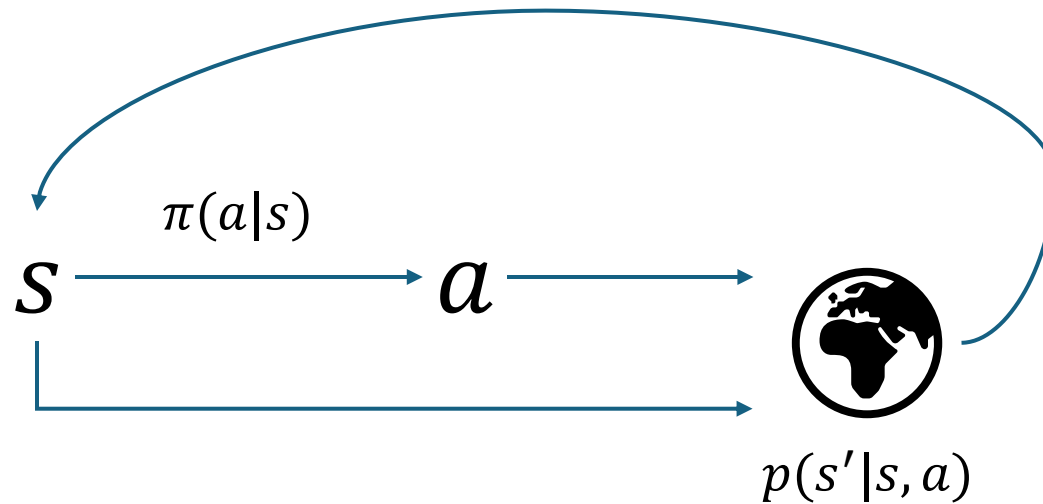


# Policy Gradient Methods

# What is the goal of Reinforcement Learning?



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

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

# Evaluating the Objective

$$\theta^* = \operatorname{argmax}_{\theta} \underbrace{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(s_t, a_t) \right]}_{J(\theta)}$$

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(s_t, a_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(s_{i,t}, a_{i,t})$$

# Direct Policy Differentiation

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

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

Reminder of key identity:

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

# Direct Policy Differentiation

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

Take Log of both sides

$$\log p_{\theta}(\tau) = \log p(s_1) + \underbrace{\sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)}_{}$$

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

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

# Evaluating the policy gradient

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

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

REINFORCE algorithm

1. Sample  $\{\tau^i\}$  using policy  $\pi_{\theta}(a_t | s_t)$  (Run the policy in the environment)
2.  $\nabla_{\theta} J(\theta) \approx \sum_i (\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)) (\sum_t r(s_t^i, a_t^i))$
3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

# What are we doing?

$$\nabla_{\theta} J_{ML}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau_i) r(\tau_i)$$

Make the good trajectories *more likely*

Make the bad trajectories *less likely*

We are formalizing the notion of “*trial and error*”

# Reducing variance

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

*Causality:* policy at time  $t'$  cannot affect reward at time  $t$  when  $t < t'$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \underbrace{\left( \sum_{t=t'}^T r(s_{i,t}, a_{i,t}) \right)}_{\text{"Reward to go"}}$$

## Baselines

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

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



# Actor-Critic Methods

$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta}[r(s_{t'}, a_{t'}) | s_t, a_t]$ : Total reward from taking  $a_t$  in  $s_t$

$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi_\theta(a_t | s_t)}[Q^\pi(s_t, a_t)]$ : Total reward from  $s_t$

$A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$ : How much better  $a_t$  is

We can fit  $Q^\pi$  and  $V^\pi$  or just  $V^\pi$

## Actor-Critic Algorithm

1. Sample  $\{\tau^i\}$  using policy  $\pi_\theta(a_t | s_t)$  (Run the policy in the environment)
2. Fit  $\hat{V}_\phi^\pi(s)$  to sampled reward sums
3. Evaluate  $\hat{A}^\pi(s_i, a_i) = r(s_i, a_i) + \hat{V}_\phi^\pi(s'_i) - \hat{V}_\phi^\pi(s_i)$
4.  $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(a_i | s_i) \hat{A}^\pi(s_i, a_i)$
5.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

# Advanced Policy Gradients

$$\theta' \leftarrow \operatorname{argmax}_{\theta'} \sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} \left[ \mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} \left[ \frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \right] \gamma^t A^{\pi_{\theta}}(s_t|a_t) \right]$$

*such that*  $|\pi_{\theta'}(a_t|s_t) - \pi_{\theta}(a_t|s_t)| \leq \epsilon$

For small enough  $\epsilon$ , this is guaranteed to improve  $J(\theta') - J(\theta)$