONTINUOUS VARIABLES

We wish to compute  $\frac{\partial}{\partial \theta} \mathbb{E}_{p_{\theta}(b)}[f(b)]$  where  $p_{\theta}(b)$  is a distribution with continuous support that admits a reparameterization  $T$  such that  $T(\theta, \epsilon) = b \sim p_{\theta}(b)$  where  $\epsilon \sim p(\epsilon)$ . We assume that  $f$  is a

black-box or involves operations such that  $\frac{\partial}{\partial \theta} f = 0$  almost everywhere or is not computable. Given the reparameterization, we can re-write this as an expectation over  $\epsilon$  as  $\frac{\partial}{\partial \theta} \mathbb{E}_{p(\epsilon)} [f(T(\epsilon, \theta))]$ . We introduce a new, differentiable function  $r_\phi$  and define our estimator below.

$$\frac{\partial}{\partial \theta} \mathbb{E}_{p_\theta(b)} [f(b)] \approx g_\phi(\epsilon) = [f(b) - r_\phi(b)] \cdot \frac{\partial}{\partial \theta} \log p(b) + \frac{\partial}{\partial \theta} r(b) \quad (7)$$

$$= [f(T(\epsilon, \theta)) - r_\phi(T(\epsilon, \theta))] \cdot \frac{\partial}{\partial \theta} \log p(b) + \frac{\partial}{\partial \theta} r(T(\epsilon, \theta)) \quad (8)$$

We note that our estimator is unbiased for all choices of  $r_\phi$ , i.e  $\mathbb{E}_\epsilon [g_\phi(\epsilon)] = \frac{\partial}{\partial \theta} \mathbb{E}_{p_\theta(b)} [f(b)]$ . A proof can be found in the appendix.

## 5.2 DISCRETE VARIABLES

In the case where  $p_\theta(b)$  is a distribution over discrete variables, we introduce some notation. We assume there exist a continuous distribution  $p_\theta(z)$  and a deterministic mapping  $H$  such that  $H(z) = b \sim p_\theta(b)$ . If  $p(b)$  is Bernoulli then  $p_\theta(z) = \text{Logistic}(z|\theta)$  and  $H(z) = \mathbb{I}(z > 0)$ . If  $b$  is Categorical, then  $p_\theta(z) = \text{Softmax}(z|\theta)$ . We also assume that  $p_\theta(z)$  is reparameterizable. Moreover, we assume the distribution  $p_\theta(z|b)$  is also reparameterizable meaning there exists  $\hat{T}$  such that  $\hat{T}(\hat{\epsilon}, b, \theta) = \hat{z} \sim p_\theta(\hat{z}|b)$  (see Appendix for details).

The estimator is based on a 2-step monte-carlo procedure. First we sample  $\epsilon \sim p(\epsilon)$  from which we produce  $z = T(\theta, \epsilon)$  and  $b = H(z)$ . We then sample  $\hat{\epsilon} \sim p(\hat{\epsilon})$  and produce  $\hat{z} = \hat{T}(\hat{\epsilon}, b, \theta)$ . We use these values to produce our estimator below.

$$\frac{\partial}{\partial \theta} \mathbb{E}_{p_\theta(b)} [f(b)] \approx g_\phi(\epsilon, \hat{\epsilon}) = [f(b) - r_\phi(\hat{z})] \cdot \frac{\partial}{\partial \theta} \log p(b) + \frac{\partial}{\partial \theta} r(z) - \frac{\partial}{\partial \theta} r(\hat{z}) \quad (9)$$

$$= [f(H(T(\epsilon, \theta))) - r_\phi(\hat{T}(\epsilon, H(T(\epsilon, \theta)), \theta))] \cdot \frac{\partial}{\partial \theta} \log p(b) \quad (10)$$

$$+ \frac{\partial}{\partial \theta} r(T(\epsilon, \theta)) - \frac{\partial}{\partial \theta} r_\phi(\hat{T}(\epsilon, H(T(\epsilon, \theta)), \theta)) \quad (11)$$

This estimator is also unbiased i.e  $\mathbb{E}_{\epsilon, \hat{\epsilon}} [g_\phi(\epsilon, \hat{\epsilon})] = \frac{\partial}{\partial \theta} \mathbb{E}_{p_\theta(b)} [f(b)]$ .

