

# **ECE/SIOC 228 Machine Learning for Physical Applications**

## **Lecture 14: Deep Reinforcement Learning I**

**Yuanyuan Shi**

Assistant Professor, ECE

University of California, San Diego

# Course Structure

UC San Diego

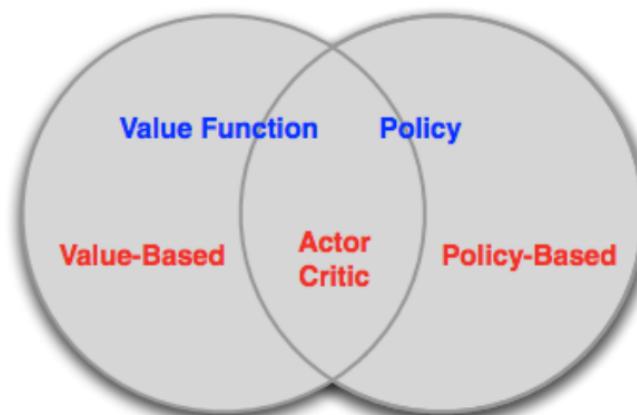
- **Planning** in Markov Decision Process (offline)
  - Have the environment model (reward, transition)
  - Find the optimal policy via value iteration and policy iteration

- **Reinforcement Learning** in Markov Decision Process (online)
  - Don't know how the environment works
  - Find the optimal policy by interacting with the environment (taking actions and collect reward)

# Reinforcement Learning Algorithms

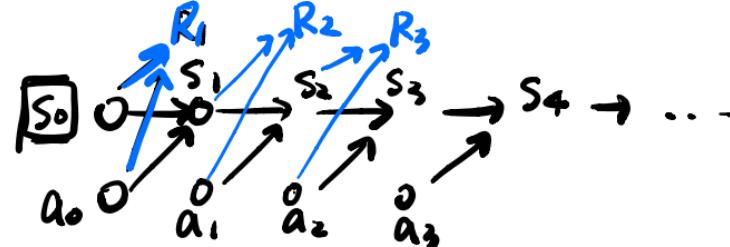
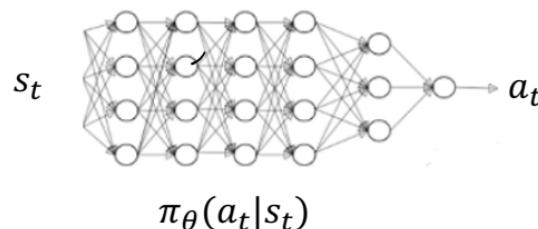
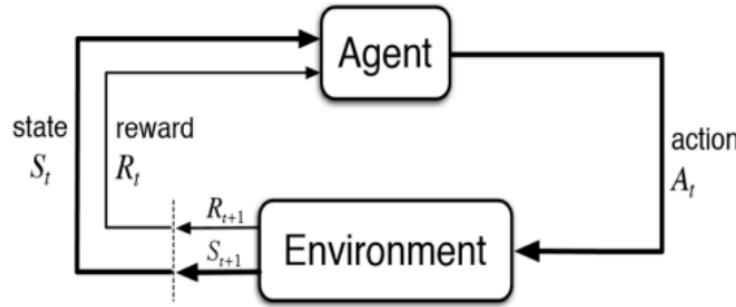
UC San Diego

- **Policy Based Reinforcement Learning**
  - Learn a policy
  - No value function / not necessarily have a value function
- Value Based Reinforcement Learning
  - Learn value function
  - Policy is implicit (e.g., greedy)
- Actor-Critic
  - Learn value function
  - Learn policy



# Policy Objective Functions

UC San Diego



Goal: given policy  $\pi_\theta(s,a)$ , with parameters  $\theta$ , find best  $\theta$

- How do we measure the Quality of policy? — **value function**:

- Episodic task, initial state  $s_0 \sim p_0$

$$J(\theta) = \mathbb{E}_{s_0 \sim p_0} [V^{\pi_\theta}(s_0)]$$

- Continuing task:  $J(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$

$d^{\pi_\theta}(s)$  is stationary distribution of Markov Chain for  $\lambda_\theta$

\* Policy based RL is an optimization problem

\* Find  $\theta^*$  that maximizes  $J(\theta)$

\* Use Gradient Ascent.

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_{\theta} J(\theta)$$

# Policy Gradient

UC San Diego

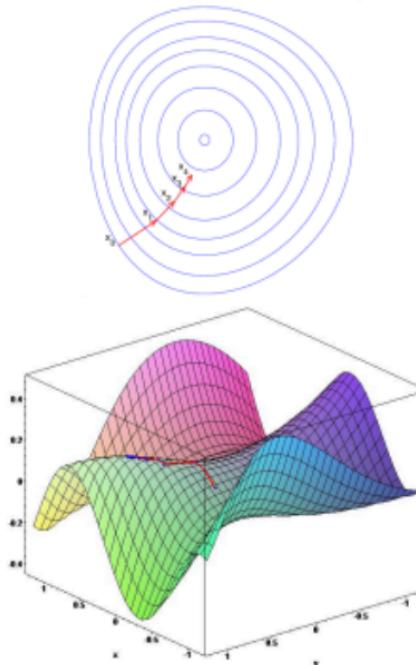
- Let  $J(\theta)$  be any policy objective function
- Policy gradient algorithms search for a *local* maximum in  $J(\theta)$  by ascending the gradient of the policy, w.r.t. parameters  $\theta$

