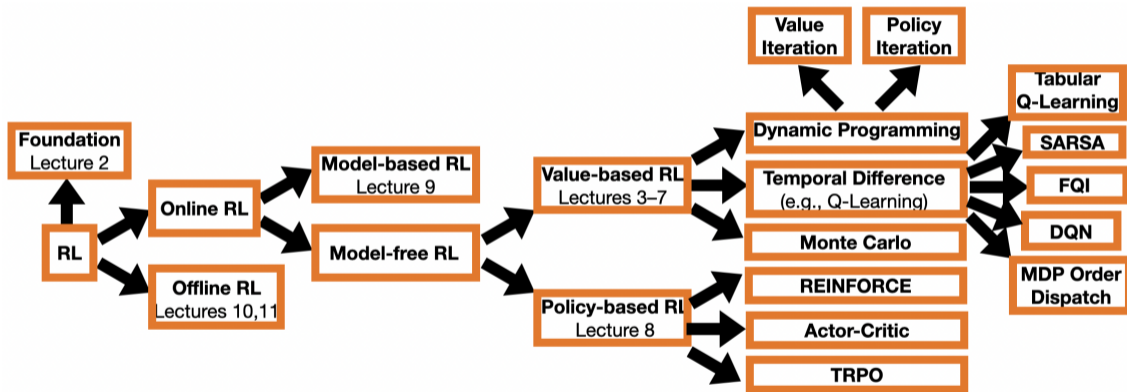


# **ST455: Reinforcement Learning**

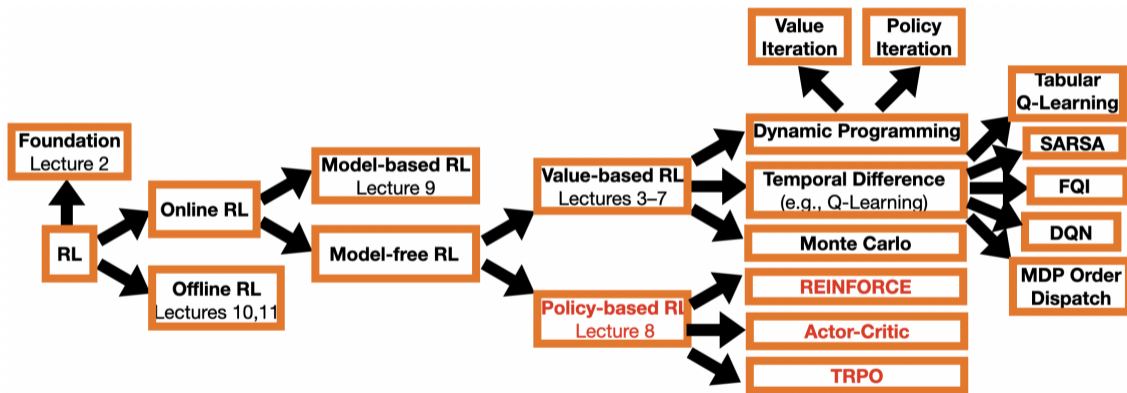
## **Lecture 8: Policy-based Learning**

Chengchun Shi

# Roadmap



# Roadmap (Cont'd)



# Lecture Outline

---

1. Introduction to Policy-based Learning
2. Policy Gradient Theorem
3. REINFORCE and Actor Critic Algorithms
4. Advantage Actor-Critic (A2C)
5. Trust Region Policy Optimization (TRPO)

# Lecture Outline

---

- 1. Introduction to Policy-based Learning**
2. Policy Gradient Theorem
3. REINFORCE and Actor Critic Algorithms
4. Advantage Actor-Critic (A2C)
5. Trust Region Policy Optimization (TRPO)

# Policy We Studied So Far

---

- Greedy policy:

$$\pi^{\text{opt}}(\mathbf{s}) = \arg \max_{\mathbf{a}} Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})$$

- $\epsilon$ -Greedy policy:

$$\begin{cases} \pi^{\text{opt}}(\mathbf{s}), & \text{with probability } \mathbf{1} - \epsilon \\ \text{random action,} & \text{with probability } \epsilon. \end{cases}$$

- **Value-based methods:** Policy Iteration, Value Iteration, SARSA, Q-Learning, etc.

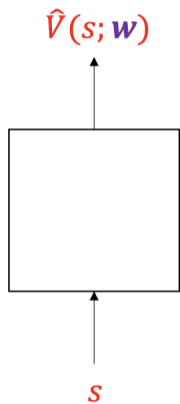
# Value-based v.s. Policy-based Methods

---

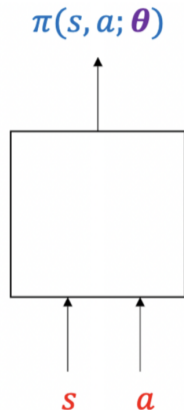
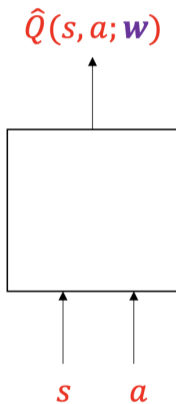
- **Value-based methods:** derive  $\pi^{\text{opt}}$  by learning an optimal Q-function (with or without function approximation)
- **Policy-based methods:** search  $\pi^{\text{opt}}$  within a restricted function class (e.g., linear, neural networks) that maximizes the value

