

# 强化学习及其应用

Reinforcement Learning and Its Applications

## 第四章 策略控制

**Policy Control**

授课人：周晓飞  
zhouxiaofei@iie.ac.cn  
2023-6-13

# 第四章 策略控制

4. 1 策略优化

4. 2 蒙特卡洛策略控制

4. 3 时序差分策略控制

4. 4 Q-Learning

4. 4 算法总结

# 第四章 策略控制

## 4.1 策略优化

## 4.2 蒙特卡洛策略控制

## 4.3 时序差分策略控制

## 4.4 Q-Learning

## 4.4 算法总结

# 策略优化

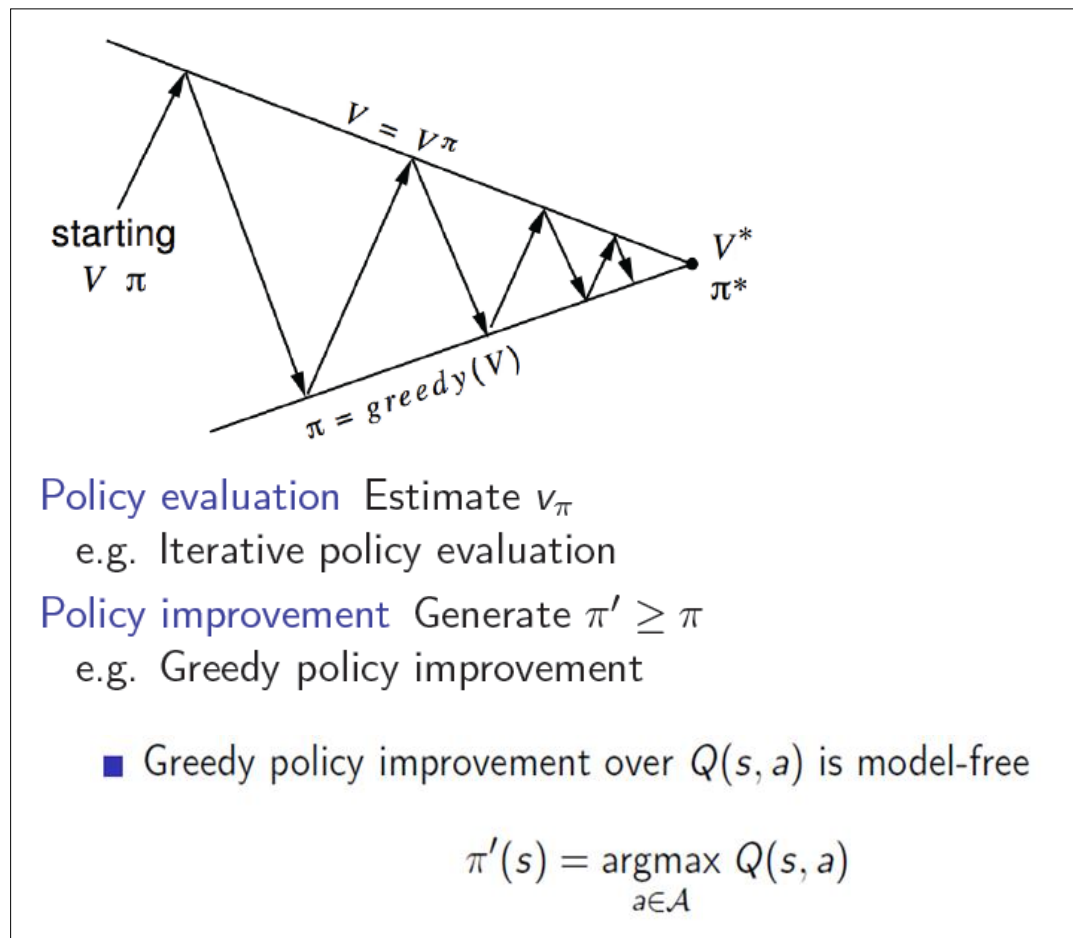
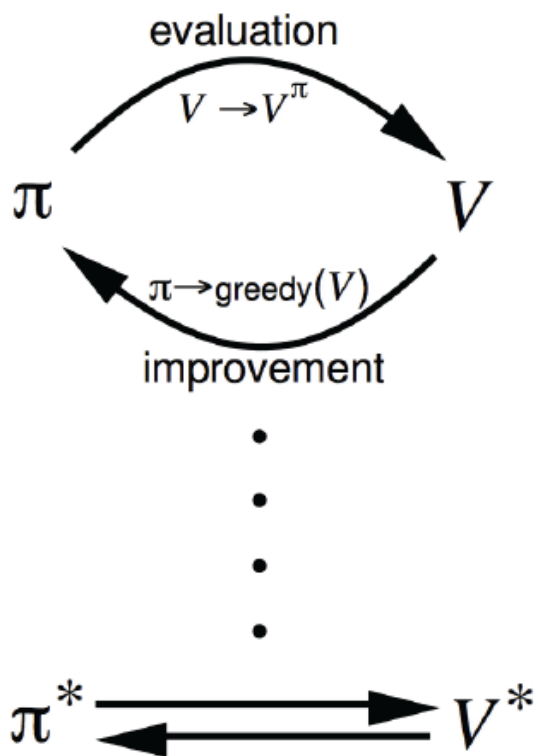
## 问题描述

