

# 11 Policy Gradient Estimation

The previous chapter discussed several ways to go about directly optimizing the parameters of a policy to maximize expected utility. In many applications, it is often useful to use the gradient of the utility with respect to the policy parameters to guide the optimization process. This chapter discusses several approaches to estimating this gradient from trajectory rollouts.<sup>1</sup> A major challenge with this approach is the variance of the estimate due to the stochastic nature of the trajectories arising from both the environment and our exploration of it. The next chapter will discuss how to use these algorithms to estimate gradients for the purpose of policy optimization.

<sup>1</sup> An additional resource on this topic is M. C. Fu, "Gradient Estimation," in *Simulation*, S. G. Henderson and B. L. Nelson, eds., Elsevier, 2006, pp. 575–616.

## 11.1 Finite Difference

*Finite difference* methods estimate the gradient of a function from small changes in its evaluation. Recall that the derivative of a univariate function  $f$  is

$$\frac{df}{dx}(x) = \lim_{\delta \rightarrow 0} \frac{f(x + \delta) - f(x)}{\delta} \quad (11.1)$$

The derivative at  $x$  can be approximated by a sufficiently small step  $\delta > 0$ :

$$\frac{df}{dx}(x) \approx \frac{f(x + \delta) - f(x)}{\delta} \quad (11.2)$$

This approximation is illustrated in figure 11.1.

The gradient of a multivariate function  $f$  with an input of length  $n$  is

$$\nabla f(\mathbf{x}) = \left[ \frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right] \quad (11.3)$$

Finite differences can be applied to each dimension to estimate the gradient.

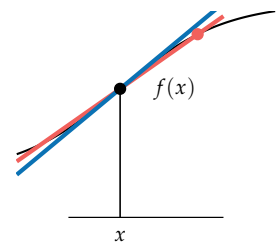


Figure 11.1. The finite difference method approximates the derivative of  $f(x)$  using an evaluation of a point near  $x$ . The finite-difference approximation, in red, is not a perfect match for the true derivative, in blue.

In the context of policy optimization, we want to estimate the gradient of the utility expected from following a policy parameterized by  $\theta$ :

$$\nabla U(\theta) = \left[ \frac{\partial U}{\partial \theta_1}(\theta), \dots, \frac{\partial U}{\partial \theta_n}(\theta) \right] \quad (11.4)$$

$$\approx \left[ \frac{U(\theta + \delta \mathbf{e}^{(1)}) - U(\theta)}{\delta}, \dots, \frac{U(\theta + \delta \mathbf{e}^{(n)}) - U(\theta)}{\delta} \right] \quad (11.5)$$

where  $\mathbf{e}^{(i)}$  is the  $i$ th *standard basis* vector, consisting of zeros except for the  $i$ th component, which is set to 1.

As discussed in section 10.1, we need to simulate policy rollouts to estimate  $U(\theta)$ . We can use algorithm 11.1 to generate trajectories. From these trajectories, we can compute their return and estimate the utility associated with the policy. Algorithm 11.2 implements the gradient estimate in equation (11.5) by simulating  $m$  rollouts for each component and averaging the returns.

```
function simulate( $\mathcal{P}$ ::MDP, s,  $\pi$ , d)
     $\tau$  = []
    for i = 1:d
        a =  $\pi$ (s)
        s', r =  $\mathcal{P}$ .TR(s, a)
        push!( $\tau$ , (s, a, r))
        s = s'
    end
    return  $\tau$ 
end
```

Algorithm 11.1. A method for generating a trajectory associated with problem  $\mathcal{P}$  starting in state  $s$  and executing policy  $\pi$  to depth  $d$ . It creates a vector  $\tau$  containing state-action-reward tuples.

A major challenge in arriving at accurate estimates of the policy gradient is the fact that the variance of the trajectory rewards can be quite high. One approach to reduce the resulting variance in the gradient estimate is to have each rollout share the same random generator seeds.<sup>2</sup> This approach can be helpful, for example, in cases where one rollout happens to hit a low-probability transition early on. Other rollouts will have the same tendency due to the shared random generator, and their rewards will tend to be biased in the same way.

Policy representations have a significant effect on the policy gradient. Example 11.1 demonstrates the sensitivity of the policy gradient to the policy parameterization. Finite differences for policy optimization can perform poorly when the parameters differ in scale.

<sup>2</sup> This random seed sharing is used in the PEGASUS algorithm. A. Y. Ng and M. Jordan, "A Policy Search Method for Large MDPs and POMDPs," in *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2000.

```

struct FiniteDifferenceGradient
     $\mathcal{P}$  # problem
     $\mathbf{b}$  # initial state distribution
     $d$  # depth
     $m$  # number of samples
     $\delta$  # step size
end

function gradient(M::FiniteDifferenceGradient,  $\pi$ ,  $\theta$ )
     $\mathcal{P}$ ,  $\mathbf{b}$ ,  $d$ ,  $m$ ,  $\delta$ ,  $\gamma$ ,  $n$  = M. $\mathcal{P}$ , M. $\mathbf{b}$ , M. $d$ , M. $m$ , M. $\delta$ , M. $\mathcal{P}$ . $\gamma$ , length( $\theta$ )
     $\Delta\theta(\mathbf{i})$  = [ $i$  ==  $k$  ?  $\delta$  : 0.0 for  $k$  in 1: $n$ ]
     $R(\tau)$  = sum( $r * \gamma^{k-1}$  for ( $k$ , ( $\mathbf{s}, \mathbf{a}, r$ )) in enumerate( $\tau$ ))
     $U(\theta)$  = mean( $R(\text{simulate}(\mathcal{P}, \text{rand}(\mathbf{b}), \mathbf{s} \rightarrow \pi(\theta, \mathbf{s}), d))$  for  $i$  in 1: $m$ )
     $\Delta U$  = [ $U(\theta + \Delta\theta(\mathbf{i})) - U(\theta)$  for  $i$  in 1: $n$ ]
    return  $\Delta U ./ \delta$ 
end

```

Algorithm 11.2. A method for estimating a policy gradient using finite differences for a problem  $\mathcal{P}$ , a parameterized policy  $\pi(\theta, \mathbf{s})$ , and a policy parameterization vector  $\theta$ . Utility estimates are made from  $m$  rollouts to depth  $d$ . The step size is given by  $\delta$ .