# Value-based v.s. Policy-based Methods (Cont'd)

---



**Value-based Methods**



**Policy-based Methods**



# Example: Linear Function Approximation

---

- Linear approximation of features  $\phi(\mathbf{s}, \mathbf{a})$
- State-action value function approximation

$$Q(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta}) = \phi^\top(\mathbf{s}, \mathbf{a})\boldsymbol{\theta}$$

- Policy function approximation

$$\pi(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta}) = \frac{\exp(\phi^\top(\mathbf{s}, \mathbf{a})\boldsymbol{\theta})}{\sum_{\mathbf{a}'} \exp(\phi^\top(\mathbf{s}, \mathbf{a}')\boldsymbol{\theta})}$$

$\phi^\top(\mathbf{s}, \mathbf{a})\boldsymbol{\theta}$  similar to the preference score in the gradient based methods in HW1

# Value-based v.s. Policy Gradient Methods (Cont'd)

---

- **Pros** of policy gradient methods:
  1. Suitable for learning general **stochastic** policies (value-based methods mainly designed for deterministic policies)
  2. More **robust** to model misspecification
  3. Scalable for **high-dimensional** or **continuous** action spaces (SARSA, Q-learning mainly designed for discrete action space)
- **Cons** of policy gradient methods:
  1. Convergence to local minima
  2. May have large variance

# Example I: Advantage of Stochastic Policy

---



- Two-player game of rock-paper-scissors
  - Scissors beats paper
  - Rock beats scissors
  - Paper beats rock
- Consider iterated rock-paper-scissors
  - A deterministic policy is easily exploited
  - A uniform random policy is optimal (Nash equilibrium)

## Example II: Robustness of Policy-based Method

---

- Q-function is more **difficult** to model compared to the optimal policy
- Example: optimal Q-function:  $Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a}) = g(\phi^\top(\mathbf{s}, \mathbf{a})\theta^*)$  for some monotonically increasing function  $g: \mathbb{R} \rightarrow \mathbb{R}$
- When  $g$  is not a **linear** function, value-based method misspecifies Q-function model

$$g(\phi^\top(\mathbf{s}, \mathbf{a})\theta^*) \neq \phi^\top(\mathbf{s}, \mathbf{a})\theta$$

- However, since  $g$  is a monotonically increasing function

$$\pi^{\text{opt}}(\mathbf{s}) = \arg \max_{\mathbf{a}} g(\phi^\top(\mathbf{s}, \mathbf{a})\theta^*) = \arg \max_{\mathbf{a}} \phi^\top(\mathbf{s}, \mathbf{a})\theta^*$$

- Policy gradient methods correctly identifies the optimal policy

$$\frac{\exp(\phi^\top(\mathbf{s}, \mathbf{a})\theta)}{\sum_{\mathbf{a}'} \exp(\phi^\top(\mathbf{s}, \mathbf{a}')\theta)} \rightarrow \mathbb{I}(\mathbf{a} = \pi^{\text{opt}}(\mathbf{s}))$$

when  $\theta = k\theta^*$  and  $k \rightarrow \infty$

# Policy Objective Functions

---

- Average rewards:

$$J(\theta) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}^{\pi(\bullet; \theta)} \left[ \sum_{t=0}^{T-1} R_t \right] = \sum_{s, a} \nu^{\pi(\bullet; \theta)}(s) \pi(s, a; \theta) \mathcal{R}_s^a$$

where  $\mathcal{R}_s^a = \mathbb{E}(R_t | A_t = a, S_t = s)$

- For each  $\pi$ , the states  $\{S_t\}_t$  forms a time-homogeneous Markov chain
- $\nu^{\pi(\bullet; \theta)}$  the stationary distribution of  $\{S_t\}_t$  under  $\pi(\bullet; \theta)$

# Policy Objective Functions (Cont'd)

---

- Discounted rewards: given a discounted factor  $\gamma \in [0, 1]$  and initial state distribution  $\nu$ , maximize the expected discounted rewards:

$$J(\theta) = \mathbb{E}^{\pi(\cdot; \theta)} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right],$$

or equivalently,

$$J(\theta) = \sum_{\mathbf{s}} \nu(\mathbf{s}) V^{\pi(\cdot; \theta)}(\mathbf{s})$$

- If  $\gamma = 1$ , the task is assumed to be episodic

# Lecture Outline

---

1. Introduction to Policy-based Learning
- 2. Policy Gradient Theorem**
3. REINFORCE and Actor Critic Algorithms
4. Advantage Actor-Critic (A2C)
5. Trust Region Policy Optimization (TRPO)

# Policy Gradient

---

- **Objective:** identify the maximizer of  $J(\theta)$
- **Method:** apply (stochastic) gradient ascent algorithm to update  $\theta$  (gradient descent to minimize  $-J(\theta)$ )

$$\theta_{t+1} = \theta_t + \alpha_t \nabla_{\theta} J(\theta_t)$$

Need to calculate the gradient  $\nabla_{\theta} J(\theta)$ !



