

Policy Gradients

CS60077: Reinforcement Learning

Abir Das

IIT Kharagpur

Nov 09, 10, 2020

Resources

- § Deep Reinforcement Learning by Sergey Levine [[Link](#)]
- § OpenAI Spinning Up [[Link](#)]

Reinforcement Learning Setting

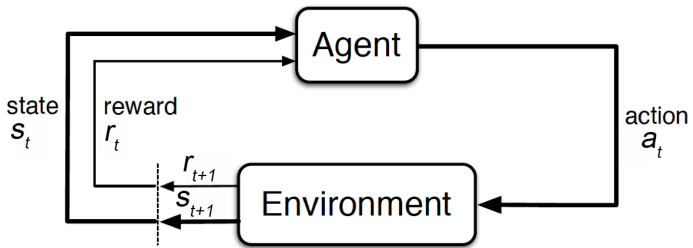


Figure credit: [SB]

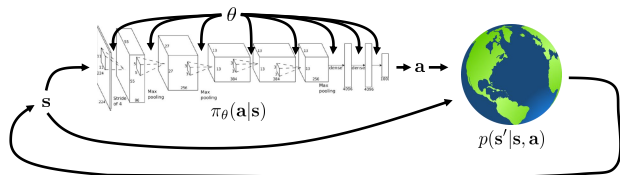


Figure credit: [Sergey Levine, UC Berkeley]

Reinforcement Learning Setting

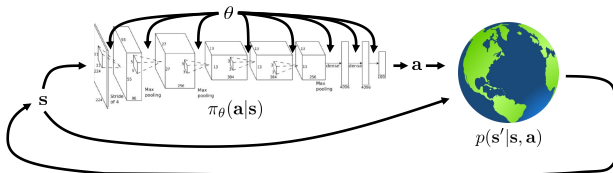


Figure credit: [Sergey Levine, UC Berkeley]

- § In the middle is the ‘policy network’ which can directly learn a parameterized policy $\pi_{\theta}(\mathbf{a}|s)$ (sometimes denoted as $\pi(\mathbf{a}|s; \theta)$) and provides the probability distribution over all actions given the state s and parameterized by θ .
- § To distinguish it from the parameter vector \mathbf{w} in value function approximator $\hat{v}(s; \mathbf{w})$, the notation θ is used.

Reinforcement Learning Setting

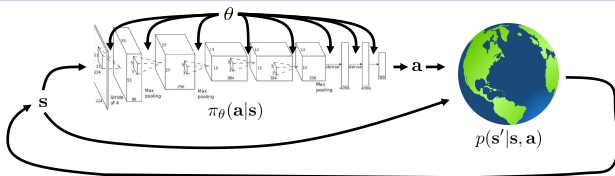


Figure credit: [Sergey Levine, UC Berkeley]

- § Goal in RL Problem is to maximize the total reward “in expectation” over long run.
- § A trajectory τ is defined as,

$$\tau = (s_1, a_1, s_2, a_2, s_3, a_3, \dots)$$

- § The probability of a trajectory is given by the joint probability of the state-action pairs.

$$p_{\theta}(s_1, a_1, s_2, a_2, \dots, s_T, a_T, s_{T+1}) = p(s_1) \prod_{t=1}^T p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t) \quad (1)$$

Reinforcement Learning Setting

§ Proof of the above relation,

$$\begin{aligned}
 & p(s_{T+1}, s_T, a_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) \\
 = & p(s_{T+1} | s_T, a_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) p(s_T, a_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) \\
 = & p(s_{T+1} | s_T, a_T) p(s_T, a_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) \\
 = & p(s_{T+1} | s_T, a_T) p(a_T | s_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) p(s_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) \\
 = & p(s_{T+1} | s_T, a_T) \pi_{\theta}(a_T | s_T) \boxed{p(s_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1)} \tag{2}
 \end{aligned}$$

§ The boxed part of the equation is very similar to the left hand side. So, using similar argument repetitively, we get,

$$\begin{aligned}
 & p(s_{T+1}, s_T, a_T, s_{T-1}, a_{T-1}, \dots, s_1, a_1) \\
 = & p(s_{T+1} | s_T, a_T) \pi_{\theta}(a_T | s_T) p(s_T | s_{T-1}, a_{T-1}) \pi_{\theta}(a_{T-1} | s_{T-1}) \\
 & p(s_{T-1}, s_{T-2}, a_{T-2}, \dots, s_1, a_1) \\
 = & p(s_1) \prod_{t=1}^T p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t) \tag{3}
 \end{aligned}$$

The Goal of Reinforcement Learning

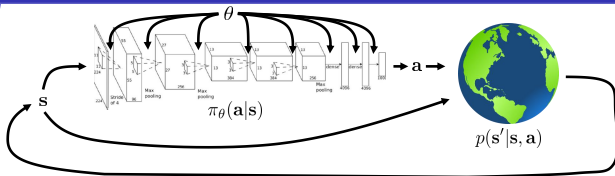


Figure credit: [Sergey Levine, UC Berkeley]

§ We will sometimes denote the probability as $p_{\theta}(\tau)$, i.e.,

$$p_{\theta}(\tau) = p_{\theta}(s_1, \mathbf{a}_1, s_2, \mathbf{a}_2, \dots, s_T, \mathbf{a}_T, s_{T+1}) = p(s_1) \prod_{t=1}^T p(s_{t+1} | s_t, \mathbf{a}_t) \pi_{\theta}(\mathbf{a}_t | s_t)$$

§ The goal can be written as,

$$\theta^* = \arg \max_{\theta} \underbrace{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, \mathbf{a}_t) \right]}_{J(\theta)}$$

§ Note that, for the time being, we are not considering discount. We will come back to that.

