

# **Brain Inspired Artificial Intelligence**

## **4: Policy Search and Actor-Critic Methods**

Eiji Uchibe

Dept. of Brain Robot Interface

ATR Computational Neuroscience Labs.

# Outline

- Policy gradient approaches to MDP problems
  - Finite Differences Method
  - REINFORCE
  - Variance reduction
- Actor-critic method
  - Stationary distribution
  - Policy Gradient Theorem
- Natural gradient
  - Kullback Leibler divergence

# Policy Based Approaches

- In the last lecture, we represent the state value or state-action value function
- A policy is computed directly from the learned value function
  - e.g. using epsilon-greedy
- In this lecture, we will directly parameterize the policy

# Three Types of RL

critic only  
(value-based RL)

learn the value  
function  $q(s, a)$



$\pi(a | s)$  is derived  
from  $q(s, a)$

Q-learning,  
double Q-learning  
SARSA, DQN

actor-critic

learn  $q(s, a)$  and  
 $\pi(a | s)$



$\pi(a | s)$  is modified  
by  $q(s, a)$

A2C, A3C, DDPG,  
soft Actor-Critic

actor only  
(policy-based RL)

learn the policy  
 $\pi(a | s)$



$q(s, a)$  is not  
estimated

REINFORCE, PGPE

# Policy-Based Approach

- Suppose that a **finite-length** trajectory  $\tau$  is generated by some stochastic policy  $\pi(a \mid s, \boldsymbol{\theta})$

$$\tau = \{S_0, A_0, R_1, S_1, A_1, R_2, S_3, \dots, S_{T-1}, A_{T-1}, R_T, S_T\}$$

- The return of a trajectory  $\tau$  is

$$G(\tau) = \sum_{t=1}^T \gamma^{t-1} R_t \quad \text{random variable!}$$

- Reminder: Under the Markovian assumption, a probability to observe  $\tau$  is

$$\Pr(\tau \mid \boldsymbol{\theta}) = p_0(S_0) \prod_{t=0}^{T-1} \pi(A_t \mid S_t, \boldsymbol{\theta}) p(S_{t+1}, R_t \mid S_t, A_t)$$

# Policy-Based Approach

- The goal is to maximize the objective function

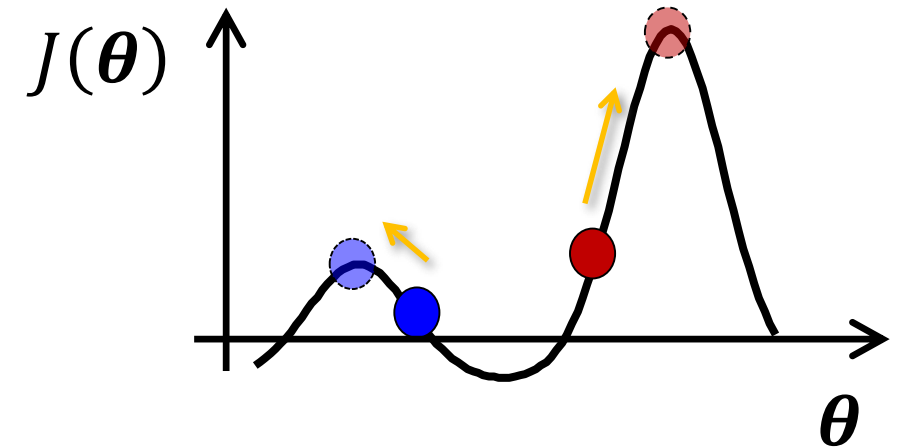
$$J(\boldsymbol{\theta}) = \mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})}[G(\tau)] = \sum_{\tau} \text{Pr}(\tau | \boldsymbol{\theta}) G(\tau)$$

- It is not feasible to compute  $\sum_{\tau}$  in practice because we have to consider all possible trajectories
- In order to maximize  $J(\boldsymbol{\theta})$ , a gradient method is applied

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

where  $\alpha$  is a positive step-size parameter

- Local optimal solution

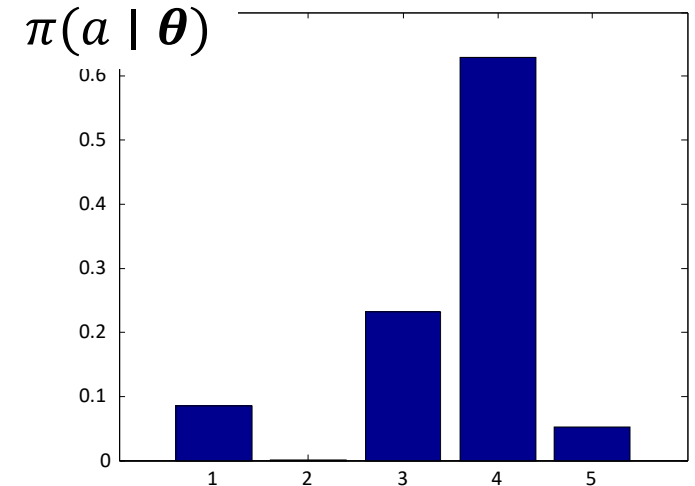
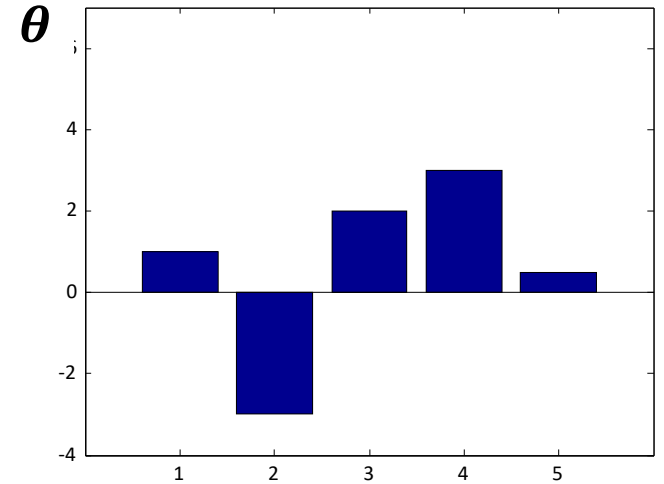


