#### Deep Reinforcement Learning (II)

Policy Gradient Methods

SUN Yinghan 2022. 03. 15

## Recap: Markov Decision Process



Markov property: Current state completely characterizes the state of the world.

A Markov decision process is defined by a tuple:  $\langle S, A, \mathcal{R}, \mathbb{P}, \gamma \rangle$ .

- S: set of possible states
- $\mathcal{A}$ : set of possible actions
- $\mathcal{R}$ : distribution of reward
- P: transition probability
- $\gamma$ : discount factor

Reward = -1 for all transitions

$$p(s', r|s, a) = \Pr\{S_t = s', R_t = r|S_{t-1} = s, A_{t-1} = a\}$$

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$
  
=  $R_{t+1} + \gamma \left( R_{t+2} + \gamma R_{t+3} + \gamma^{2} R_{t+4} + \cdots \right)$   
=  $R_{t+1} + \gamma G_{t+1}$ 

# Recap: Value-Based Approach

**Policy:** ways of acting. It is a mapping from states to probabilities of selecting each possible action.

**Value Function:** the expected return when starting in *s* and following  $\pi$  thereafter.

State-value function for policy  $\pi$ :  $v_{\pi}(s) = \mathbb{E}_{\pi}[G_t|S_t = s]$ Action-value function for policy  $\pi$ :  $q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t|S_t = s, A_t = a]$ 

#### **Relationship:**

$$v_{\pi}(s) = \sum_{a} \pi(a|s)q_{\pi}(s,a)$$
$$q_{\pi}(s,a) = \sum_{s'} \sum_{r} p(s',r|s,a)[r + \gamma v_{\pi}(s')]$$

**Objective:** find the optimal policy, which maximizes cumulative discounted reward (value function).

**Q-Learning:** a temporal-difference learning method.  $Q(s,a) \leftarrow Q(s,a) + \alpha[R + \gamma \max_{a} Q(s',a) - Q(s,a)]$ 

What if there are large number of states, or even infinitely many states? DQN

## Policy-Based Methods

**RL Goal:** Learning a policy to maximize the expectation of the accumulate returns (rewards) when an agent interacts with the environment.

$$\theta^* = \arg \max_{\theta} \underbrace{\mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(s_t, a_t) \right]}_{J(\theta)}$$

**Value-Based Methods:** learning action-value firstly, then select actions based on the learned action-value (policy).

**Policy-Based Methods:** compute the gradient of  $J(\theta)$  w.r.t.  $\theta$  directly, then update the parameter  $\theta$  of the policy according to the gradient.

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(s_{t}, a_{t}) \right] = \int p_{\theta}(\tau) r(\tau) d\tau$$
$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau$$
$$= \int p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} r(\tau) d\tau$$
$$= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) d\tau$$
$$= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) \right]$$

### Compute the Gradient

**Question:** How to compute  $\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) \right]$ ?

$$p_{\theta}(\tau) = p_{\theta}(s_{1}, a_{1}, s_{2}, a_{2}, s_{3}, a_{3}, \cdots, s_{T}, a_{T})$$

$$= p(s_{1})p(a_{1}|s_{1})p(s_{2}|s_{1}, a_{1})p(a_{2}|s_{2})p(s_{3}|s_{1}, a_{1}, s_{2}, a_{2})p(a_{3}|s_{3}) \cdots$$

$$= p(s_{1})\pi(a_{1}|s_{1})p(s_{2}|s_{1}, a_{1})\pi(a_{2}|s_{2})p(s_{3}|s_{2}, a_{2})\pi(a_{3}|s_{3}) \cdots$$

$$= p(s_{1})\prod_{t=1}^{T} [\pi_{\theta}(a_{t}|s_{t})p(s_{t+1}|s_{t}, a_{t})]$$

$$\log p_{\theta}(\tau) = \log p(s_{1}) + \sum_{t=1}^{T} [\log \pi_{\theta}(a_{t}|s_{t}) + \log p(s_{t+1}|s_{t}, a_{t})]$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} \log p(s_{1}) + \sum_{t=1}^{T} [\nabla_{\theta} \log \pi_{\theta}(a_{t}|s_{t}) + \nabla_{\theta} \log p(s_{t+1}|s_{t}, a_{t})]$$

$$= \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(a_{t}|s_{t})$$

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

**Monte-Carlo Sampling:**  $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left[ \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(a_{i,t}|s_{i,t}) \right) \left( \sum_{t=1}^{T} r(s_{i,t}|a_{i,t}) \right) \right]$ 

# Algorithm: REINFORCE

1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(a_t|s_t)$ . (i.e. run the policy)

2. compute gradient of the cost function.

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

3. update policy parameters:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ 

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

We need to choose some proper policies with parameter  $\theta$ , or use a neural network to approximate this mapping. We need to design some proper reward functions.

### Example: Gaussian Policies (Continuous Case)



$$\pi_{\theta}(a_t|s_t) = \mathcal{N}(g_{\theta}(s_t), \sigma^2)$$
$$= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\frac{(a_t - g_{\theta}(s_t))^2}{\sigma^2}\right\}$$
$$\log \pi_{\theta}(a_t|s_t) = -\frac{1}{2\sigma^2}\left(a_t - g_{\theta}(s_t)\right)^2 + \log\frac{1}{\sqrt{2\pi\sigma}}$$
$$\nabla_{\theta}\log \pi_{\theta}(a_t|s_t) = \frac{1}{\sigma^2}\left(a_t - g_{\theta}(s_t)\right)\frac{\partial g_{\theta}(s_t)}{\partial \theta}$$

# Summary: RL Algorithms (So Far)



# Case Study: Application on Robot Locomotion



Hwangbo J, Lee J, Dosovitskiy A, et al. Learning Agile and Dynamic Motor Skills for Legged Robots[J]. Science Robotics, 2019, 4(26): eaau5872.

# Case Study: Application on Robot Locomotion

#### Base Velocity Tracking Performance of the Learned Controller While Following Random Commands



(A) Forward velocity, (B) Lateral velocity, (C) yaw rate. For all graphs, the dotted lines represent the commanded velocity and the solid lines represent the measured velocity. All commands are followed with a reasonable accuracy even when the commands are given in a random fashion.

Hwangbo J, Lee J, Dosovitskiy A, et al. Learning Agile and Dynamic Motor Skills for Legged Robots[J]. Science Robotics, 2019, 4(26): eaau5872.



[1] Sutton R S, Barto A G. Reinforcement Learning: An Introduction[M]. MIT press, 2018.

[2] Williams R J. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning[J]. Machine learning, 1992, 8(3): 229-256.

[3] Schulman J, Levine S, Abbeel P, et al. Trust Region Policy Optimization[C]//International Conference on Machine Learning. PMLR, 2015: 1889-1897.

[4] Schulman J, Wolski F, Dhariwal P, et al. Proximal Policy Optimization Algorithms[J]. arXiv preprint arXiv:1707.06347, 2017.

[5] Hwangbo J, Lee J, Dosovitskiy A, et al. Learning Agile and Dynamic Motor Skills for Legged Robots[J]. Science Robotics, 2019, 4(26): eaau5872.

# **Thanks for Listening**