# Policy Gradient Theorem

## Theorem

For any differentiable policy  $\pi(\mathbf{s}, \mathbf{a}; \theta)$  with respect to parameter  $\theta$ , the policy gradient for average reward and discounted expected rewards objective is

$$\nabla_{\theta} J(\theta) = \sum_{\mathbf{s}, \mathbf{a}} \mu^{\pi(\bullet; \theta)}(\mathbf{s}, \mathbf{a}) \nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta)) Q^{\pi(\bullet; \theta)}(\mathbf{s}, \mathbf{a})$$

- For average reward objective:  
 $\mu^{\pi(\bullet; \theta)}$  is the stationary distribution of  $\{(\mathbf{S}_t, \mathbf{A}_t)\}_t$  under  $\pi(\bullet; \theta)$
- For discounted expected rewards objective:

$$\mu^{\pi(\bullet; \theta)}(\mathbf{s}, \mathbf{a}) = \sum_{t \geq 0} \gamma^t \Pr^{\pi(\bullet; \theta)}(\mathbf{S}_t = \mathbf{s}, \mathbf{A}_t = \mathbf{a})$$

Discounted state-action visitation probability

# Policy Gradient Theorem (Cont'd)

## Theorem

For any differentiable policy  $\pi(\mathbf{s}, \mathbf{a}; \theta)$  with respect to parameter  $\theta$ , the policy gradient for average reward and discounted expected rewards objective is

$$\nabla_{\theta} J(\theta) = \sum_{\mathbf{s}, \mathbf{a}} \mu^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a}) \nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta)) Q^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a})$$

- For average reward objective:

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) = \mathbb{E}^{\pi} \left[ \sum_{t \geq 0} (R_t - J(\theta)) \mid \mathbf{S}_0 = \mathbf{s}, \mathbf{A}_0 = \mathbf{a} \right]$$

- For discounted expected rewards objective: Q-function defined as usual.
- Proof given in the appendix

# Policy Score

---

- For any state-action pair  $(\mathbf{s}, \mathbf{a})$ , the term

$$\nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta))$$

is referred as the **policy score**

# Example 1: Softmax Policy Gradient

---

- State-action pairs weighted by linear combination of features

$$\pi(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\phi}^\top(\mathbf{s}, \mathbf{a})\boldsymbol{\theta})}{\sum_{\mathbf{a}'} \exp(\boldsymbol{\phi}^\top(\mathbf{s}, \mathbf{a}')\boldsymbol{\theta})}$$

- The score function

$$\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta}) = \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}) - \frac{\sum_{\mathbf{a}'} \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}') \exp(\boldsymbol{\phi}^\top(\mathbf{s}, \mathbf{a}')\boldsymbol{\theta})}{\sum_{\mathbf{a}'} \exp(\boldsymbol{\phi}^\top(\mathbf{s}, \mathbf{a}')\boldsymbol{\theta})}$$

or equivalently,

$$\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta}) = \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{s}, \bullet; \boldsymbol{\theta})} \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}')$$

## Example 2: Continuous Action Space

---

- Action space: set of real numbers  $\mathcal{A} = \mathbb{R}$
- Policy approximator:

$$\pi(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi}\sigma(\mathbf{s}; \boldsymbol{\theta})} \exp\left(-\frac{(\mathbf{a} - \boldsymbol{\mu}(\mathbf{s}; \boldsymbol{\theta}))^2}{2\sigma^2(\mathbf{s}; \boldsymbol{\theta})}\right),$$

where  $\boldsymbol{\mu}$  and  $\sigma$  are mean and deviation function approximators

- Linear function approximator with feature vectors  $\boldsymbol{\phi}_\mu(\mathbf{s})$  and  $\boldsymbol{\phi}_\sigma(\mathbf{s})$ 
  - $\boldsymbol{\mu}(\mathbf{s}; \boldsymbol{\theta}) = \boldsymbol{\phi}_\mu^\top(\mathbf{s})\boldsymbol{\theta}_\mu$  and  $\sigma(\mathbf{s}; \boldsymbol{\theta}) = \boldsymbol{\phi}_\sigma^\top(\mathbf{s})\boldsymbol{\theta}_\sigma$
  - $\nabla_{\boldsymbol{\theta}_\mu} \log \pi(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) = \frac{\mathbf{a} - \boldsymbol{\mu}(\mathbf{s}; \boldsymbol{\theta})}{\sigma^2(\mathbf{s}; \boldsymbol{\theta})} \boldsymbol{\phi}_\mu(\mathbf{s})$
  - $\nabla_{\boldsymbol{\theta}_\sigma} \log \pi(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) = \frac{(\mathbf{a} - \boldsymbol{\mu}(\mathbf{s}; \boldsymbol{\theta}))^2 - \sigma^2(\mathbf{s}; \boldsymbol{\theta})}{\sigma^2(\mathbf{s}; \boldsymbol{\theta})} \boldsymbol{\phi}_\sigma(\mathbf{s})$

## Example 3: Bernoulli, Logistic Example

---

- Actions space: binary,  $\{0, 1\}$
- Policy approximator:

$$\pi(\mathbf{1}, \mathbf{s}; \theta) = 1 - \pi(\mathbf{0}, \mathbf{s}; \theta) = \mathbf{p}(\mathbf{s}; \theta)$$

where  $\mathbf{p}(\mathbf{s}; \theta)$  is a function approximator

- Linear function approximator with feature vectors  $\phi(\mathbf{s})$ 
  - Logistic function  $\sigma(\mathbf{x}) = [\mathbf{1} + \exp(-\mathbf{x})]^{-1}$
  - For exponential soft-max policy  $\mathbf{p}(\mathbf{s}; \theta) = \sigma(\phi^\top(\mathbf{s})\theta)$
  - $\nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta)) = (\mathbf{a} - \sigma(\phi^\top(\mathbf{s})\theta))\phi(\mathbf{s})$

# Lecture Outline

---

1. Introduction to Policy-based Learning
2. Policy Gradient Theorem
- 3. REINFORCE and Actor Critic Algorithms**
4. Advantage Actor-Critic (A2C)
5. Trust Region Policy Optimization (TRPO)

