

# 11-695: Competitive Engineering Reinforcement Learning

Spring 2018

## ① Settings

Notations

Trajectory

Reward

Formulation of RL Objective

## ② Policy Gradients

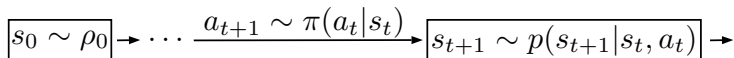
REINFORCE

Actor-Critic

---

Name	Notation	Meaning
Environment	$\mathcal{S}$	Set of states
	$\mathcal{A}$	Set of actions
	$p(s' s, a)$	State transition probabilities
	$s_0 \sim \rho_0$	Initial state and initial distribution
	$r(s, a)$	Reward function
	$\gamma \in [0, 1]$	Discount factor
Policy	$\pi(a \in \mathcal{A} s)$	Distribution of actions given a state

---



## Trajectory

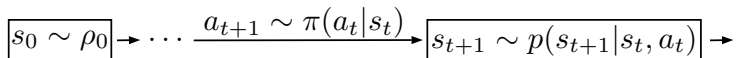
Infinite horizon trajectory  $\tau \sim \pi$

$\tau = (s_0, a_0, s_1, a_1, \dots)$ , where:

$$s_0 \sim \rho_0$$

$$a_t \sim \pi(s_t)$$

$$s_{t+1} \sim p(\cdot | s_t, a_t)$$



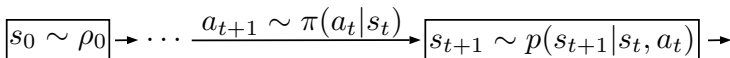
## Trajectory

Infinite horizon trajectory  $\tau \sim \pi$

$\tau = (s_0, a_0, s_1, a_1, \dots)$ , where:

- $s_0 \sim \rho_0$
- $a_t \sim \pi(s_t)$
- $s_{t+1} \sim p(\cdot | s_t, a_t)$

**From the probability perspective:**  $\tau$  is a random variable.



## Trajectory

Infinite horizon trajectory  $\tau \sim \pi$

$\tau = (s_0, a_0, s_1, a_1, \dots)$ , where:

- $s_0 \sim \rho_0$
- $a_t \sim \pi(s_t)$
- $s_{t+1} \sim p(\cdot | s_t, a_t)$

**From the probability perspective:**  $\tau$  is a random variable.

**Overloading notations:** We will use  $s_t(\tau)$  and  $a_t(\tau)$  to denote the  $t^{\text{th}}$  state of trajectory  $\tau$  and the  $t^{\text{th}}$  action of trajectory  $\tau$ , respectively. In that case,  $s_t(\tau)$  and  $a_t(\tau)$  serve as random variables, just like  $\tau$

## Probability and Likelihood of a Trajectory

Given *environment* and *policy*  $\pi$ , a trajectory  $\tau = (s_0, a_0, s_1, a_1, \dots)$  has probability

$$p_{\pi}(\tau) = p(s_0) \cdot \prod_{t=0}^{\infty} (\pi(a_t | s_t) \cdot p(s_{t+1} | s_t, a_t)) \quad (1)$$

Equivalently,  $\tau$  has the likelihood

$$\log p_{\pi}(\tau) = \log p(s_0) + \sum_{t=0}^{\infty} (\log \pi(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)) \quad (2)$$

- In theory, a trajectory  $\tau = (s_0, a_0, s_1, a_1, \dots)$  can last forever
- In practice, we can only simulate finitely many steps of  $\tau$

$$\tau_T(s_0, a_0, s_1, a_1, \dots) = (s_0, a_0, s_1, a_1, \dots, s_T, a_T) \quad (3)$$

- Probability

$$p_\pi(\tau_T) = p(s_0) \cdot \prod_{t=0}^{T-1} (\pi(a_t | s_t) \cdot p(s_{t+1} | s_t, a_t)) \quad (4)$$

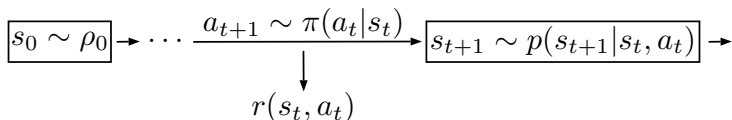
- Likelihood

$$\log p_\pi(\tau_T) = \log p(s_0) + \sum_{t=0}^{T-1} (\log \pi(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)) \quad (5)$$