## 5.3 OPTIMIZING THE CONTROL VARIATE

For both estimators, we would like to optimize the parameters  $\phi$  to produce the lowest variance. Following ?, we note that for any unbiased estimator  $g_\phi$  we have

$$\frac{\partial}{\partial \phi} \text{Var}(g_\phi) = \frac{\partial}{\partial \phi} \mathbb{E}[g_\phi^2] - \frac{\partial}{\partial \phi} \mathbb{E}[g_\phi]^2 = \frac{\partial}{\partial \phi} \mathbb{E}[g_\phi^2] - \frac{\partial}{\partial \phi} \mathbb{E}_{p_\theta(b)} [f(b)] = \frac{\partial}{\partial \phi} \mathbb{E}[g_\phi^2] = \mathbb{E} \left[ \frac{\partial}{\partial \phi} g_\phi^2 \right]. \quad (12)$$

Therefore, we can minimize the variance of the estimator by training to minimize the square of the gradients it produces. [TODO: talk about how this is probably better than minimizing the squared prediction error]

## 6 APPLICATIONS

We demonstrate the effectiveness of our estimator on a number of challenging optimization problems. Following ? we begin with a simple toy example to illuminate the potential of our method and then continue to the more relevant problems of optimize binary VAE's and reinforcement learning.

## 6.1 TOY EXPERIMENT

We seek to minimize  $\mathbb{E}_{p(b|\theta)}[(b-t)^2]$  as a function of the parameter  $\theta$  where  $p(b|\theta) = \text{Bern}(b|\theta)$ . We set the target  $t = .45$ . We focus on the more challenging case where  $t = .499$ . With this setting of the target, REBAR and competing methods suffer from high variance and are unable to discover the optimal solution  $\theta = 0$ .

The fixed Concrete relaxation of REBAR is unable to produce a gradient whose signal outweighs the sample noise and is therefore unable to solve this problem noticeably faster than REINFORCE. Figure ?? plots the learned relaxations for a fixed value of  $\theta$ .

It can be seen that RELAX learns a relaxation whose derivative points in the direction of decreased loss for all values of reparameterization noise  $u$ , whereas REBAR’s fixed relaxation only does so for values of  $u > t$ .

## 6.2 DISCRETE VARIATIONAL AUTOENCODER

As in (?), we benchmark the RELAX estimator on the task of training a variational autoencoder (??) where all random variables are Bernoulli taking values in  $\{-1, 1\}$ . As in ?, we compare training the variational lower-bound across the MNIST and Omniglot (?) datasets. As in ? we test models with 1 and 2 of Bernoulli random variables with linear mappings between them and a model with 1 layer of Bernoulli random variables with non-linear mappings between layers.

We found that due to the complicated structure of the loss function, the RELAX estimator performed worse than REBAR. Instead we add a learned relaxation to REBAR’s control variate which we denote relaxed-REBAR. Our estimator takes the form of (??) with

$$\bar{r}(z) = r(z) + f(\sigma_\lambda(z))$$

where  $r(z)$  is a learned neural network and  $f(\sigma_\lambda(z))$  is the Concrete relaxation of REBAR with temperature parameter  $\lambda$ . In all experiments, adding the learned  $r(z)$  reduced the variance of the gradients and improved the final results.

	NVIL	MuProp	REBAR ?	REBAR ours	RELAX
<b>MNIST</b>					
Nonlinear	-102.2	-101.5	-101.1	-83.02	<b>-79.49</b>
Linear 1 layer	-112.5	-111.7	-111.6	-111.66	<b>-111.22</b>
Linear 2 Layer	-99.6	-99.07	-98.8	-98.23	<b>-98.04</b>
<b>Omniglot</b>					
Nonlinear	-110.4	-109.58	-108.72	-62.28	<b>-58.55</b>
Linear 1 layer	-117.44	-117.09	-116.83	-116.75	<b>-116.62</b>
Linear 2 Layer	-109.98	-109.55	-108.99	-108.74	<b>-108.59</b>

Table 1: Training variational lower bound after training.