The Goal of Reinforcement Learning

§ Goal for a finite horizon setting:

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim p_{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)]$$

§ The same for the infinite horizon setting

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim p_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a})]$$

§ We will consider only finite horizon case in this topic.

Evaluating the Objective

- § We will see how we can optimize this objective - the expected value of the total reward under the trajectory distribution induced by the policy θ .
- § But before that let us see how we can evaluate the objective in model free setting.

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad (4)$$

Evaluating the Objective

- § We will see how we can optimize this objective - the expected value of the total reward under the trajectory distribution induced by the policy θ .
- § But before that let us see how we can evaluate the objective in model free setting.

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \quad (4)$$

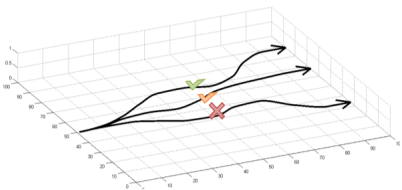


Figure credit: [Sergey Levine, UC Berkeley]

Maximizing the Objective

- § Now that we have seen how to evaluate the objective, the next step is to maximize it.

Maximizing the Objective

- § Now that we have seen how to evaluate the objective, the next step is to maximize it.
- § Compute the gradient and take steps in the direction of the gradient.

Maximizing the Objective

- § Now that we have seen how to evaluate the objective, the next step is to maximize it.
- § Compute the gradient and take steps in the direction of the gradient.

$$\theta^* = \arg \max_{\theta} \underbrace{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\overbrace{\sum_t r(\mathbf{s}_t, \mathbf{a}_t)}^{r(\tau)} \right]}_{J(\theta)}$$
$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau$$

Maximizing the Objective

- § Now that we have seen how to evaluate the objective, the next step is to maximize it.
- § Compute the gradient and take steps in the direction of the gradient.

$$\theta^* = \arg \max_{\theta} \underbrace{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\underbrace{\sum_t \overbrace{r(\mathbf{s}_t, \mathbf{a}_t)}^{r(\tau)}} \right]}_{J(\theta)}$$

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau$$

- § How to compute this complicated looking gradient!

Maximizing the Objective

- § Now that we have seen how to evaluate the objective, the next step is to maximize it.
- § Compute the gradient and take steps in the direction of the gradient.

$$\theta^* = \arg \max_{\theta} \underbrace{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\overbrace{\sum_t r(\mathbf{s}_t, \mathbf{a}_t)}^{r(\tau)} \right]}_{J(\theta)}$$

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau$$

- § How to compute this complicated looking gradient! [The log-derivative trick is our rescue.](#)

Log Derivative Trick

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(\tau) &= \frac{\partial \log p_{\theta}(\tau)}{\partial p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) = \frac{1}{p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) \\ \implies \nabla_{\theta} p_{\theta}(\tau) &= p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)\end{aligned}\tag{5}$$

Log Derivative Trick

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(\tau) &= \frac{\partial \log p_{\theta}(\tau)}{\partial p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) = \frac{1}{p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) \\ \implies \nabla_{\theta} p_{\theta}(\tau) &= p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)\end{aligned}\tag{5}$$

§ So using eqn. (5) we get the gradient of the objective as,

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]\end{aligned}\tag{6}$$

Log Derivative Trick

$$\begin{aligned}\nabla_{\theta} \log p_{\theta}(\tau) &= \frac{\partial \log p_{\theta}(\tau)}{\partial p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) = \frac{1}{p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) \\ \implies \nabla_{\theta} p_{\theta}(\tau) &= p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)\end{aligned}\tag{5}$$

§ So using eqn. (5) we get the gradient of the objective as,

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]\end{aligned}\tag{6}$$

§ Remember that

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau$$

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

§ We have also seen,

$$p_{\boldsymbol{\theta}}(\tau) = p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{s}_{T+1}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

§ We have also seen,

$$p_{\boldsymbol{\theta}}(\tau) = p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{s}_{T+1}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking log both sides,

$$\log p_{\boldsymbol{\theta}}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t=1}^T \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

§ We have also seen,

$$p_{\boldsymbol{\theta}}(\tau) = p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{s}_{T+1}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking log both sides,

$$\log p_{\boldsymbol{\theta}}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t=1}^T \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking $\nabla_{\boldsymbol{\theta}}$ both sides,

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

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

§ We have also seen,

$$p_{\boldsymbol{\theta}}(\tau) = p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{s}_{T+1}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking log both sides,

$$\log p_{\boldsymbol{\theta}}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t=1}^T \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking $\nabla_{\boldsymbol{\theta}}$ both sides,

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) = \cancel{\log p(\mathbf{s}_1)} \rightarrow 0$$

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

§ We have also seen,

$$p_{\boldsymbol{\theta}}(\tau) = p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{s}_{T+1}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking log both sides,

$$\log p_{\boldsymbol{\theta}}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t=1}^T \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking $\nabla_{\boldsymbol{\theta}}$ both sides,

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) = \cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \xrightarrow{0}$$

Log Derivative Trick

§ Till now we have the following,

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} J(\boldsymbol{\theta}); \quad J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [r(\tau)]$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)]$$

§ We have also seen,

$$p_{\boldsymbol{\theta}}(\tau) = p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{s}_{T+1}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking log both sides,

$$\log p_{\boldsymbol{\theta}}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t=1}^T \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

§ Taking $\nabla_{\boldsymbol{\theta}}$ both sides,

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) = \cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \cancel{\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} + \sum_{t=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t)$$

Log Derivative Trick

§ Thus,

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

Log Derivative Trick

§ Thus,

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\tau) r(\tau)] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

§ So, to get the estimate of the gradient we take samples and average not only the sum of rewards but also average the sum of the gradients of the policy values.