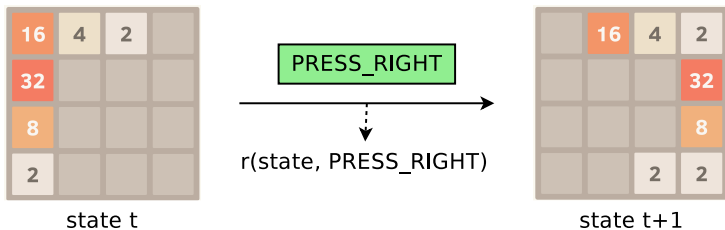


## Reward

Instant reward	$r(s, a)$	Received upon taking action $a$ at state $s$
Discount factor	$\gamma$	Further rewards count less
Discounted reward at $T$	$R_T(\tau)$	$\sum_{t=0}^{\infty} \gamma^t r(s_{T+t}, a_{T+t})$



# Reward: An Example



## Objective

Find a policy  $\pi^*$  that maximizes the expected discounted reward of all trajectories

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (6)$$

where

$$R_0(\tau = (s_0, a_0, s_1, a_1, \dots)) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \quad (7)$$

Three forms of RL expected discounted reward

Trajectory form	Discounted reward form
$\mathbb{E}_{\tau \sim \pi} [R_0(\tau)]$	$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$

## ① Settings

Notations

Trajectory

Reward

Formulation of RL Objective

## ② Policy Gradients

REINFORCE

Actor-Critic

REINFORCE equation; episodic version [Williams, 1992]

Let  $\pi(a|s; \theta)$  be a parameterized policy and let  $J(\theta)$  be the RL objective

$$J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (8)$$

Then one can compute the gradient  $\nabla_{\theta} J$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau)] \quad (9)$$

REINFORCE equation; episodic version [Williams, 1992]

Let  $\pi(a|s; \theta)$  be a parameterized policy and let  $J(\theta)$  be the RL objective

$$J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (8)$$

Then one can compute the gradient  $\nabla_{\theta} J$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau)] \quad (9)$$

**Why *episodic*?** Because it has to see the whole *episode*, aka, *trajectory*  $\tau$

$$\nabla_{\theta} \log p_{\pi}(\tau) = \nabla_{\theta} \left[ \log \rho_0(s_0) + \sum_{t=0}^{\infty} (\log p(s_{t+1}|s_t, a_t) + \log \pi(a_t|s_t; \theta)) \right] \quad (10)$$

$$= \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi(a_t|s_t; \theta) \quad (11)$$

$$\nabla_{\theta} J(\theta) \tag{12}$$

$$\tag{13}$$

$$\tag{14}$$

$$\tag{15}$$

$$\tag{16}$$

$$\tag{17}$$



$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \tag{12}$$

$$\tag{13}$$

$$\tag{14}$$

$$\tag{15}$$

$$\tag{16}$$

$$\tag{17}$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (12)$$

$$= \nabla_{\theta} \sum_{\tau} R_0(\tau) p_{\pi}(\tau) \quad (13)$$

$$(14)$$

$$(15)$$

$$(16)$$

$$(17)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (12)$$

$$= \nabla_{\theta} \sum_{\tau} R_0(\tau) p_{\pi}(\tau) \quad (13)$$

$$= \sum_{\tau} R_0(\tau) \nabla_{\theta} p_{\pi}(\tau) \quad (14)$$

$$(15)$$

$$(16)$$

$$(17)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (12)$$

$$= \nabla_{\theta} \sum_{\tau} R_0(\tau) p_{\pi}(\tau) \quad (13)$$

$$= \sum_{\tau} R_0(\tau) \nabla_{\theta} p_{\pi}(\tau) \quad (14)$$

$$= \sum_{\tau} R_0(\tau) \cdot \frac{\nabla_{\theta} p_{\pi}(\tau)}{p_{\pi}(\tau)} \cdot p_{\pi}(\tau) \quad (15)$$

$$(16)$$

$$(17)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (12)$$

$$= \nabla_{\theta} \sum_{\tau} R_0(\tau) p_{\pi}(\tau) \quad (13)$$

$$= \sum_{\tau} R_0(\tau) \nabla_{\theta} p_{\pi}(\tau) \quad (14)$$

$$= \sum_{\tau} R_0(\tau) \cdot \frac{\nabla_{\theta} p_{\pi}(\tau)}{p_{\pi}(\tau)} \cdot p_{\pi}(\tau) \quad (15)$$

