

Some notes on Policy Gradients for Reinforcement Learning

David Meyer

dmm@{1-4-5.net,uoregon.edu,...}

Last update: September 11, 2017

1 Introduction

Policy gradients for Reinforcement Learning (RL) are really beautiful in form and function. In these notes I'll try to point out a few of these amazing features¹.

Here we consider methods that learn a parameterized policy that can select actions without consulting a value function. A value function may still be used to learn the policy weights, but that is not required for action selection. The notation $\theta \in \mathbb{R}^n$ is typically used for the primary learned weight vector, and $\pi(a | s, \theta) = Pr(A_t = a | s_t = s, \theta_t = \theta)$ is the probability that action a is taken at time t given that the agent is in state s at time t with weight vector θ . If a method uses a learned value function as well, then the value function's weight vector is denoted w to distinguish it from θ , such as in $\hat{v}(s, w)$.

We wish to consider methods for learning the policy weights θ based on the gradient of some performance measure $\eta(\theta)$ with respect to the policy weights. These methods seek to maximize performance, so their updates approximate gradient ascent in η . As is typical for gradient descent/ascent,

$$\theta_{t+1} := \theta_t + \alpha \widehat{\nabla \eta(\theta_t)} \tag{1}$$

where, as is typical with gradient ascent/descent methods, $\widehat{\nabla \eta(\theta_t)}$ is a stochastic estimate whose expectation approximates the gradient of the performance measure η with respect to its argument θ_t . As we will see, $\widehat{\nabla \eta(\theta_t)}$ will turn out to be something like $\nabla_{\theta} \mathbb{E}_{\tau} [R(\tau)]$, where τ is a *trajectory* and $R(\tau)$ is the return for path τ (see Section 2).

¹Here we follow the notation used in [1].

2 Score Function/Likelihood Ratio Estimators

First notices that Likelihood Ratio methods only change the probabilities of experienced paths, and further, these methods do not try to change the actions taken in a given path. That said, recall that our problem here is to solve the following:

$$\max_{\theta} \mathbb{E} \left[\sum_{t=0}^H R(s_t) \mid \pi_{\theta} \right] \quad (2)$$

where $\pi_{\theta}(u \mid s)$ is the probability of action u in state s . Now, because ultimately we'll want to compute the gradient $\nabla_{\theta} \mathbb{E}[R \mid \pi_{\theta}]$, we need to prove one identity, namely (imagine that $f(x) = R$)

$$\nabla_{\theta} \mathbb{E}_{x \sim p(x|\theta)}[f(x)] = \mathbb{E}_x[f(x) \nabla_{\theta} \log p(x \mid \theta)] \quad (3)$$

So the proof of this is amazingly simple given the log derivative trick². First, consider a function $f(x)$ for which we wish to find $\nabla_{\theta} \mathbb{E}_{x \sim p(x|\theta)}[f(x)]$. Then

$$\nabla_{\theta} \mathbb{E}_{x \sim p(x|\theta)}[f(x)] = \nabla_{\theta} \int_x p(x \mid \theta) f(x) dx \quad \# \text{ definition} \quad (4)$$

$$= \int_x \nabla_{\theta} p(x \mid \theta) f(x) dx \quad \# \text{ Lebesgue Integral} \quad (5)$$

$$= \int_x p(x \mid \theta) \frac{\nabla_{\theta} p(x \mid \theta)}{p(x \mid \theta)} f(x) dx \quad \# \text{ multiply by } \frac{p(x \mid \theta)}{p(x \mid \theta)} \quad (6)$$

$$= \int_x p(x \mid \theta) \nabla_{\theta} \log p(x \mid \theta) f(x) \quad \# \text{ log-derivative trick} \quad (7)$$

$$= \mathbb{E}_{x \sim p(x|\theta)}[f(x) \nabla_{\theta} \log p(x \mid \theta)] \quad \# \text{ defn expectation} \quad (8)$$

Define $\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i \mid \theta)$. An empirical estimate the gradient for m samples is

$$\nabla_{\theta} \mathbb{E}_{x \sim p(x|\theta)}[f(x)] = g \approx \frac{1}{m} \sum_{i=1}^m \hat{g}_i = \frac{1}{m} \sum_{i=1}^m f(x_i) \nabla_{\theta} \log p(x_i \mid \theta) \quad (9)$$