# REINFORCE: MC Policy Gradient Algorithm

---

- To maximize  $J(\theta)$ , we apply (stochastic) gradient ascent algorithm

$$\theta_{t+1} = \theta_t + \alpha_t \nabla_{\theta} J(\theta_t)$$

- According to the policy gradient theorem,

$$\nabla_{\theta} J(\theta) = \sum_{\mathbf{s}, \mathbf{a}} \mu^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a}) \nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta)) Q^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a})$$

- Focus on the average reward setting
- $\mu^{\pi}$  (stationary state-action distribution) is unknown: use empirical state-action distribution  $\{(\mathbf{S}_t, \mathbf{A}_t)\}_t$  as an approx
- $Q^{\pi}$  is unknown: use empirical return  $\mathbf{G}_t = \sum_{j=t}^T \mathbf{R}_j$  as an approx



# REINFORCE: Pseudocode

---

- **Initialization:**  $\theta$  arbitrary
- **For each** episode  $(S_0, A_0, R_0, \dots, S_T, A_T, R_T)$  generated using policy  $\pi(\bullet; \theta)$

**For**  $t = 0, 1, 2, \dots, T$  **do:**

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log(\pi(S_t, A_t; \theta)) G_t$$

**end for**

**return**  $\theta$

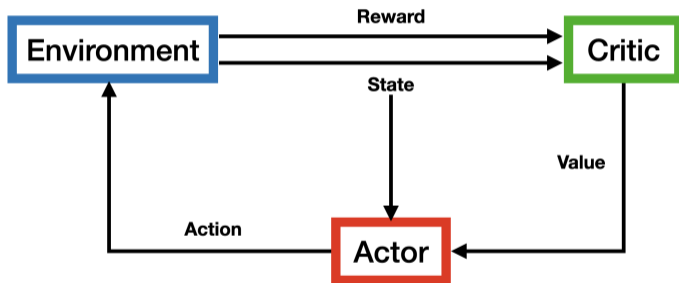
# Actor-Critic Algorithm

---

- MC policy gradient algorithm may have a large **variance**
  - Return involves many state transitions, many actions and many rewards
- Solution sought by using **actor-critic algorithms**
- Actor-critic algorithms combine **policy gradient** with **value function estimation**

# Actor-Critic Algorithm (Cont'd)

---



- **Critic** uses function approximator to learn value function
- **Actor** uses policy approximator to learn optimal policy

# Actor-Critic Control

---

- **Critic:** estimates  $Q^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a})$  by a function approximator  $\widehat{Q}(\mathbf{s}, \mathbf{a}; \omega)$ 
  - The critic performs **policy evaluation**
  - Standard methods can be applied: MC, TD(0), TD( $\lambda$ ), gradient-based methods
- **Actor:** updates policy parameter  $\theta$ 
  - The actor performs control using approximate policy gradient

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mu} \nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta)) \widehat{Q}(\mathbf{s}, \mathbf{a}; \omega)$$

- Parameter update
  - Average reward setting

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t; \theta)) \widehat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega)$$

- Discounted reward setting

$$\theta \leftarrow \theta + \alpha \gamma^t \nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t; \theta)) \widehat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega)$$

# Example: Actor-Critic with Linear Value Function

---

- Linear value function approximator

$$\hat{Q}(s, a; \omega) = \phi^\top(s, a)\omega$$

- Focus on the discounted reward setting
- **Critic:** updates  $\omega$  by linear TD(0)

$$\omega_{t+1} = \omega_t + \eta \phi(S_t, A_t) (R_t + \gamma \phi^\top(S_{t+1}, A_{t+1})\omega_t - \phi^\top(S_t, A_t)\omega_t)$$

# Pseudocode

---

- **Initialization:**  $s, \theta, \omega$

- **For each** episode:

**Initialize**  $t = 0$

Sample action  $a$  from  $\pi(\bullet, s; \theta)$

**Repeat** until  $s$  is terminal

Receive reward  $r$  and next state  $s'$

Sample action  $a'$  from  $\pi(\bullet, s; \theta)$

$$\theta \leftarrow \theta + \alpha \gamma^t \log(\pi(s, a; \theta)) \phi^\top(s, a) \omega$$

$$\omega \leftarrow \omega + \eta \phi(s, a) [r + \gamma \phi^\top(s', a') \omega - \phi^\top(s, a) \omega]$$

$a \leftarrow a'$  and  $s \leftarrow s'$

$t \leftarrow t + 1$

# Bias-Variance Tradeoff

---

- **REINFORCE** uses Return  $G_t$ , an unbiased estimate of  $Q^{\pi(\cdot;\theta)}(s, a)$
- **Actor-critic** uses  $\hat{Q}(s, a; \omega)$ , a biased estimate of  $Q^{\pi(\cdot;\theta)}(s, a)$
- REINFORCE gradient has **high variance** and **zero bias**
- Actor-critic gradient has **low variance** and **some bias**
- Similar to Pros & Cons of MC vs TD (Lecture 4, p13)
- Perhaps surprisingly, actor-critic gradient can be **unbiased** under certain conditions (see appendix)

# Lecture Outline

---

1. Introduction to Policy-based Learning
2. Policy Gradient Theorem
3. REINFORCE and Actor Critic Algorithms
- 4. Advantage Actor-Critic (A2C)**
5. Trust Region Policy Optimization (TRPO)



# Variance Reduction Using a Baseline

---

- Recall that policy parameter update

$$\theta \leftarrow \theta + \alpha \gamma^t \nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t; \theta)) \widehat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega)$$