$$\Delta\theta = \alpha \nabla_{\theta} J(\theta)$$

- Where  $\nabla_{\theta} J(\theta)$  is the **policy gradient**

$$\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix}$$

- and  $\alpha$  is a step-size parameter



# Computing Gradients By Finite Differences

UC San Diego

- To evaluate policy gradient of  $\pi_\theta(s, a)$
- For each dimension  $k \in [1, n]$ 
  - Estimate  $k$ th partial derivative of objective function w.r.t.  $\theta$
  - By perturbing  $\theta$  by small amount  $\epsilon$  in  $k$ th dimension

$$\theta \in \mathbb{R}^N, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \xrightarrow{\theta} \underbrace{\begin{bmatrix} \theta_1 + \epsilon u_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}}_{\theta + \epsilon u_k}$$

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon u_k) - J(\theta)}{\epsilon}$$

where  $u_k$  is unit vector with 1 in  $k$ th component, 0 elsewhere

- Uses  $n$  evaluations to compute policy gradient in  $n$  dimensions
- Simple, noisy, inefficient - but sometimes effective
- Works for arbitrary policies, even if policy is not differentiable

- We now compute the policy gradient analytically
- Assume policy  $\pi_\theta(s, a)$  is differentiable

Episodic task:  $J(\theta) = E_{S_0 \sim p_0} [V^{\pi_\theta}(S_0)]$

$$= E_{S_0 \sim p_0} [E_{\pi_\theta} [G_0 | S_0]] = E_{z \sim p_\theta(z)} [\underbrace{\sum_{t=0}^{T-1} \gamma^t R(S_t, A_t)}_{\gamma(z)}]$$

$\gamma(z)$  is the return (total discounted reward) of an episode

$$z \doteq S_0, A_0, S_1, A_1, S_2, A_2, \dots, S_{T-1}, A_{T-1}, S_T$$

## Direct Policy Differentiation

UC San Diego

$$P_{\theta}(z) = P_{\theta}(s_0) \prod_{t=0}^{T-1} \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

Find the best  $\theta$  that optimize  $J(\theta)$ :

$$\theta^* = \underset{\theta}{\operatorname{argmax}} J(\theta) = \underset{\theta}{\operatorname{argmax}} E_{z \sim P_{\theta}(z)}[r(z)]$$

Policy gradient

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} E_z[r(z)] = \underline{\nabla_{\theta}} \int P_{\theta}(z) r(z) dz = \underline{\int} \nabla_{\theta} P_{\theta}(z) r(z) dz \\ &= \int P_{\theta}(z) \frac{\nabla_{\theta} P_{\theta}(z)}{P_{\theta}(z)} r(z) dz = \int P_{\theta}(z) \nabla_{\theta} \log P_{\theta}(z) r(z) dz \\ &= E_{z \sim P_{\theta}(z)} [\nabla_{\theta} \log P_{\theta}(z) r(z)] \end{aligned}$$

Leibniz integral rule

## Direct Policy Differentiation

UC San Diego

(From previous page)  $\nabla_{\theta} J(\theta) = E_{z \sim p_{\theta}(z)} \left[ \frac{\nabla_{\theta} \log p_{\theta}(z)}{?} \gamma(z) \right]$

$$p_{\theta}(z) = p_{\theta}(s_0) \prod_{t=0}^{T-1} \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$\log p_{\theta}(z) = \log p_{\theta}(s_0) + \sum_{t=0}^{T-1} \log \pi_{\theta}(a_t | s_t) + \sum_{t=0}^{T-1} \log p(s_{t+1} | s_t, a_t)$$

$$\nabla_{\theta} \log p_{\theta}(z) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

$$\Rightarrow \nabla_{\theta} J(\theta) = E_{z \sim p_{\theta}(z)} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \gamma(z) \right]$$

Return of  
episode z

- This result is elegant, since to compute  $\nabla_{\theta} J(\theta)$ , we don't need to know the transition matrix  $P(S_{t+1}|S_t, A_t)$ .  $\Rightarrow$  "model-free"
- Evaluating the expectation  $E_{z \sim p_{\theta}(z)}[\dots]$  still remains a challenge  $\Rightarrow$  use sampled episodes

# Direct Policy Differentiation

UC San Diego

$$\nabla_{\theta} J(\theta) = E_{z \sim p_{\theta}(z)} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r(z) \right]$$

$\approx \frac{1}{N} \sum_{i=1}^N \underbrace{\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)}_{\text{Do log } p_{\theta}(z)} \left( \underbrace{\sum_{t=0}^{T-1} Y^t r(s_t^i, a_t^i)}_{\text{"return"-weighted likelihood}} \right)$

- high return episode is given a larger weight  $\Rightarrow$  become more likely
- low return episode is given a smaller weight  $\Rightarrow$  become less likely

