

Chapter 3 Policy-based Methods

3.1 Policy Gradient

3.1.1 Direct Policy Differentiation

During previous discussions, we learned that trying to find θ^* would help us get more reward.

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad (3.1.1)$$

$$p_{\theta}(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \quad (3.1.2)$$

We can rewrite the objective function as:

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \quad (3.1.3)$$

where \sum_i means the sum over samples from π_{θ} . If we abbreviate $\sum_t r(\mathbf{s}_t, \mathbf{a}_t)$ as $r(\tau)$, then we have:

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)}[r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau \quad (3.1.4)$$

To find the optimal θ , we calculate the policy differentiation:

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) d\tau = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)] \quad (3.1.5)$$

The derivation from the second term to the third term of the equation might be a bit confusing. Here we use a convenient identity:

$$\frac{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)}{p_{\theta}(\tau)} = \nabla_{\theta} p_{\theta}(\tau) \quad (3.1.6)$$

Also, note that $p_{\theta}(\tau)$ is given in (4.2), and take \log on both sides, then we have:

$$\log p_{\theta}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \quad (3.1.7)$$

Substituting into equation (4.5), we obtain the *direct policy differentiation*.

Theorem 3.1 (Direct policy differentiation).

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

Algorithm 15 REINFORCE algorithm**Require:** arbitrarily initialized π_θ

- 1: **repeat**
- 2: sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
- 3: $\nabla_\theta J(\theta) \approx \sum_i (\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
- 4: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
- 5: **until** convergence

Note. Markov property is not actually used in policy gradient. Therefore, you can use it in partially observed MDPs without modification.

3.1.2 Comparison to Maximum Likelihood

In policy gradient:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \quad (3.1.8)$$

In maximum likelihood:

$$\nabla_\theta J_{\text{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \quad (3.1.9)$$

Insight. Good stuff is made more likely; bad stuff is made less likely. What we just did in policy gradient simply formalize the notion of “through trial and error”!

3.1.3 The High Variance of Policy Gradient

One significant drawback of policy gradient is that the gradient is often noisy (i.e. the variance is high). Here are some methods of reducing variance.

Causality

Firstly, causality tells us that policy at time t' should not affect the reward at time t when $t \leq t'$. The modified gradient function is:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) \quad (3.1.10)$$

The expression in brackets means “reward to go”. And sometimes we also write it as $\hat{Q}_{i,t}^\pi$: estimate of expected reward if we take action $\mathbf{a}_{i,t}$ in state $\mathbf{s}_{i,t}$ ($\hat{Q}_{i,t}^\pi = \sum_{t'=1}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$)

Baseline

The second approach is to introduce **baseline** to the reward part:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log p_\theta(\tau) [r(\tau) - b], \quad b = \frac{1}{N} \sum_{i=1}^N r(\tau) \quad (3.1.11)$$

We are allowed to do that because what we really care about is *expectation* ($E[\nabla_\theta \log p_\theta(\tau)b]$), and subtracting a baseline is *unbiased in expectation*:

$$E[\nabla_\theta \log p_\theta(\tau)b] = \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) b d\tau = \int \nabla_\theta p_\theta(\tau) b d\tau = b \nabla_\theta \int p_\theta(\tau) d\tau = b \nabla_\theta 1 = 0 \quad (3.1.12)$$

Here we use the average reward as the baseline. It is not the best baseline (as to reducing variance), but it works pretty good. If you are not satisfied with it, we can try to find the best baseline:

$$\text{Var}[x] = E[x^2] - E[x]^2 \quad (3.1.13)$$

$$\nabla_\theta J(\theta) = E_{\tau \sim p_\theta(\tau)} [\nabla_\theta \log p_\theta(\tau) (r(\tau) - b)] \quad (3.1.14)$$

$$\text{Var} = E_{\tau \sim p_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau)(r(\tau) - b))^2 \right] - E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau)(r(\tau) - b)]^2 \quad (3.1.15)$$

$$\frac{d \text{Var}}{db} = \frac{d}{db} E [g(\tau)^2 (r(\tau) - b)^2] = -2E [g(\tau)^2 r(\tau)] + 2bE [g(\tau)^2] = 0 \quad (3.1.16)$$

Solve this equation and you'll find the optimal baseline:

$$b = \frac{E [g(\tau)^2 r(\tau)]}{E [g(\tau)^2]} \quad (3.1.17)$$

Note that this is just expected reward, but weighted by gradient magnitudes!