- For any  $\theta$ , when  $\mathbf{A}_t \sim \pi(\mathbf{S}_t, \bullet, \theta)$

$$\mathbb{E}[\nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t, \theta)) | \mathbf{S}_t] = \mathbf{0}$$

- For any baseline function  $\mathbf{B}(s)$ , consider the update

$$\theta \leftarrow \theta + \alpha \gamma^t \nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t; \theta)) [\widehat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega) - \mathbf{B}(\mathbf{S}_t)]$$

- The **mean** of gradient is the same without baseline
- However, the **variance** of the gradient would be smaller with a properly chosen  $\mathbf{B}$

## Variance Reduction Using a Baseline (Cont'd)

---

- Consider the baseline that minimizes the variance of the gradient
- For any random variable  $\mathbf{Z}$ , the mean  $\mathbb{E}\mathbf{Z}$  minimizes  $\arg \min_z \mathbb{E}(\mathbf{Z} - z)^2$
- To minimize variance of the gradient  $\nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t; \theta))[\widehat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega) - \mathbf{B}(\mathbf{S}_t)]$ , the baseline is set to the conditional mean of Q-function given the state
- i.e.,  $\mathbf{B}(\mathbf{s}) = \sum_{\mathbf{a}} \pi(\mathbf{s}, \mathbf{a}; \theta) \widehat{Q}(\mathbf{s}, \mathbf{a}; \omega)$ , e.g., the estimated state-value
- Similar ideas have been employed in gradient-based algorithms in HW1

# Policy Gradient Using Advantage Function

---

- Advantage function:  $\mathbf{A}^{\pi(\cdot;\theta)}(\mathbf{s}, \mathbf{a}) = \mathbf{Q}^{\pi(\cdot;\theta)}(\mathbf{s}, \mathbf{a}) - \mathbf{V}^{\pi(\cdot;\theta)}(\mathbf{s})$
- Policy gradient based on advantage function

$$\nabla_{\theta} \mathbf{J}(\theta) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mu^{\pi(\cdot;\theta)}} \nabla_{\theta} \log(\pi(\mathbf{s}, \mathbf{a}; \theta)) \mathbf{A}^{\pi(\cdot;\theta)}(\mathbf{s}, \mathbf{a})$$

- The advantage function reduces the variance of policy gradient

# An Approach for Estimating Advantage Function

---

- The critic may compute estimators of both value functions

$$\hat{Q}(\mathbf{s}, \mathbf{a}; \omega) \text{ for } Q^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a})$$

and

$$\hat{V}(\mathbf{s}; \omega) \text{ for } V^{\pi(\cdot; \theta)}(\mathbf{s})$$

which can be done by standard methods such as TD learning

- The estimator of the advantage function

$$\hat{A}(\mathbf{s}, \mathbf{a}; \omega) = \hat{Q}(\mathbf{s}, \mathbf{a}; \omega) - \hat{V}(\mathbf{s}; \omega)$$

# Another Approach

---

- $r + \gamma V^{\pi(\cdot;\theta)}(s') - V^{\pi(\cdot;\theta)}(s)$  is unbiased to  $A^{\pi(\cdot;\theta)}(s, a)$

$$\begin{aligned} & \mathbb{E}[r + \gamma V^{\pi(\cdot;\theta)}(s') - V^{\pi(\cdot;\theta)}(s) | a, s] \\ &= \mathbb{E}[r + \gamma V^{\pi(\cdot;\theta)}(s') - Q^{\pi(\cdot;\theta)}(s, a) + Q^{\pi(\cdot;\theta)}(s, a) - V^{\pi(\cdot;\theta)}(s) | a, s] \\ &= Q^{\pi(\cdot;\theta)}(s, a) - V^{\pi(\cdot;\theta)}(s) = A^{\pi(\cdot;\theta)}(s, a) \end{aligned}$$

- As such,

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{(s,a) \sim \mu^{\pi(\cdot;\theta)}} \nabla_{\theta} \log(\pi(s, a; \theta)) [r + \gamma V^{\pi(\cdot;\theta)}(s') - V^{\pi(\cdot;\theta)}(s)]$$

- No need to estimate the advantage. It suffices to estimate the state-value and use the estimator to compute the policy gradient

# Critic Policy Evaluation Methods

---

- When specialized to linear methods  $\hat{V}(\mathbf{s}; \boldsymbol{\omega}) = \boldsymbol{\phi}^\top(\mathbf{s})\boldsymbol{\omega}$ , the critic can use different targets to evaluate

$$\boldsymbol{\omega}_{t+1} \leftarrow \boldsymbol{\omega}_t + \eta_t[\mathbf{v}_t - \boldsymbol{\phi}^\top(\mathbf{S}_t)\boldsymbol{\omega}_t]\boldsymbol{\phi}(\mathbf{S}_t)$$

- The target is defined differently for different methods
  - MC:  $\mathbf{v}_t = G_t$
  - TD:  $\mathbf{v}_t = R_t + \gamma \hat{V}(\mathbf{S}_{t+1})$
  - TD( $\lambda$ ):  $\mathbf{v}_t = G_t^\lambda$

# Actor Policy Gradient Methods

---

- The policy gradient

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{(s, a) \sim \mu^{\pi(\cdot; \theta)}} \nabla_{\theta} \log(\pi(s, a; \theta)) \mathbf{A}^{\pi(\cdot; \theta)}(s, a)$$

- Gradient-based method

$$\theta \leftarrow \theta + \alpha \gamma^t \nabla_{\theta} \log(\pi(\mathbf{S}_t, \mathbf{A}_t; \theta)) \widehat{\mathbf{A}}(\mathbf{S}_t, \mathbf{A}_t; \omega)$$