Consider a single-state, single-step MDP with a one-dimensional continuous action space and a reward function  $R(s, a) = a$ . In this case, larger actions produce higher rewards.

Suppose we have a stochastic policy  $\pi_\theta$  that samples its action according to a uniform distribution between  $\theta_1$  and  $\theta_2$  for  $\theta_2 > \theta_1$ . The expected value is

$$U(\theta) = \mathbb{E}[a] = \int_{\theta_1}^{\theta_2} a \frac{1}{\theta_2 - \theta_1} da = \frac{\theta_1 + \theta_2}{2}$$

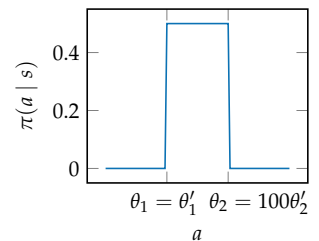
The policy gradient is

$$\nabla U(\theta) = [1/2, 1/2]$$

The policy could be reparameterized to draw actions from a uniform distribution between  $\theta'_1$  and  $100\theta'_2$ , for  $100\theta'_2 > \theta'_1$ . Now the expected reward is  $(\theta'_1 + 100\theta'_2)/2$  and the policy gradient is  $[1/2, 50]$ .

The two parameterizations can represent the same policies, but they have very different gradients. Finding a suitable perturbation scalar for the second policy is much more difficult because the parameters vary widely in scale.

Example 11.1. An example of how policy parameterization has a significant impact on the policy gradient.



## 11.2 Regression Gradient

Instead of estimating the gradient at  $\theta$  by taking a fixed step along each coordinate axis, as done in the previous section, we can use *linear regression*<sup>3</sup> to estimate the gradient from the results of random perturbations from  $\theta$ . These perturbations are stored in a matrix as follows:<sup>4</sup>

$$\Delta\Theta = \begin{bmatrix} (\Delta\theta^{(1)})^\top \\ \vdots \\ (\Delta\theta^{(m)})^\top \end{bmatrix} \quad (11.6)$$

More policy parameter perturbations will tend to produce better gradient estimates.<sup>5</sup>

For each of these perturbations, we perform a rollout and estimate the change in utility:<sup>6</sup>

$$\Delta\mathbf{U} = \left[ U(\theta + \Delta\theta^{(1)}) - U(\theta), \dots, U(\theta + \Delta\theta^{(m)}) - U(\theta) \right] \quad (11.7)$$

The policy gradient estimate using linear regression is then<sup>7</sup>

$$\nabla U(\theta) \approx \Delta\Theta^+ \Delta\mathbf{U} \quad (11.8)$$

Algorithm 11.3 provides an implementation of this approach in which the perturbations are drawn uniformly from a hypersphere with radius  $\delta$ . Example 11.2 demonstrates this approach with a simple function.

## 11.3 Likelihood Ratio

The *likelihood ratio* approach<sup>8</sup> to gradient estimation uses an analytical form of  $\nabla\pi_\theta$  to improve our estimate of  $\nabla U(\theta)$ . Recall from equation (10.2) that

$$U(\theta) = \int p_\theta(\tau) R(\tau) d\tau \quad (11.9)$$

<sup>3</sup>Linear regression is covered in section 8.6.

<sup>4</sup>This general approach is sometimes referred to as *simultaneous perturbation stochastic approximation* by J. C. Spall, *Introduction to Stochastic Search and Optimization*. Wiley, 2003. The general connection to linear regression is provided by J. Peters and S. Schaal, “Reinforcement Learning of Motor Skills with Policy Gradients,” *Neural Networks*, vol. 21, no. 4, pp. 682–697, 2008.

<sup>5</sup>A recommended rule of thumb is to use about twice as many perturbations as the number of parameters.

<sup>6</sup>This equation shows the *forward difference*. Other finite-difference formulations, such as the central difference, can also be used.

<sup>7</sup>As discussed in section 8.6,  $\mathbf{X}^+$  denotes the pseudoinverse of  $\mathbf{X}$ .

<sup>8</sup>P. W. Glynn, “Likelihood Ratio Gradient Estimation for Stochastic Systems,” *Communications of the ACM*, vol. 33, no. 10, pp. 75–84, 1990.

```

struct RegressionGradient
   $\mathcal{P}$  # problem
   $\mathbf{b}$  # initial state distribution
   $d$  # depth
   $m$  # number of samples
   $\delta$  # step size
end

function gradient(M::RegressionGradient,  $\pi$ ,  $\theta$ )
   $\mathcal{P}$ ,  $\mathbf{b}$ ,  $d$ ,  $m$ ,  $\delta$ ,  $\gamma$  = M. $\mathcal{P}$ , M. $\mathbf{b}$ , M. $d$ , M. $m$ , M. $\delta$ , M. $\gamma$ 
   $\Delta\theta$  = [ $\delta$ .*normalize(randn(length( $\theta$ )), 2) for  $i$  = 1: $m$ ]
   $R(\tau)$  = sum( $r$ * $\gamma^{k-1}$  for ( $k$ , ( $s,a,\tau$ )) in enumerate( $\tau$ ))
   $U(\theta)$  = R(simulate( $\mathcal{P}$ , rand( $\mathbf{b}$ ),  $s \rightarrow \pi(\theta,s)$ ,  $d$ ))
   $\Delta U$  = [ $U(\theta + \Delta\theta)$  -  $U(\theta)$  for  $\Delta\theta$  in  $\Delta\theta$ ]
  return pinv(reduce(hcat,  $\Delta\theta$ )') *  $\Delta U$ 
end

```

Algorithm 11.3. A method for estimating a policy gradient using finite differences for an MDP  $\mathcal{P}$ , a stochastic parameterized policy  $\pi(\theta, s)$ , and a policy parameterization vector  $\theta$ . Policy variation vectors are generated by normalizing normally distributed samples and scaling by a perturbation scalar  $\delta$ . A total of  $m$  parameter perturbations are generated, and each is evaluated in a rollout from an initial state drawn from  $\mathbf{b}$  to depth  $d$  and compared to the original policy parameterization.