- For prediction:
  - Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and policy  $\pi$
  - or: MRP  $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$
  - Output: value function  $v_\pi$
- Or for control:
  - Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
  - Output: optimal value function  $v_*$
  - and: optimal policy  $\pi_*$

# 策略优化

## 策略迭代

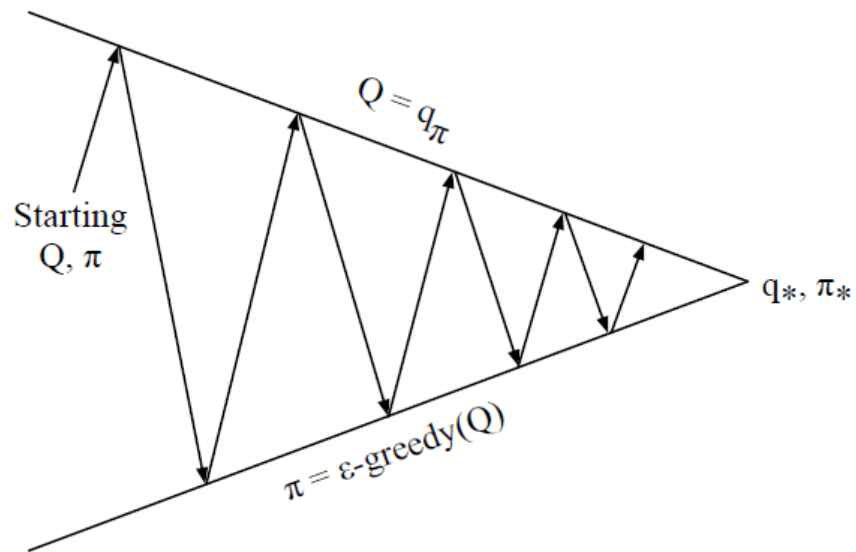
### 两个步骤: Evaluation & Improvement



# 策略优化

## 策略迭代

Q替代V值，进行迭代



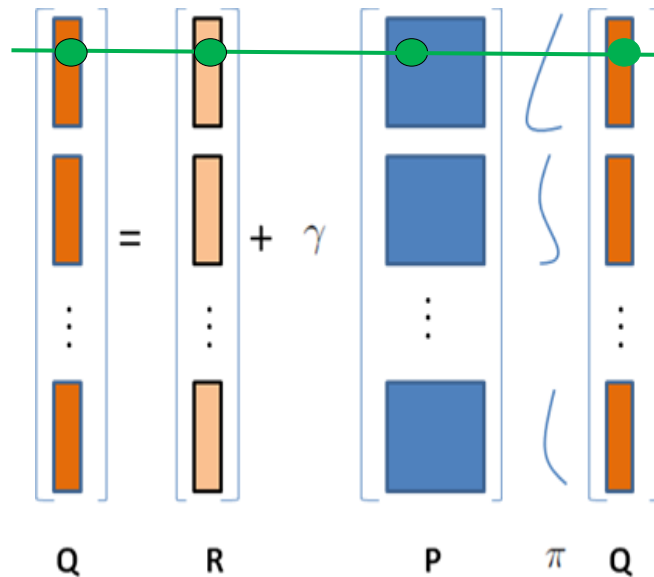
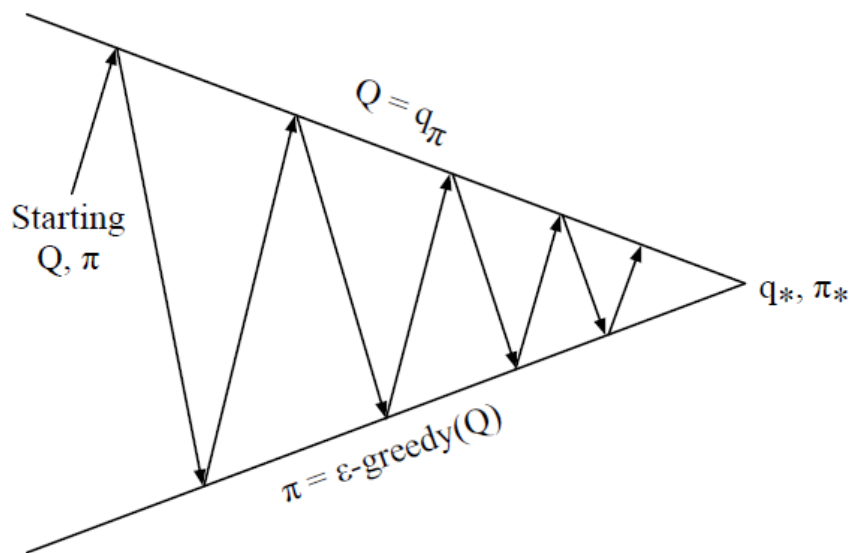
$$\begin{bmatrix} \text{orange bar} \\ \text{orange bar} \\ \vdots \\ \text{orange bar} \end{bmatrix} = \begin{bmatrix} \text{light orange bar} \\ \text{light orange bar} \\ \vdots \\ \text{light orange bar} \end{bmatrix} + \gamma \left( \begin{bmatrix} \text{blue square} \\ \text{blue square} \\ \vdots \\ \text{blue square} \end{bmatrix} \begin{bmatrix} \text{blue line} \\ \text{blue line} \\ \vdots \\ \text{blue line} \end{bmatrix} \begin{bmatrix} \text{orange bar} \\ \text{orange bar} \\ \vdots \\ \text{orange bar} \end{bmatrix} \right)$$

$Q \quad R \quad P \quad \pi \quad Q$

# 策略优化

## 策略迭代

Q替代V值，进行迭代



# 策略优化

## 策略迭代

### 如何改善 Policy?

$\epsilon$ -greedy policy improvement

- Simplest idea for ensuring continual exploration
- All  $m$  actions are tried with non-zero probability
- With probability  $1 - \epsilon$  choose the greedy action
- With probability  $\epsilon$  choose an action at random

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{if } a^* = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a) \\ \epsilon/m & \text{otherwise} \end{cases}$$