- Examples:

- MC:  $\widehat{\mathbf{A}}(\mathbf{S}_t, \mathbf{A}_t; \omega) = \mathbf{G}_t - \widehat{\mathbf{V}}(\mathbf{S}_t; \omega)$
- TD:  $\widehat{\mathbf{A}}(\mathbf{S}_t, \mathbf{A}_t; \omega) = \mathbf{R}_t + \gamma \widehat{\mathbf{V}}(\mathbf{S}_{t+1}; \omega) - \widehat{\mathbf{V}}(\mathbf{S}_t; \omega)$

# Lecture Outline

---

1. Introduction to Policy-based Learning
2. Policy Gradient Theorem
3. REINFORCE and Actor Critic Algorithms
4. Advantage Actor-Critic (A2C)
- 5. Trust Region Policy Optimization (TRPO)**



# TRPO: Introduction

---

- Limitations of policy gradient methods (REINFORCE & Actor Critic):
  - Convergence to local minima
  - Derivative-based, cannot use **non-smooth** policy classes, e.g.,  
 $\mathbf{a} = \arg \max_{\mathbf{a}'} \phi(\mathbf{s}, \mathbf{a}')^\top \boldsymbol{\theta}$  (used in application such as deep brain stimulation due to device constraints)
- Trust region policy optimization:
  - **Similarities** to policy gradient methods: an **iterative** algorithm
  - **Difference** from policy gradient methods: **derivative-free**

# TRPO: Foundation

## Theorem

For any two policies with parameters  $\theta_1$  and  $\theta_0$ ,

$$J(\theta_1) - J(\theta_0) = \sum_{\mathbf{s}, \mathbf{a}} \mathbf{A}^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \pi(\mathbf{s}, \mathbf{a}; \theta_1) \mu^{\pi(\cdot; \theta_1)}(\mathbf{s}, \mathbf{a})$$

- Considers discounted expected rewards objective

$$\mu^{\pi(\cdot; \theta)}(\mathbf{s}, \mathbf{a}) = \sum_{t \geq 0} \gamma^t \Pr^{\pi(\cdot; \theta)}(\mathbf{S}_t = \mathbf{s}, \mathbf{A}_t = \mathbf{a})$$

Discounted state-action visitation probability

- Proof given in Kakade and Langford [2002]

# TRPO: Challenge

## Theorem

For any two policies with parameters  $\theta_1$  and  $\theta_0$ ,

$$J(\theta_1) - J(\theta_0) = \sum_{\mathbf{s}, \mathbf{a}} \mathbf{A}^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \pi(\mathbf{s}, \mathbf{a}; \theta_1) \mu^{\pi(\cdot; \theta_1)}(\mathbf{s}, \mathbf{a})$$

- For a given  $\theta_0$ , suffices to search  $\theta_1$  that maximizes RHS
- Not feasible due to the complex dependence of  $\mu^{\pi(\cdot; \theta_1)}$  on  $\theta_1$
- Consider the following approximation with  $\mu^{\pi(\cdot; \theta_0)}$ :

$$\sum_{\mathbf{s}, \mathbf{a}} \mathbf{A}^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \pi(\mathbf{s}, \mathbf{a}; \theta_1) \mu^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a})$$

- The resulting maximizer is not guaranteed to improve the value function

# TRPO: Idea

---

$$\begin{aligned} J(\theta_1) - J(\theta_0) &= \sum_{\mathbf{s}, \mathbf{a}} \mathbf{A}^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \pi(\mathbf{s}, \mathbf{a}; \theta_1) \mu^{\pi(\cdot; \theta_1)}(\mathbf{s}, \mathbf{a}) \\ &= \sum_{\mathbf{s}, \mathbf{a}} \mathbf{A}^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \pi(\mathbf{s}, \mathbf{a}; \theta_1) \mu^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) + \text{Remainder} \end{aligned}$$

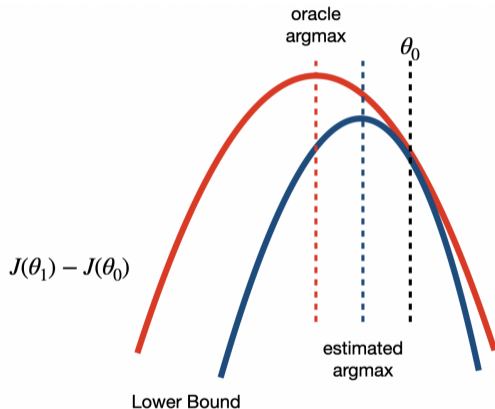
- The remainder term can be upper bounded by  $\mathbf{CKL}(\theta_0, \theta_1)$  for some  $\mathbf{C} > 0$ .
- Searching  $\theta_1$  that maximizes the leading term is not guaranteed to improve the value
- TRPO searches  $\theta_1$  that the lower bound: The leading term  $-\mathbf{CKL}(\theta_0, \theta_1)$
- The resulting maximizer satisfies

$$J(\theta_1) - J(\theta_0) \geq \text{leading}(\theta_0, \theta_1) - \mathbf{CKL}(\theta_0, \theta_1) \geq \text{leading}(\theta_0, \theta_0) - \mathbf{CKL}(\theta_0, \theta_0) = 0$$

leading to guaranteed monotonic improvement

# TRPO: Idea (Cont'd)

---



Idea similar to the MM algorithm for solving generic optimization problems

# TRPO: Implementation

---

- Given an initial  $\theta_0$
- For  $k = 1, \dots, K$ 
  1. Solve an optimization with a trust-region constraint:

$$\mathbf{A}^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \pi(\mathbf{s}, \mathbf{a}; \theta_1) \mu^{\pi(\cdot; \theta_0)}(\mathbf{s}, \mathbf{a}) \quad \text{subject to } \text{KL}(\theta_1, \theta_0) \leq \delta,$$

for some small  $\delta$

2. Set  $\theta_0$  to  $\theta_1$
- $\mathbf{A}^\pi$  and  $\mu^\pi$  can be similarly estimated as in actor-critic methods