Hence,

$$\nabla U(\theta) = \nabla_{\theta} \int p_{\theta}(\tau) R(\tau) d\tau \quad (11.10)$$

$$= \int \nabla_{\theta} p_{\theta}(\tau) R(\tau) d\tau \quad (11.11)$$

$$= \int p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} R(\tau) d\tau \quad (11.12)$$

$$= \mathbb{E}_{\tau} \left[ \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} R(\tau) \right] \quad (11.13)$$

The name for this method comes from this trajectory likelihood ratio. This likelihood ratio can be seen as a weight in likelihood weighted sampling (section 3.7) over trajectory rewards.

Applying the log derivative trick,<sup>9</sup> we have

$$\nabla U(\theta) = \mathbb{E}_{\tau} [\nabla_{\theta} \log p_{\theta}(\tau) R(\tau)] \quad (11.14)$$

We can estimate this expectation using trajectory rollouts. For each trajectory  $\tau$ , we need to compute the product  $\nabla_{\theta} \log p_{\theta}(\tau) R(\tau)$ . Recall that  $R(\tau)$  is the return associated with trajectory  $\tau$ . If we have a stochastic policy,<sup>10</sup> the gradient  $\nabla_{\theta} \log p_{\theta}(\tau)$  is

$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{k=1}^d \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \quad (11.15)$$

<sup>9</sup> The log derivative trick was introduced in section 10.5. It uses the following equality:

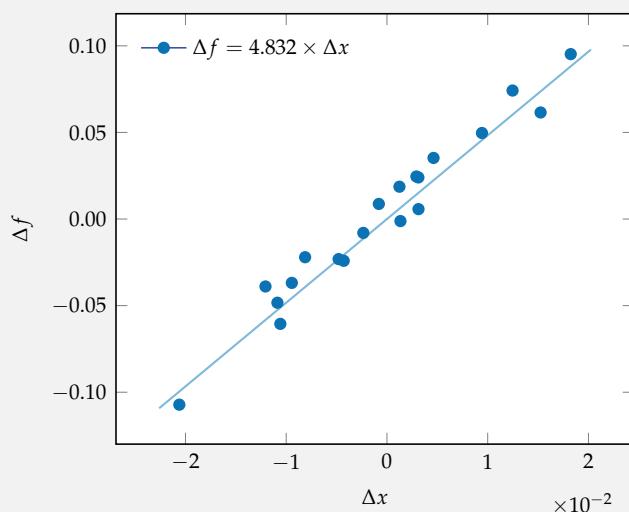
$$\nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} p_{\theta}(\tau) / p_{\theta}(\tau)$$

<sup>10</sup> We use  $\pi_{\theta}(a | s)$  to represent the probability (either density or mass) the policy  $\pi_{\theta}$  assigns to taking action  $a$  from state  $s$ .

We would like to apply the regression gradient to estimate the gradient of a simple, one-dimensional function  $f(x) = x^2$ , evaluated at  $x_0 = 2$  from  $m = 20$  samples. To imitate the stochasticity inherent in policy evaluation, we add noise to the function evaluations. We generate a set of disturbances  $\Delta X$ , sampled from  $\mathcal{N}(0, \delta^2)$ , and evaluate  $f(x_0 + \Delta x) - f(x_0)$  for each disturbance  $\Delta x$  in  $\Delta X$ . We can then estimate the one-dimensional gradient (or derivative)  $\Delta X^+ \Delta F$  with this code:

```
f(x) = x^2 + 1e-2*randn()
m = 20
delta = 1e-2
DeltaX = [delta.*randn() for i = 1:m]
x0 = 2.0
DeltaF = [f(x0 + DeltaX) - f(x0) for DeltaX in DeltaX]
pinv(DeltaX) * DeltaF
```

The samples and linear regression are shown here. The slope of the regression line is close to the exact solution of 4:



Example 11.2. Using the regression gradient method on a one-dimensional function.

because  $p_{\theta}(\tau)$  takes the form

$$p_{\theta}(\tau) = p(s^{(1)}) \prod_{k=1}^d T(s^{(k+1)} | s^{(k)}, a^{(k)}) \pi_{\theta}(a^{(k)} | s^{(k)}) \quad (11.16)$$

where  $s^{(k)}$  and  $a^{(k)}$  are the  $k$ th state and action, respectively, in trajectory  $\tau$ . Algorithm 11.4 provides an implementation in which  $m$  trajectories are sampled to arrive at a gradient estimate. Example 11.3 illustrates the process.

If we have a deterministic policy, the gradient requires computing:<sup>11</sup>

$$\nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} \log \left[ p(s^{(1)}) \prod_{k=1}^d T(s^{(k+1)} | s^{(k)}, \pi_{\theta}(s^{(k)})) \right] \quad (11.17)$$

$$= \sum_{k=1}^d \nabla_{\theta} \pi_{\theta}(s^{(k)}) \frac{\partial}{\partial a^{(k)}} \log T(s^{(k+1)} | s^{(k)}, a^{(k)}) \quad (11.18)$$

Equations (11.17) and (11.18) require knowing the transition likelihood, which is in contrast with equation (11.15) for stochastic policies.

```

struct LikelihoodRatioGradient
   $\mathcal{P}$  # problem
  b # initial state distribution
  d # depth
  m # number of samples
   $\nabla \log \pi$  # gradient of log likelihood
end

function gradient(M::LikelihoodRatioGradient,  $\pi$ ,  $\theta$ )
   $\mathcal{P}$ , b, d, m,  $\nabla \log \pi$ ,  $\gamma$  = M. $\mathcal{P}$ , M.b, M.d, M.m, M. $\nabla \log \pi$ , M. $\mathcal{P}$ . $\gamma$ 
   $\pi_{\theta}(s)$  =  $\pi(\theta, s)$ 
   $R(\tau)$  = sum( $r * \gamma^{k-1}$  for (k, (s, a, r)) in enumerate( $\tau$ ))
   $\nabla U(\tau)$  = sum( $\nabla \log \pi(\theta, a, s)$  for (s, a) in  $\tau$ ) *  $R(\tau)$ 
  return mean( $\nabla U(\text{simulate}(\mathcal{P}, \text{rand}(\mathbf{b}), \pi_{\theta}, \mathbf{d}))$  for i in 1:m)
end