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$

Log Derivative Trick

§ Thus,

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

§ So, to get the estimate of the gradient we take samples and average not only the sum of rewards but also average the sum of the gradients of the policy values.

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$

§ And the last bit is to update θ along the gradient direction.

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta) \quad (7)$$

Fitting in Generic RL Pipeline

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$
$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

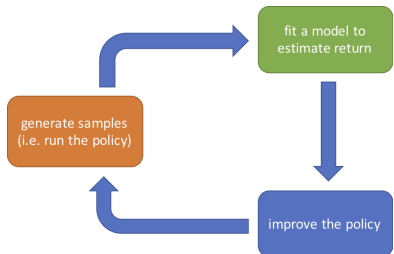


Figure credit: [Sergey Levine, UC Berkeley]

Fitting in Generic RL Pipeline

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

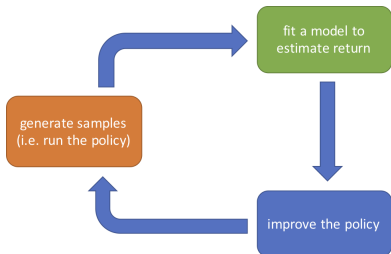


Figure credit: [Sergey Levine, UC Berkeley]

REINFORCE Algorithm

- 1 Sample $\{r^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
- 2 $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$
- 3 $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$
- 4 Repeat

Taking a Closer Look

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$

- § What is given by $\log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$? - It is log of the probability of action $\mathbf{a}_{i,t}$ at state $\mathbf{s}_{i,t}$ **under the distribution parameterized by θ** .
- § This gives the **likelihood, i.e., how likely**, we are to see $\mathbf{a}_{i,t}$ as the action, if our policy is defined by the current θ that we have.
- § Computing the gradient and taking a step along the direction of the gradient, changes θ in such a way that the likelihood of the action $\mathbf{a}_{i,t}$ increases.

Taking a Closer Look

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$

- § What is given by $\log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$? - It is log of the probability of action $\mathbf{a}_{i,t}$ at state $\mathbf{s}_{i,t}$ **under the distribution parameterized by θ** .
- § This gives the **likelihood, i.e., how likely**, we are to see $\mathbf{a}_{i,t}$ as the action, if our policy is defined by the current θ that we have.
- § Computing the gradient and taking a step along the direction of the gradient, changes θ in such a way that the likelihood of the action $\mathbf{a}_{i,t}$ increases.

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right]$$

- § Now consider the case, when it is getting multiplied by $\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$.
- § Those actions with high rewards are getting more likely.

Taking a Closer Look

- § Good stuff is made more likely.
- § Bad stuff is made less likely.
- § Formalizes the 'trial and error' learning.

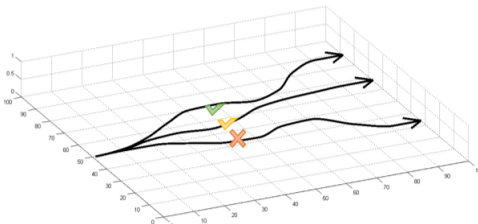


Figure credit: [Sergey Levine, UC Berkeley]

Taking a Closer Look

- § Good stuff is made more likely.
- § Bad stuff is made less likely.
- § Formalizes the 'trial and error' learning.

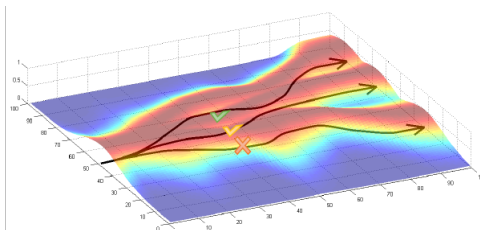


Figure credit: [Sergey Levine, UC Berkeley]

Taking a Closer Look

- § Good stuff is made more likely.
- § Bad stuff is made less likely.
- § Formalizes the 'trial and error' learning.

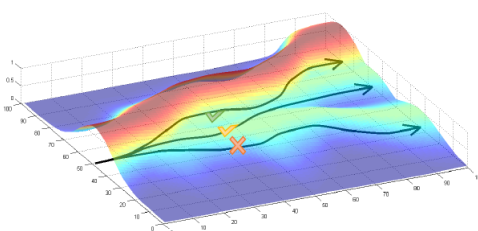


Figure credit: [Sergey Levine, UC Berkeley]

Taking a Closer Look

- § Good stuff is made more likely.
- § Bad stuff is made less likely.
- § Formalizes the 'trial and error' learning.

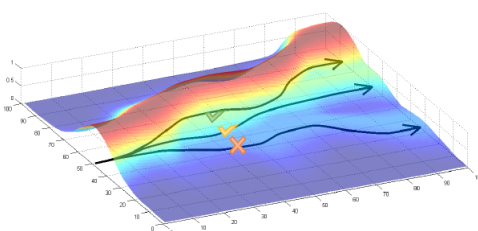


Figure credit: [Sergey Levine, UC Berkeley]

Bias and Variance in Estimation

- § One way to work with values we do not know is to estimate them by experimenting repeatedly.
- § Monte-Carlo methods provide the estimate of the true value and we have used Monte-Carlo methods to estimate the value functions and to estimate the gradient of the expected return.
- § The estimator is a function of the data which itself are random variables. So the estimated value is subject to many possible outcomes if employed repeatedly, *i.e.*, if you conduct the experiment multiple times, in general, the estimator will provide different values.
- § An estimator is good if,
 - ▶ On average the estimated values are close to the true value for different trials - (Bias)
 - ▶ The estimates do not vary much in each trial - (variance)