# Summary

## Policy Function Approximation

No

Yes

Value Function Approximation

No

**Value-based  
(tabular)**

**REINFORCE**

Yes

**Value-based**

**Actor-Critic  
TRPO**

- **Value-based**

- Tabular (Lectures 3 & 4)
- Function approx (Lectures 5 & 7)

- **REINFORCE**

- No value function
- Learn policy

- **Actor-critic**

- Learn value
- Learn policy

- **Advantage actor-critic**

- Variance reduction

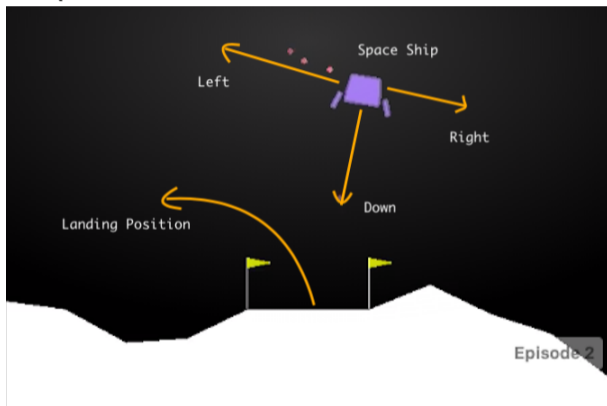
- **TRPO**

- Guaranteed monotonic improvement

# Seminar Exercise

---

- Solution to HW7 (Deadline: Wed 12pm)
- Implementation of DQN on LunarLander



Taken from <https://shiva-verma.medium.com/solving-lunar-lander-openaigym-reinforcement-learning-785675066197>



# References I

---

- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 267–274, 2002.
- Hamid Reza Maei, Csaba Szepesvári, Shalabh Bhatnagar, and Richard S Sutton. Toward off-policy learning control with function approximation. In *ICML*, 2010.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.

# Questions

# Appendix: Proof of Policy Gradient Theorem

---

- We focus on the discounted reward setting. Proofs in the average reward setting can be found in Sutton et al. [1999]
- Basic identities

$$(A) \quad \mathbf{V}^\pi(\mathbf{s}) = \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) Q^\pi(\mathbf{s}, \mathbf{a})$$

$$(B) \quad Q^\pi(\mathbf{s}, \mathbf{a}) = R_s^{\mathbf{a}} + \gamma \sum_{\mathbf{s}'} P_{\mathbf{s}, \mathbf{s}'}^{\mathbf{a}} V^\pi(\mathbf{s}')$$

$$(C) \quad \nabla_\theta \mathbf{V}^\pi(\mathbf{s}) = \sum_{\mathbf{a}} [\nabla_\theta \pi(\mathbf{a}|\mathbf{s})] Q^\pi(\mathbf{s}, \mathbf{a}) + \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) [\nabla_\theta Q^\pi(\mathbf{s}, \mathbf{a})]$$

$$(D) \quad \nabla_\theta Q^\pi(\mathbf{s}, \mathbf{a}) = \gamma \sum_{\mathbf{s}'} P_{\mathbf{s}, \mathbf{s}'}^{\mathbf{a}} \nabla_\theta V^\pi(\mathbf{s}')$$

## Appendix: Proof (Cont'd)

---

$$\begin{aligned}\nabla_{\theta} \mathbf{V}^{\pi}(\mathbf{s}) &\stackrel{(C)}{=} \sum_{\mathbf{a}} [\nabla_{\theta} \pi(\mathbf{a}|\mathbf{s})] \mathbf{Q}^{\pi}(\mathbf{s}, \mathbf{a}) + \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) [\nabla_{\theta} \mathbf{Q}^{\pi}(\mathbf{s}, \mathbf{a})] \\ &\stackrel{(D)}{=} \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) [\nabla_{\theta} \log(\pi(\mathbf{a}|\mathbf{s}))] \mathbf{Q}^{\pi}(\mathbf{s}, \mathbf{a}) + \underbrace{\gamma \sum_{\mathbf{a}, \mathbf{s}'} \pi(\mathbf{a}|\mathbf{s}) \mathbf{P}_{\mathbf{s}, \mathbf{s}'}^{\mathbf{a}} \nabla_{\theta} \mathbf{V}^{\pi}(\mathbf{s}')}_I\end{aligned}$$

Now, consider  $I$ . Similarly, we have

$$\begin{aligned}I &= \sum_{\mathbf{a}, \mathbf{s}', \mathbf{a}'} \pi(\mathbf{a}|\mathbf{s}) \mathbf{P}_{\mathbf{s}, \mathbf{s}'}^{\mathbf{a}} \pi(\mathbf{a}'|\mathbf{s}') [\nabla_{\theta} \log(\pi(\mathbf{a}'|\mathbf{s}'))] \mathbf{Q}^{\pi}(\mathbf{s}', \mathbf{a}') \\ &\quad + \gamma \sum_{\mathbf{a}, \mathbf{s}', \mathbf{a}', \mathbf{s}''} \pi(\mathbf{a}|\mathbf{s}) \mathbf{P}_{\mathbf{s}, \mathbf{s}'}^{\mathbf{a}} \pi(\mathbf{a}'|\mathbf{s}') \mathbf{P}_{\mathbf{s}', \mathbf{s}''}^{\mathbf{a}'} \nabla_{\theta} \mathbf{V}^{\pi}(\mathbf{s}'')\end{aligned}$$

