

# 强化学习2024

## 第4节

涉及知识点：

SARSA、Q学习算法及其收敛性、多步自助法

# 无模型控制方法

张伟楠 – [上海交通大学](#)

# 2024年上海交通大学ACM班强化学习课程大纲

## 强化学习基础部分

(中文课件)

1. 强化学习、探索与利用
2. MDP和动态规划
3. 值函数估计
4. 无模型控制方法
5. 参数化的值函数和策略
6. 规划与学习
7. 深度强化学习价值方法
8. 深度强化学习策略方法

## 强化学习前沿部分

(英文课件)

9. 基于模型的深度强化学习
10. 离线强化学习
11. 模仿学习
12. 多智能体强化学习基础
13. 多智能体强化学习前沿
14. 基于扩散模型的强化学习
15. AI Agent与决策大模型
16. 技术与交流与回顾

# 从知道什么是好的，到如何做好行动

- 从知道什么是好的：估计 $V^\pi(s_t)$

$$V(s_t) \leftarrow V(s_t) + \alpha(g_t - V(s_t))$$

$$V(s_t) \leftarrow V(s_t) + \alpha(r_t + \gamma V(s_{t+1}) - V(s_t))$$

- 基于 $V$ 函数，如何选择好的行动？

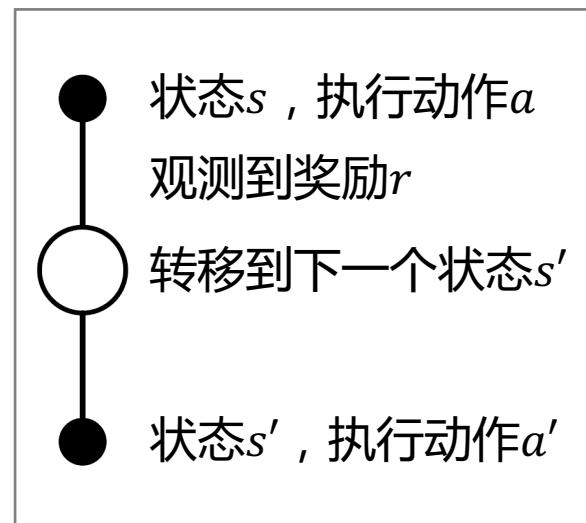
$$\pi(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$$

需要知道环境模型

- 基于 $Q$ 函数，如何选择好的行动？

$$\pi(s) = \arg \max_{a \in A} Q(s, a)$$

- 因此，估计 $Q$ 函数对直接做行动（控制）有直接的作用



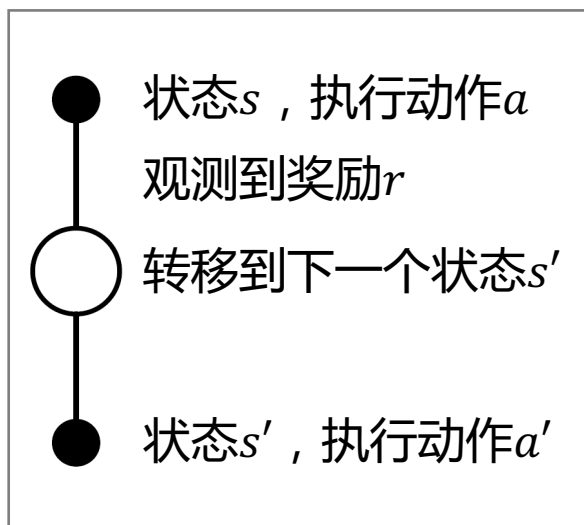


# SARSA

讲师：张伟楠 - [上海交通大学](#)