# Reminder: Stochastic Policy in the Bandit Problem

- We considered the policy represented by the Boltzmann distribution

$$\pi(a_i | \boldsymbol{\theta}) = \frac{\exp(\theta_i)}{\sum_{j=1}^k \exp(\theta_j)}$$

- $\theta_i$  is interpreted as a preference value of action  $a_i$
- We are going to consider  $\pi(a | s, \boldsymbol{\theta})$  that is a conditional probability on state  $s$



# Linear-Exponential Policy

- We introduce a state-action feature vector  $\boldsymbol{\phi}(s, a)$  to represent the preference for action  $a$  in state  $s$

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

- The derivative of the log policy is given by

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta}) &= \boldsymbol{\phi}(s, a) - \mathbb{E}_{\pi}[\boldsymbol{\phi}(s, a)] \\ &= \boldsymbol{\phi}(s, a) - \sum_{a'} \pi(a' \mid s, \boldsymbol{\theta}) \boldsymbol{\phi}(s, a') \end{aligned}$$

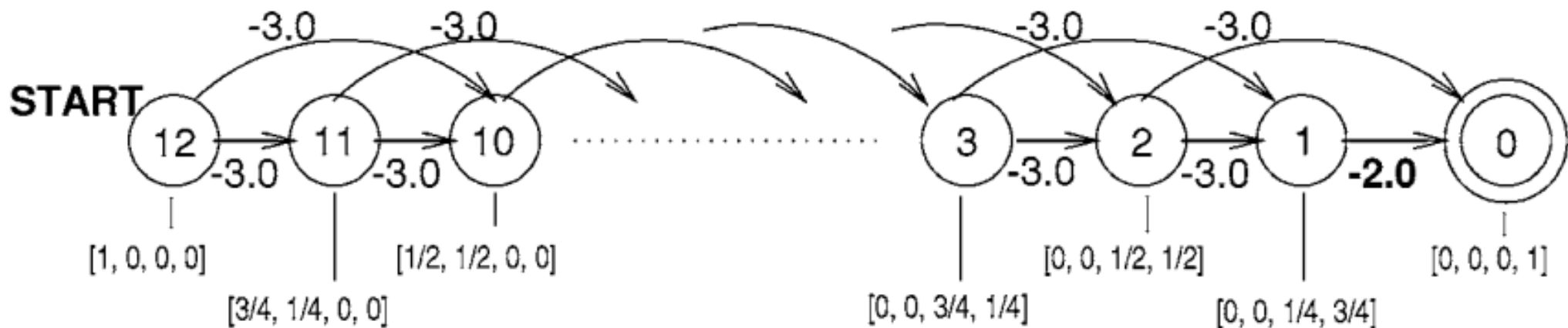


# Policy Parameterization

- State-dependent feature vector is widely used for simplicity

$$\pi(a | s, \theta) = \frac{\exp(\theta_a^\top \phi(s))}{\sum_{b \in \mathcal{A}} \exp(\theta_b^\top \phi(s))}$$

- Example of the feature vector  $\phi(s)$  [Boyan 2002]
  - 12 discrete states
  - 4 dimensional feature vector



# Policy Parameterization (Continuous Actions)

- Gaussian policy

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

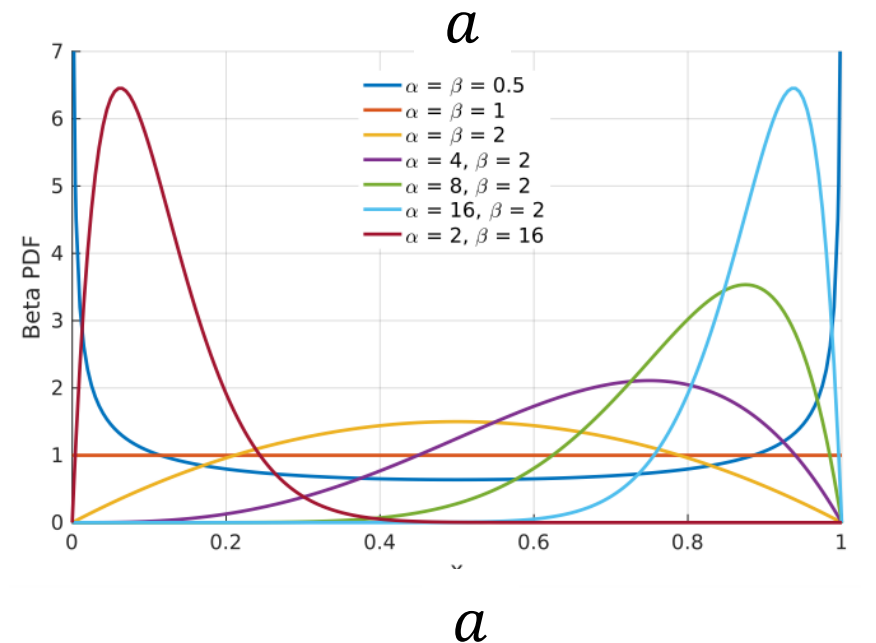
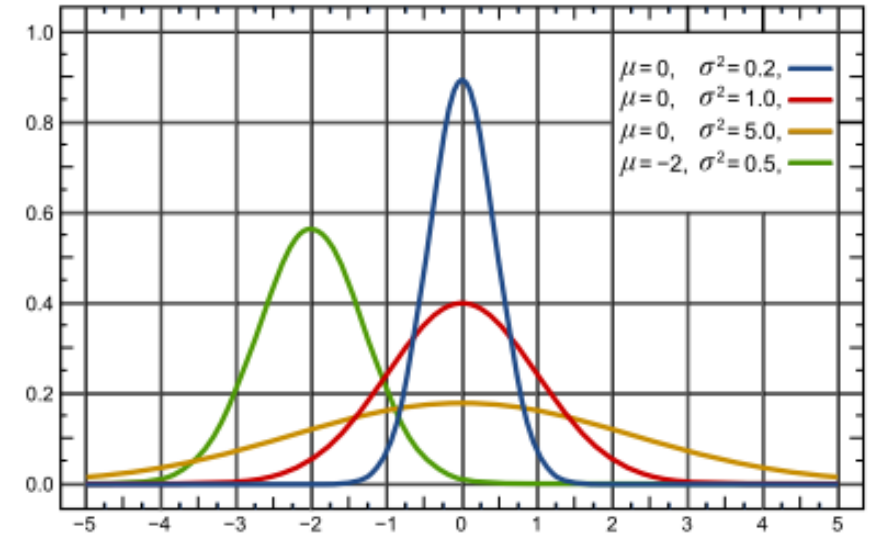
where  $\mu(s, \boldsymbol{\theta}) = \boldsymbol{\phi}(s)^\top \boldsymbol{\theta}$