始终让最优的  $a$  具有最大的概率，策略分布近似 one-hot 分布，同时实现 Exploration。



# 第四章 策略控制

4. 1 策略优化

**4. 2 蒙特卡洛策略控制**

4. 3 时序差分策略控制

4. 4 Q-Learning

4. 4 算法总结

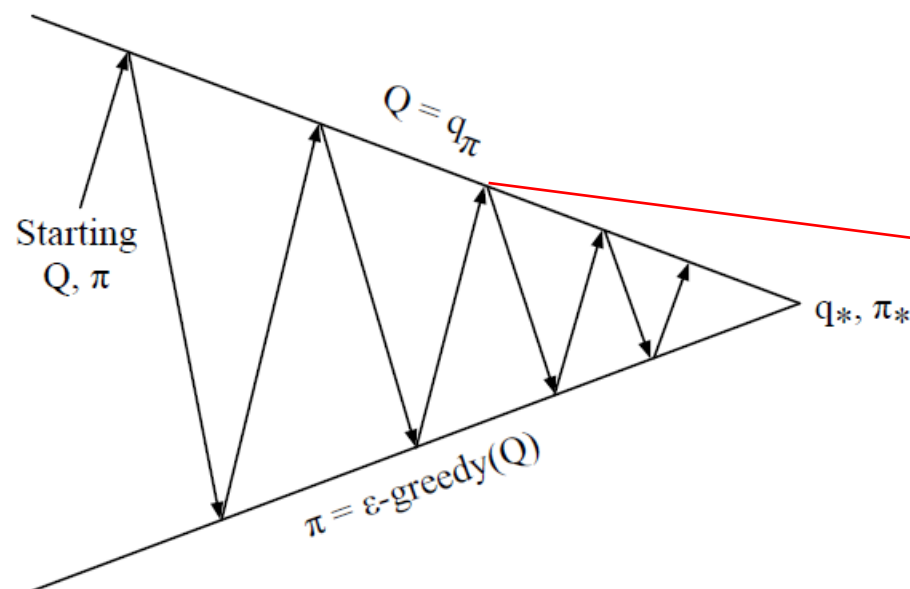
# 蒙特卡洛策略控制

## MC Policy Iteration

两个部分:

Policy evaluation Monte-Carlo policy evaluation,  $Q = q_\pi$

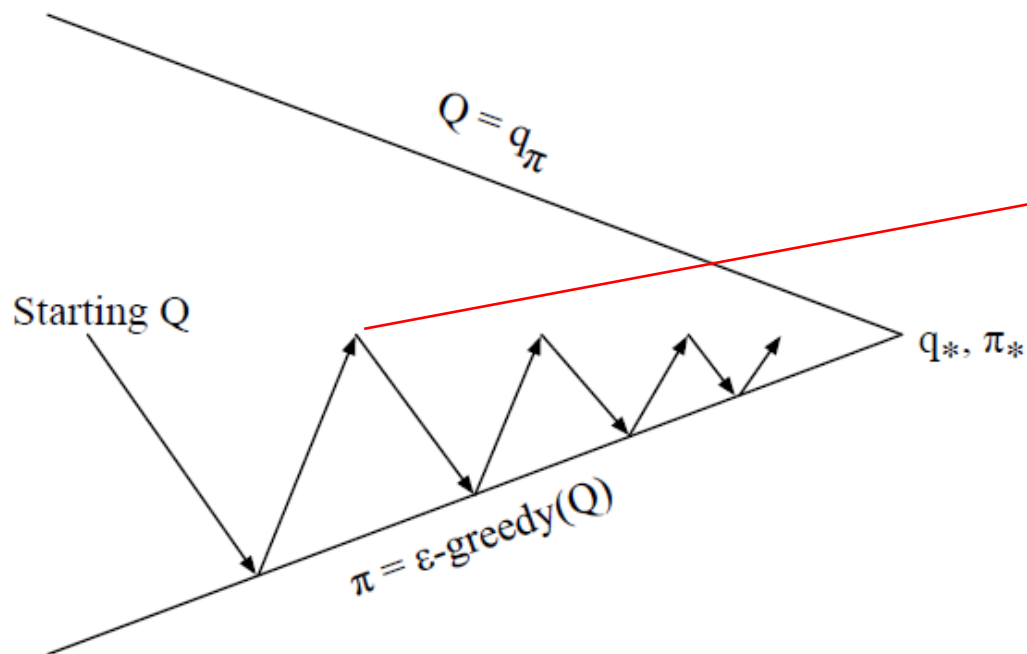
Policy improvement  $\epsilon$ -greedy policy improvement



每次 policy 的提升，都充分的采样了

# 蒙特卡洛策略控制

## MC Control



每次 **Policy** 的提升，没有充分采样，都是随机的逼近

Every episode:

Policy evaluation Monte-Carlo policy evaluation,  $Q \approx q_\pi$

Policy improvement  $\epsilon$ -greedy policy improvement

## GLIE MC Control

- Sample  $k$ th episode using  $\pi$ :  $\{S_1, A_1, R_2, \dots, S_T\} \sim \pi$
- For each state  $S_t$  and action  $A_t$  in the episode,

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$$

- Improve policy based on new action-value function

$$\epsilon \leftarrow 1/k$$

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

### Theorem

*GLIE Monte-Carlo control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$*

# 蒙特卡洛策略控制

## GLIE MC Control

- Sample  $k$ th episode using  $\pi$ :  $\{S_1, A_1, R_2, \dots, S_T\} \sim \pi$
- For each state  $S_t$  and action  $A_t$  in the episode,

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$$

MC Evaluation

- Improve policy based on new action-value function

$$\epsilon \leftarrow 1/k \longrightarrow \text{收敛技巧}$$

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

$\epsilon$ -greedy policy improvement

### Theorem

*GLIE Monte-Carlo control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$*

# 第四章 策略控制

4. 1 策略优化

4. 2 蒙特卡洛策略控制

**4. 3 时序差分策略控制**

4. 4 Q-Learning

4. 4 算法总结

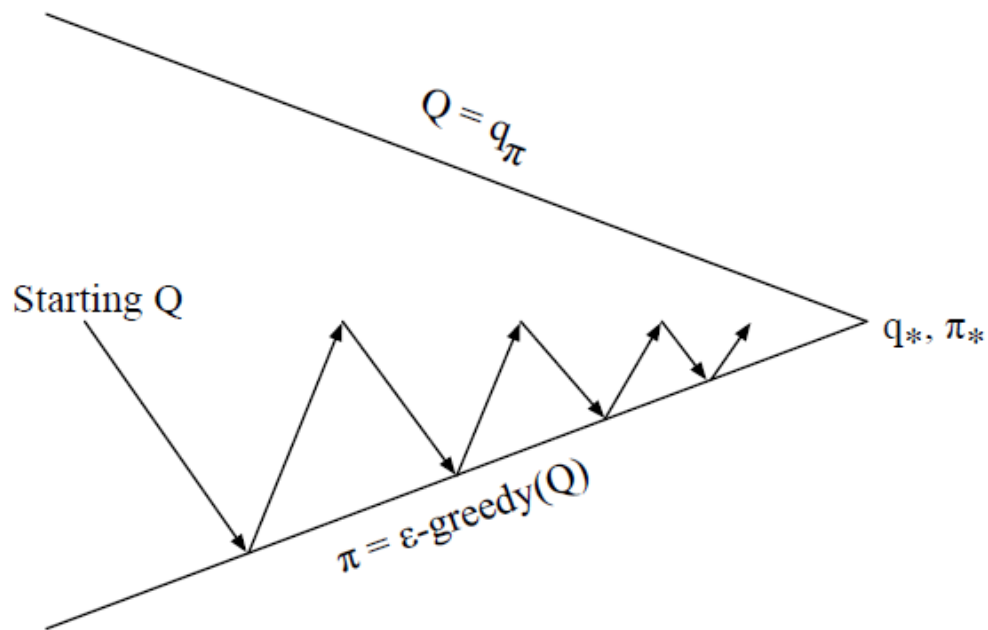
# 时序差分策略控制

## TD for Policy Iteration

- Temporal-difference (TD) learning has several advantages over Monte-Carlo (MC)
  - Lower variance
  - Online
  - Incomplete sequences
- Natural idea: use TD instead of MC in our control loop
  - Apply TD to  $Q(S, A)$
  - Use  $\epsilon$ -greedy policy improvement
  - Update every time-step

# 时序差分策略控制

## Sarsa Control



Every **time-step**:

Policy evaluation **Sarsa**,  $Q \approx q_\pi$

Policy improvement  $\epsilon$ -greedy policy improvement



## Sarsa Control

### ■ Action-Value Evaluation

$$Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma Q(S', A') - Q(S, A))$$

### ■ Policy Improvement

$\epsilon$ -greedy policy improvement

# 时序差分策略控制

## Sarsa Control

### ■ Sarsa Algorithm

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal

*Have a break !*

# 时序差分策略控制

## Sarsa ( $\lambda$ )

### ■ n-step returns

- Consider the following  $n$ -step returns for  $n = 1, 2, \infty$ :

$$\begin{array}{ll} n = 1 & \text{(Sarsa)} \quad q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1}) \\ n = 2 & q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2}) \\ & \vdots \\ n = \infty & \text{(MC)} \quad q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T \end{array}$$

- Define the  $n$ -step Q-return

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

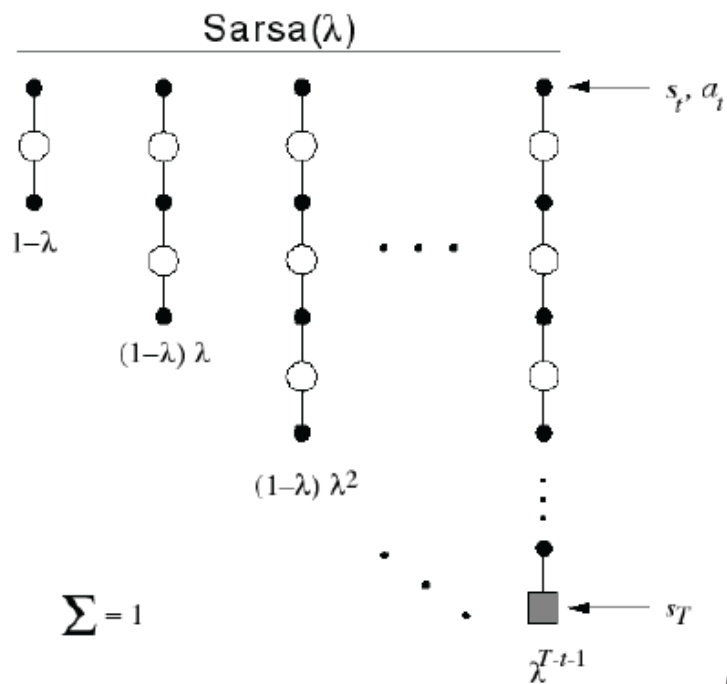
- $n$ -step Sarsa updates  $Q(s, a)$  towards the  $n$ -step Q-return

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( q_t^{(n)} - Q(S_t, A_t) \right)$$

# 时序差分策略控制

## Sarsa ( $\lambda$ )

### Forward Sarsa( $\lambda$ )



- The  $q^\lambda$  return combines all  $n$ -step Q-returns  $q_t^{(n)}$

- Using weight  $(1-\lambda)\lambda^{n-1}$

$$q_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$$

- Forward-view Sarsa( $\lambda$ )

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^\lambda - Q(S_t, A_t))$$

# 时序差分策略控制

## Sarsa ( $\lambda$ )

### ■ Backward Sarsa( $\lambda$ )

- Just like TD( $\lambda$ ), we use **eligibility traces** in an online algorithm
- But Sarsa( $\lambda$ ) has one eligibility trace for each state-action pair

$$E_0(s, a) = 0$$

$$E_t(s, a) = \gamma\lambda E_{t-1}(s, a) + \mathbf{1}(S_t = s, A_t = a)$$

- $Q(s, a)$  is updated for every state  $s$  and action  $a$
- In proportion to TD-error  $\delta_t$  and eligibility trace  $E_t(s, a)$

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t E_t(s, a)$$

# 时序差分策略控制

## Sarsa ( $\lambda$ )

### ■ Backward Sarsa( $\lambda$ ) Algorithm

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat (for each episode):

$E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

    Initialize  $S, A$

    Repeat (for each step of episode):

        Take action  $A$ , observe  $R, S'$

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$

$E(S, A) \leftarrow E(S, A) + \delta$

        For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :

$Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$

$E(s, a) \leftarrow \gamma \lambda E(s, a)$

$S \leftarrow S'; A \leftarrow A'$

    until  $S$  is terminal

# 第四章 策略控制

4. 1 策略优化

4. 2 蒙特卡洛策略控制

4. 3 时序差分策略控制

**4. 4 Q-Learning**

4. 4 算法总结



## Off-Policy Learning

- Evaluate target policy  $\pi(a|s)$  to compute  $v_\pi(s)$  or  $q_\pi(s, a)$
- While following behaviour policy  $\mu(a|s)$

$$\{S_1, A_1, R_2, \dots, S_T\} \sim \mu$$

- Why is this important?
- Learn from observing humans or other agents
- Re-use experience generated from old policies  $\pi_1, \pi_2, \dots, \pi_{t-1}$
- Learn about *optimal* policy while following *exploratory* policy
- Learn about *multiple* policies while following *one* policy

## 重要性采样

- Estimate the expectation of a different distribution

$$\begin{aligned}\mathbb{E}_{X \sim P}[f(X)] &= \sum P(X)f(X) \\ &= \sum Q(X) \frac{P(X)}{Q(X)} f(X) \\ &= \mathbb{E}_{X \sim Q} \left[ \frac{P(X)}{Q(X)} f(X) \right]\end{aligned}$$

## 重要性采样

- Estimate the expectation of a different distribution

$$\begin{aligned}\mathbb{E}_{X \sim P}[f(X)] &= \sum P(X)f(X) \\ &= \sum Q(X) \frac{P(X)}{Q(X)} f(X) \\ &= \mathbb{E}_{X \sim Q} \left[ \frac{P(X)}{Q(X)} f(X) \right]\end{aligned}$$

随机过程中的样本

## Off-Policy Monte-Carlo

- Use returns generated from  $\mu$  to evaluate  $\pi$
- Weight return  $G_t$  according to similarity between policies
- Multiply importance sampling corrections along whole episode

$$G_t^{\pi/\mu} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t$$

- Update value towards *corrected* return

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{\pi/\mu} - V(S_t) \right)$$

- Cannot use if  $\mu$  is zero when  $\pi$  is non-zero
- Importance sampling can dramatically increase variance

## Off-Policy TD

- Use TD targets generated from  $\mu$  to evaluate  $\pi$
- Weight TD target  $R + \gamma V(S')$  by importance sampling
- Only need a single importance sampling correction

$$V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

- Much lower variance than Monte-Carlo importance sampling
- Policies only need to be similar over a single step

## Q-Learning Control

**策略控制中，Q 值的随机更新公式：**

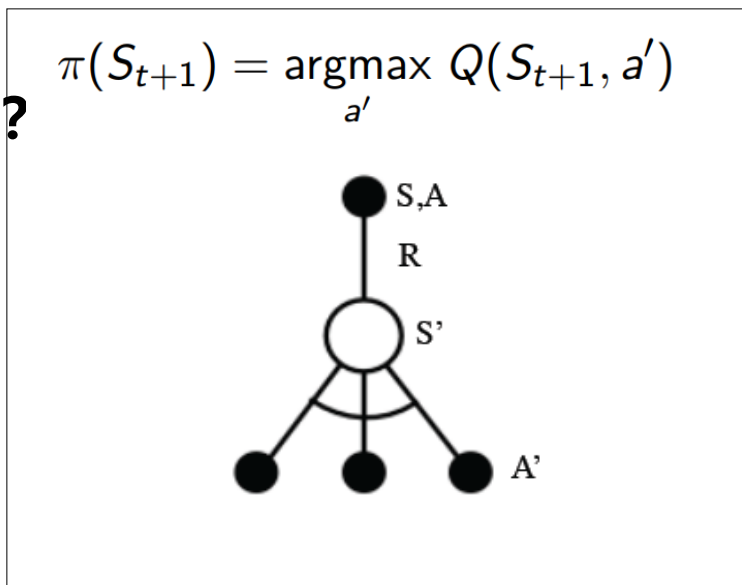
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t))$$