# SARSA

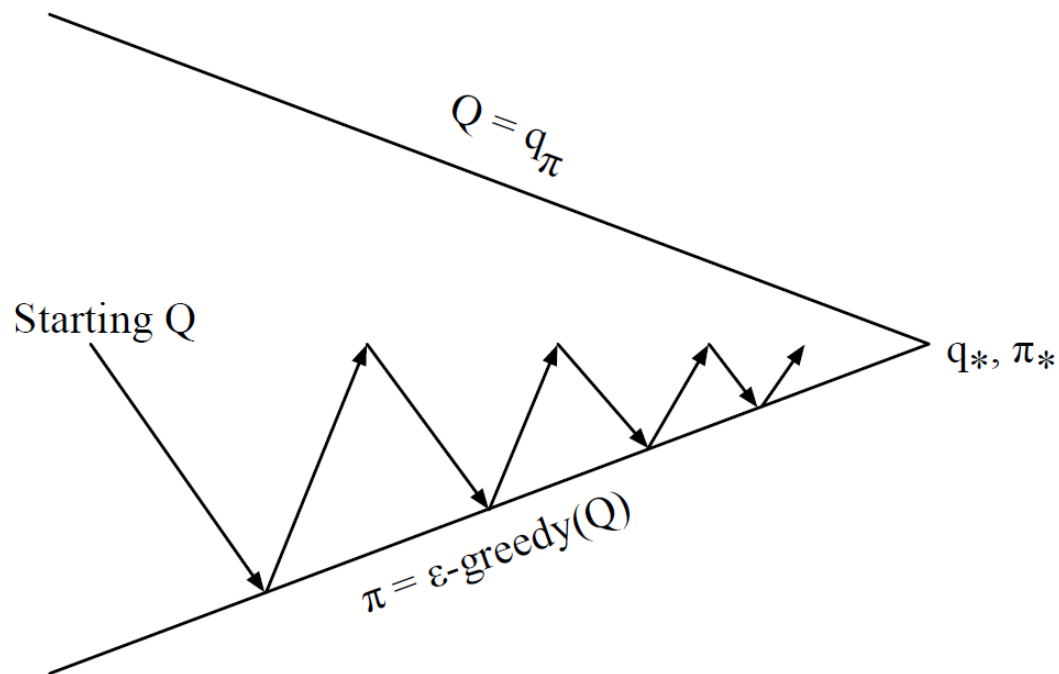
- 对于当前策略执行的每个（状态-动作-奖励-状态-动作）元组



- SARSA更新状态-动作值函数为

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$$

# 使用SARSA的在线策略控制



## □ 每个时间步长：

- 策略评估：SARSA  $Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$
- 策略改进： $\epsilon$ -greedy策略改进

# SARSA算法

## Sarsa: An on-policy TD control 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.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

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

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

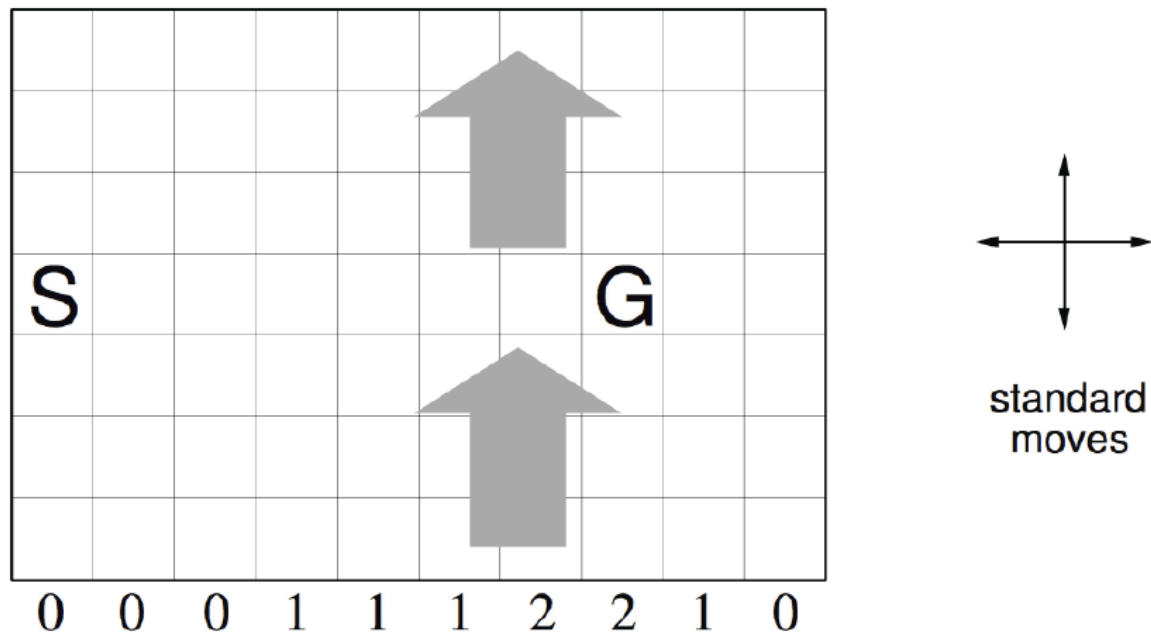
$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