Unbiased Estimators

§ An unbiased estimator is the one that yields the true value of the variable being estimated on average. With θ denoting the true value and $\hat{\theta}$ denoting the estimated value, and unbiased estimator is one with,

$$\mathbb{E}[\hat{\theta}] = \theta$$

§ Naturally bias is defined as,

$$b = \mathbb{E}[\hat{\theta}] - \theta$$

Unbiased Estimators

§ An unbiased estimator is the one that yields the true value of the variable being estimated on average. With θ denoting the true value and $\hat{\theta}$ denoting the estimated value, and unbiased estimator is one with,

$$\mathbb{E}[\hat{\theta}] = \theta$$

§ Naturally bias is defined as,

$$b = \mathbb{E}[\hat{\theta}] - \theta$$

§ Let us consider estimating a constant value (say temperature of this room) by some sensors which are not perfect. Consider the observations.

$$x[n] = \theta + w[n] \quad n = 0, 1, \dots, N-1. \quad w[n] \text{ is WGN with variance} = \sigma^2.$$

Unbiased Estimators

§ An unbiased estimator is the one that yields the true value of the variable being estimated on average. With θ denoting the true value and $\hat{\theta}$ denoting the estimated value, and unbiased estimator is one with,

$$\mathbb{E}[\hat{\theta}] = \theta$$

§ Naturally bias is defined as,

$$b = \mathbb{E}[\hat{\theta}] - \theta$$

§ Let us consider estimating a constant value (say temperature of this room) by some sensors which are not perfect. Consider the observations.

$$x[n] = \theta + w[n] \quad n = 0, 1, \dots, N-1. \quad w[n] \text{ is WGN with variance} = \sigma^2.$$

§ A reasonable estimator is the average value of $x[n]$ *i.e.*,

$$\hat{\theta} = \frac{1}{N} \sum_{n=0}^{N-1} x[n]$$

Estimator Bias

§ The sample mean estimator is unbiased.

$$\begin{aligned}\mathbb{E}[\hat{\theta}] &= \mathbb{E}\left[\frac{1}{N} \sum_{n=0}^{N-1} x[n]\right] = \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}[x[n]] \\ &= \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}([\theta + w[n]]) = \frac{1}{N} \sum_{n=0}^{N-1} (\mathbb{E}[\theta] + \mathbb{E}[w[n]]) \\ &= \frac{1}{N} \sum_{n=0}^{N-1} = (\theta + 0) = \theta\end{aligned}$$

Estimator Bias

§ The sample mean estimator is unbiased.

$$\begin{aligned}\mathbb{E}[\hat{\theta}] &= \mathbb{E}\left[\frac{1}{N} \sum_{n=0}^{N-1} x[n]\right] = \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}[x[n]] \\ &= \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}([\theta + w[n]]) = \frac{1}{N} \sum_{n=0}^{N-1} (\mathbb{E}[\theta] + \mathbb{E}[w[n]]) \\ &= \frac{1}{N} \sum_{n=0}^{N-1} = (\theta + 0) = \theta\end{aligned}$$

§ Let us see what happens with a modified estimator, $x[n]$ *i.e.*,

$$\check{\theta} = \frac{1}{2N} \sum_{n=0}^{N-1} x[n]$$

Estimator Bias

§ The sample mean estimator is unbiased.

$$\begin{aligned}\mathbb{E}[\hat{\theta}] &= \mathbb{E}\left[\frac{1}{N} \sum_{n=0}^{N-1} x[n]\right] = \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}[x[n]] \\ &= \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}([\theta + w[n]]) = \frac{1}{N} \sum_{n=0}^{N-1} (\mathbb{E}[\theta] + \mathbb{E}[w[n]]) \\ &= \frac{1}{N} \sum_{n=0}^{N-1} = (\theta + 0) = \theta\end{aligned}$$

§ Let us see what happens with a modified estimator, $x[n]$ *i.e.*,

$$\check{\theta} = \frac{1}{2N} \sum_{n=0}^{N-1} x[n]$$

§ It is easy to see that $\mathbb{E}[\check{\theta}] = \frac{1}{2}\theta$.

§ So the bias is $b = \mathbb{E}[\check{\theta}] - \theta = -\frac{1}{2}\theta$

Estimator Variance

- § That an estimator is unbiased does not necessarily mean that it is a good estimator. It is reasonable to check by repeating the experiment how the results differ in successive trials.
- § Thus the variance of the estimate is another measure of goodness of the estimator. And the aim will be to see how small we can make $\text{var}(\hat{\theta})$.
- § Let us take the following 3 estimators for θ and see the variances of all these.

$$\begin{array}{lll}
 \hat{\theta}_a = 0 & \hat{\theta}_b = x[0] & \hat{\theta}_c = \frac{1}{N} \sum_{n=0}^{N-1} x[n] \\
 \mathbb{E}(\hat{\theta}_a) = 0 & \mathbb{E}(\hat{\theta}_b) = \mathbb{E}(x[0]) & \mathbb{E}(\hat{\theta}_c) = \theta \quad (\text{already seen}) \\
 \text{var}(\hat{\theta}_a) = 0 & = \mathbb{E}(\theta + w[0]) & \text{var}(\hat{\theta}_c) = \mathbb{E}[(\hat{\theta}_c - \mathbb{E}[\hat{\theta}_c])^2] \\
 & = \theta + 0 = \theta & \text{(Continued on next slide.)} \\
 & \text{var}(\hat{\theta}_b) = \text{var}(x[0]) = \sigma^2 &
 \end{array}$$