3.1.4 More about reducing the variance

Before moving on, let's take a moment to review what it actually means to reduce the variance.

In a nutshell, in REINFORCE algorithm we are trying to optimize some function (i.e. the objective function $J(\theta)$) on a random variable (the return of some stochastic policy). What you need to do is sample a bunch of trajectories from the current policy, estimate the current expected value (remember, you are **sampling** trajectories so you don't know the expected value you can only **estimate** it), and update the parameters of the current policy in the direction that optimizes your estimate.

But you only have a finite amount of time to sample (at some point you have to update the parameters to make progress), so your estimate will be off of the true value and you will making updates based on that estimate. So let's say you are only going to sample 10 times, if the distribution of the values of that random variable has high variance, then sampling only 10 times will be a pretty poor estimate. If it has low variance, then 10 samples might be enough to approximate the value you want. Therefore we introduce the baseline to reduce the variance of the samples.

Another way of understanding *reducing the variance* is *reducing the aggressiveness of the updates*. Here is a concrete example: assuming we have a reward function that looks like

$$f(x) = \begin{cases} -x, & x < 0 \\ 0, & x \geq 0 \end{cases} \quad (3.1.18)$$

This can have a practical interpretation. For example, if you are in a fire, the only way to receive a reward of 0 is by escaping the fire in the positive x direction; otherwise, if you move in the negative x direction, you will get hurt and the reward will be negative.

In this scenario, if the baseline is not introduced, the gradient at 0 will be infinite, so one would not know if this policy would lead to some strange places. However, by introducing the baseline, a reasonable gradient will be obtained. This is why adding a baseline can also be interpreted as reducing the policy's aggressiveness.

3.2 Off-Policy Policy Gradients

If we revisit the derivation of policy differentiation, it's obvious that policy gradient is on-policy. The underlined part would be a trouble, because the neural networks change only a little bit with each gradient step, but you need to do the sampling all over again. Therefore, on-policy learning can be extremely inefficient!

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

The question here is: what if we skip the sampling step? What if we don't have samples from $p_{\theta}(\tau)$ (we have samples from some $\bar{p}(\tau)$ instead)?

Before delving into this issue, let's talk about some maths first: *importance sampling*

Theorem 3.2 (Importance Sampling).

$$\begin{aligned} E_{x \sim p(x)} [f(x)] &= \int p(x) f(x) dx \\ &= \int \frac{q(x)}{q(x)} p(x) f(x) dx \\ &= \int q(x) \frac{p(x)}{q(x)} f(x) dx \\ &= E_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

Its mathematical essence lies in using a *known distribution* (typically an easy-to-sample distribution) to estimate the expectation of another distribution (usually a difficult-to-sample distribution). Back to off-policy learning,

here we can turn the original objective function $J(\theta) = E_{\tau \sim p_\theta(\tau)}[r(\tau)]$, into a new objective function using importance sampling:

$$J(\theta) = E_{\tau \sim \bar{p}(\tau)} \left[\frac{p_\theta(\tau)}{\bar{p}(\tau)} r(\tau) \right] \quad (3.2.1)$$

Here we have:

$$p_\theta(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \quad (3.2.2)$$

$$\frac{p_\theta(\tau)}{\bar{p}(\tau)} = \frac{p(\mathbf{s}_1) \prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_t, \mathbf{a}_t)}{p(\mathbf{s}_1) \prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_t, \mathbf{a}_t)} = \frac{\prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t)} \quad (3.2.3)$$

for some new parameters θ' , similarly we have:

$$\frac{p_{\theta'}(\tau)}{p_\theta(\tau)} = \frac{\prod_{t=1}^T \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \quad (3.2.4)$$

therefore,

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim p_\theta(\tau)} \left[\frac{\nabla_{\theta'} p_{\theta'}(\tau)}{p_\theta(\tau)} r(\tau) \right] = E_{\tau \sim p_\theta(\tau)} \left[\frac{p_{\theta'}(\tau)}{p_\theta(\tau)} \nabla_{\theta'} \log p_{\theta'}(\tau) r(\tau) \right] \quad (3.2.5)$$