$$= \sum_{\tau} R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau) \cdot p_{\pi}(\tau) \quad (16)$$

$$(17)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (12)$$

$$= \nabla_{\theta} \sum_{\tau} R_0(\tau) p_{\pi}(\tau) \quad (13)$$

$$= \sum_{\tau} R_0(\tau) \nabla_{\theta} p_{\pi}(\tau) \quad (14)$$

$$= \sum_{\tau} R_0(\tau) \cdot \frac{\nabla_{\theta} p_{\pi}(\tau)}{p_{\pi}(\tau)} \cdot p_{\pi}(\tau) \quad (15)$$

$$= \sum_{\tau} R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau) \cdot p_{\pi}(\tau) \quad (16)$$

$$= \mathbb{E}_{\tau \sim \pi} [R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau)] \quad (17)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (12)$$

$$= \nabla_{\theta} \sum_{\tau} R_0(\tau) p_{\pi}(\tau) \quad (13)$$

$$= \sum_{\tau} R_0(\tau) \nabla_{\theta} p_{\pi}(\tau) \quad (14)$$

$$= \sum_{\tau} R_0(\tau) \cdot \frac{\nabla_{\theta} p_{\pi}(\tau)}{p_{\pi}(\tau)} \cdot p_{\pi}(\tau) \quad (15)$$

$$= \sum_{\tau} R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau) \cdot p_{\pi}(\tau) \quad (16)$$

$$= \mathbb{E}_{\tau \sim \pi} [R_0(\tau) \cdot \nabla_{\theta} \log p_{\pi}(\tau)] \quad (17)$$

REINFORCE doesn't care about **the reward function  $R_0(\tau)$** !

$$\nabla_{\theta} \log p_{\pi}(\tau) \tag{18}$$

$$\tag{19}$$
$$\tag{20}$$
$$\tag{21}$$
$$\tag{22}$$



$$\nabla_{\theta} \log p_{\pi}(\tau) = \nabla_{\theta} \log \left( p(s_0) \prod_{t=0}^{\infty} \pi(a_t | s_t; \theta) \cdot p(s_{t+1} | s_t, a_t) \right) \quad (18)$$

(19)

(20)

(21)

(22)

$$\nabla_{\theta} \log p_{\pi}(\tau) = \nabla_{\theta} \log \left( p(s_0) \prod_{t=0}^{\infty} \pi(a_t | s_t; \theta) \cdot p(s_{t+1} | s_t, a_t) \right) \quad (18)$$

$$= \nabla_{\theta} \left[ \log p(s_0) + \sum_{t=0}^{\infty} \left( \log \pi(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t) \right) \right] \quad (19)$$

(20)

(21)

(22)

$$\nabla_{\theta} \log p_{\pi}(\tau) = \nabla_{\theta} \log \left( p(s_0) \prod_{t=0}^{\infty} \pi(a_t | s_t; \theta) \cdot p(s_{t+1} | s_t, a_t) \right) \quad (18)$$

$$= \nabla_{\theta} \left[ \log p(s_0) + \sum_{t=0}^{\infty} \left( \log \pi(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t) \right) \right] \quad (19)$$

$$= \nabla_{\theta} \log p(s_0) + \sum_{t=0}^{\infty} \left( \nabla_{\theta} \log \pi(a_t | s_t; \theta) + \nabla_{\theta} \log p(s_{t+1} | s_t, a_t) \right) \quad (20)$$

$$(21)$$

$$(22)$$

$$\nabla_{\theta} \log p_{\pi}(\tau) = \nabla_{\theta} \log \left( p(s_0) \prod_{t=0}^{\infty} \pi(a_t | s_t; \theta) \cdot p(s_{t+1} | s_t, a_t) \right) \quad (18)$$

$$= \nabla_{\theta} \left[ \log p(s_0) + \sum_{t=0}^{\infty} \left( \log \pi(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t) \right) \right] \quad (19)$$

$$= \nabla_{\theta} \log p(s_0) + \sum_{t=0}^{\infty} \left( \nabla_{\theta} \log \pi(a_t | s_t; \theta) + \nabla_{\theta} \log p(s_{t+1} | s_t, a_t) \right) \quad (20)$$