## Appendix: Proof (Cont'd)

---

Recursively applying the first identity, we obtain

$$\nabla_{\theta} \mathbf{V}^{\pi}(s) = \mu^{\pi(\cdot; \theta)}(s', a'; s) \nabla_{\theta} \log(\pi(s', a')) \mathbf{Q}^{\pi}(s', a')$$

where

$$\mu^{\pi(\cdot; \theta)}(s', a'; s) = \sum_{t \geq 0} \gamma^t \pi(s', a') \Pr^{\pi(\cdot; \theta)}(S_t = s' | S_0 = s)$$

# Compatible Function Approximation Theorem

## Theorem

Assume the following two conditions:

(C1) Compatibility of value function approximator and the policy

$$\nabla_{\omega} \hat{Q}(s, \mathbf{a}; \omega) = \nabla_{\theta} \log \pi(s, \mathbf{a}; \theta)$$

(C2) Value function approximator minimizes the mean-squared error:

$$\mathbb{E}_{(s, \mathbf{a}) \sim \mu^{\pi(\cdot; \theta)}} \left[ Q^{\pi(\cdot; \omega)}(s, \mathbf{a}) - \hat{Q}(s, \mathbf{a}; \omega) \right]^2$$

Then the gradient is unbiased

$$\nabla_{\theta} \mathbf{J}(\theta) = \mathbb{E}_{(s, \mathbf{a}) \sim \mu^{\pi(\cdot; \theta)}} \nabla_{\theta} (\log \pi(s, \mathbf{a}; \theta)) \hat{Q}(s, \mathbf{a}; \omega)$$

# Compatible Linear Function Approximation

---

- Consider the soft-max policy, for a given state-action feature vector  $\phi(\mathbf{s}, \mathbf{a})$ :

$$\pi(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\phi}^\top(\mathbf{s}, \mathbf{a})\boldsymbol{\theta})}{\sum_{\mathbf{a}'} \exp(\boldsymbol{\phi}^\top(\mathbf{s}, \mathbf{a}')\boldsymbol{\theta})}$$

- Compatibility condition requires that

$$\nabla_{\boldsymbol{\omega}} \widehat{Q}(\mathbf{s}, \mathbf{a}'; \boldsymbol{\omega}) = \nabla_{\boldsymbol{\theta}} \log(\pi(\mathbf{s}, \mathbf{a}; \boldsymbol{\theta})) = \underbrace{\boldsymbol{\phi}(\mathbf{s}, \mathbf{a}) - \sum_{\mathbf{a}'} \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}')\pi(\mathbf{s}, \mathbf{a}'; \boldsymbol{\theta})}_{\text{centered state-action features vectors}}$$

which leads to a linear approximation for the value function

$$\widehat{Q}(\mathbf{s}, \mathbf{a}'; \boldsymbol{\omega}) = \left[ \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}) - \sum_{\mathbf{a}'} \boldsymbol{\phi}(\mathbf{s}, \mathbf{a}')\pi(\mathbf{s}, \mathbf{a}'; \boldsymbol{\theta}) \right]^\top \boldsymbol{\omega}$$

# Convergence Theorem

## Theorem

Assume

- $\pi(\bullet; \theta)$  and  $\hat{Q}(\bullet; \omega)$  are differentiable functions
- Compatibility assumption holds
- The Hessian matrix  $\nabla_{\theta}^2 \pi(\mathbf{s}, \mathbf{a}; \theta)$  are uniformly bounded away from infinity
- Step sizes are such that  $\sum_t \alpha_t = \infty$  and  $\sum_t \alpha_t^2 < \infty$
- At each step,  $\omega_t$  is chosen to be the solution of

$$\mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mu} \pi(\mathbf{s}, \mathbf{a}; \theta_t) [Q^{\pi(\bullet; \theta_t)}(\mathbf{s}, \mathbf{a}) - \hat{Q}(\mathbf{s}, \mathbf{a}; \omega)] \nabla_{\omega} \hat{Q}(\mathbf{s}, \mathbf{a}; \omega) = \mathbf{0}$$

Then  $\{\theta_t\}_t$  are convergent in the sense that  $\lim_{t \rightarrow \infty} \|\nabla_{\theta} J(\theta_t)\| \rightarrow \mathbf{0}$ .



# Separation of Timescales

---

- The last condition defines  $\omega_t$  as a solution of a fixed-point equation which has the policy's parameter vector  $\theta_t$  as a parameter
- In practice, we update  $\omega_t$  using stochastic gradient descent algorithm. SGD would update  $\omega_t$  in a similar manner with a larger step size than  $\alpha_t$ . It ensures  $\omega$  converges faster than  $\theta$ , thus closer to the solution of the fixed-point equation at each time
- This can be seen as a **separation of timescales**:
  - Critic updates the value function approximator at a **faster** timescale trying to evaluate the current policy chosen by the actor
  - Actor varies the policy's parameter more **slowly** to allow the critic to evaluate the current policy
- Similar assumptions are imposed in gradient Q-learning algorithms [Maei et al., 2010]