Now estimate locally, at $\theta = \theta'$, we have:

$$\nabla_\theta J(\theta) = E_{\tau \sim p_\theta(\tau)} [\nabla_\theta \log p_\theta(\tau) r(\tau)] \quad (3.2.6)$$

When $\theta \neq \theta'$, then:

$$\begin{aligned} \nabla_{\theta'} J(\theta') &= E_{\tau \sim p_\theta(\tau)} \left[\frac{p_{\theta'}(\tau)}{p_\theta(\tau)} \nabla_{\theta'} \log p_{\theta'}(\tau) r(\tau) \right] \\ &= E_{\tau \sim p_\theta(\tau)} \left[\left(\prod_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \right) \left(\sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \end{aligned} \quad (3.2.7)$$

Also, take causality into consideration, future actions don't affect current weight:

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim p_\theta(\tau)} \left[\sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \left(\prod_{t'=1}^t \frac{\pi_{\theta'}(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{\pi_\theta(\mathbf{a}_{t'} | \mathbf{s}_{t'})} \right) \left(\sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \left(\prod_{t'' \neq t}^t \frac{\pi_{\theta'}(\mathbf{a}_{t''} | \mathbf{s}_{t''})}{\pi_\theta(\mathbf{a}_{t''} | \mathbf{s}_{t''})} \right) \right) \right] \quad (3.2.8)$$

The underlined part in [3.2.8](#) is distribution mismatch between different policies. If we ignore this part (later we'll see why this is reasonable), we'll get the following:

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim p_\theta(\tau)} \left[\sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \left(\prod_{t'=1}^t \frac{\pi_{\theta'}(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{\pi_\theta(\mathbf{a}_{t'} | \mathbf{s}_{t'})} \right) \left(\sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right] \quad (3.2.9)$$

In the next section, we'll see that [3.2.9](#) is equivalent to policy iteration.

3.3 (Optional) Policy Gradient as Policy Iteration

Now we give further analysis of policy gradient as policy iteration. The trick here is to express the expected value under policy π_θ in terms of the expected value with respect to the trajectory τ under policy $\pi_{\theta'}$:

$$\begin{aligned}
J(\theta') - J(\theta) &= J(\theta') - E_{\mathbf{s}_0 \sim p(\mathbf{s}_0)} [V^{\pi_{\theta}}(\mathbf{s}_0)] \\
&= J(\theta') - E_{\tau \sim p_{\theta'}(\tau)} [V^{\pi_{\theta}}(\mathbf{s}_0)] \\
&= J(\theta') - E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t V^{\pi_{\theta}}(\mathbf{s}_t) - \sum_{t=1}^{\infty} \gamma^t V^{\pi_{\theta}}(\mathbf{s}_t) \right] \\
&= J(\theta') + E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_{\theta}}(\mathbf{s}_{t+1}) - V^{\pi_{\theta}}(\mathbf{s}_t)) \right] \\
&= E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] + E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_{\theta}}(\mathbf{s}_{t+1}) - V^{\pi_{\theta}}(\mathbf{s}_t)) \right] \\
&= E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (r(\mathbf{s}_t, \mathbf{a}_t) + \gamma V^{\pi_{\theta}}(\mathbf{s}_{t+1}) - V^{\pi_{\theta}}(\mathbf{s}_t)) \right] \\
&= E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right]
\end{aligned} \tag{3.3.1}$$