```

<sup>11</sup> Many problems have vector-valued actions  $\mathbf{a} \in \mathbb{R}^n$ . In this case,  $\nabla_{\theta} \pi_{\theta}(s^{(k)})$  is replaced with a Jacobian matrix whose  $j$ th column is the gradient with respect to the  $j$ th action component, and the  $\frac{\partial}{\partial a^{(k)}}$   $\log T(s^{(k+1)} | s^{(k)}, a^{(k)})$  is replaced with an action gradient.

Algorithm 11.4. A method for estimating a policy gradient of a policy  $\pi(s)$  for an MDP  $\mathcal{P}$  with initial state distribution  $\mathbf{b}$  using the likelihood ratio trick. The gradient with respect to the parameterization vector  $\theta$  is estimated from  $m$  rollouts to depth  $d$  using the log policy gradients  $\nabla \log \pi$ .

## 11.4 Reward-to-Go

The likelihood ratio policy gradient method is unbiased but has high variance. Example 11.4 reviews bias and variance. The variance generally increases significantly with rollout depth due to the correlation between actions, states, and

Consider the single-step, single-state problem from example 11.1. Suppose we have a stochastic policy  $\pi_{\theta}$  that samples its action according to a Gaussian distribution  $\mathcal{N}(\theta_1, \theta_2^2)$ , where  $\theta_2^2$  is the variance.

$$\begin{aligned}\log \pi_{\theta}(a | s) &= \log \left( \frac{1}{\sqrt{2\pi\theta_2^2}} \exp \left( -\frac{(a - \theta_1)^2}{2\theta_2^2} \right) \right) \\ &= -\frac{(a - \theta_1)^2}{2\theta_2^2} - \frac{1}{2} \log(2\pi\theta_2^2)\end{aligned}$$

The gradient of the log policy likelihood is

$$\begin{aligned}\frac{\partial}{\partial \theta_1} \log \pi_{\theta}(a | s) &= \frac{a - \theta_1}{\theta_2^2} \\ \frac{\partial}{\partial \theta_2} \log \pi_{\theta}(a | s) &= \frac{(a - \theta_1)^2 - \theta_2^2}{\theta_2^3}\end{aligned}$$

Suppose we run three rollouts with  $\theta = [0, 1]$ , taking actions  $\{0.5, -1, 0.7\}$  and receiving the same rewards ( $R(s, a) = a$ ). The estimated policy gradient is

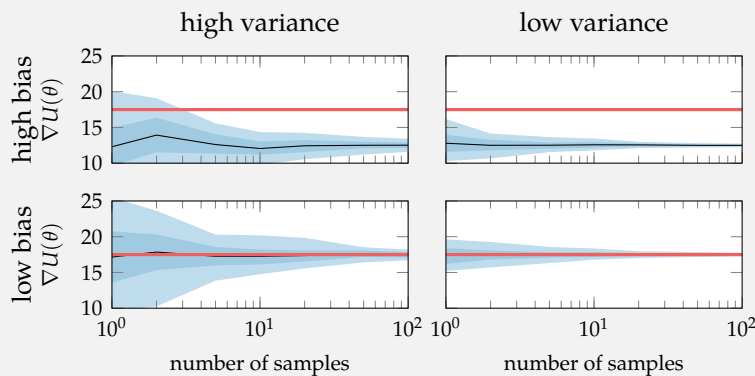
$$\begin{aligned}\nabla U(\theta) &\approx \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log p_{\theta}(\tau^{(i)}) R(\tau^{(i)}) \\ &= \frac{1}{3} \left( \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix} 0.5 + \begin{bmatrix} -1.0 \\ 0.0 \end{bmatrix} (-1) + \begin{bmatrix} 0.7 \\ -0.51 \end{bmatrix} 0.7 \right) \\ &= [0.58, -0.244]\end{aligned}$$

Example 11.3. Applying the likelihood ratio trick to estimate a policy gradient in a simple problem.



When estimating a quantity of interest from a collection of simulations, we generally want to use a scheme that has both low *bias* and low *variance*. In this chapter, we want to estimate  $\nabla U(\theta)$ . Generally, with more simulation samples, we can arrive at a better estimate. Some methods can lead to bias, where—even with infinitely many samples—it does not lead to an accurate estimate. Sometimes methods with nonzero bias may still be attractive if they also have low variance, meaning that they require fewer samples to converge.

Here are plots of the estimates from four notional methods for estimating  $\nabla U(\theta)$ . The true value is 17.5, as indicated by the red lines. We ran 100 simulations 100 times for each method. The variance decreases as the number of samples increases. The blue regions indicate the 5% to 95% and 25% to 75% empirical quantiles of the estimates.



Example 11.4. An empirical demonstration of bias and variance when estimating  $\nabla U(\theta)$ .

rewards across time steps. The *reward-to-go* approach attempts to reduce the variance in the estimate.