In (?), a separate REBAR estimator was used to estimate the gradients of each model parameter (each weight matrix and bias vector). To apply our estimator to this formulation, we would need to learn a separate relaxation for each model parameter. To get around this, we use our gradient estimator to approximate  $g_\phi \approx \frac{\partial}{\partial \theta} \mathbb{E}_{q(b|\theta)}[f(b)]$  where  $x \cdot W = \theta$  is the parameters of the Bernoulli latent variables,  $W$  is our layer’s weight matrix. We then obtain an estimate of  $\frac{\partial}{\partial W} \mathbb{E}_{q(b|\theta)}[f(b)] = g_\phi \cdot \frac{\partial \theta}{\partial W}$ . We note this gives us unbiased gradients because

$$\mathbb{E}_\epsilon[g_\phi(\epsilon) \cdot \frac{\partial \theta}{\partial W}] = \mathbb{E}_\epsilon[g_\phi(\epsilon)] \cdot \frac{\partial \theta}{\partial W} = \frac{\partial}{\partial \theta} \mathbb{E}_{q(b|\theta)}[f(b)] \cdot \frac{\partial \theta}{\partial W} = \frac{\partial}{\partial W} \mathbb{E}_{q(b|\theta)}[f(b)] \quad (13)$$

To provide a fair comparison, we re-implemented REBAR in this way (denoted REBAR-ours in table 6.2). We believe this explains the large difference in performance between our implementation and that of (?) for the nonlinear models since there are 3 layers of parameters that all share the same gradient estimator. In the linear models, each layer has its own gradient estimator making our implementation closer to that of (?).

### 6.3 REINFORCEMENT LEARNING

Reinforcement learning is strong motivating problem for methods similar to ours. In reinforcement learning we seek to optimize the parameters of a policy distribution  $\pi(a|s; \phi)$  to maximize the (often discounted) sum of future rewards given that policy  $\mathbb{E}_{\pi(\phi)}[\sum_{t=1}^{\infty} r_t]$ . We can view the sum of future rewards as a black-box function of our policy's actions  $f(a)$ . Thus, as before we have reduced the problem to that of estimating  $\frac{\partial \mathbb{E}_{\pi(a|\phi)}[f(a)]}{\partial \phi}$  which is the standard policy gradient algorithm ?.

We test our approach on simple reinforcement learning environments with discrete actions. We use the RELAX estimator and compare with the advantage actor critic algorithm (A2C ?) as a baseline. In all experiments we utilized the same learning rate for the policy network for RELAX and A2C to so differences in performance depended solely on the control variate used.

To compare these approaches in the most illustrative setting possible we run these algorithms on one environment at a time, running each episode to completion. After completion, we generate the discounted reward for each timestep, treat the episode as a single batch of data, and perform one step of gradient decent. We test our algorithm on the Cart-Pole and Lunar-Lander environments from the OpenAI Gym ?. We run the Cart-Pole and Lunar-Lander environments for 250 and 1000 episodes, respectively and plot reward and the log-variance of the policy gradients in figure X.

### 6.4 RL INTRODUCTION

We seek to compute

$$\frac{\partial \mathbb{E}_{\tau}[R]}{\partial \theta} = \mathbb{E} \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \sum_{t'=t}^T r_{t'} \right]$$

but the estimator on the right hand side can have potentially high variance. Instead, we typically compute

$$\frac{\partial \mathbb{E}_{\tau}[R]}{\partial \theta} = \mathbb{E}_{\tau} \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \left[ \left( \sum_{t'=t}^T r_{t'} \right) - b(s_t) \right] \right]$$

This is unbiased because

$$E_{a_{1:T}, s_{1:T}} \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \cdot b(s_t) \right] = E_{a_{1:t-1}, s_{1:t}} \left[ E_{a_{t:T}, s_{t+1:T}} \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \cdot b(s_t) \right] \right] \quad (14)$$

$$= E_{a_{1:t-1}, s_{1:t}} \left[ b(s_t) \cdot E_{a_{t:T}, s_{t+1:T}} \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \right] \right] \quad (15)$$

$$= E_{a_{1:t-1}, s_{1:t}} \left[ b(s_t) \cdot \frac{\partial}{\partial \theta} E_{a_{t:T}, s_{t+1:T}} [1] \right] \quad (16)$$

$$= 0 \quad (17)$$

$$= 0 \quad (18)$$

Where  $b(s_t)$  is trained to minimize  $(b(s_t) - \sum_{i=t}^T r_i)^2$ . This is unbiased for any choice of  $b$  since  $b$  does not depend on  $a_t$ .

#### 6.4.1 VARIANCE REDUCTION WITH RELAX

We begin in the case where actions are continuous and any expectation over an action is replaced with an expectation over reparameterization variables  $\epsilon$

$$\mathbb{E}_a[f(a)] = \mathbb{E}_{\epsilon}[f(a(\epsilon))]$$

. We are interested in replacing  $b(s_t)$  with a function  $m(a_t, s_t)$  which is a function of both actions and states.