– infinite support  $a \in (-\infty, +\infty)$

- Beta policy

$$\pi(a | s, \boldsymbol{\theta}) = \mathcal{B}\left(\frac{a + h}{2h} \mid \alpha(s, \boldsymbol{\theta}), \beta(s, \boldsymbol{\theta})\right)$$

– finite support  $a \in [-h, h]$



# Finite Differences Method

- For each dimension
  - Estimate  $k$ -th partial derivative of objective function with respect to  $\theta$
  - By perturbing  $\theta$  by small amount  $\epsilon$  in  $k$ -th dimension

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon \mathbf{u}_k) - J(\theta)}{\epsilon}$$

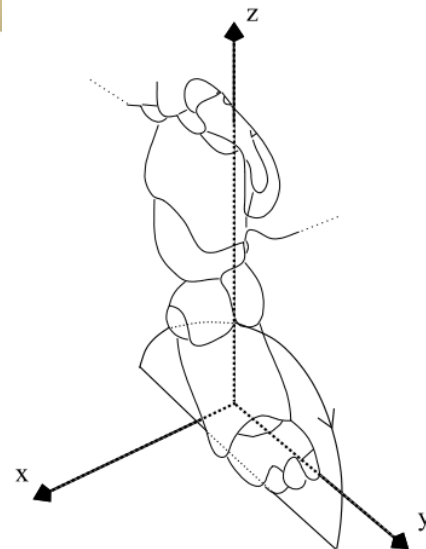
$\swarrow$   
 $k$ -th element

where  $\mathbf{u}_k = [0, \dots, 0, 1, 0, \dots, 0]^\top$

- Use  $n$  evaluations to compute policy gradient in  $n$  dimensions
- Sometimes effective

# Training AIBO to Walk

- $J$  is given by field traversal time



Parameter	Initial Value	$\epsilon$	Best Value
Front locus:			
(height)	4.2	0.35	4.081
(x offset)	2.8	0.35	0.574
(y offset)	4.9	0.35	5.152
Rear locus:			
(height)	5.6	0.35	6.02
(x offset)	0.0	0.35	0.217
(y offset)	-2.8	0.35	-2.982
Locus length	4.893	0.35	5.285
Locus skew multiplier	0.035	0.175	0.049
Front height	7.7	0.35	7.483
Rear height	11.2	0.35	10.843
Time to move through locus	0.704	0.016	0.679
Time on ground	0.5	0.05	0.430

# Training AIBO to Walk



initial behavior



during training



learned behaviors



# Computing Policy Gradient Analytically

- Finite differences method works for arbitrary policies even if a policy is not differentiable
- However, it often estimates noisy policy gradient
- We now compute the policy gradient analytically


$$\begin{aligned}\bullet \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \frac{\partial}{\partial \boldsymbol{\theta}} \left( \sum_{\tau} G(\tau) \Pr(\tau \mid \boldsymbol{\theta}) \right) \\ &= \sum_{\tau} G(\tau) \frac{\partial \Pr(\tau \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{\tau} \Pr(\tau \mid \boldsymbol{\theta}) G(\tau) \frac{\partial \ln \Pr(\tau \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \\ &= \mathbb{E}_{\Pr(\tau \mid \boldsymbol{\theta})} [G(\tau) \nabla_{\boldsymbol{\theta}} \ln \Pr(\tau \mid \boldsymbol{\theta})]\end{aligned}$$

# Derivation

- $$\ln \Pr(\tau \mid \boldsymbol{\theta}) = \ln p_0(S_0) + \sum_{t=0}^{T-1} \ln \pi(A_t \mid S_t, \boldsymbol{\theta}) + \sum_{t=0}^{T-1} \ln p(S_{t+1}, R_t \mid S_t, A_t)$$

- Since the first and third terms do not depend on  $\boldsymbol{\theta}$ , the derivative is given by

$$\nabla_{\boldsymbol{\theta}} \ln \Pr(\tau \mid \boldsymbol{\theta}) = \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \ln \pi(A_t \mid S_t, \boldsymbol{\theta})$$

- We do not need to know the environmental dynamics  
 model free

# Baseline

- As we studied in the bandit problems, the policy gradient remains unchanged even if a constant baseline is subtracted

$$\begin{aligned}\mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})}[G(\tau)\nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})] \\ = \mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})}[(G(\tau) - b)\nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})]\end{aligned}$$

where  $b$  is a constant baseline



# Reminder: Proof Sketch

- $\mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})}[(G(\tau) - b)\nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})]$   
 $= \mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})} [G(\tau)\nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})] - b\mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})]$
- $\mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})] = \sum_{\tau} \text{Pr}(\tau | \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ln \text{Pr}(\tau | \boldsymbol{\theta})$   
 $= \sum_{\tau} \frac{\partial \text{Pr}(\tau | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{\partial}{\partial \boldsymbol{\theta}} \left[ \sum_{\tau} \text{Pr}(\tau | \boldsymbol{\theta}) \right]$   
 $= \frac{\partial}{\partial \boldsymbol{\theta}} [1] = \mathbf{0}$

# Variance Reduction

- We can tune the baseline  $b$  to reduce the variance of the policy gradient

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})} \left[ (G(\tau) - b) \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \ln \pi(A_t | S_t, \boldsymbol{\theta}) \right]$$

- For each  $\theta_i$ , the baseline can be optimized
  - For example, the value-based baseline is used

$$b = \mathbb{E}_{\text{Pr}(\tau|\boldsymbol{\theta})} [G(\tau)]$$

- Note that this does not minimize the variance of the gradient

# REINFORCE algorithm: Monte Carlo Policy Gradient

Input: a differential policy parameterization  $\pi(a | s, \theta)$

Algorithm parameter: step size  $\alpha > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^d$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot | \cdot, \theta)$

$$b \leftarrow T^{-1} \sum_{t=1}^T R_t$$

    Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla_{\theta} \ln \pi(A_t | S_t, \theta)$$