To derive this approach, we begin by expanding equation (11.14):

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \left( \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \right) \left( \sum_{k=1}^d r^{(k)} \gamma^{k-1} \right) \right] \quad (11.19)$$

Let  $f^{(k)}$  replace  $\nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)})$  for convenience. We then expand as follows:

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \left( \sum_{k=1}^d f^{(k)} \right) \left( \sum_{k=1}^d r^{(k)} \gamma^{k-1} \right) \right] \quad (11.20)$$

$$= \mathbb{E}_\tau \left[ \left( f^{(1)} + f^{(2)} + f^{(3)} + \dots + f^{(d)} \right) \left( r^{(1)} + r^{(2)}\gamma + r^{(3)}\gamma^2 + \dots + r^{(d)}\gamma^{d-1} \right) \right] \quad (11.21)$$

$$= \mathbb{E}_\tau \begin{bmatrix} f^{(1)}r^{(1)} + f^{(1)}r^{(2)}\gamma + f^{(1)}r^{(3)}\gamma^2 + \dots + f^{(1)}r^{(d)}\gamma^{d-1} \\ + f^{(2)}r^{(1)} + f^{(2)}r^{(2)}\gamma + f^{(2)}r^{(3)}\gamma^2 + \dots + f^{(2)}r^{(d)}\gamma^{d-1} \\ + f^{(3)}r^{(1)} + f^{(3)}r^{(2)}\gamma + f^{(3)}r^{(3)}\gamma^2 + \dots + f^{(3)}r^{(d)}\gamma^{d-1} \\ \vdots \\ + f^{(d)}r^{(1)} + f^{(d)}r^{(2)}\gamma + f^{(d)}r^{(3)}\gamma^2 + \dots + f^{(d)}r^{(d)}\gamma^{d-1} \end{bmatrix} \quad (11.22)$$

The first reward,  $r^{(1)}$ , is affected only by the first action. Thus, its contribution to the policy gradient should not depend on subsequent time steps. We can remove other such causality-violating terms as follows:<sup>12</sup>

<sup>12</sup> The term  $\sum_{\ell=k}^d r^{(\ell)} \gamma^{\ell-k}$  is often called the *reward-to-go* from step  $k$ .

$$\nabla U(\theta) = \mathbb{E}_\tau \begin{bmatrix} f^{(1)}r^{(1)} + f^{(1)}r^{(2)}\gamma + f^{(1)}r^{(3)}\gamma^2 + \dots + f^{(1)}r^{(d)}\gamma^{d-1} \\ + f^{(2)}r^{(2)}\gamma + f^{(2)}r^{(3)}\gamma^2 + \dots + f^{(2)}r^{(d)}\gamma^{d-1} \\ + f^{(3)}r^{(3)}\gamma^2 + \dots + f^{(3)}r^{(d)}\gamma^{d-1} \\ \vdots \\ + f^{(d)}r^{(d)}\gamma^{d-1} \end{bmatrix} \quad (11.23)$$

$$= \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \left( \sum_{\ell=k}^d r^{(\ell)} \gamma^{\ell-1} \right) \right] \quad (11.24)$$

$$= \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \left( \gamma^{k-1} \sum_{\ell=k}^d r^{(\ell)} \gamma^{\ell-k} \right) \right] \quad (11.25)$$

$$= \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{to-go}}^{(k)} \right] \quad (11.26)$$

Algorithm 11.5 provides an implementation of this.

Notice that the reward-to-go for a state-action pair  $(s, a)$  under a policy parameterized by  $\theta$  is really an approximation of the state-action value from that state,  $Q_\theta(s, a)$ . The action value function, if known, can be used to obtain the policy gradient:

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} Q_\theta(s^{(k)}, a^{(k)}) \right] \quad (11.27)$$

```

struct RewardToGoGradient
   $\mathcal{P}$  # problem
  b # initial state distribution
  d # depth
  m # number of samples
   $\nabla \log \pi$  # gradient of log likelihood
end

function gradient(M::RewardToGoGradient,  $\pi$ ,  $\theta$ )
   $\mathcal{P}$ , b, d, m,  $\nabla \log \pi$ ,  $\gamma$  = M. $\mathcal{P}$ , M.b, M.d, M.m, M. $\nabla \log \pi$ , M. $\mathcal{P}$ . $\gamma$ 
   $\pi_\theta(s)$  =  $\pi(\theta, s)$ 
   $R(\tau, j)$  = sum( $r * \gamma^{k-1}$  for (k, (s, a, r)) in zip(j:d,  $\tau$ [j:end]))
   $\nabla U(\tau)$  = sum( $\nabla \log \pi(\theta, a, s) * R(\tau, j)$  for (j, (s, a, r)) in enumerate( $\tau$ ))
  return mean( $\nabla U(\text{simulate}(\mathcal{P}, \text{rand}(\mathbf{b}), \pi_\theta, \mathbf{d}))$  for i in 1:m)
end

```

Algorithm 11.5. A method that uses reward-to-go for estimating a policy gradient of a policy  $\pi(s)$  for an MDP  $\mathcal{P}$  with initial state distribution **b**. The gradient with respect to the parameterization vector  $\theta$  is estimated from **m** rollouts to depth **d** using the log policy gradient  $\nabla \log \pi$ .

## 11.5 Baseline Subtraction

We can further build on the approach presented in the previous section by subtracting a *baseline* value from the reward-to-go<sup>13</sup> to reduce the variance of the gradient estimate. This subtraction does not bias the gradient.

We now subtract a baseline  $r_{\text{base}}(s^{(k)})$ :

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} (r_{\text{to-go}}^{(k)} - r_{\text{base}}(s^{(k)})) \right] \quad (11.28)$$

To show that baseline subtraction does not bias the gradient, we first expand:

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{to-go}}^{(k)} - \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{base}}(s^{(k)}) \right] \quad (11.29)$$

The *linearity of expectation* states that  $\mathbb{E}[a + b] = \mathbb{E}[a] + \mathbb{E}[b]$ , so it is sufficient to prove that equation (11.29) is equivalent to equation (11.26), if for each step  $k$ , the expected associated baseline term is  $\mathbf{0}$ :

$$\mathbb{E}_{\tau} \left[ \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{base}}(s^{(k)}) \right] = \mathbf{0} \quad (11.30)$$

We begin by converting the expectation into nested expectations, as illustrated in figure 11.2:

$$\mathbb{E}_{\tau} \left[ \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{base}}(s^{(k)}) \right] = \mathbb{E}_{\tau_{1:k}} \left[ \mathbb{E}_{\tau_{k+1:d}} \left[ \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{base}}(s^{(k)}) \right] \right] \quad (11.31)$$

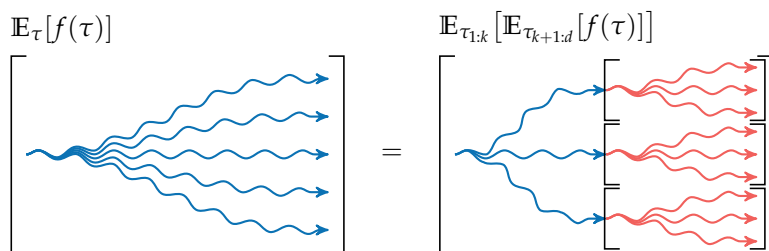


Figure 11.2. The expectation of a function of trajectories sampled from a policy can be viewed as an expectation over a nested expectation of subtrajectories. For a mathematical derivation, see exercise 11.4.