$$= \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (21)$$

$$(22)$$

$$\nabla_{\theta} \log p_{\pi}(\tau) = \nabla_{\theta} \log \left( p(s_0) \prod_{t=0}^{\infty} \pi(a_t | s_t; \theta) \cdot p(s_{t+1} | s_t, a_t) \right) \quad (18)$$

$$= \nabla_{\theta} \left[ \log p(s_0) + \sum_{t=0}^{\infty} \left( \log \pi(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t) \right) \right] \quad (19)$$

$$= \nabla_{\theta} \log p(s_0) + \sum_{t=0}^{\infty} \left( \nabla_{\theta} \log \pi(a_t | s_t; \theta) + \nabla_{\theta} \log p(s_{t+1} | s_t, a_t) \right) \quad (20)$$

$$= \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (21)$$

$$(22)$$

REINFORCE doesn't care about **the environment!**

**Temporal Structure:** Stuff that involves multiple steps ( $t$  in equations)

**Temporal Structure:** Stuff that involves multiple steps ( $t$  in equations)

We have the following equation (proof skipped for now)

$$\nabla_{\theta} \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (23)$$

**Temporal Structure:** Stuff that involves multiple steps ( $t$  in equations)

We have the following equation (proof skipped for now)

$$\nabla_{\theta} \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (23)$$



Put together so far

$$\nabla_{\theta} \mathbb{E}_{\tau \sim \pi} [R_0(\tau)] \quad (24)$$

$$= \mathbb{E}_{\tau \sim \pi} \left[ R_0(\tau) \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (25)$$

$$= \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (26)$$

## Monte Carlo Episodic Policy Gradient

REINFORCE (Eqn 26) induces a method to estimate policy gradient

$$\nabla_{\theta} \hat{J}(\theta) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (27)$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{\infty} \gamma^t R_t(\tau_i) \nabla_{\theta} \log \pi(a_t(\tau_i) | s_t(\tau_i); \theta) \quad (28)$$

Eqn 28, albeit useful, has **high variance** and **high sample complexity**.

To reduce the high variance in the previous slide.

(29)

To reduce the high variance in the previous slide.

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (29)$$

To reduce the high variance in the previous slide.

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (29)$$

**Proof.** skip for now :(

To reduce the high variance in the previous slide.

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (29)$$

**Proof.** skip for now :(

Therefore

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{\infty} \gamma^t (R_t(\tau_i) - b(s_t(\tau_i))) \nabla_{\theta} \log \pi(a_t(\tau_i) | s_t(\tau_i); \theta) \quad (30)$$

## REINFORCE algorithm

- 1 Initialize  $\pi(a|s; \theta)$
- 2 Repeat until converge
  - Sample  $N$  trajectories  $\tau_i \sim \pi$
  - Compute  $\nabla_{\theta} \hat{J}(\theta)$

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{\infty} \gamma^t (R_t(\tau_i) - b(s_t(\tau_i))) \nabla_{\theta} \log \pi(a_t(\tau_i) | s_t(\tau_i); \theta)$$

- Update  $b(s_t(\tau_i))$  using  $\ell_2$  loss  $\|b(s_t(\tau_i)) - R_t(\tau_i)\|^2$  or moving averages
- Update  $\pi(a|s; \theta)$

---

Name	Notation	Meaning
Value	$V_\pi(s)$	Expected discounted reward from state $s$ $V_\pi(s) = \mathbb{E}_{a \sim \pi(a s)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim p(s' s, a)} [V_\pi(s')] \right]$
$Q$ -value	$Q_\pi(s, a)$	Expected discounted reward if take action $a$ at state $s$ $Q_\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s' s, a)} [V_\pi(s')]$
Advantage	$A_\pi(s, a)$	How better is the action than average $A_\pi(s, a) = Q_\pi(s, a) - V_\pi(s)$

---



Recall from REINFORCE. For each trajectory  $\tau$ , we do this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (31)$$

Recall from REINFORCE. For each trajectory  $\tau$ , we do this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (31)$$

Problems:

Recall from REINFORCE. For each trajectory  $\tau$ , we do this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (31)$$

Problems:

- The  $\infty$ . No update until  $\tau$  finishes.

Recall from REINFORCE. For each trajectory  $\tau$ , we do this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (31)$$

Problems:

- The  $\infty$ . No update until  $\tau$  finishes.
- The rumors have it: [high variance](#).

Recall from REINFORCE. For each trajectory  $\tau$ , we do this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (31)$$

Problems:

- The  $\infty$ . No update until  $\tau$  finishes.
- The rumors have it: [high variance](#).
- The rumors have it: [high sample complexity](#).