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

until  $S$  is terminal

注：在线策略时序差分控制（**on-policy TD control**）使用**当前策略进行动作采样**。即，SARSA算法中的两个“A”都是由**当前策略**选择的

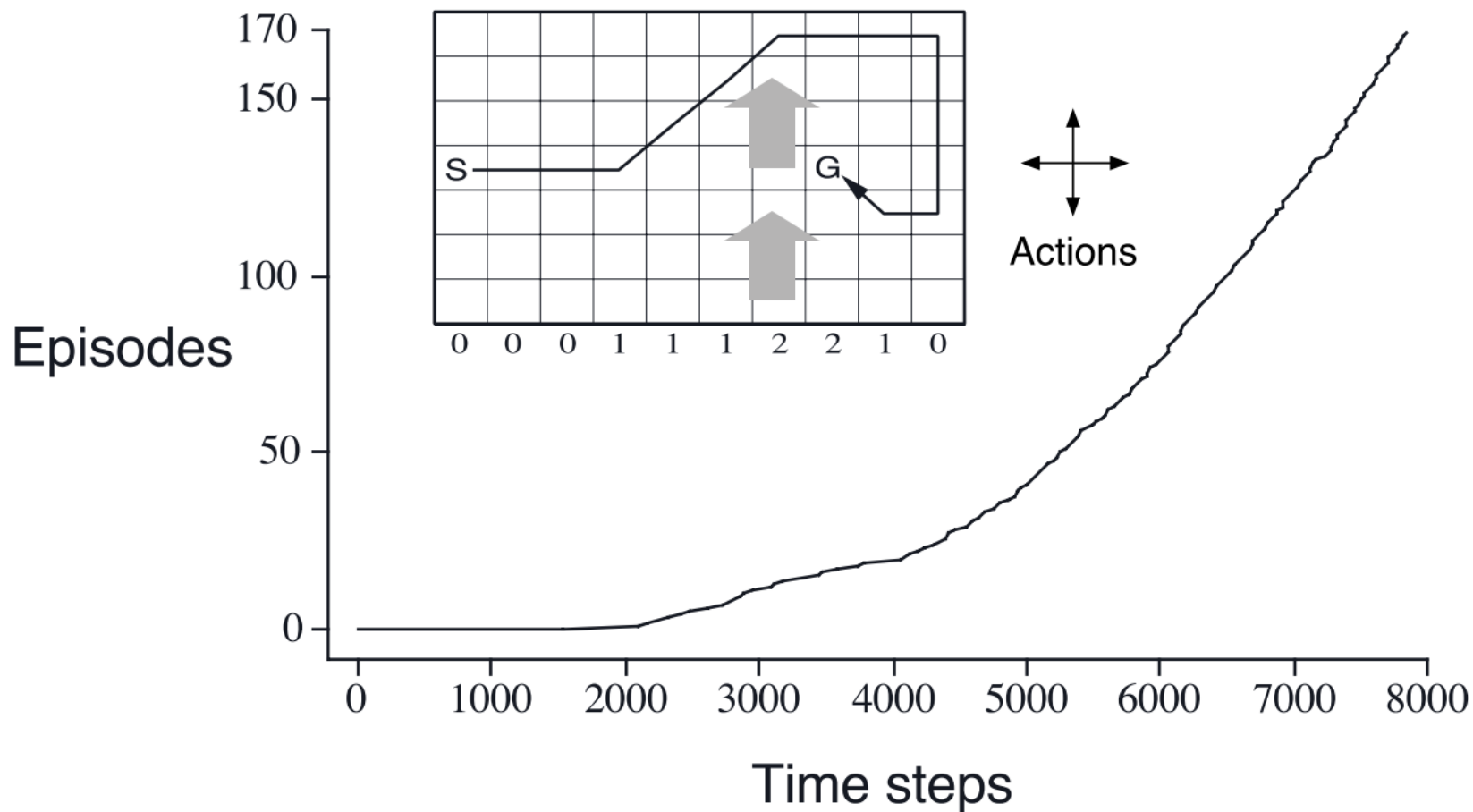
# SARSA示例：Windy Gridworld



- 每步的奖励 = -1，直到智能体抵达目标网格
- 无折扣因子



# SARSA示例：Windy Gridworld



注意：随着训练的进行，SARSA策略越来越快速地抵达目标



# Q 学习算法及其收敛性

讲师：张伟楠 – [上海交通大学](#)

# 目录

Contents

01 Q 学习

02 收敛性证明



01

Q 学习

# Q 学习

- 学习状态-动作值函数  $Q(s, a) \in \mathbb{R}$ ，不直接优化策略
- 一种离线策略 ( **off-policy** ) 学习方法

策略函数，一般是给定的策略， $\mu(\cdot | s_t) \in \mathbb{R}^{|A|}$

$$Q(s_t, a_t) = \sum_{t=0}^T \gamma^t r(s_t, a_t), a_t \sim \mu(s_t)$$

动作空间， $a \sim A$

奖励函数， $r(s_t, a_t) \in \mathbb{R}$

迭代式： $Q(s_t, a_t) = r(s_t, a_t) + \gamma Q(s_{t+1}, a_{t+1})$