We continue with our derivation, using the same log derivative trick from section 11.3:

$$\begin{aligned} \mathbb{E}_{\tau_{1:k}} \left[ \mathbb{E}_{\tau_{k+1:d}} \left[ \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \gamma^{k-1} r_{\text{base}}(s^{(k)}) \right] \right] \\ = \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \mathbb{E}_{\tau_{k+1:d}} \left[ \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \right] \right] \end{aligned} \quad (11.32)$$

$$= \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \mathbb{E}_{a^{(k)}} \left[ \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \right] \right] \quad (11.33)$$

$$= \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \int \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \pi_{\theta}(a^{(k)} | s^{(k)}) \mathbf{d}a^{(k)} \right] \quad (11.34)$$

$$= \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \int \frac{\nabla_{\theta} \pi_{\theta}(a^{(k)} | s^{(k)})}{\pi_{\theta}(a^{(k)} | s^{(k)})} \pi_{\theta}(a^{(k)} | s^{(k)}) \mathbf{d}a^{(k)} \right] \quad (11.35)$$

$$= \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \nabla_{\theta} \int \pi_{\theta}(a^{(k)} | s^{(k)}) \mathbf{d}a^{(k)} \right] \quad (11.36)$$

$$= \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \nabla_{\theta} \mathbf{1} \right] \quad (11.37)$$

$$= \mathbb{E}_{\tau_{1:k}} \left[ \gamma^{k-1} r_{\text{base}}(s^{(k)}) \mathbf{0} \right] \quad (11.38)$$

Therefore, subtracting a term  $r_{\text{base}}(s^{(k)})$  does not bias the estimate. This derivation assumed continuous state and action spaces. The same result applies to discrete spaces.

We can choose a different  $r_{\text{base}}(s)$  for every component of the gradient, and we will select them to minimize the variance. For simplicity, we will drop the dependence on  $s$  and treat each baseline component as constant.<sup>14</sup> For compactness in writing the equations in our derivation, we define

$$\ell_i(a, s, k) = \gamma^{k-1} \frac{\partial}{\partial \theta_i} \log \pi_{\theta}(a | s) \quad (11.39)$$

The variance of the  $i$ th component of our gradient estimate in equation (11.28) is

$$\mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ (\ell_i(a, s, k) (r_{\text{to-go}} - r_{\text{base},i}))^2 \right] - \mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ \ell_i(a, s, k) (r_{\text{to-go}} - r_{\text{base},i}) \right]^2 \quad (11.40)$$

where the expectation is over the  $(a, s, r_{\text{to-go}})$  tuples in our trajectory samples, and  $k$  is each tuple's depth.

We have just shown that the second term is zero. Hence, we can focus on choosing  $r_{\text{base},i}$  to minimize the first term by taking the derivative with respect to the baseline and setting it to zero:

$$\begin{aligned} & \frac{\partial}{\partial r_{\text{base},i}} \mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ (\ell_i(a, s, k) (r_{\text{to-go}} - r_{\text{base},i}))^2 \right] \\ &= \frac{\partial}{\partial r_{\text{base},i}} \left( \mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ \ell_i(a, s, k)^2 r_{\text{to-go}}^2 \right] - 2 \mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ \ell_i(a, s, k)^2 r_{\text{to-go}} r_{\text{base},i} \right] + r_{\text{base},i}^2 \mathbb{E}_{a,s,k} \left[ \ell_i(a, s, k)^2 \right] \right) \end{aligned} \quad (11.41)$$

$$= -2 \mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ \ell_i(a, s, k)^2 r_{\text{to-go}} \right] + 2 r_{\text{base},i} \mathbb{E}_{a,s,k} \left[ \ell_i(a, s, k)^2 \right] = 0 \quad (11.42)$$

Solving for  $r_{\text{base},i}$  yields the baseline component that minimizes the variance:

$$r_{\text{base},i} = \frac{\mathbb{E}_{a,s,r_{\text{to-go}},k} \left[ \ell_i(a, s, k)^2 r_{\text{to-go}} \right]}{\mathbb{E}_{a,s,k} \left[ \ell_i(a, s, k)^2 \right]} \quad (11.43)$$