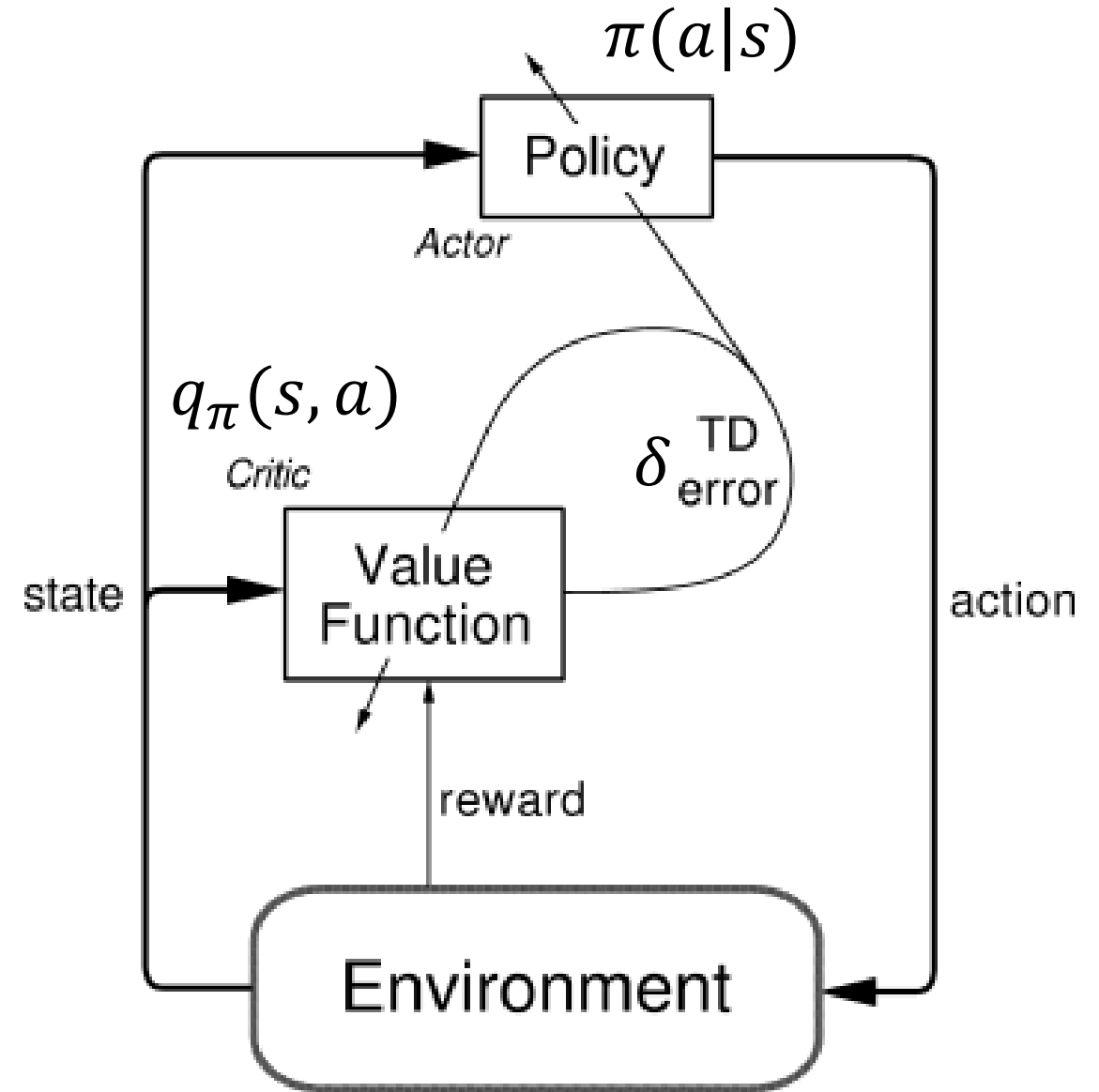
# Pros and cons

	Value-based	Pure Policy-based
learning speed	slow	much slower
convergence	global optimal	local optimal
Markovian	necessary	not necessary
action	discrete	discrete/continuous

- Actor-critic approaches combines the value- and the policy-based approaches

# Actor-Critic Methods

- Two learning components
  - actor:  $\pi(a|s)$
  - critic:  $v_{\pi}(s)$  or  $q_{\pi}(s, a)$
- TD error is used to train  $\pi(a|s)$  as well as  $v_{\pi}(s)$
- Appropriate if the action is continuous



# Policy-Gradient theorem

- Reminder: the update rule of the policy gradient method

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\Pr(\tau|\theta)} [(G(\tau) - b) \nabla_{\theta} \ln \Pr(\tau | \theta)]$$

- Since  $G(\tau)$  is an actual return, we have to wait until an episode ends in order to know  $G(\tau)$
- We will replace the actual return  $G(\tau)$  with a state-action value function

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{d_{\pi}, \pi} [(q_{\pi}(s, a) - b(s)) \nabla_{\theta} \ln \pi(a | s, \theta)] \\ &= \sum_s d_{\pi}(s) \sum_a \pi(a | s, \theta) (q_{\pi}(s, a) - b(s)) \nabla_{\theta} \ln \pi(a | s, \theta) \end{aligned}$$