Estimator Variance

$$\begin{aligned}\text{var}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta}_c - \mathbb{E}[\hat{\theta}_c])^2] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} x[n] - \mathbb{E}[\hat{\theta}_c]\right)^2\right] \\ &= \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} \theta + w[n] - \theta\right)^2\right] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} w[n]\right)^2\right] = \frac{1}{N^2} \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n]\right)^2\right]\end{aligned}\tag{8}$$

Estimator Variance

$$\begin{aligned}\text{var}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta}_c - \mathbb{E}[\hat{\theta}_c])^2] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} x[n] - \mathbb{E}[\hat{\theta}_c]\right)^2\right] \\ &= \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} \theta + w[n] - \theta\right)^2\right] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} w[n]\right)^2\right] = \frac{1}{N^2} \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n]\right)^2\right]\end{aligned}\quad (8)$$

§ Now,

$$\begin{aligned}\text{var}\left(\sum_{n=0}^{N-1} w[n]\right) &= \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n] - \mathbb{E}\left[\sum_{n=0}^{N-1} w[n]\right]\right)^2\right] \\ &= \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n] - \sum_{n=0}^{N-1} \overbrace{\mathbb{E}[w[n]]}^0\right)^2\right] = \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n]\right)^2\right]\end{aligned}$$

Estimator Variance

$$\begin{aligned} \text{var}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta}_c - \mathbb{E}[\hat{\theta}_c])^2] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} x[n] - \mathbb{E}[\hat{\theta}_c]\right)^2\right] \\ &= \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} \theta + w[n] - \theta\right)^2\right] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{n=0}^{N-1} w[n]\right)^2\right] = \frac{1}{N^2} \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n]\right)^2\right] \end{aligned} \quad (8)$$

§ Now,

$$\begin{aligned} \text{var}\left(\sum_{n=0}^{N-1} w[n]\right) &= \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n] - \mathbb{E}\left[\sum_{n=0}^{N-1} w[n]\right]\right)^2\right] \\ &= \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n] - \sum_{n=0}^{N-1} \overbrace{\mathbb{E}[w[n]]}^0\right)^2\right] = \mathbb{E}\left[\left(\sum_{n=0}^{N-1} w[n]\right)^2\right] \end{aligned}$$

§ Using the above in eqn. (8)

$$\begin{aligned} \text{var}(\hat{\theta}_c) &= \frac{1}{N^2} \text{var}\left(\sum_{n=0}^{N-1} w[n]\right) = \frac{1}{N^2} \left(\sum_{n=0}^{N-1} \text{var}(w[n])\right) \quad (\text{WGN}) \\ &= \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N} \end{aligned}$$

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}\text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)]\end{aligned}$$

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}\text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])] \\ &\quad (\text{why?}) - (\text{Hint: What is random here?})\end{aligned}$$

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}
 \text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])] \\
 &\quad (\text{why?}) - (\text{Hint: What is random here?}) \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}]) \xrightarrow{0}
 \end{aligned}$$

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}\text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])] \\ &\quad (\text{why?}) - (\text{Hint: What is random here?}) \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}]) \xrightarrow{0} \\ &= \text{var}(\hat{\theta}) + \text{bias}^2(\hat{\theta})\end{aligned}$$

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}
 \text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])] \\
 &\quad (\text{why?}) - (\text{Hint: What is random here?}) \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}]) \xrightarrow{0} \\
 &= \text{var}(\hat{\theta}) + \text{bias}^2(\hat{\theta})
 \end{aligned}$$

§ So the mean square error in estimation, is composed of errors due to the variance of the estimator as well as the bias.

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}
 \text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])] \\
 &\quad (\text{why?}) - (\text{Hint: What is random here?}) \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}]) \\
 &= \text{var}(\hat{\theta}) + \text{bias}^2(\hat{\theta})
 \end{aligned}$$

§ So the mean square error in estimation, is composed of errors due to the variance of the estimator as well as the bias.

§ Recall MC evaluation

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T \quad \text{and} \quad v_\pi(s) = \mathbb{E}[G_t | S_t = s]$$

$$\hat{v}_\pi(s) = \frac{1}{N} \sum_{i=1}^N G_t^{(i)}(S_t = s)$$

Estimator Mean Square Error

§ The mean of the square error of estimation is,

$$\begin{aligned}
 \text{mse}(\hat{\theta}_c) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + \mathbb{E}[(\mathbb{E}[\hat{\theta}] - \theta)^2] + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta)] \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)\mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])] \\
 &\quad (\text{why?}) - (\text{Hint: What is random here?}) \\
 &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 2(\mathbb{E}[\hat{\theta}] - \theta)(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}]) \xrightarrow{0} \\
 &= \text{var}(\hat{\theta}) + \text{bias}^2(\hat{\theta})
 \end{aligned}$$

§ So the mean square error in estimation, is composed of errors due to the variance of the estimator as well as the bias.

§ Recall MC evaluation

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T \quad \text{and} \quad v_\pi(s) = \mathbb{E}[G_t | S_t = s]$$

$$\hat{v}_\pi(s) = \frac{1}{N} \sum_{i=1}^N G_t^{(i)}(S_t = s)$$

§ So $\hat{v}_\pi(s)$ is an unbiased estimator but with variance (inversely proportional to number of samples N .)