This changes the first equation to

$$\frac{\partial \mathbb{E}_{\tau}[R]}{\partial \theta} = \mathbb{E}_{\tau} \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \left[ \left( \sum_{t'=t}^T r_{t'} \right) - m(a_t, s_t) \right] \right]$$

but this makes the estimator biased as  $\mathbb{E}_\tau \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} m(a_t, s_t) \right] \neq 0$ . Because of the dependence on  $a_t$  we cannot pull this out of the expectation as we could in line 2 of the above equation.

Instead we subtract out a term whose expectation should be the same as the score function version of the control variate changing our estimator to

$$\frac{\partial \mathbb{E}_\tau[R]}{\partial \theta} = \mathbb{E}_\tau \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \left[ \left( \sum_{t'=t}^T r_{t'} \right) - m(a_t, s_t) \right] + \frac{\partial m(a_t, s_t)}{\partial \theta} \right].$$

For this estimator to be unbiased, we must have

$$\mathbb{E}_\tau \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} m(a_t, s_t) - \frac{\partial m(a_t, s_t)}{\partial \theta} \right] = 0$$

and we claim that  $\forall t$ , we have

$$\mathbb{E}_\tau \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} m(a_t, s_t) - \frac{\partial m(a_t, s_t)}{\partial \theta} \right] = 0$$

We begin with

$$\mathbb{E}_\tau \left[ \frac{\partial m(a_t, s_t)}{\partial \theta} \right] = \mathbb{E}_{a_{1:t-1}, s_{1:t}} \left[ E_{a_{t:T}, s_{t+1:T}} \left[ \frac{\partial m(a_t, s_t)}{\partial \theta} \right] \right] \quad (19)$$

$$= \mathbb{E}_{a_{1:t-1}, s_{1:t}} \left[ E_{a_t} \left[ \frac{\partial m(a_t, s_t)}{\partial \theta} \right] \right] \quad (20)$$

$$= \mathbb{E}_{a_{1:t-1}, s_{1:t}} \left[ \frac{\partial E_{a_t} [m(a_t, s_t)]}{\partial \theta} \right] \quad (21)$$

$$= \mathbb{E}_{a_{1:t-1}, s_{1:t}} \left[ E_{a_t} \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} m(a_t, s_t) \right] \right] \quad (22)$$

$$= E_\tau \left[ \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} m(a_t, s_t) \right] \quad (23)$$

which completes our proof.

#### 6.4.2 DISCRETE CASE

In the discrete action setting our policy parameterizes a soft-max distribution which we use to sample actions. In the discrete case we define  $z_t = f(\pi, u) = \sigma(\log \pi - \log(-\log(u)))$  where  $u \sim \text{Unif}[0, 1]$ ,  $a_t = \text{argmax}(z_t)$ ,  $\sigma$  is the soft-max function.

We also define  $\tilde{z}_t \sim p(z_t|a_t)$  so if the  $z_t$  are gumbel softmax samples, then  $\tilde{z}_t$  are gumbel softmax samples constrained so that the largest value is that of the index of  $a_t$ . See appendix for how to efficiently generate samples  $\tilde{z}_t$ .

In the discrete case, we use  $m(\tilde{z}_t, s_t) \cdot \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta}$  as our control variate giving us

$$\frac{\partial \mathbb{E}_\tau[R]}{\partial \theta} = \mathbb{E}_\tau \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \left[ \left( \sum_{t'=t}^T r_{t'} \right) - m(\tilde{z}_t, s_t) \right] \right]$$

which is biased, so we must add a term to remove this bias which will be  $\frac{\partial m(z_t, s_t)}{\partial \theta} - \frac{\partial m(\tilde{z}_t, s_t)}{\partial \theta}$  making the full estimator

$$\frac{\partial \mathbb{E}_\tau[R]}{\partial \theta} = \mathbb{E}_\tau \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} \left[ \left( \sum_{t'=t}^T r_{t'} \right) - m(\tilde{z}_t, s_t) \right] + \frac{\partial m(z_t, s_t)}{\partial \theta} - \frac{\partial m(\tilde{z}_t, s_t)}{\partial \theta} \right].$$