Writing out the expectation with respect to the trajectory τ explicitly:

$$\begin{aligned}
E_{\tau \sim p_{\theta'}(\tau)} \left[\sum_t \gamma^t A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right] &= \sum_t E_{\mathbf{s}_t \sim p_{\theta'}(\mathbf{s}_t)} [E_{\mathbf{a}_t \sim \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)} [\gamma^t A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t)]] \\
&= \sum_t E_{\mathbf{s}_t \sim p_{\theta'}(\mathbf{s}_t)} \left[E_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)} \left[\frac{\pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)} \gamma^t A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right] \right]
\end{aligned} \tag{3.3.2}$$

where the second equation is obtained by using the importance sampling formula [3.2](#). Now the remaining question is: can we use $p_{\theta'}$ instead of p_{θ} in the expectation? The answer is yes, because we can actually bound the distribution gap between p_{θ} and $p_{\theta'}$ if the two policies π_{θ} and $\pi_{\theta'}$ are close enough.

3.4 (Optional) Bounding the Distribution Gap

3.5 (Optional) Advanced Policy Gradient Methods

In this section, we'll introduce the ideas of two advanced policy gradient methods. Feel free to skip this section if you are not interested in mathy details.

3.5.1 Trust Region Policy Optimization (TRPO)

Add TRPO algorithm

3.5.2 Proximal Policy Optimization (PPO)

Add PPO algorithm

3.6 Implementation Tips

There are some tips in implementing policy gradient:

1. Remember that the gradient has high variance. This isn't the same as supervised learning, actually the gradients would be really noisy.
2. Consider using much larger batches.
3. Tweaking learning rates is very hard. (We'll learn about policy gradient-specific learning rate adjustment methods later)

Chapter 4 Actor-Critic Methods

4.1 Introducing the Actor-Critic Methods

4.1.1 Recap

We mention Q function and V function here again as a recap.

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \mid \mathbf{s}_t, \mathbf{a}_t] : \text{total reward from taking } \mathbf{a}_t \text{ in } \mathbf{s}_t$$

$$V^\pi(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t \mid \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)] : \text{total reward from } \mathbf{s}_t$$

$$A^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t) : \text{how much better } \mathbf{a}_t \text{ is than average}$$

Policy gradient:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} \mid \mathbf{s}_{i,t}) Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \quad (4.1.1)$$

Adding a baseline:

$$\begin{aligned} \nabla_\theta J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} \mid \mathbf{s}_{i,t}) (Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - V(\mathbf{s}_{i,t})) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} \mid \mathbf{s}_{i,t}) A^\pi(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \end{aligned} \quad (4.1.2)$$

4.1.2 Policy Evaluation

Let's go back to three basics steps of RL: **generate samples** (i.e. run the policy), **fit a model to estimate return**, and **improve the policy**. In policy gradient, we figure out how to calculate $\nabla_\theta J(\theta)$, which tells us how to improve the policy. But we still need to fit the model, which quantitatively tell us how good the policy is. We have three possible quantities, Q^π , V^π or A^π , so the question is: *what* should we fit to *what*? Since the current reward of taking \mathbf{a}_t at \mathbf{s}_t is fixed, so we can rewrite $Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ and $A^\pi(\mathbf{s}_t, \mathbf{a}_t)$:

$$\begin{aligned} Q^\pi(\mathbf{s}_t, \mathbf{a}_t) &= r(\mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=t+1}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \mid \mathbf{s}_t, \mathbf{a}_t] \\ &\approx r(\mathbf{s}_t, \mathbf{a}_t) + V^\pi(\mathbf{s}_{t+1}) \end{aligned} \quad (4.1.3)$$

$$A^\pi(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^\pi(\mathbf{s}_{t+1}) - V^\pi(\mathbf{s}_t) \quad (4.1.4)$$

Therefore V^π would be a nice choice, since Q^π and A^π depend on both actions and states, while V^π depends on solely states. Of course this is not the only option in actor-critic algorithm. We'll talk about that later.

Actually calculating the value function is essentially calculating the expectation of reward under a given policy. That's why the process of fitting value function is also called policy evaluation.