策略是随时被贪婪策略改善，即产生  $A_{t+1}$  的策略比产生  $A_t$  的要好；  
因此  $Q(S_t, A_t)$  被改善后的策略引导修正。

# Q-Learning

## Q-Learning Control

通过重要性采样思想，可否让  $Q(S_t, A_t)$  被**最优的策略**引导？



# Q-Learning

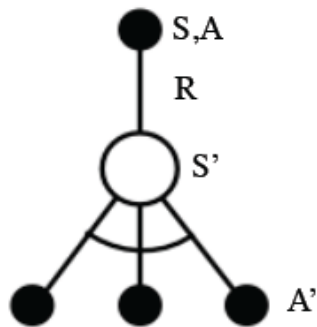
## Q-Learning Control

通过重要性采样思想，可否让  $Q(S_t, A_t)$  被**最优的策略**引导？

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$

以目标的最优分布 ( $\max_a$ ) 的  $Q$  值，修正当前分布的  $Q$  值。

$$\pi(S_{t+1}) = \operatorname{argmax}_{a'} Q(S_{t+1}, a')$$



### Theorem

*Q-learning control converges to the optimal action-value function,*  
 $Q(s, a) \rightarrow q_*(s, a)$



## Q-Learning

### ■ Q-Learning Algorithm

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$S \leftarrow S'$ ;

until  $S$  is terminal

# 第四章 策略控制

4. 1 策略优化

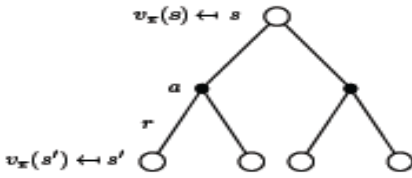

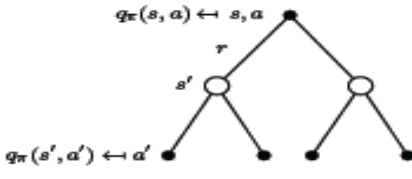

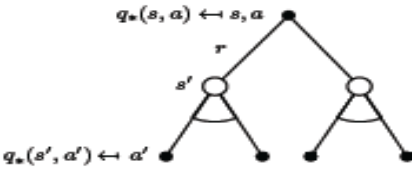
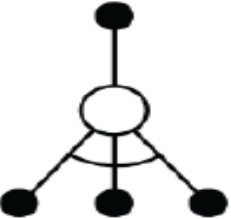
4. 2 蒙特卡洛策略控制

4. 3 时序差分策略控制

4. 4 Q-Learning

**4. 4 算法总结**

# TD V.S. DP

	Full Backup (DP)	Sample Backup (TD)
Bellman Expectation Equation for $v_{\pi}(s)$	 <p>Iterative Policy Evaluation</p>	 <p>TD Learning</p>
Bellman Expectation Equation for $q_{\pi}(s, a)$	 <p>Q-Policy Iteration</p>	 <p>Sarsa</p>
Bellman Optimality Equation for $q_{*}(s, a)$	 <p>Q-Value Iteration</p>	 <p>Q-Learning</p>

## TD V.S. DP

<i>Full Backup (DP)</i>	<i>Sample Backup (TD)</i>
Iterative Policy Evaluation	TD Learning
$V(s) \leftarrow \mathbb{E}[R + \gamma V(S') \mid s]$	$V(S) \stackrel{\alpha}{\leftarrow} R + \gamma V(S')$
Q-Policy Iteration	Sarsa
$Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A') \mid s, a]$	$Q(S, A) \stackrel{\alpha}{\leftarrow} R + \gamma Q(S', A')$
Q-Value Iteration	Q-Learning
$Q(s, a) \leftarrow \mathbb{E}\left[R + \gamma \max_{a' \in \mathcal{A}} Q(S', a') \mid s, a\right]$	$Q(S, A) \stackrel{\alpha}{\leftarrow} R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')$

# 本讲参考文献

1. Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. (Second edition, in progress, draft).
2. David Silver, Slides@ 《Reinforcement Learning: An Introduction》, 2016.