# What is $d_\pi(s)$

- For the average-reward formulation,

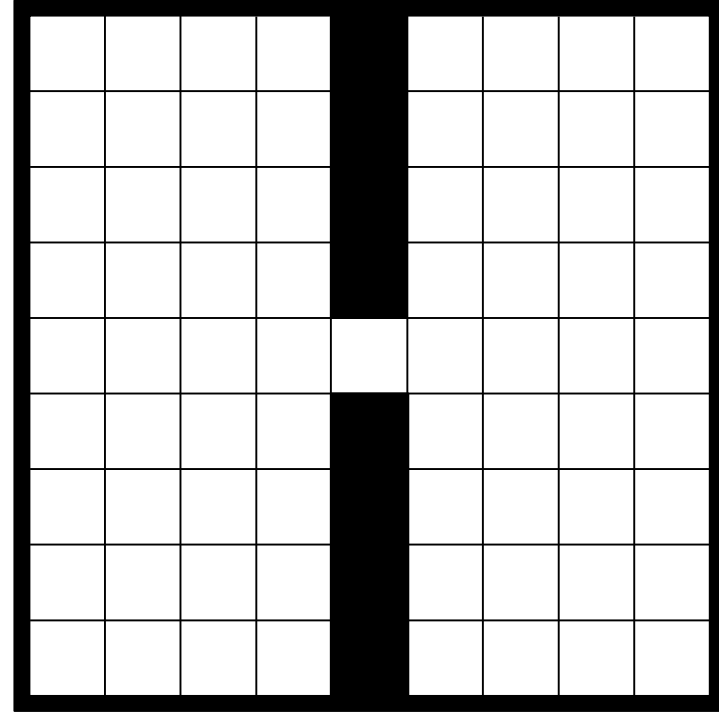
$$d_\pi(s) \triangleq \lim_{t \rightarrow \infty} \Pr(S_t = s \mid A_{0:t} \sim \pi)$$

- $d_\pi(s)$  is called the stationary distribution
  - Under some assumptions, it converges to the stationary distribution of states under  $\pi$
  - The stationary distribution does not depend on the starting state (ergodicity)
- The following equation satisfies

$$\sum_s d_\pi(s) \sum_a \pi(a \mid s) p(s' \mid s, a) = d_\pi(s')$$

# Example: Stationary Distribution

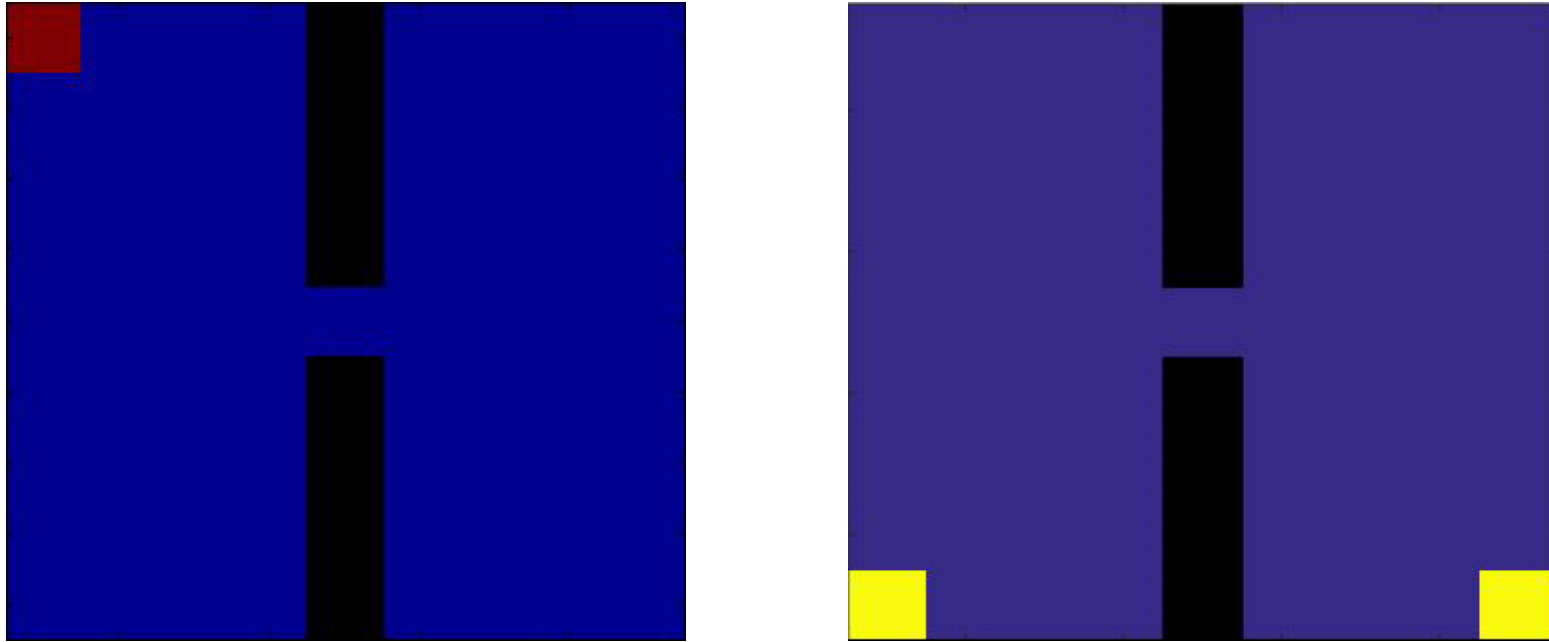
- Consider the 9x9 grid world
- 5 actions (stop, left, right, up, and down)
- Random policy
$$\pi(a | s) = 1/5$$
- Move to the intended direction with probability 0.9.





# Example: Stationary Distribution

- Converge to the same stationary distribution



- For the infinite horizon average reward problem, the stationary distribution plays an important role to define the objective function

# Average reward objective function

- So far, we considered the discounted sum of rewards to formulate the objective function
- Here, we consider the average rate of reward per time step because Policy Gradient Theorem for the discounted problems is complicated to derive
- The new objective function is given by

$$\begin{aligned} J(\boldsymbol{\theta}) &\triangleq \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[R_t \mid A_{0:t-1} \sim \pi] \\ &= \sum_s d_{\pi}(s) \sum_a \pi(a \mid s) \sum_{s', r} r p(s', r \mid s, a) \end{aligned}$$

# Value functions for the new setting

- Value functions are re-defined as follows

$$v_{\pi}(s) \triangleq \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

$$q_{\pi}(s, a) \triangleq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a]$$