Theorem 4.1 (Monte Carlo policy evaluation). Monte Carlo policy evaluation estimates the value of state-action pairs based on random sampling of experiences. It involves averaging the returns observed from multiple episodes to approximate the true value function for a given policy.

$$V^\pi(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

Note that in traditional Monte Carlo policy evaluation you need to generate samples from the same initial state, which means constantly resetting the simulator. Another way of doing that is training a neural network to estimate value function:

$$\text{training data: } \left\{ \left(\mathbf{s}_{i,t}, \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) \right\}$$

$$\text{supervised regression: } \mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$

In fact, the value function is very intuitive. For example, when training an agent to play a chess-like game, if we set the probability of winning the game to 1 and the probability of losing to 0, the value function typically indicates the probability of eventually winning the game in the current state.

4.1.3 From Evaluation to Actor Critic

Now that we have learned policy evaluation and policy gradients, now we are able to put the two pieces together. And that makes an actor-critic algorithm:

Algorithm 16 Batch actor-critic algorithm

- 1: **repeat**
 - 2: sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
 - 3: fit $\hat{V}_\phi^\pi(\mathbf{s})$ to sampled reward sums
 - 4: evaluate $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \hat{V}_\phi^\pi(\mathbf{s}'_i) - \hat{V}_\phi^\pi(\mathbf{s}_i)$
 - 5: $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
 - 6: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
 - 7: **until** convergence
-

Note that while fitting the value function, we are fitting the sum of reward in a given episode length T . What if T is ∞ ? In many cases \hat{V}_ϕ^π can get infinitely large. A simple trick here is to all discount factor γ , which would tell the policy that its better to get rewards sooner than later.

Without discount factors, the target function and loss function of the neural network is:

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

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2 \quad (4.1.5)$$

Now adding the discount factors, we'll have:

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1}), \quad \gamma \in [0, 1] (0.99 \text{ works well}) \quad (4.1.6)$$

Insight. One way of understanding discount factors is that we change the MDP by adding an extra state of death. Once we enter the death state we never leave, and the reward is zero. In each state the agent has the probability of $1 - \gamma$ of falling into the death state.

4.1.4 Aside: the Discount Factor

Then how do we introduce discount factors to (Monte Carlo) policy gradients? Actually there seems to be two ways of doing this.

$$\text{option 1: } \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$

$$\text{option 2: } \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T \gamma^{t-1} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \quad (4.1.7)$$

The only difference is whether we introduce the discount factors before or after we use the causality rules.

To showcase the differences more effectively, we can rewrite the expression of option 2 in the form of option 1.

$$\begin{aligned}
\nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} \mid \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T \gamma^{t-1} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} \mid \mathbf{s}_{i,t}) \left(\sum_{t'=t}^T \gamma^{t'} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \gamma^{t-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} \mid \mathbf{s}_{i,t}) \left(\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)
\end{aligned} \tag{4.1.8}$$

Note that γ^{t-1} comes before $\nabla_{\theta} \log \pi_{\theta}$. Then it becomes much clearer. Option 2 implies that we not only care less about the rewards in the future, we also care less about the decisions in the future. As a result you are discounting the gradients. In other word, making right decision in the first time step is more important than making right decision in future time steps. If you are solving a real discount problem, this is exactly what we are trying to do. Because remember the modified MDP settings, later steps don't matter when you are dead. **But**, in reality, this is often not quite what we want. The version that we actually use is **option 1**.

Take some time to review why we introduced the discount factor. For tasks with particularly long time episodes, we introduced the discount factor to artificially reduce the rewards for future actions in order to compute the value function. However, in practical tasks, we do not want to do this. For example, if we want a robot to run steadily forward, we actually want it to keep running. We want to learn a policy that can do the right thing for a long time, rather than just at nearby timesteps. Therefore, option 1 becomes more reasonable.