Bias and Variance of MC and TD

- § One key contribution of variance in MC evaluation comes from the randomness at each timestep.
- § This is not the case in TD as the G_t is estimated by bootstrapping,

$$\hat{G}_t = R_{t+1} + \gamma \hat{V}(S_{t+1})$$

- § This makes the estimator suffer less from variance as randomness comes from only one random step taken. The rest is deterministic.
- § But this introduces bias. The estimate always have the deterministic additive component $\gamma \hat{V}(S_{t+1})$

Reducing Variance in Policy Gradient Estimate

§ We have seen,

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

§ Inside each trajectory, a lot of randomness is there.

§ We can derive versions of this formula that eliminate terms to reduce variance.

§ Let us apply the log derivative trick ($\nabla_{\theta} \log p_{\theta}(\tau) = \sum \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$) to compute the gradient for a single reward term.

$$\nabla_{\theta} \mathbb{E}_{\tau} [r(\mathbf{s}_t, \mathbf{a}_t)] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t'=1}^t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \right) r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad (9)$$

§ Note that the sum goes up to t . Why?

Reducing Variance in Policy Gradient Estimate

§ We have seen,

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

§ Inside each trajectory, a lot of randomness is there.

§ We can derive versions of this formula that eliminate terms to reduce variance.

§ Let us apply the log derivative trick ($\nabla_{\theta} \log p_{\theta}(\tau) = \sum \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$) to compute the gradient for a single reward term.

$$\nabla_{\theta} \mathbb{E}_{\tau} [r(\mathbf{s}_t, \mathbf{a}_t)] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t'=1}^t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \right) r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad (9)$$

§ Note that the sum goes up to t . Why? - The reward at timestep t depends on actions till $t' \leq t$. - **Causality**

Reducing Variance in Policy Gradient Estimate

§ Summing over time we get (with some reordering of the sums, last)

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau} [r(\tau)] &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \sum_{t'=1}^t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \quad (10)\end{aligned}$$

§ With less randomness inside each trajectory the variance is less, but what about bias?

Reducing Variance in Policy Gradient Estimate

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t=1}^T \left(r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t'=1}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \sum_{t'=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right] \\ &= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\overbrace{\left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]}^{f(t, t')} \right] \tag{11}\end{aligned}$$

Reducing Variance in Policy Gradient Estimate

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t=1}^T \left(r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \\
 &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t'=1}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right] \\
 &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \sum_{t'=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right] \\
 &= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\overbrace{\left[\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]}^{f(t, t')} \right] \tag{11}
 \end{aligned}$$

§ Let us consider the term,

$$\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \tag{12}$$

§ We will show that for the case of $t' < t$ (reward coming before the action is performed) the above term is zero.

Reducing Variance in Policy Gradient Estimate

$$\begin{aligned}
 \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] &= \int p(\tau) f(t, t') d(\tau) \\
 &= \int p(s_1, a_1, \dots, s_t, a_t, \dots, s_{t'}, a_{t'}, \dots) f(t, t') \\
 &\quad d(s_1, a_1, \dots, s_t, a_t, \dots, s_{t'}, a_{t'}, \dots) \\
 &= \int p(s_t, a_t, s_{t'}, a_{t'}) f(t, t') d(s_t, a_t, s_{t'}, a_{t'}) \quad (13)
 \end{aligned}$$

§ The above comes from the property below.

$$\begin{aligned}
 \int_X \int_Y f(X) P(X, Y) dY dX &= \int_X \int_Y f(X) P(X) P(Y|X) dY dX \\
 &= \int_X f(X) P(X) dX \int_Y P(Y|X) dY \\
 &= \int_X f(X) P(X) dX \quad (14)
 \end{aligned}$$

§ Taking $X = \{s_t, a_t, s_{t'}, a_{t'}\}$ and Y the rest.

Reducing Variance in Policy Gradient Estimate

§ Till now we have,

$$\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] = \int p(s_t, a_t, s_{t'}, a_{t'}) f(t, t') d(s_t, a_t, s_{t'}, a_{t'}) \quad (15)$$

§ We will now use a variation of iterated expectation.

$$\begin{aligned} \mathbb{E}_{A,B}[f(A, B)] &= \int P(A, B) f(A, B) dB dA \\ &= \int P(B|A) P(A) f(A, B) dB dA \\ &= \int P(A) \int P(B|A) f(A, B) dB dA \\ &= \int P(A) \mathbb{E}_B [f(A, B) | A] dA \\ &= \mathbb{E}_A [\mathbb{E}_B [f(A, B) | A]] \end{aligned}$$

§ Taking $A = s_{t'}, a_{t'}$ and $B = s_t, a_t$, eqn. (15) can be written as,

$$\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] = \mathbb{E}_{s_{t'}, a_{t'}} \left[\mathbb{E}_{s_t, a_t} [f(t, t') | s_{t'}, a_{t'}] \right] \quad (16)$$

Reducing Variance in Policy Gradient Estimate

§ Putting the value of $f(t, t')$ back in eqn. (16), we get,

$$\begin{aligned} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] &= \mathbb{E}_{s_{t'}, a_{t'}} \left[\mathbb{E}_{s_t, a_t} [f(t, t') | s_{t'}, a_{t'}] \right] \\ &= \mathbb{E}_{s_{t'}, a_{t'}} \left[\mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) r(s_{t'}, \mathbf{a}_{t'}) | s_{t'}, \mathbf{a}_{t'}] \right] \\ &= \mathbb{E}_{s_{t'}, a_{t'}} [r(s_{t'}, \mathbf{a}_{t'}) \mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) | s_{t'}, \mathbf{a}_{t'}]] \end{aligned} \quad (17)$$

§ Let us take a closer look at the inner expectation,

$$\mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) | s_{t'}, \mathbf{a}_{t'}] = \int P(s_t, \mathbf{a}_t | s_{t'}, \mathbf{a}_{t'}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) d(\mathbf{a}_t, s_t) \quad (18)$$