# 离线策略学习

## 什么是离线策略学习

- 目标策略  $\pi(a|s)$  进行值函数评估 (  $V^\pi(s)$ 或 $Q^\pi(s, a)$  )
- 行为策略  $\mu(a|s)$  收集数据 :  $\{s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T\} \sim \mu$

## 为什么使用离线策略学习

- 平衡探索 ( [exploration](#) ) 和利用 ( [exploitation](#) )
- 通过观察人类或其他智能体学习策略
- 重用旧策略所产生的经验
- 遵循探索策略时学习最优策略
- 遵循一个策略时学习多个策略
- 在剑桥MSR研究时的一个xbox music的例子
  - Collective Noise Contrastive Estimation for Policy Transfer Learning. AAI 2016

# Q 学习

- 使用数据片段 $(s_t, a_t, r_t, s_{t+1})$ 无需重要性采样（为什么？）
- 根据行为策略选择动作  $a_t \sim \mu(\cdot | s_t)$
- 根据目标策略选择后续动作  $a'_{t+1} \sim \pi(\cdot | s_{t+1})$

- 目标  $Q^*(s_t, a_t) = r_t + \gamma Q(s_{t+1}, a'_{t+1})$

- 更新  $Q(s_t, a_t)$  的值以逼近目标状态-动作值

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (r_t + \gamma Q(s_{t+1}, a'_{t+1}) - Q(s_t, a_t))$$