Algorithm 17 Batch actor-critic algorithm (with discount)

- 1: **repeat**
 - 2: sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t)$ (run the policy)
 - 3: fit $\hat{V}_{\phi}^{\pi}(\mathbf{s})$ to sampled reward sums
 - 4: evaluate $\hat{A}^{\pi}(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}'_i) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_i)$
 - 5: $\nabla_{\theta} J(\theta) \approx \sum_i \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_i \mid \mathbf{s}_i) \hat{A}^{\pi}(\mathbf{s}_i, \mathbf{a}_i)$
 - 6: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$
 - 7: **until** convergence
-

The only difference that in step 3 we introduce the discount factor. In previous discussions we have been using policy gradients in episodic batch mode settings. If we modify a little bit, we'll get the online actor-critic algorithm:

Algorithm 18 Online actor-critic algorithm (with discount)

- 1: **repeat**
 - 2: take action $\mathbf{a} \sim \pi_{\theta}(\mathbf{a} \mid \mathbf{s})$, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$
 - 3: update $\hat{V}_{\phi}^{\pi}(\mathbf{s})$ using target $r + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}')$
 - 4: evaluate $\hat{A}^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}') - \hat{V}_{\phi}^{\pi}(\mathbf{s})$
 - 5: $\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \pi_{\theta}(\mathbf{a} \mid \mathbf{s}) \hat{A}^{\pi}(\mathbf{s}, \mathbf{a})$
 - 6: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$
 - 7: **until** convergence
-

Insight. In actor-critic algorithm, **actor** is the policy, and **critic** is the value function. It can be viewed as a improved version of policy gradient, with reduced variance.

4.2 Further Analysis

4.2.1 Critics as Baselines

Let's make a comparison here between the actor-critic algorithm and the policy gradient algorithm:

$$\begin{aligned} \text{Actor-critic: } \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t}) \right) \\ \text{Policy gradient: } \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\left(\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right) \end{aligned} \quad (4.2.1)$$

In actor-critic, we have lower variance since we have the critic, but it is not unbiased if the critic is not perfect. While in policy gradient, it is unbiased but the variance would be high due to the single-sample estimate.

So intuitively a nature step we can take is to use \hat{V}_{ϕ}^{π} while still keep the estimator unbiased:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\left(\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t}) \right) \quad (4.2.2)$$

Up to now we are all using the state-dependent functions as baselines. So can we use functions that involves both action and state as baselines? Here we come to control variates, which means baselines that are action-dependent. Here we are going to talk about that.

Baselines that use both actions and states are also called *control variates*. The true advantage function is:

$$\begin{aligned} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) &= Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi}(\mathbf{s}_t) \\ &= \sum_{t'=t}^T E_{\pi_{\theta}}[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t] - E_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}[Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)] \end{aligned} \quad (4.2.3)$$

While the approximate advantage function we use in policy gradient is:

$$\hat{A}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - V_{\phi}^{\pi}(\mathbf{s}_t) \quad (4.2.4)$$

This is unbiased and has a lower variance (but still a higher variance than the actor-critic because of the single sample estimate). But we can make the variance even lower if we subtract the Q value.

$$\hat{A}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \quad (4.2.5)$$

This version has a nice property that it goes to zero in expectation if critic is correct. Unfortunately, it does not work if you simply plug it into the policy gradient, since there is an error term you have to compensate for. Now taking that error term into accounts, we have:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\hat{Q}_{i,t} - Q_{\phi}^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) + \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} E_{\mathbf{a} \sim \pi_{\theta}(\mathbf{a} | \mathbf{s}_{i,t})} [Q_{\phi}^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_t)] \quad (4.2.6)$$

The second term represents the gradient of expectation value under the policy of the baseline. Note that this equation is valid even when the baseline depends merely on the state functions, but in that case the second term equals to zero. This kind of trick can provide for a very low variance policy gradient.