Aside: The "log derivative trick", which is also called the "likelihood ratio trick" since $\frac{\nabla_{\theta} p(x|\theta)}{p(x|\theta)}$ is called the "likelihood ratio" and $\frac{\nabla_{\theta} p(x|\theta)}{p(x|\theta)} = \nabla_{\theta} \log p(x \mid \theta)$, is sometimes framed

²<http://blog.shakirm.com/2015/11/machine-learning-trick-of-the-day-5-log-derivative-trick>

up like this: The gradient of something divided by something is the gradient of log something. Also frequently pronounced "grad something divided by something is grad log something".

So $\nabla_{\theta} E_{x \sim p(x|\theta)}[f(x)] = E_{x \sim p(x|\theta)}[f(x) \nabla_{\theta} \log p(x|\theta)]$. This gives us an unbiased gradient estimator; to compute the gradient estimate, just sample $x_i \sim p(x|\theta)$ and then compute the gradient estimate $\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i|\theta)$. Now, let the trajectory (sometimes path) τ be the state-action sequence $(s_0, a_0, r_0, s_1, a_1, r_1, \dots, a_{T-1}, a_{T-1}, r_{T-1})$ and suppose $f(x) = R(\tau)$, the total return for path τ . Then

$$\nabla_{\theta} \mathbb{E}_{\tau}[R(\tau)] = E_{\tau}[\nabla_{\theta} \log p(\tau|\theta) R(\tau)] \quad (10)$$

so that all we really need to do is see that we can compute $p(\tau|\theta)$:

$$p(\tau|\theta) = \mu_0(s_0) \prod_{t=0}^{T-1} \underbrace{\pi(a_t | s_t, \theta)}_{\text{policy}} \underbrace{P(s_{t+1}, r_t | s_t, a_t)}_{\text{dynamics}} \quad \# \text{definition} \quad (11)$$

$$\log p(\tau|\theta) = \log \mu_0(s_0) + \sum_{t=0}^{T-1} \log \pi(a_t | s_t, \theta) + \log P(s_{t+1}, r_t | s_t, a_t) \quad (12)$$

$$\nabla_{\theta} \log p(\tau|\theta) = \nabla_{\theta} \sum_{t=0}^{T-1} \log \pi(a_t | s_t, \theta) \quad (13)$$

$$= \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t, \theta) \quad (14)$$

So $\nabla_{\theta} \mathbb{E}_{\tau}[R] = \mathbb{E}_{\tau}[R(\tau) \nabla_{\theta} \sum_{t=0}^{T-1} \log \pi(a_t | s_t, \theta)]$, which gives a direct way of computing a gradient we can use with standard Stochastic Gradient Ascent (or Descent), where the parameter update rule would be something like $\theta := \theta + \alpha \nabla_{\theta} \mathbb{E}_{\tau}[R(\tau)]$. Essentially move θ in the direction of better expected return. Note also that the model dynamics, $P(s_{t+1}, r_t | s_t, a_t)$, conveniently drops out of the gradient (good thing given that we're trying to learn model-free control).

3 Derivation of the Gradient via Importance Sampling

It turns out that almost all on-policy control methods utilize importance sampling, which is a general technique for estimating expected values under one distribution given samples

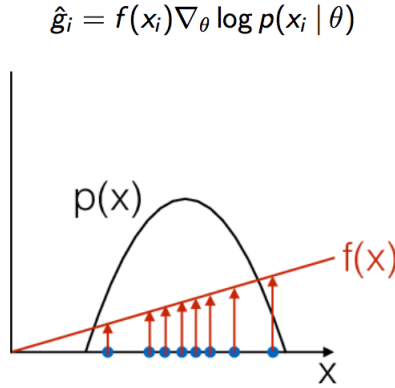


Figure 1: Score/likelihood Gradient Estimator Intuition 1. Image courtesy John Schulman <https://learning.mpi-sws.org/mlss2016/slides/2016-MLSS-RL.pdf>