§ Now, let us consider the timestep t be greater than t' , i.e., the action occurs after the reward. In such a case, $P(s_t, \mathbf{a}_t | s_{t'}, \mathbf{a}_{t'})$ can be broken down to $P(\mathbf{a}_t | s_t) P(s_t | s_{t'}, \mathbf{a}_{t'})$. Thus eqn. (18) becomes,

$$\begin{aligned} \mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) | s_{t'}, \mathbf{a}_{t'}] &= \int \int P(\mathbf{a}_t | s_t) P(s_t | s_{t'}, \mathbf{a}_{t'}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) d\mathbf{a}_t ds_t \\ &= \int P(s_t | s_{t'}, \mathbf{a}_{t'}) \int P(\mathbf{a}_t | s_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) d\mathbf{a}_t ds_t \\ &= \mathbb{E}_{s_t} \left[\mathbb{E}_{a_t} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t) | s_t] | s_{t'}, \mathbf{a}_{t'} \right] \end{aligned} \quad (19)$$

Reducing Variance in Policy Gradient Estimate

§ Now we will use a neat trick known as ‘*Expected Grad Log Probability*’ (EGLP) lemma which says $\mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] = 0$.

$$\begin{aligned}\mathbb{E}_{x \sim p_{\theta}(x)}[\nabla_{\theta} \log p_{\theta}(x)] &= \int p_{\theta}(x) \nabla_{\theta} \log p_{\theta}(x) dx = \int p_{\theta}(x) \frac{\nabla_{\theta} p_{\theta}(x)}{p_{\theta}(x)} dx \\ &= \int \nabla_{\theta} p_{\theta}(x) dx = \nabla_{\theta} \int p_{\theta}(x) dx = \nabla_{\theta} 1 = 0\end{aligned}$$

§ Thus the inner expectation in eqn. (19) is 0. This, in turn, means eqn. (17), (16) and (15) are all 0.

§ That is, $\mathbb{E}_{\tau \sim p_{\theta}(\tau)}[f(t, t')] = 0$ for $t > t'$.

§ Now for $t \leq t'$, $P(\mathbf{s}_t, \mathbf{a}_t | \mathbf{s}_{t'}, \mathbf{a}_{t'})$ can **not** be broken down to $P(\mathbf{a}_t | \mathbf{s}_t)P(\mathbf{s}_t | \mathbf{s}_{t'}, \mathbf{a}_{t'})$, as past state (\mathbf{s}_t) will get conditioned on future state and actions ($\mathbf{s}_{t'}, \mathbf{a}_{t'}$) violating the Markov property.

§ So, $\mathbb{E}_{\tau \sim p_{\theta}(\tau)}[f(t, t')] \neq 0$ for $t \leq t'$.

Reducing Variance in Policy Gradient Estimate

§ So we began with,

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] \quad (20)$$

and have shown that

$$\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] \begin{cases} = 0 & \text{if } t' < t \\ \neq 0 & \text{if } t' \geq t \end{cases}$$

Reducing Variance in Policy Gradient Estimate

§ So we began with,

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] \quad (20)$$

and have shown that

$$\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] \begin{cases} = 0 & \text{if } t' < t \\ \neq 0 & \text{if } t' \geq t \end{cases}$$

§ So, the gradient of the total expected return can be written as,

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \sum_{t=1}^T \sum_{t'=t}^T \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \sum_{t'=t}^T f(t, t') \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t'=t}^T \left(r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \quad (21) \end{aligned}$$

Reducing Variance in Policy Gradient Estimate

§ So we began with,

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] \quad (20)$$

and have shown that

$$\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] \begin{cases} = 0 & \text{if } t' < t \\ \neq 0 & \text{if } t' \geq t \end{cases}$$

§ So, the gradient of the total expected return can be written as,

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \sum_{t=1}^T \sum_{t'=t}^T \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [f(t, t')] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \sum_{t'=t}^T f(t, t') \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t'=t}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right] \quad (21) \end{aligned}$$

§ This is the 'reward to go' formulation we have seen earlier and which has less variance. But this also is same as the total expected reward expression which is unbiased. So this is unbiased and less variance estimator of the total expected reward.

Baselines

- § Good stuff is made more likely.
- § Bad stuff is made less likely.

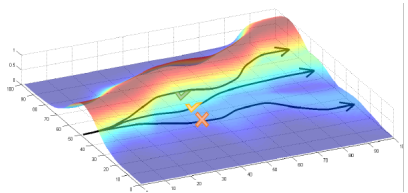


Figure credit: [Sergey Levine, UC Berkeley]

Baselines

- § Good stuff is made more likely.
- § Bad stuff is made less likely.
- § What if all have high reward?

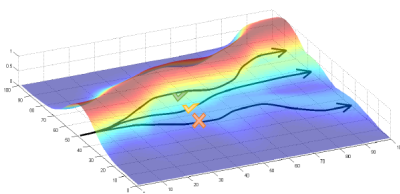


Figure credit: [Sergey Levine, UC Berkeley]

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

Baselines

- § Good stuff is made more likely.
- § Bad stuff is made less likely.
- § What if all have high reward?

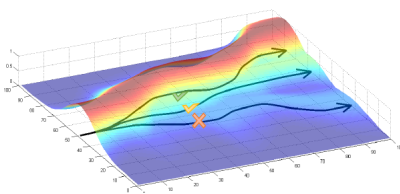


Figure credit: [Sergey Levine, UC Berkeley]

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)] = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b]]$$

- § Will it remain unbiased?
- § Only if $\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) b] = b \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau)] = 0$
- § And $\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau)] = 0$ by EGLP Lemma.