<sup>14</sup>Some methods approximate a state-dependent baseline using  $r_{\text{base}}(s^{(k)}) = \Phi(s^{(k)})^\top \mathbf{w}$ . Selecting appropriate baseline functions tends to be difficult. J. Peters and S. Schaal, "Reinforcement Learning of Motor Skills with Policy Gradients," *Neural Networks*, vol. 21, no. 4, pp. 682–697, 2008.

```

struct BaselineSubtractionGradient
   $\mathcal{P}$  # problem
   $\mathbf{b}$  # initial state distribution
   $d$  # depth
   $m$  # number of samples
   $\nabla \log \pi$  # gradient of log likelihood
end

function gradient(M::BaselineSubtractionGradient,  $\pi$ ,  $\theta$ )
   $\mathcal{P}$ ,  $\mathbf{b}$ ,  $d$ ,  $m$ ,  $\nabla \log \pi$ ,  $\gamma$  = M. $\mathcal{P}$ , M. $\mathbf{b}$ , M. $d$ , M. $m$ , M. $\nabla \log \pi$ , M. $\mathcal{P}$ . $\gamma$ 
   $\pi\theta(\mathbf{s}) = \pi(\theta, \mathbf{s})$ 
   $\ell(\mathbf{a}, \mathbf{s}, k) = \nabla \log \pi(\theta, \mathbf{a}, \mathbf{s}) * \gamma^{k-1}$ 
   $R(\tau, k) = \text{sum}(r * \gamma^{j-1} \text{ for } (j, (\mathbf{s}, \mathbf{a}, r)) \text{ in enumerate}(\tau[k:\text{end}])))$ 
   $\text{numer}(\tau) = \text{sum}(\ell(\mathbf{a}, \mathbf{s}, k) .^2 * R(\tau, k) \text{ for } (k, (\mathbf{s}, \mathbf{a}, r)) \text{ in enumerate}(\tau))$ 
   $\text{denom}(\tau) = \text{sum}(\ell(\mathbf{a}, \mathbf{s}, k) .^2 \text{ for } (k, (\mathbf{s}, \mathbf{a})) \text{ in enumerate}(\tau))$ 
   $\text{base}(\tau) = \text{numer}(\tau) ./ \text{denom}(\tau)$ 
   $\text{trajs} = [\text{simulate}(\mathcal{P}, \text{rand}(\mathbf{b}), \pi\theta, d) \text{ for } i \text{ in } 1:m]$ 
   $\text{rbase} = \text{mean}(\text{base}(\tau) \text{ for } \tau \text{ in trajs})$ 
   $\nabla U(\tau) = \text{sum}(\ell(\mathbf{a}, \mathbf{s}, k) .* (R(\tau, k) - \text{rbase}) \text{ for } (k, (\mathbf{s}, \mathbf{a}, r)) \text{ in enumerate}(\tau))$ 
  return  $\text{mean}(\nabla U(\tau) \text{ for } \tau \text{ in trajs})$ 
end

```

Algorithm 11.6. Likelihood ratio gradient estimation with reward-to-go and baseline subtraction for an MDP  $\mathcal{P}$ , policy  $\pi$ , and initial state distribution  $\mathbf{b}$ . The gradient with respect to the parameterization vector  $\theta$  is estimated from  $m$  rollouts to depth  $d$  using the log policy gradients  $\nabla \log \pi$ .

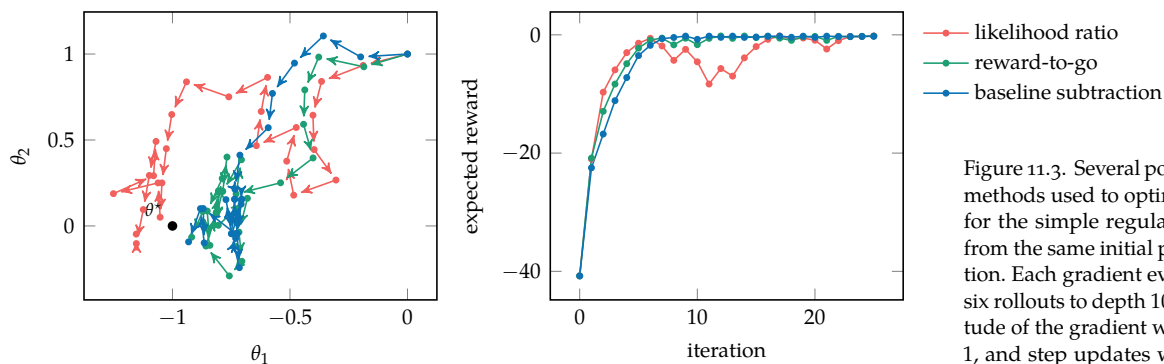


Figure 11.3. Several policy gradient methods used to optimize policies for the simple regulator problem from the same initial parameterization. Each gradient evaluation ran six rollouts to depth 10. The magnitude of the gradient was limited to 1, and step updates were applied with step size 0.2. The optimal policy parameterization is shown in black.

It is common to use likelihood ratio policy gradient estimation with this baseline subtraction (algorithm 11.6).<sup>15</sup> Figure 11.3 compares the methods discussed here.

Qualitatively, when considering the gradient contribution of state-action pairs, what we really care about is the relative value of one action over another. If all actions in a particular state produce the same high value, there is no real signal in the gradient, and baseline subtraction can zero that out. We want to identify the actions that produce a higher value than others, regardless of the mean value across actions.

An alternative to the action value is the *advantage*,  $A(s, a) = Q(s, a) - U(s)$ . Using the state value function in baseline subtraction produces the advantage. The policy gradient using the advantage is unbiased and typically has much lower variance. The gradient computation takes the following form:

$$\nabla U(\theta) = \mathbb{E}_\tau \left[ \sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} A_\theta(s^{(k)}, a^{(k)}) \right] \quad (11.44)$$

As with the state and action value functions, the advantage function is typically unknown. Other methods, covered in chapter 13, are needed to approximate it.

## 11.6 Summary

- A gradient can be estimated using finite differences.
- Linear regression can also be used to provide more robust estimates of the policy gradient.
- The likelihood ratio can be used to derive a form of the policy gradient that does not depend on the transition model for stochastic policies.
- The variance of the policy gradient can be significantly reduced using the reward-to-go and baseline subtraction.

<sup>15</sup>This combination is used in the class of algorithms called *REINFORCE* as introduced by R. J. Williams, “Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning,” *Machine Learning*, vol. 8, pp. 229–256, 1992.

## 11.7 Exercises

**Exercise 11.1.** If we estimate the expected discounted return of a given parameterized policy  $\pi_{\theta}$  defined by an  $n$ -dimensional vector of parameters  $\theta$  using  $m$  rollouts, how many total rollouts do we need to perform to compute the policy gradient using a finite difference approach?

*Solution:* In order to estimate the policy gradient using a finite difference approach, we need to estimate the utility of the policy given the current parameter vector  $U(\theta)$ , as well as all  $n$  variations of the current parameter vector  $U(\theta + \delta e^{(i)})$  for  $i = 1 : n$ . Since we estimate each of these using  $m$  rollouts, we need to perform a total of  $m(n + 1)$  rollouts.