A solution: fix this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (32)$$

(33)

A solution: fix this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (32)$$

$$\approx \sum_{t=0}^T \gamma^t Q_{\pi}(s_t, a_t) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (33)$$

A solution: fix this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (32)$$

$$\approx \sum_{t=0}^T \gamma^t Q_{\pi}(s_t, a_t) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (33)$$

**How does this help?**  $Q_{\pi}(s_t, a_t)$  does **not** need  $\tau$  to finish.



A solution: fix this

$$\sum_{t=0}^{\infty} \gamma^t R_t(\tau) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (32)$$

$$\approx \sum_{t=0}^T \gamma^t Q_{\pi}(s_t, a_t) \nabla_{\theta} \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (33)$$

**How does this help?**  $Q_{\pi}(s_t, a_t)$  does **not** need  $\tau$  to finish.

→ Can update  $\theta$  after simulating  $\tau$  for  $T$  steps.

Something's wrong! How do we know  $Q_\pi(s_t, a_t)$ ?

**Something's wrong!** How do we know  $Q_\pi(s_t, a_t)$ ?

We can approximate  $Q_\pi(s, a)$  with

$$\hat{Q}(s, a; \theta_Q) \approx Q_\pi(s, a) \quad (34)$$

**Something's wrong!** How do we know  $Q_\pi(s_t, a_t)$ ?

We can approximate  $Q_\pi(s, a)$  with

$$\hat{Q}(s, a; \theta_Q) \approx Q_\pi(s, a) \quad (34)$$

**Plug it in**

$$\sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_\theta \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (35)$$

**Something's wrong!** How do we know  $Q_\pi(s_t, a_t)$ ?

We can approximate  $Q_\pi(s, a)$  with

$$\hat{Q}(s, a; \theta_Q) \approx Q_\pi(s, a) \quad (34)$$

Plug it in

$$\sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_\theta \log \pi(a_t(\tau) | s_t(\tau); \theta) \quad (35)$$

**How do we estimate  $\theta_Q$ ?** Simulate  $\tau$  and use  $\ell_2$  error;  $Q$ -learning, etc.

## Actor-Critic Algorithm

- 1 Initialize  $\pi(a, s; \theta)$  and  $Q(s, a; \theta_Q)$
- 2 Repeat until convergence
  - Sample  $N$  trajectories  $\tau_i$ , each for  $T$  steps
  - Compute

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (36)$$

- Update  $\theta$  using  $\nabla_{\theta} \hat{J}(\theta)$
- Update  $\theta_Q$  using  $\sum_{i=1}^N \sum_{t=0}^T \|Q(s_t, a_t; \theta_Q) - R_t(\tau_i)\|^2$

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Actor-Critic estimates  $\nabla_{\theta}$  for each trajectory  $\tau$  as follows

$$\sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (37)$$

Problems:

- ~~The  $\infty$ . No update until  $\tau$  finishes.~~
- The rumors have it: **high variance.**
- The rumors have it: **high sample complexity.**

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Solution: back to the baseline form of REINFORCE

(38)

(39)



# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Solution: back to the baseline form of REINFORCE

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (38)$$

(39)

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Solution: back to the baseline form of REINFORCE

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (38)$$

$$\approx \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t (\hat{R}_t - V_{\pi}(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (39)$$

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Solution: back to the baseline form of REINFORCE

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (38)$$

$$\approx \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t (\hat{R}_t - V_{\pi}(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (39)$$

Same as Actor-Critic, we will parameterize something. This time

$$\hat{V}(s; \theta_V) \approx V_{\pi}(s) \quad (40)$$

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Solution: back to the baseline form of REINFORCE

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (38)$$

$$\approx \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t (\hat{R}_t - V_{\pi}(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (39)$$

Same as Actor-Critic, we will parameterize something. This time

$$\hat{V}(s; \theta_V) \approx V_{\pi}(s) \quad (40)$$

With  $V_{\pi}(s)$  known, one can estimate  $\hat{R}_t$  the following quantity with  $n$ -steps sampling from  $t$

$$\hat{R}_t = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^{n-1} r(s_{t+n-1}, a_{t+n-1}) + \gamma^{t+n} V(s_{t+n}; \theta_V) \quad (41)$$

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

---

Solution: back to the baseline form of REINFORCE

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (38)$$

$$\approx \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t (\hat{R}_t - V_{\pi}(s_t)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \right] \quad (39)$$