For this estimator to be unbiased, we must have

$$\mathbb{E}_\tau \left[ \sum_{t=1}^T \frac{\log \pi(a_t|s_t, \theta)}{\partial \theta} m(\tilde{z}_t, s_t) - \left( \frac{\partial m(z_t, s_t)}{\partial \theta} - \frac{\partial m(\tilde{z}_t, s_t)}{\partial \theta} \right) \right] = 0$$

and we claim that  $\forall t$

$$\mathbb{E}_\tau \left[ \frac{\log \pi(a_t | s_t, \theta)}{\partial \theta} m(\tilde{z}_t, s_t) - \left( \frac{\partial m(z_t, s_t)}{\partial \theta} - \frac{\partial m(\tilde{z}_t, s_t)}{\partial \theta} \right) \right] = 0$$

We introduce some notation that  $a(z) = \operatorname{argmax}(z)$

$$0 = \frac{\partial}{\partial \theta} \mathbb{E}_\tau \left[ m(z_t, s_t) - m(\tilde{z}_t, s_t) \right] = \mathbb{E}_{a_{<t}, s_{\leq t}} \left[ \frac{\partial}{\partial \theta} \mathbb{E}_{z_t | s_t} \left[ m(z_t, s_t) \right] - \frac{\partial}{\partial \theta} \mathbb{E}_{a_t | s_t} \left[ \mathbb{E}_{z_t | a_t, s_t} \left[ m(z_t, s_t) \right] \right] \right] \quad (24)$$

We drop the outer expectation for brevity and continue

$$\frac{\partial}{\partial \theta} \mathbb{E}_{z_t | s_t} \left[ m(z_t, s_t) \right] - \frac{\partial}{\partial \theta} \mathbb{E}_{a_t | s_t} \left[ \mathbb{E}_{z_t | a_t, s_t} \left[ m(z_t, s_t) \right] \right] \quad (25)$$

$$= \frac{\partial}{\partial \theta} \mathbb{E}_{z_t | s_t} \left[ m(z_t, s_t) \right] - \mathbb{E}_{a_t | s_t} \left[ \mathbb{E}_{z_t | a_t, s_t} \left[ m(z_t, s_t) \frac{\log \pi(a_t | s_t, \theta)}{\partial \theta} - \frac{\partial}{\partial \theta} \mathbb{E}_{z_t | a_t, s_t} \left[ m(z_t, s_t) \right] \right] \right] \quad (26)$$

$$= \mathbb{E}_{z_t | s_t} \left[ \frac{\partial}{\partial \theta} m(z_t, s_t) \right] - \mathbb{E}_{a_t | s_t} \left[ \mathbb{E}_{z_t | a_t, s_t} \left[ m(z_t, s_t) \frac{\log \pi(a_t | s_t, \theta)}{\partial \theta} - \mathbb{E}_{z_t | a_t, s_t} \left[ \frac{\partial}{\partial \theta} m(z_t, s_t) \right] \right] \right] \quad (27)$$

Thus due to the markov property of the MDP, then we can wrap all of these expectations into the expectation over trajectories giving us

$$\mathbb{E}_{a_{<t}, s_{\leq t}} \left[ \mathbb{E}_{z_t | s_t} \left[ \frac{\partial}{\partial \theta} m(z_t, s_t) \right] - \mathbb{E}_{a_t | s_t} \left[ \mathbb{E}_{z_t | a_t, s_t} \left[ m(z_t, s_t) \frac{\log \pi(a_t | s_t, \theta)}{\partial \theta} - \mathbb{E}_{z_t | a_t, s_t} \left[ \frac{\partial}{\partial \theta} m(z_t, s_t) \right] \right] \right] \right] \quad (28)$$

$$= \mathbb{E}_\tau \left[ \frac{\log \pi(a_t | s_t, \theta)}{\partial \theta} m(\tilde{z}_t, s_t) - \left( \frac{\partial m(z_t, s_t)}{\partial \theta} - \frac{\partial m(\tilde{z}_t, s_t)}{\partial \theta} \right) \right] = 0 \quad (29)$$

#### 6.4.3 EXPERIMENTS

We run experiments on some small reinforcement learning tasks with both discrete and continuous actions. We compare our method against a baseline of advantage actor-critic. We do not bootstrap the rewards with the value function. We estimate the policy gradient using a single monte-carlo sample produced from a one episode roll-out of the current policy. We note that improved performance could be achieved by using more samples, but we intended to test our model in the highest-possible variance setting. In all experiments the model architectures were identical between our model and the baseline and the learning rates were also constant across tests making it that all improvements derive from having a lower variance estimate of the policy gradient.