- $G_t$  is also re-defined as the differential return

$$G_t \triangleq \underbrace{R_{t+1} - J(\boldsymbol{\theta}) + R_{t+2} - J(\boldsymbol{\theta}) + R_{t+3} - J(\boldsymbol{\theta}) + \dots}_{\text{deviation from the average}}$$

deviation from the average

# Derivation 1/5 (Modified Bellman equation)

- Bellman expectation equation for average reward formulation is

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r - J(\boldsymbol{\theta}) + v_{\pi}(s')]$$

where  $v_{\pi}(s) = \sum_a \pi(a \mid s) q_{\pi}(s, a)$

- Reminder: Bellman expectation equation for discounted problems

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


## Derivation 2/5 (Derivative w.r.t. policy parameter)

- $\nabla_{\theta} v_{\pi}(s) = \nabla_{\theta} \left( \sum_a \pi(a | s) q_{\pi}(s, a) \right)$
- $= \sum_a [\nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \nabla_{\theta} q_{\pi}(s, a)]$
- $= \sum_a \left[ \nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \nabla_{\theta} \left( \sum_{s', r} p(s', r | s, a) (r - J(\theta) + v_{\pi}(s')) \right) \right]$ 

Bellman expectation equation
- $= \sum_a \left\{ \nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \left[ -\nabla_{\theta} J(\theta) + \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right] \right\}$ 

$r$  does not depend on  $\theta$

## Derivation 3/5

- $\nabla_{\theta} v_{\pi}(s) = \sum_a \left\{ \nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \left[ -\nabla_{\theta} J(\theta) + \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right] \right\}$   

- $\sum_a \pi(a | s) \nabla_{\theta} J(\theta) = \sum_a \left\{ \nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \left[ \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right] \right\} - \nabla_{\theta} v_{\pi}(s)$
- Note that the left hand side can be simplified as follows

$$\sum_a \pi(a | s) \nabla_{\theta} J(\theta) = \nabla_{\theta} J(\theta)$$

## Derivation 4/5

- $\nabla_{\theta} J(\theta) = \sum_a \left\{ \nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \left[ \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right] \right\} - \nabla_{\theta} v_{\pi}(s)$
- The left hand side of the obtained equation does not depend on state  $s$
- Thus, the right hand side does not depend on  $s$  either, and we can safely sum it over all  $s \in \mathcal{S}$  weighted by  $d_{\pi}(s)$

$$\sum_s d_{\pi}(s) \nabla_{\theta} J(\theta) = \sum_s d_{\pi}(s) \sum_a \left\{ \nabla_{\theta} \pi(a | s) q_{\pi}(s, a) + \pi(a | s) \left[ \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right] \right\} - \sum_s d_{\pi}(s) \nabla_{\theta} v_{\pi}(s)$$

## Derivation 5/5

- $$\sum_s d_\pi(s) \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \sum_s d_\pi(s) \sum_a \left\{ \nabla_{\boldsymbol{\theta}} \pi(a | s) q_\pi(s, a) + \pi(a | s) \left[ \sum_{s'} p(s' | s, a) \nabla_{\boldsymbol{\theta}} v_\pi(s') \right] \right\} - \sum_s d_\pi(s) \nabla_{\boldsymbol{\theta}} v_\pi(s)$$

- Left hand side

$$\sum_s d_\pi(s) \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

- Second term of the right hand side

$$\sum_s d_\pi(s) \sum_a \pi(a | s) \left[ \sum_{s'} p(s' | s, a) \nabla_{\boldsymbol{\theta}} v_\pi(s') \right] = \sum_{s'} d_\pi(s') \nabla_{\boldsymbol{\theta}} v_\pi(s')$$

– This is identical to the third term



# Policy gradient theorem

- Finally the policy gradient is given by

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \sum_s d_{\pi}(s) \sum_a q_{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \pi(a | s, \boldsymbol{\theta}) \\ &= \sum_s d_{\pi}(s) \sum_a \pi(a | s, \boldsymbol{\theta}) q_{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \ln \pi(a | s, \boldsymbol{\theta}) \\ &= \mathbb{E}_{d_{\pi}, \pi} [q_{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \ln \pi(a | s, \boldsymbol{\theta})]\end{aligned}$$

- This is called the policy gradient theorem  
(Sutton et al., 2000; Konda and Tsitsiklis, 2000)
- Reminder: Gradient of the bandit problems

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi} [(r - b) \nabla_{\boldsymbol{\theta}} \ln \pi(a | \boldsymbol{\theta})]$$

# Action-value actor-critic

- Prepare a **critic**  $\hat{q}(s, a, \mathbf{w})$  and an **actor**  $\pi(a \mid s, \boldsymbol{\theta})$
- Sample  $A \sim \pi(a \mid S, \boldsymbol{\theta})$
- For each step do
  - Sample reward  $r$  and transition  $S', R \sim p(s, r \mid S, A)$
  - Sample action  $A' \sim \pi(\cdot \mid S', \boldsymbol{\theta})$
  - Compute the TD error  $\delta = r + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$
  - Update the policy parameter
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta}) \hat{q}(s, a, \mathbf{w})$$
  - Update the value parameter
$$\mathbf{w} \leftarrow \mathbf{w} + \alpha_{\mathbf{w}} \delta \nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w})$$
  - $s \leftarrow s', a \leftarrow a'$
- end for

# Approximating $q_\pi$

- Let  $\hat{q}(s, a, \mathbf{w})$  be an approximation of  $q_\pi(s, a)$ , parameterized by  $\mathbf{w}$
- How should we approximate?
- The policy gradient is exact when the following two conditions are satisfied

1. Value function approximator is compatible to the policy

$$\nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w}) = \nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta})$$

2. The parameter  $\mathbf{w}$  minimizes the mean-squared error  $\varepsilon$

$$\varepsilon = \mathbb{E}_{d_{\pi, \pi}} \left[ \left( q_\pi(s, a) - \hat{q}(s, a, \mathbf{w}) \right)^2 \right]$$

# Approximating $q_\pi$

- Since  $\mathbf{w}$  minimizes  $\varepsilon$ ,  $\nabla_{\mathbf{w}}\varepsilon = \mathbf{0}$
- $$\nabla_{\mathbf{w}}\varepsilon = -\mathbb{E}_{d_{\pi,\pi}}[(q_\pi(s, a) - \hat{q}(s, a, \mathbf{w}))\nabla_{\mathbf{w}}\hat{q}(s, a, \mathbf{w})]$$
- Using the first assumption, we obtain

$$\mathbb{E}_{d_{\pi,\pi}}[(q_\pi(s, a) - \hat{q}(s, a, \mathbf{w}))\nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta})] = \mathbf{0}$$