Same as Actor-Critic, we will parameterize something. This time

$$\hat{V}(s; \theta_V) \approx V_{\pi}(s) \quad (40)$$

With  $V_{\pi}(s)$  known, one can estimate  $\hat{R}_t$  the following quantity with  $n$ -steps sampling from  $t$

$$\hat{R}_t = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^{n-1} r(s_{t+n-1}, a_{t+n-1}) + \gamma^{t+n} V(s_{t+n}; \theta_V) \quad (41)$$

**How do we train  $\theta_V$ ?** Small  $\ell_2$  estimate error in Eqn 41 to  $V(s_t; \theta_V)$ .

# Advantage Actor-Critic [Mnih et al., 2016]

Carnegie Mellon

## Advantage Actor-Critic Algorithm

- 1 Initialize  $\pi(a, s; \theta)$  and  $V(s, a; \theta_V)$
- 2 Repeat until convergence
  - Sample  $N$  trajectories  $\tau_i$
  - For each state  $s_t$  in any trajectory  $\tau_i$ , estimate using  $n$ -step sampling

$$\hat{R}_t(\tau_i) = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^{n-1} r(s_{t+n-1}, a_{t+n-1}) + \gamma^{t+n} V(s_{t+n}; \theta_V) \quad (42)$$

- Compute

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t (\hat{R}_t - V(s; \theta_V)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (43)$$

- Update  $\theta$  using  $\nabla_{\theta} \hat{J}(\theta)$
- Update  $\theta_V$  using  $\sum_{i=1}^N \sum_{t=0}^T \left\| \hat{R}_t(\tau_i) - V(s_t(\tau_i); \theta_V) \right\|^2$

Actor-Critic [Sutton et al., 1999]

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (44)$$

$$\mathcal{L}(\theta_Q) = \sum_{i=1}^N \sum_{t=0}^T \|Q(s_t, a_t; \theta_Q) - R_t(\tau_i)\|^2 \quad (45)$$

**Actor-Critic** [Sutton et al., 1999]

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (44)$$

$$\mathcal{L}(\theta_Q) = \sum_{i=1}^N \sum_{t=0}^T \|Q(s_t, a_t; \theta_Q) - R_t(\tau_i)\|^2 \quad (45)$$

**Advantage Actor-Critic** [Mnih et al., 2016]

$$\hat{R}_t(\tau_i) = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^{n-1} r(s_{t+n-1}, a_{t+n-1}) + \gamma^{t+n} V(s_{t+n}; \theta_V) \quad (46)$$

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t (\hat{R}_t(\tau_i) - V(s; \theta_V)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (47)$$

$$\mathcal{L}(\theta_V) = \sum_{i=1}^N \sum_{t=0}^T \|\hat{R}_t(\tau_i) - V(s_t(\tau_i); \theta_V)\|^2 \quad (48)$$



**Actor-Critic** [Sutton et al., 1999]

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t \hat{Q}(s_t, a_t; \theta_Q) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (44)$$

$$\mathcal{L}(\theta_Q) = \sum_{i=1}^N \sum_{t=0}^T \|Q(s_t, a_t; \theta_Q) - R_t(\tau_i)\|^2 \quad (45)$$

**Advantage Actor-Critic** [Mnih et al., 2016]

$$\hat{R}_t(\tau_i) = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^{n-1} r(s_{t+n-1}, a_{t+n-1}) + \gamma^{t+n} V(s_{t+n}; \theta_V) \quad (46)$$

$$\nabla_{\theta} \hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t (\hat{R}_t(\tau_i) - V(s; \theta_V)) \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (47)$$

$$\mathcal{L}(\theta_V) = \sum_{i=1}^N \sum_{t=0}^T \|\hat{R}_t(\tau_i) - V(s_t(\tau_i); \theta_V)\|^2 \quad (48)$$

	Name	Meaning	Role
<b>Etymology:</b>	Actor	The policy $\pi$	Samples trajectories
	Critic	$Q(s, a; \theta_Q)$ or $V(s; \theta_V)$	Scales the temporal gradient $\nabla_{\theta} \log \pi(a_t   s_t; \theta)$
	Advantage	$\hat{R}_t(\tau) - V(s; \theta_V)$	Estimate $A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(s)$

Volodymyr Mnih, Adri Puigdomnech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, 2016.

Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, 1999.

Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 1992.