策略  $\pi$  的动作，而非策略  $\mu$

## 使用 Q 学习的离线策略控制

- 允许行为策略和目标策略都进行改进
- 目标策略 $\pi$ 是关于 $Q(s, a)$ 的贪心策略

$$\pi(s_{t+1}) = \arg \max_{a'} Q(s_{t+1}, a')$$

- 行为策略 $\mu$ 是关于 $Q(s, a)$ 的 $\epsilon$ -greedy策略
- Q-学习目标函数可以简化为

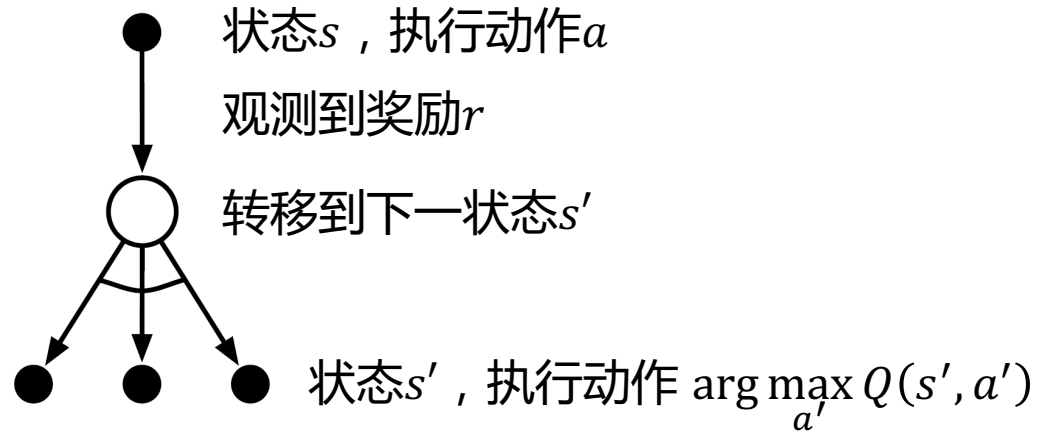
$$\begin{aligned} r_t + \gamma Q(s_{t+1}, a'_{t+1}) &= r_t + \gamma Q(s_{t+1}, \arg \max_{a'_{t+1}} Q(s_{t+1}, a'_{t+1})) \\ &= r_t + \gamma \max_{a'_{t+1}} Q(s_{t+1}, a'_{t+1}) \end{aligned}$$

- Q-学习更新方式

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a'_{t+1}} Q(s_{t+1}, a'_{t+1}) - Q(s_t, a_t))$$



# Q 学习控制算法

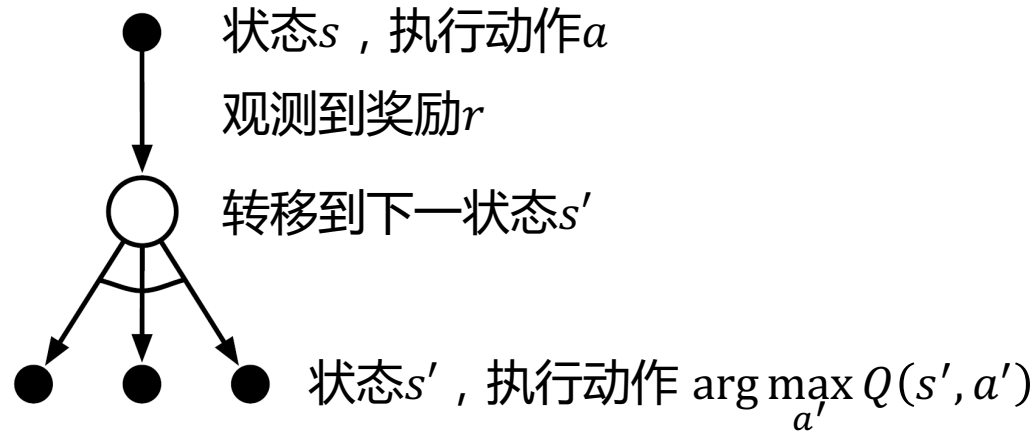


$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$$

□ 定理：Q-学习控制收敛到最优状态-动作值函数

$$Q(s, a) \rightarrow Q^*(s, a)$$

# Q 学习控制算法



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$$

## □ 为什么不需要重要性采样？

- 使用了状态-动作值函数而不是使用状态值函数
- 使用数据片段 $(s_t, a_t, r_t, s_{t+1})$

02

# 收敛性证明

(思路)

# Q 学习的收敛性

## 收缩算子 ( contraction operator )

- $Q(s, a) = r(s, a) + \gamma \max_{a'} Q(s', a')$
- 定义  $H$  算子 :  $HQ = r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} [ \max_{a'} Q(s', a') ]$
- 最优值函数  $Q^*$  是  $H$  的不动点 , 意味着 :  $Q^* = HQ^*$