- $$\begin{aligned} \mathbb{E}_{d_{\pi,\pi}}[\hat{q}(s, a, \mathbf{w})\nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta})] \\ = \mathbb{E}_{d_{\pi,\pi}}[q_\pi(s, a)\nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta})] = \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \end{aligned}$$

# Theorem

- Policy gradient  $\nabla_{\theta} J$  is calculated by

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{d_{\pi, \pi}} [\nabla_{\theta} \ln \pi(a | s, \theta) \hat{q}(s, a, \mathbf{w})]$$

where

$$\nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w}) = \nabla_{\theta} \ln \pi(a | s, \theta)$$

- The equation of the first assumption implies that  $\hat{q}$  is represented by

$$\hat{q}(s, a, \mathbf{w}) = \mathbf{w}^{\top} \underbrace{\nabla_{\theta} \ln \pi(a | s, \theta)}$$

- Then, we can write

basis function vector

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{d_{\pi, \pi}} [\nabla_{\theta} \ln \pi(a | s, \theta) \nabla_{\theta} \ln \pi(a | s, \theta)^{\top} \mathbf{w}]$$

## Note on $\hat{q}(s, a, \mathbf{w})$

- Since the projected value function is given by

$$\hat{q}(s, a, \mathbf{w}) = \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s, \boldsymbol{\theta})$$

- It satisfies

$$\begin{aligned} \sum_a \pi(a \mid s; \boldsymbol{\theta}) \hat{q}(s, a, \mathbf{w}) &= \sum_a \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \pi(a \mid s, \boldsymbol{\theta}) \\ &= \mathbf{w}^\top \sum_a \nabla_{\boldsymbol{\theta}} \pi(a \mid s, \boldsymbol{\theta}) = \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \left[ \sum_a \pi(a \mid s, \boldsymbol{\theta}) \right] \\ &= \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} (1) = 0 \end{aligned}$$

- It suggests that

$$\mathbb{E}_{\pi}[\hat{q}(s, a, \mathbf{w})] = 0$$

## Note on $\hat{q}(s, a, w)$

- Interestingly, it is enough to use a linear function approximator to update the actor

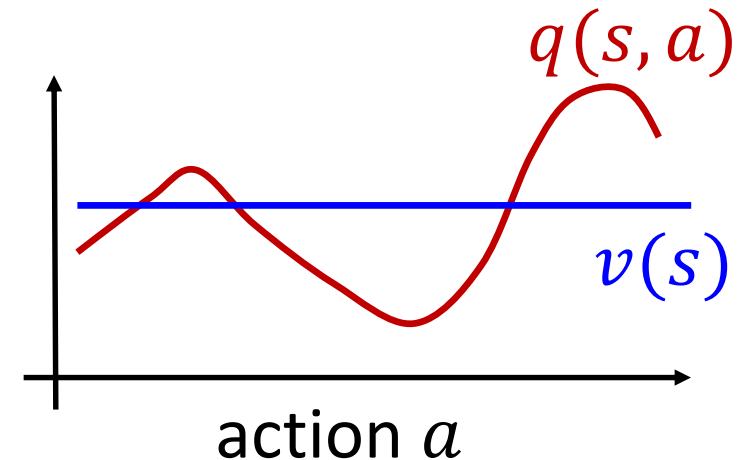
$$q(s, a; \mathbf{w}) = \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s; \boldsymbol{\theta})$$

- Choice of  $b(s)$  is important. When  $b(s) = V(s)$ ,

$$q(s, a) - v(s) = A(s, a)$$

is usually called the advantage function

- The TD error can be used as a target value to train the advantage function [Morimura et al., 2005]



# Summary

- REINFORCE algorithm is the most classical policy gradient, and it is applicable even though the environment is not Markovian
- The Actor-Critic has many equivalent forms:

$$\begin{aligned}\nabla_{\theta} J(\boldsymbol{\theta}) &= \mathbb{E}[\nabla_{\theta} \ln \pi(a | s, \theta) Q^w(s, a)] && \text{Q Actor-Critic} \\ &= \mathbb{E}[\nabla_{\theta} \ln \pi(a | s, \theta) A^w(s, a)] && \text{Advantage Actor-Critic} \\ &= \mathbb{E}[\nabla_{\theta} \ln \pi(a | s, \theta) \delta] && \text{TD Actor-Critic}\end{aligned}$$

- Each leads a stochastic gradient ascent algorithm
- Critic uses policy evaluation to estimate  $q_{\pi}$ ,  $A_{\pi}$ , or  $v_{\pi}$



# References

- Choromanski, K., Iscen, A., Sindhvani, V., Tan, J. and Coumans, E. Optimizing Simulations with Noise-Tolerant Structured Exploration. In Proc. of ICRA 2018.
- Deisenroth, M., Neumann, G., and Peters, J. (2013). A survey on policy search for robotics. Foundations and Trends in Robotics, 2(1-2): 1-142.
- Kober, J. and Peters, J. Policy Search for Motor Primitives in Robotics. NIPS 2009.
- Kohl, N. and Stone, P. (2004). [Policy Gradient Reinforcement Learning for Fast Quadrupedal Locomotion](#). In Proc. of IEEE International Conference on Robotics and Automation, 2619-2624.

# References

- Konda, V. and Tsitsiklis, J. (2000). Actor-Critic Algorithms. NIPS 12.
- Lagoudakis, M. and Parr, R. (2003). Least-squares policy iteration. Machine Learning, 4, 1107-1149.
- Peters, J. and Schaal, S. (2008). Natural Actor-Critic. Neurocomputing, 71(7-9), 1180-1190.
- Peters, J. and Schaal, S. (2008). Reinforcement learning of motor skills with policy gradients. Neural Networks, 21(4), 682-697.
- Sehnke, F. et al. Parameter-exploring policy gradients. Neural Networks, 23. 551-559, 2010.
- Sutton, R., McAllester, D. Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. NIPS 12.