from another[1]. The basic approach to applying importance sampling to on-policy learning is by weighting returns according to the relative probability of their trajectories occurring under the target and behavior policies; this is called the importance-sampling ratio. Now given a starting state S_t , the probability of the subsequent state-action trajectory $\tau = A_t, S_{t+1}, A_{t+1} \dots, S_T$ occurring under any policy π is

$$p_t(\tau) = \prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k) \tag{15}$$

where p is the state-transition probability function. The relative probability of the trajectory under the target and behavior policies (the importance-sampling ratio) is then

$$\rho_t^T = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} u(A_k | S_k) p(S_{k+1} | S_k, A_k)} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k)}{\prod_{k=t}^{T-1} u(A_k | S_k)} \tag{16}$$

One of the many interesting things here is that it turns out that ρ_t^T doesn't depend on the MDP's transition probability function (which is what we want since we want to solve the model-free problem). In the following example, consider $\rho_t^T = \frac{p(x|\theta)}{p(x|\theta_{old})}$.

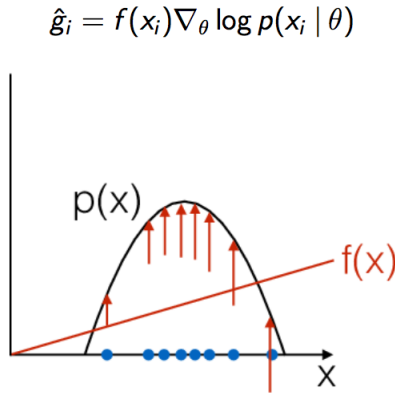


Figure 2: Score/likelihood Gradient Estimator Intuition 2. Image courtesy John Schulman <https://learning.mpi-sws.org/mlss2016/slides/2016-MLSS-RL.pdf>

$$\mathbb{E}_{x \sim \theta} [f(x)] = \mathbb{E}_{x \sim \theta_{old}} \left[\frac{p(x | \theta)}{p(x | \theta_{old})} f(x) \right] \quad (17)$$

$$\nabla_{\theta} \mathbb{E}_{x \sim \theta} [f(x)] = \mathbb{E}_{x \sim \theta_{old}} \left[\frac{\nabla_{\theta} p(x | \theta)}{p(x | \theta_{old})} f(x) \right] \quad (18)$$

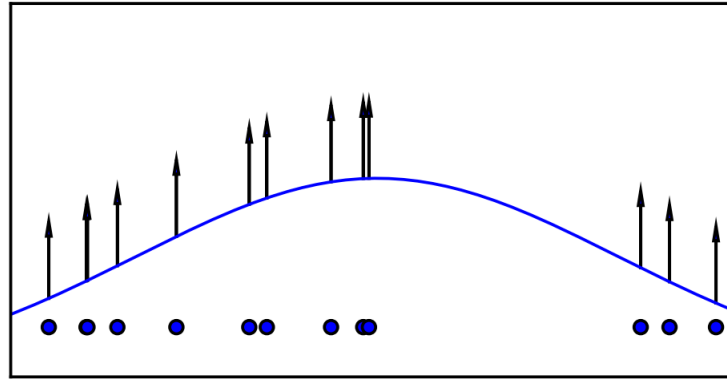
$$\nabla_{\theta} \mathbb{E}_{x \sim \theta} [f(x)] \Big|_{\theta = \theta_{old}} = \mathbb{E}_{x \sim \theta_{old}} \left[\frac{\nabla_{\theta} p(x | \theta) \Big|_{\theta = \theta_{old}}}{p(x | \theta_{old})} f(x) \right] \quad (19)$$

$$= \mathbb{E}_{x \sim \theta_{old}} \left[\nabla_{\theta} \log p(x | \theta) \Big|_{\theta = \theta_{old}} f(x) \right] \quad (20)$$

4 A Bit of Intuition

What we've learned so far is that we can compute the an estimate of the gradient for sample i is $\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i | \theta)$. Now, suppose $f(x)$ somehow measures how good the sample x is. Then we can conclude a few things about \hat{g}_i :

- Moving in the direction \hat{g}_i pushes up the log probability ("logprob") of the sample in proportion to how good it is (i.e., $f(x_i)$), remembering that the SGD update rule is



$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} \log p_{\text{model}}(x | \theta)$$

Figure 3: Maximum Likelihood Estimation. Image courtesy Ian Goodfellow [2].

something like $\theta := \theta + \alpha \nabla_{\theta} \mathbb{E}_{\tau} [R(\tau)]$. This pushing up is depicted in Figures 1 and 2.