We test our approach in the discrete action domains of the CartPole and LunarLander problems as provided by the OpenAI gym ?.

?

[Do we use ADAM (?) for optimization?]

## 7 RELATED WORK

? further reduce the variance of reparameterization gradients in an orthogonal way.

As gradient estimators become more complex, checking their unbiasedness numerically becomes difficult. The automatic theorem-proving-based unbiasedness checker developed by ? may become relevant to this line of research.

NVIL (?), VIMCO (?)



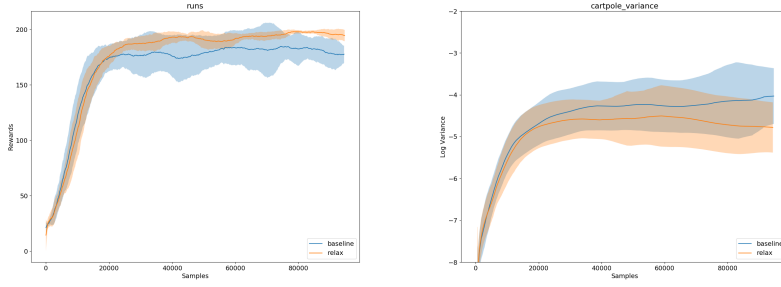


Figure 2: Cartpole Experiments. *Left: Rewards Right: Log Variances.*

? address the general problem of developing gradient estimators for deterministic black-box functions or discrete optimization. They introduce a sampling distribution, and optimize an objective similar to ours.

? also introduce a sampling distribution to build a gradient estimator, and consider optimizing the sampling distribution.

? reduce the variance of actor-critic gradient estimates by simply summing over all possible actions.

? estimate gradients using a form of finite differences, evaluating hundreds of different parameter values in parallel to construct a gradient estimator. In contrast, our method is a simple-sample estimator.

Generalized Reparameterization Gradients REBAR and the generalization in this paper uses a mixture of score function and reparameterization gradients. A recent paper by ? unifies these two gradient estimators as the generalized reparameterization gradient (GRG). This framework can help disentangle the various components of generalized REBAR.

REBAR innovation as further decomposition the correction term into secondary reparameterization components note this is a recursive application of the principles of GRG observe that the GRG suggests this recursive application to components of an estimator propose that other estimators could be similarly recursively decomposed?

## 8 CONCLUSIONS AND FUTURE WORK

Other possible applications:

GANs (?) that generate text or other discrete objects.

Learning to parse (?)

VAEs with continuous latent variables but non-differentiable likelihood functions.

## 9 APPENDIX A: CONTROL VARIATES

Generalizing the reparameterization trick

Write sample from distribution  $s(\epsilon)$  as  $\epsilon = \mathcal{T}^{-1}(\mathbf{z}; \nu)$  for some invertible transform  $\mathcal{T}$  with variational parameters  $\nu$ . write out transformed density example: normal with standard normal  $s$  example: inverse CDF of Gaussian with uniform  $s$  write out expected gradient under transformation show decomposition of expected gradient into reparameterization and correction terms

## 10 APPENDIX B: CATEGORICAL VARIABLES

Let  $G_{1:k} = -\log -\log(U_{i:k})$  be samples from the Gumbel distribution, and learnable parameters  $(\alpha_1, \dots, \alpha_k)$  be interpreted as some unnormalized parameterization of the discrete distribution under consideration. Then, consider the following sampling procedure: for each  $k$ , find the  $k$  that maximizes  $\log \alpha_k - G_k$ , and then set  $D_k = 1$  and  $D_{i \neq k} = 0$ . The Gumbel-Max trick states that sampling from the discrete distribution is equivalent to taking this argmax, that is,  $p(D_k = 1) = \alpha_k / \sum_{i=1}^n \alpha_i$ .

Since taking an argmax is still a discontinuous operation, ? and ? proposed further relaxing the argmax operator through the softmax function with an additional temperature parameter  $\lambda$ :

$$x_k = \frac{\exp\{(\log \alpha_k + G_k)/\lambda\}}{\sum_{i=1}^n \exp\{(\log \alpha_i + G_i)/\lambda\}} \quad (30)$$

This relaxation allows values within the simplex, but in the low temperature limit, it becomes exactly the discrete argmax. One limitation of the concrete distribution is that it is a biased estimator except in limiting temperature. In other words, a small amount of bias is present for a non-zero temperature.