# References

- Theodorou, E. A. et al. A generalized path integral control approach to reinforcement learning. *Journal of Machine Learning Research*. 11, 3153-3197, 2010.
- Vlassis, N. et al. Learning model-free robot control by a Monte Carlo EM. *Autonomous Robotics*. 27, 123-130, 2009.
- Williams, R. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*. 8(3): 229-256, 1992.

# What is $d_\pi(s)$

- For the discounted-reward formulation,  $d^\pi(s)$  is defined by

$$d_\pi(s) = \sum_{t=0}^{\infty} \gamma^t \Pr(S_t = s \mid S_0, \pi)$$

- Dependent on the starting state
- For the average-reward formulation,

$$d_\pi(s) = \lim_{t \rightarrow \infty} \Pr(S_t = s \mid S_0, \pi)$$

- Under some assumptions, it converges to the stationary distribution of states under  $\pi$
  - The stationary distribution does not depend on the starting state  $S_0$

# Function approximation

- Consider a linear function approximator

$$G(\tau) \approx \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \ln p(\tau \mid \boldsymbol{\theta}) + \mathbf{v}^\top \boldsymbol{\varphi}(S_0)$$

where  $\mathbf{w} \in \mathbb{R}^n$  is the parameter to be tuned.

- $\boldsymbol{\varphi}(S_0) \in \mathbb{R}^{n'}$  and  $\mathbf{v} \in \mathbb{R}^{n'}$  denote the basis function and the weight vector, respectively.
- The parameter vectors  $\mathbf{w}$  and  $\mathbf{v}$  can be estimated by the least squares method from samples
- $\mathbf{v}^\top \boldsymbol{\varphi}(S_0)$  can be interpreted as an approximator of the baseline  $b$

# How to estimate $w$ ?

- The dataset  $\mathcal{D} = \{\tau_i\}_{i=1}^N$  is given, where

$$\tau_i = \{S_0^i, A_0^i, R_1^i, \dots, S_{T-1}^i, A_{T-1}^i, R_{T-1}^i, S_T^i\}$$

- The error of the  $i$ -th data is computed by

$$e_i = [\nabla_{\boldsymbol{\theta}} \ln p(\tau_i \mid \boldsymbol{\theta})^\top \quad \boldsymbol{\varphi}^\top(S_0^i)] \begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} - G_i$$

where  $G_i = \sum_{t=0}^{T-1} R_t^i$

- Since this is the ordinary least squares problem, the objective function is constructed by

$$E = \sum_{i=1}^N e_i^2$$

# Algorithm of Episodic NAC

---

## Algorithm 8 Episodic Natural Actor Critic

---

Input: Policy parametrization  $\theta$ ,

data set  $\mathcal{D} = \left\{ \mathbf{x}_{1:T}^{[i]}, \mathbf{u}_{1:T-1}^{[i]}, r_{1:T}^{[i]} \right\}_{i=1 \dots N}$

**for** each sample  $i = 1 \dots N$  **do**

  Compute returns:  $R^{[i]} = \sum_{t=0}^T r_t^{[i]}$

  Compute features:  $\boldsymbol{\psi}^{[i]} = \begin{bmatrix} \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}} \left( \mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right) \\ \boldsymbol{\varphi}(\mathbf{x}_0^{[i]}) \end{bmatrix}$

**end for**

Fit advantage function and initial value function

$$\mathbf{R} = \left[ R^{[1]}, \dots, R^{[N]} \right]^T, \quad \boldsymbol{\Psi} = \left[ \boldsymbol{\psi}^{[1]}, \dots, \boldsymbol{\psi}^{[N]} \right]^T$$

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} = (\boldsymbol{\Psi}^T \boldsymbol{\Psi})^{-1} \boldsymbol{\Psi}^T \mathbf{R}$$

return  $\nabla_{\boldsymbol{\theta}}^{\text{eNAC}} J_{\boldsymbol{\theta}} = \mathbf{w}$

---

note

$$\begin{aligned} & \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}_t^i \mid \mathbf{x}_t^i) \\ &= \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\tau}_i) \end{aligned}$$

# Reminder: Policy Gradient Theorem

- The vanilla gradient for the infinite horizon can be derived from the Policy Gradient Theorem:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{d_{\pi, \pi}} \left[ \nabla_{\boldsymbol{\theta}} \ln \pi(a | s, \boldsymbol{\theta}) \left( \hat{q}(s, a, \boldsymbol{w}) - \hat{b}(s) \right) \right]$$

- This theorem is independently proved in 2000 by Sutton et al. and Konda and Tsitsiklis
- Then, a linear function estimator is constructed to approximate  $\hat{q}(s, a, \boldsymbol{w}) - \hat{b}(s)$



# Advantage function

- If  $\hat{b}(s) = \hat{v}(s, \boldsymbol{v})$ , the following function is called the advantage function

$$A^\pi(s, a) = \hat{q}(s, a, \boldsymbol{w}) - \hat{v}(s, \boldsymbol{v})$$

where

$$\hat{v}(s, \boldsymbol{v}) \approx \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

$$\hat{q}(s, a, \boldsymbol{w}) \approx \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

- Then, consider the following approximators

$$\hat{A}(s, a, \boldsymbol{w}) \triangleq \boldsymbol{w}^\top \nabla_{\boldsymbol{\theta}} \ln \pi(a \mid s)$$

$$\hat{v}(s, \boldsymbol{v}) \triangleq \boldsymbol{v}^\top \boldsymbol{\varphi}(s)$$

- Important features

$$\mathbb{E}_{\pi}\{A^{\pi}(s, a)\} = \sum_a \pi(a | s) A^{\pi}(s, a) = 0$$

$$\mathbb{E}_{\pi}\{\mathbf{w}^{\top} \nabla_{\theta} \ln \pi(a | s)\} = 0$$

- Using the approximators in the previous slide,  $\hat{q}(s, a, \mathbf{w}, \mathbf{v})$  can be approximated by

$$\hat{q}(s, a, \mathbf{w}, \mathbf{v}) = \mathbf{w}^{\top} \nabla_{\theta} \ln \pi(a | s) + \mathbf{v}^{\top} \varphi(s)$$

- We can apply any policy evaluation algorithm