We can also take a look at the two different versions of advantage function:

$$\begin{aligned} \hat{A}_{\text{C}}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) &= r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t) && \text{+ lower variance} \\ &&& \text{- higher bias if value is wrong (it always is)} \\ \hat{A}_{\text{MC}}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t) && \text{+ no bias} \\ &&& \text{- higher variance (because single-sample estimate)} \end{aligned} \quad (4.2.7)$$

So can we find something in the middle, combine these two, to control bias/variance tradeoff? The trick here is that instead of using an infinite time horizon, you cut it off before the variance goes too high, given that the variance is generally small in the near future and goes higher in the far future. Here is what an n-step return estimator does:

$$\hat{A}_n^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{t+n} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_\phi^\pi(\mathbf{s}_t) + \gamma^n \hat{V}_\phi^\pi(\mathbf{s}_{t+n}) \quad (4.2.8)$$

Generally the larger n you choose, the bias goes smaller and the variance goes larger. The first and the third term contributes to variance, and the second term contributes to bias. The n here we use in n-step return estimator is fixed. Actually we do not have to choose a single n . Instead, we can take all possible n at one time, which is the generalized advantage estimation.

$$\hat{A}_{\text{GAE}}^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{n=1}^{\infty} w_n \hat{A}_n^\pi(\mathbf{s}_t, \mathbf{a}_t) \quad (4.2.9)$$

The generalized advantage estimation is actually a weighted combination of all n-step returns. In terms of determining the weights, we mostly prefer cutting earlier because it brings less variance. Therefore, a descent choice is to use exponential falloff ($w_n \propto \lambda^{n-1}$).

$$\hat{A}_{\text{GAE}}^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{\infty} (\gamma\lambda)^{t'-t} \delta_{t'} \quad \delta_{t'} = r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) + \gamma \hat{V}_\phi^\pi(\mathbf{s}_{t'+1}) - \hat{V}_\phi^\pi(\mathbf{s}_{t'}) \quad (4.2.10)$$

4.3 (Optional) Advanced Actor-Critic Methods

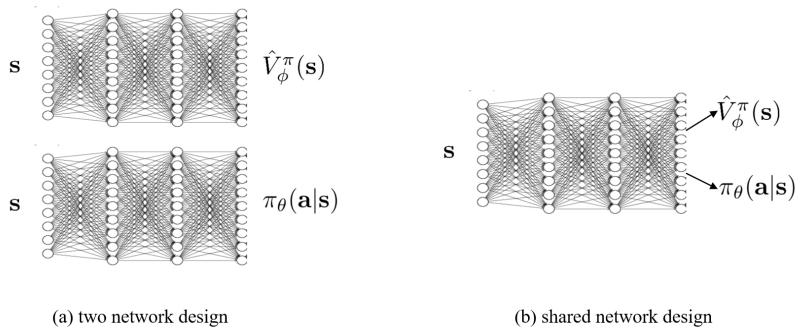
4.3.1 Soft Actor-Critic (SAC)

Add SAC

4.4 Implementations Tips

4.4.1 Architecture Design

There are two options of the neural network architecture design: two network design and shared network design.



The two-network design is simple and stable, but lacks efficient feature sharing between the actor and critic. In contrast, the shared network design allows for potential feature sharing efficiency, but it has more hyperparameters and can be more complex and less stable during training.

4.4.2 Parallel Actor-Critic

In a typical online actor-critic algorithm, step2 and step4 work best with a batch (e.g., parallel workers). There are also two types of parallel actor-critic design: **synchronized parallel actor-critic** and **asynchronous parallel actor-critic**.

Synchronized parallel actor-critic: multiple agents learn simultaneously, but they must synchronously update parameters at each step to ensure consistent behavior across all agents. However different agents would use different random seeds so their actions would be a little bit different.

Asynchronous parallel actor-critic: different agents independently update parameters during learning without the need for synchronous updates, leading to improved learning efficiency. However, each agent independently update parameters at their own pace, therefore the parameters update may be inconsistent.

4.4.3 Off-Policy Actor-Critic Algorithm

The problem of sample efficiency gets us thinking about another problem: can we remove the on-policy assumption entirely? A primitive way is to use a replay buffer where we store the transitions we saw in prior time steps. Based on this we have an immature off policy actor-critic algorithm:

Algorithm 19 (Fake) Off-policy actor-critic algorithm

- 1: **repeat**
 - 2: take action $\mathbf{a} \sim \pi_\theta(\mathbf{a} | \mathbf{s})$, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$, store in \mathcal{R}
 - 3: sample a batch $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$ from buffer \mathcal{R}
 - 4: update $\hat{V}_\phi^\pi(\mathbf{s})$ using target $y_i = r_i + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i)$ for each \mathbf{s}_i , $\mathcal{L}(\phi) = \frac{1}{N} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$
 - 5: evaluate $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) - \hat{V}_\phi^\pi(\mathbf{s}_i)$
 - 6: $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
 - 7: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
 - 8: **until** convergence
-

There are two fallacies lying in this immature algorithm. The first problem comes from the target value. When you are loading transitions from the replay buffer, you are actually loading actions from the older version of policy, not the latest one. Therefore \mathbf{s}' comes from the older actions, which is not we actually want.