Baselines

§ So subtracting a constant baseline keeps the estimate unbiased.

§ A reasonable choice of baseline is average reward across the

trajectories, $b = \frac{1}{N} \sum_{i=1}^N r(\tau)$

§ What about variance?

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b]]$$

$$\text{var} = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b] \right)^2 \right] - \left(\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b]] \right)^2$$

$$= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b] \right)^2 \right] - \left(\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)] \right)^2$$

$$\frac{\partial \text{var}}{\partial b} = \frac{\partial \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b] \right)^2 \right]}{\partial b} - 0$$

$$= \frac{\partial \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\nabla_{\theta} \log p_{\theta}(\tau) \right)^2 [r^2(\tau) - 2r(\tau)b + b^2] \right]}{\partial b}$$

Baselines

$$\begin{aligned}\frac{\partial \text{var}}{\partial b} &= \frac{\partial \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 [r^2(\tau) - 2r(\tau)b + b^2] \right]}{\partial b} \\ &= 0 - 2 \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 r(\tau) \right] + 2b \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 \right]\end{aligned}$$

§ For minimum variance,

$$\begin{aligned}\frac{\partial \text{var}}{\partial b} &= 0 \\ - \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 r(\tau) \right] + b \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 \right] &= 0 \\ b &= \frac{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 r(\tau) \right]}{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 \right]}\end{aligned}$$

Advantage Function

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \underbrace{\sum_{t'=t}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right)}_{\widehat{Q}^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t)} \right] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \widehat{Q}^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right] \end{aligned}$$

Advantage Function

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \underbrace{\sum_{t'=t}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right)}_{\widehat{Q}^{\theta}(\mathbf{s}_t, \mathbf{a}_t)} \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \widehat{Q}^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

§ It would be good to have the true value of Q to be used in the equation.

Advantage Function

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \underbrace{\sum_{t'=t}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right)}_{\widehat{Q}^{\theta}(\mathbf{s}_t, \mathbf{a}_t)} \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \widehat{Q}^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

- § It would be good to have the true value of Q to be used in the equation.
- § But that is not available to us.

Advantage Function

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \underbrace{\sum_{t'=t}^T \left(r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right)}_{\widehat{Q}^{\theta}(\mathbf{s}_t, \mathbf{a}_t)} \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \widehat{Q}^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

- § It would be good to have the true value of Q to be used in the equation.
- § But that is not available to us.
- § Other alternatives are to estimate this value using methods that we have seen earlier - MC evaluation, Bootstrapped evaluation (TD), using function approximation for these.

Advantage Function

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\theta}(\mathbf{s}_t, \mathbf{a}_t) - \mathbb{E}_{\mathbf{a}_t} [Q^{\theta}(\mathbf{s}_t, \mathbf{a}_t)] \right) \right]$$

§ We can also use a baseline version of this.

Advantage Function

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - \mathbb{E}_{\mathbf{a}_t} [Q^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t)] \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - V^{\boldsymbol{\theta}}(\mathbf{s}_t) \right) \right]\end{aligned}$$

§ We can also use a baseline version of this.

Advantage Function

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - \mathbb{E}_{\mathbf{a}_t} [Q^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t)] \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - V^{\boldsymbol{\theta}}(\mathbf{s}_t) \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{s}_t) \right) A^{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

§ We can also use a baseline version of this.

§ This is called the 'Advantage function'.

Advantage Function

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\theta}(\mathbf{s}_t, \mathbf{a}_t) - \mathbb{E}_{\mathbf{a}_t} [Q^{\theta}(\mathbf{s}_t, \mathbf{a}_t)] \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(Q^{\theta}(\mathbf{s}_t, \mathbf{a}_t) - V^{\theta}(\mathbf{s}_t) \right) \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) A^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]\end{aligned}$$

- § We can also use a baseline version of this.
- § This is called the 'Advantage function'.
- § $A(\mathbf{s}_t, \mathbf{a}_t)$ can be approximated following the methods we used earlier (single sample backup or bootstrapping)

Advantage Function

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) A^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left(\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) A^{\theta}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})\end{aligned}$$

§ $Q^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^{\theta}(\mathbf{s}_t)$

§ $A^{\theta}(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^{\theta}(\mathbf{s}_{t+1}) - V^{\theta}(\mathbf{s}_t)$

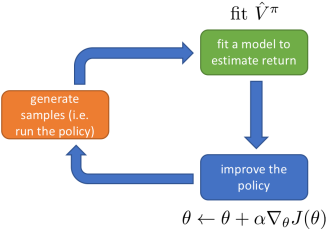
§ So we can use a neural network which learns to produce $V(\mathbf{s})$

Actor-Critic

An actor-critic algorithm

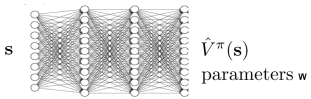
batch actor-critic algorithm:

1. sample $\{\mathbf{s}_i, \mathbf{a}_i\}$ from $\pi_\theta(\mathbf{a}|\mathbf{s})$ (run it on the robot)
2. fit $\hat{V}_w^\pi(\mathbf{s})$ to sampled reward sums
3. evaluate $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \hat{V}_w^\pi(\mathbf{s}'_i) - \hat{V}_w^\pi(\mathbf{s}_i)$
4. $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i|\mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
5. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$



$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_w^\pi(\mathbf{s}_{i,t+1})$$

$$\mathcal{L}(w) = \frac{1}{2} \sum_i \left\| \hat{V}_w^\pi(\mathbf{s}_i) - y_i \right\|^2$$



$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

Figure credit: [Sergey Levine, UC Berkeley]