- Note that this is true even if $f(x)$ is discontinuous, and unknown, or the sample space (containing x) is a discrete set!
- More generally, Maximum Likelihood Estimation (MLE) consists of taking several samples from the Data Generating Distribution (DGD) to form a training set, then pushing up on the probability the model assigns to those points, thus maximizing the likelihood of the training data. Figure 3 shows how different data points push up on different parts of the density function for a Gaussian model applied to 1-D data. The fact that the density function must sum to one means that we cannot simply assign infinite likelihood to all points; as one point pushes up in one place it inevitably pulls down in other places. The resulting density function balances out the upward forces from all the data points in different locations.

Basically, the gradient tries to increase probability of paths with positive return R and decrease the probability of paths with negative return.

4.1 Reorganizing Equation 10

Equation 10 can be rewritten as follows:

$$\nabla_{\theta} \mathbb{E}_{\tau} [R(\tau)] = E_{\tau} [\nabla_{\theta} \log p(\tau | \theta) R(\tau)] \quad (21)$$

$$= \mathbb{E} \left[\left(\sum_{t=0}^{T-1} r_t \right) \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t, \theta) \right) \right] \quad (22)$$

But notice that we can limit the trajectory to t steps and derive gradient estimator for one reward term $r_{t'}$. In particular

$$\nabla_{\theta} \mathbb{E} [r_{t'}] = \mathbb{E} \left[r_{t'} \sum_{t=0}^t \nabla_{\theta} \log \pi(a_t | s_t, \theta) \right] \quad (23)$$

Now, we can sum over t to get

$$\nabla_{\theta} \mathbb{E} [R] = \mathbb{E} \left[\sum_{t'=0}^{T-1} r_{t'} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t, \theta) \right] \quad (24)$$

$$= \mathbb{E} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t, \theta) \sum_{t'=t}^{T-1} r_{t'} \right] \quad (25)$$

5 Adding a Baseline: Deriving an Unbiased Estimator

Here's the problem: suppose that $f(x) \geq 0 \forall x$. Then for every x_i , the gradient estimator \hat{g}_i tries to push up its density (recall that $\theta_{t+1} := \theta_t + \alpha \widehat{\nabla_{\theta} \eta(\theta_t)}$). Better would be to have a (new) unbiased estimator that avoids this problem; this estimator would only push up the density for better-than-average x_i . What would this look like?

$$\nabla_{\theta} \mathbb{E}_x [f(x)] = \nabla_{\theta} \mathbb{E}_x [(f(x) - b)] \quad (26)$$

$$= \mathbb{E}_x [\nabla_{\theta} \log p(x | \theta) (f(x) - b)] \quad (27)$$

Note that a near optimal choice for b is always $\mathbb{E}[f(x)]$; however $\mathbb{E}[f(x)]$ must itself be estimated.

To derive our estimate with baseline, recall that

$$\nabla_{\theta} \mathbb{E}_{\tau} [R] = \left[\sum_{t'=0}^{T-1} r_{t'} \sum_{t=t'}^{t'} \nabla_{\theta} \log \pi(a_t | s_t, \theta) \right] \quad (28)$$

$$= \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \log \pi(a_t | s_t, \theta) \left(\sum_{t'=t}^{T-1} r_{t'} - b(s_t) \right) \right] \quad (29)$$

for any "baseline" function $b : \mathcal{S} \rightarrow \mathbb{R}$. See Schulman, J. , et al. [3] for details of the proof.

References

- [1] R. S. Sutton and A. G. Barto., *Reinforcement learning: An introduction*, vol. 1. MIT Press, 1998.
- [2] I. Goodfellow, “NIPS 2016 tutorial: Generative adversarial networks,” 01 2017.
- [3] J. Schulman, N. Heess, T. Weber, and P. Abbeel, “Gradient estimation using stochastic computation graphs,” 06 2015.