The second issue comes from the same reason. Because the action \mathbf{a}_i does not come from the latest policy, you cannot calculate policy gradient.

To fix the first problem, the method we take here is to substitute V-function with Q-function. Note the difference between these two functions:

$$\begin{aligned} V^\pi(\mathbf{s}_t) &= \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t] \\ Q^\pi(\mathbf{s}_t, \mathbf{a}_t) &= \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t] \end{aligned} \quad (4.4.1)$$

Q-function cares about the total reward from taking \mathbf{a}_t in \mathbf{s}_t , and then following the policy π_θ . However we do not restrict that action \mathbf{a}_t must come from π_θ . It is a valid function for any action.

So here in the third step in Algorithm 6, instead of using $\hat{V}_\phi^\pi(\mathbf{s})$, we update $\hat{Q}_\phi^\pi(\mathbf{s})$ using target y_i . This time we have:

$$\begin{aligned} y_i &= r_i + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) \\ &= r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i) \end{aligned} \quad (4.4.2)$$

The trick here is to use:

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t] = E_{\mathbf{a} \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)] \quad (4.4.3)$$

Note that action \mathbf{a}'_i in (5.10) does not come from the replay buffer \mathcal{R} . Instead, it is the action an agent would have taken following policy π_θ under state \mathbf{s}'_i . This procedure does not require state \mathbf{s}'_i to actually happens in a simulator, because all you need is to plug state \mathbf{s}'_i into the neural network and get the output action.

We can use a similar approach to resolve the policy gradient issue. Instead of using \mathbf{a}_i which comes from the replay buffer, we use $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a} | \mathbf{s}_i)$ that comes from the latest policy.

So now we have step 5 in Algorithm 6 as:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi | \mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi) \quad (4.4.4)$$

One more thing to say. In practice, we don't actually use the advantage function in policy gradient. We simply use the Q-function here:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_i^{\pi} | \mathbf{s}_i) \hat{Q}^{\pi}(\mathbf{s}_i, \mathbf{a}_i^{\pi}) \quad (4.4.5)$$

This would increase the variance because it does not involve a baseline. However, a high variance here is OK since we don't need to interact with the simulators to sample actions. So it's easy to lower the variance by generating more samples of actions \mathbf{a}_i^{π} , without generating states \mathbf{s}_i .

So we have the fixed version of off-policy actor-critic algorithm here:

Algorithm 20 Off-policy actor-critic algorithm

- 1: **repeat**
 - 2: take action $\mathbf{a} \sim \pi_{\theta}(\mathbf{a} | \mathbf{s})$, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$, store in \mathcal{R}
 - 3: sample a batch $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$ from buffer \mathcal{R}
 - 4: update $\hat{Q}_{\phi}^{\pi}(\mathbf{s})$ using target $y_i = r_i + \gamma \hat{Q}_{\phi}^{\pi}(\mathbf{s}'_i, \mathbf{a}'_i)$ for each \mathbf{s}_i , $\mathcal{L}(\phi) = \frac{1}{N} \sum_i \left\| \hat{Q}_{\phi}^{\pi}(\mathbf{s}_i) - y_i \right\|^2$
 - 5: evaluate $\hat{A}^{\pi}(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}'_i) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_i)$
 - 6: $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_i^{\pi} | \mathbf{s}_i) \hat{Q}^{\pi}(\mathbf{s}_i, \mathbf{a}_i^{\pi})$
 - 7: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$
 - 8: **until** convergence
-