**Exercise 11.2.** Suppose we have a robotic arm with which we are able to run experiments manipulating a wide variety of objects. We would like to use the likelihood ratio policy gradient or one of its extensions to train a policy that is efficient at picking up and moving these objects. Would it be more straightforward to use a deterministic or a stochastic policy, and why?

*Solution:* The likelihood ratio policy gradient requires an explicit representation of the transition likelihood when used with deterministic policies. Specifying an accurate explicit transition model for a real-world robotic arm manipulation task would be challenging. Computing the policy gradient for a stochastic policy does not require having an explicit representation of the transition likelihood, making the use of a stochastic policy more straightforward.

**Exercise 11.3.** Consider policy gradients of the form

$$\nabla_{\theta} U(\theta) = \mathbb{E}_{\tau} \left[ \sum_{k=1}^d \gamma^{k-1} y \nabla_{\theta} \log \pi_{\theta}(a^{(k)} | s^{(k)}) \right]$$

Which of the following values of  $y$  result in a valid policy gradient? Explain why.

- (a)  $\gamma^{1-k} \sum_{\ell=1}^{\infty} r^{(\ell)} \gamma^{\ell-1}$
- (b)  $\sum_{\ell=k}^{\infty} r^{(\ell)} \gamma^{\ell-k}$
- (c)  $\left( \sum_{\ell=k}^{\infty} r^{(\ell)} \gamma^{\ell-k} \right) - r_{\text{base}}(s^{(k)})$
- (d)  $U(s^{(k)})$
- (e)  $Q(s^{(k)}, a^{(k)})$
- (f)  $A(s^{(k)}, a^{(k)})$
- (g)  $r^{(k)} + \gamma U(s^{(k+1)}) - U(s^{(k)})$



*Solution:*

- (a)  $\sum_{\ell=1}^{\infty} r^{(\ell)}$  results in the total discounted reward, as

$$\gamma^{k-1} \gamma^{1-k} \sum_{\ell=1}^{\infty} r^{(\ell)} \gamma^{\ell-1} = \sum_{\ell=1}^{\infty} r^{(\ell)} \gamma^{\ell-1}$$

and produces a valid policy gradient, as given in equation (11.19).

- (b)  $\sum_{\ell=k}^{\infty} r^{(\ell)} \gamma^{\ell-k}$  is the reward-to-go and produces a valid policy gradient, as given in equation (11.26).
- (c)  $\left(\sum_{\ell=k}^{\infty} r^{(\ell)}\right) - r_{\text{base}}(s^{(k)})$  is the baseline subtracted reward-to-go and produces a valid policy gradient, as given in equation (11.28).
- (d)  $U(s^{(k)})$  is the state value function and does not produce a valid policy gradient.
- (e)  $Q(s^{(k)}, a^{(k)})$  is the state-action value function and produces a valid policy gradient, as given in equation (11.27).
- (f)  $A(s^{(k)}, a^{(k)})$  is the advantage function and produces a valid policy gradient, as given in equation (11.44).
- (g)  $r^{(k)} + \gamma U(s^{(k+1)}) - U(s^{(k)})$  is the temporal difference residual (to be discussed further in chapter 13) and produces a valid policy gradient because it is an unbiased approximation of the advantage function.

**Exercise 11.4.** Show that  $\mathbb{E}_{\tau \sim \pi}[f(\tau)] = \mathbb{E}_{\tau_{1:k} \sim \pi}[\mathbb{E}_{\tau_{k:d} \sim \pi}[f(\tau)]]$  for step  $k$ .

*Solution:* The nested expectations can be proven by writing the expectation in integral form and then converting back:

$$\begin{aligned} \mathbb{E}_{\tau \sim \pi}[f(\tau)] &= \\ &= \int p(\tau) f(\tau) \, d\tau \\ &= \int \left( p(s^{(1)}) \prod_{k=1}^d p(s^{(k+1)} \mid s^{(k)}, a^{(k)}) \pi(a^{(k)} \mid s^{(k)}) \right) f(\tau) \, d\tau \\ &= \int \int \int \int \cdots \int \left( p(s^{(1)}) \prod_{k=1}^d p(s^{(k+1)} \mid s^{(k)}, a^{(k)}) \pi(a^{(k)} \mid s^{(k)}) \right) f(\tau) \, ds^{(d)} \cdots da^{(2)} \, ds^{(2)} \, da^{(1)} \, ds^{(1)} \\ &= \mathbb{E}_{\tau_{1:k} \sim \pi} \left[ \int \int \int \int \cdots \int \left( \prod_{q=k}^d p(s^{(q+1)} \mid s^{(q)}, a^{(q)}) \pi(a^{(q)} \mid s^{(q)}) \right) f(\tau) \, ds^{(d)} \cdots da^{(k+1)} \, ds^{(k+1)} \, da^{(k)} \, ds^{(k)} \right] \\ &= \mathbb{E}_{\tau_{1:k} \sim \pi}[\mathbb{E}_{\tau_{k:d} \sim \pi}[f(\tau)]] \end{aligned}$$

**Exercise 11.5.** Our implementation of the regression gradient (algorithm 11.3) fits a linear mapping from perturbations to the difference in returns,  $U(\theta + \Delta\theta^{(i)}) - U(\theta)$ . We evaluate  $U(\theta + \Delta\theta^{(i)})$  and  $U(\theta)$  for each of the  $m$  perturbations, thus reevaluating  $U(\theta)$  a total of  $m$  times. How might we reallocate the samples in a more effective manner?

*Solution:* One approach is to evaluate  $U(\theta)$  once and use the same value for each perturbation, thereby conducting only  $m + 1$  evaluations. Having an accurate estimate of  $U(\theta)$  is particularly important for an accurate regression gradient estimate. An alternative is to still compute  $U(\theta)$  once, but use  $m$  rollouts, thus preserving the total number of rollouts per iteration. This approach uses the same amount of computation as algorithm 11.3, but it can produce a more reliable gradient estimate.