# Q 学习的收敛性

## □ 直接从Q函数证明

$$(\mathbf{H}Q)(s, a) = \sum_{s' \in S} p(s'|s, a) [r(s, a) + \gamma \max_{a' \in A} Q(s', a')]$$

$$\|\mathbf{H}Q_1 - \mathbf{H}Q_2\|_\infty$$

$$= \max_{s, a} \left| \sum_{s' \in S} p(s'|s, a) [r(s, a) + \gamma \max_{a' \in A} Q_1(s', a') - r(s, a) - \gamma \max_{a' \in A} Q_2(s', a')] \right|$$

$$= \max_{s, a} \gamma \left| \sum_{s' \in S} p(s'|s, a) [\max_{a' \in A} Q_1(s', a') - \max_{a' \in A} Q_2(s', a')] \right|$$

$$\leq \max_{s, a} \gamma \sum_{s' \in S} p(s'|s, a) \left| \max_{a' \in A} Q_1(s', a') - \max_{a' \in A} Q_2(s', a') \right|$$

$$\leq \max_{s, a} \gamma \sum_{s' \in S} p(s'|s, a) \max_{s'', a''} |Q_1(s'', a'') - Q_2(s'', a'')|$$

$$= \max_{s, a} \gamma \sum_{s' \in S} p(s'|s, a) \|Q_1 - Q_2\|_\infty$$

$$= \gamma \|Q_1 - Q_2\|_\infty$$



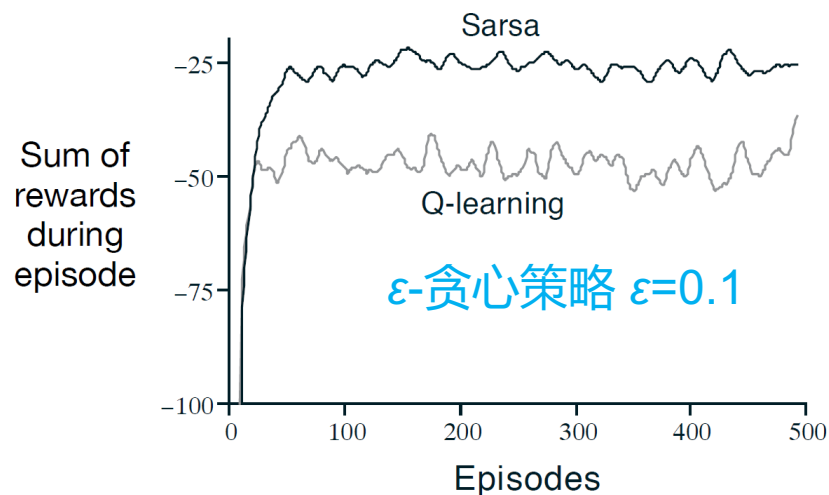
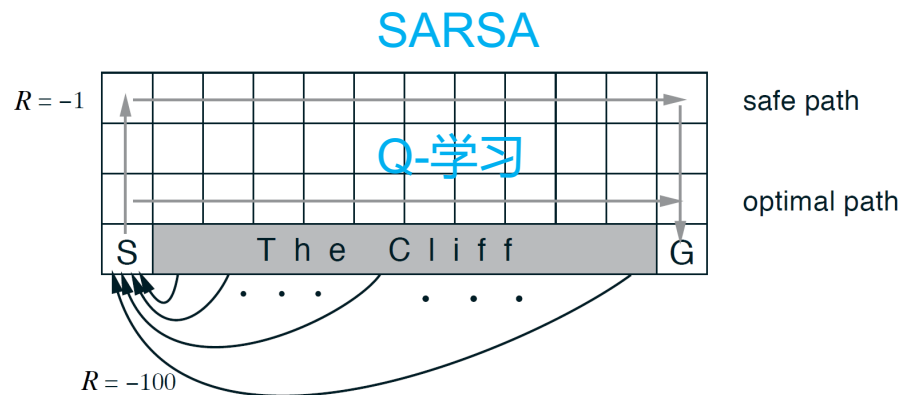
$$\|\mathbf{H}Q_1 - \mathbf{H}Q_2\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty$$

# SARSA与Q 学习对比实验

## □ 悬崖边行走 (Cliff-walking)

- 无折扣的奖励
- 片段式的任务
- 所有移动奖励 = -1
- 踏入悬崖区域会产生-100奖励并将智能体送回开始处

## □ 为什么会有图示结果？





# 多步自助法

讲师：张伟楠 - [上海交通大学](#)

# 目录

## Contents

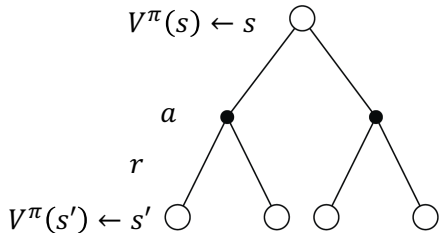
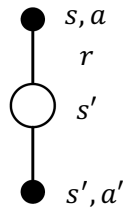
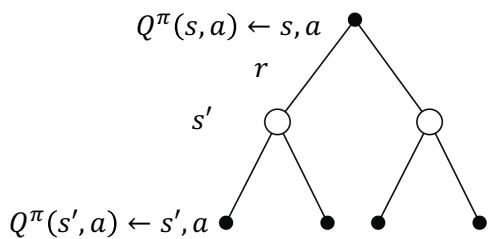
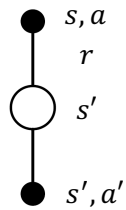
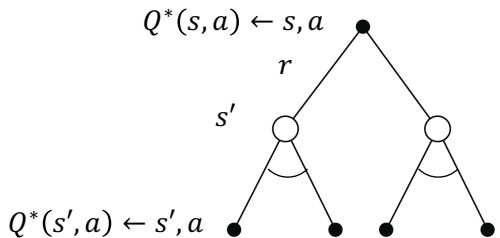
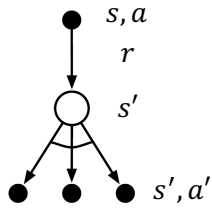
- 01 多步时序差分预测
- 02 多步Sarsa
- 03 使用重要性采样的多步离线学习法
- 04 多步备份树算法



01

**多步时序差分  
预测**

# 回顾：动态规划 (DP) 和时序差分 (TD) 的关系

	完全反向传播(DP)	采样反向传播(TD)
状态值函数的 贝尔曼 期望方程 $V^\pi(s)$	 <p>迭代的策略评估</p>	 <p>时序差分学习</p>
状态-动作值函 数的贝尔曼期 望方程 $Q^\pi(s, a)$	 <p>Q - 策略迭代</p>	 <p>SARSA</p>
状态-动作值函 数的贝尔曼 最优方程 $Q^*(s, a)$	 <p>Q - 价值迭代</p>	 <p>Q学习</p>

# 回顾：动态规划 (DP) 和时序差分 (TD) 的关系

完全反向传播(DP)	采样反向传播(TD)
迭代的策略评估 $V(s) \leftarrow \mathbb{E}[R + \gamma V(S')   s]$	时序差分学习 $V(s) \stackrel{\alpha}{\leftarrow} r + \gamma V(s')$
$Q$ - 策略迭代 $Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A')   s, a]$	SARSA $Q(s, a) \stackrel{\alpha}{\leftarrow} r + \gamma Q(s', a')$
$Q$ - 价值迭代 $Q(s, a) \leftarrow \mathbb{E} \left[ R + \gamma \max_{a'} Q(S', a') \mid s, a \right]$	$Q$ 学习 $Q(s, a) \stackrel{\alpha}{\leftarrow} r + \gamma \max_{a'} Q(s', a')$

其中  $x \stackrel{\alpha}{\leftarrow} y \equiv x \leftarrow x + \alpha(y - x)$

# 回顾：蒙特卡洛方法&时序差分法

## 蒙特卡洛方法

$$V(s_t) \leftarrow V(s_t) + \alpha(g_t - V(s_t))$$

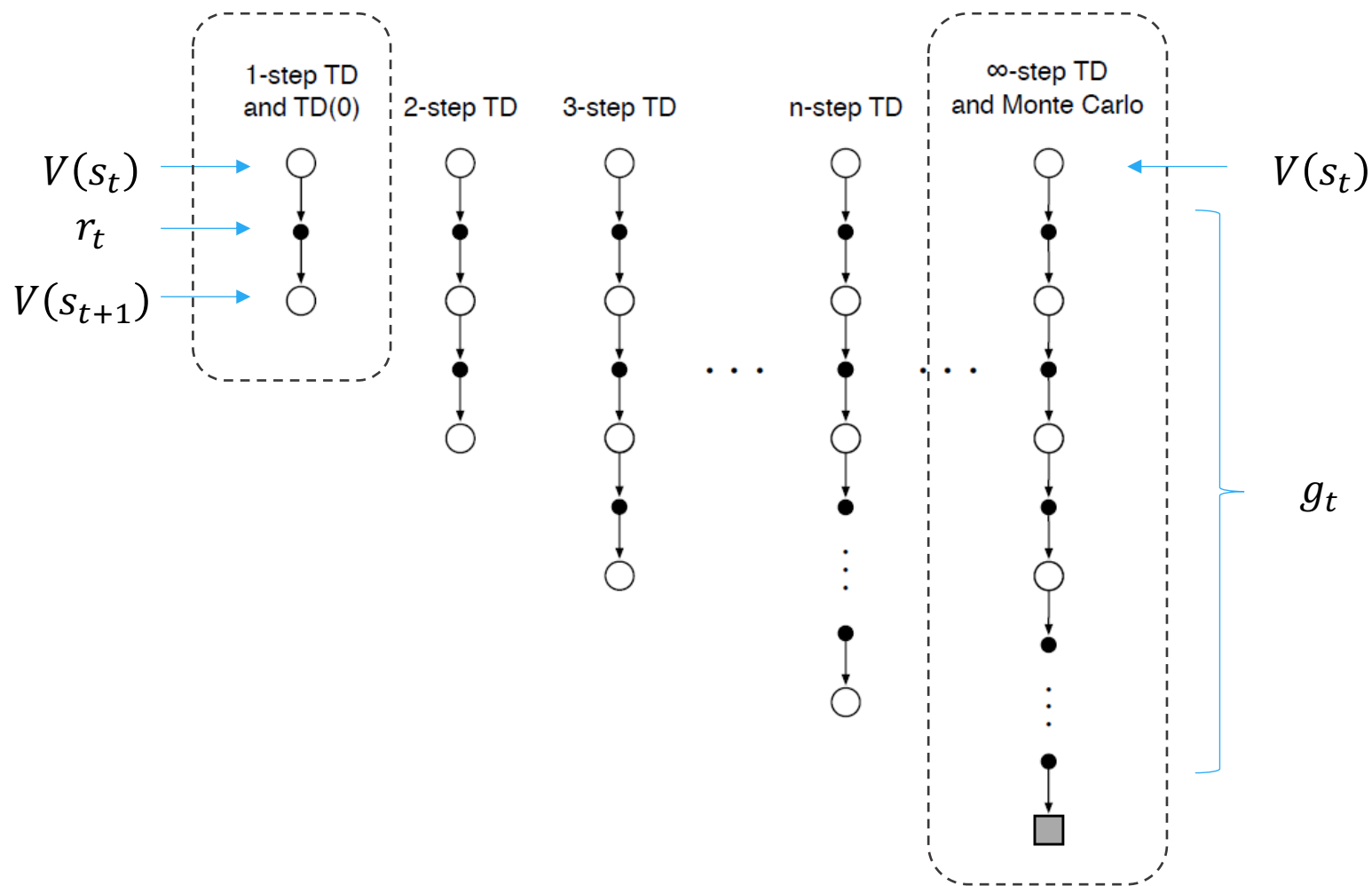
$$g_t = r_t + \gamma r_{t+1} + \dots \gamma^{T-t} r_T$$

## 时序差分法

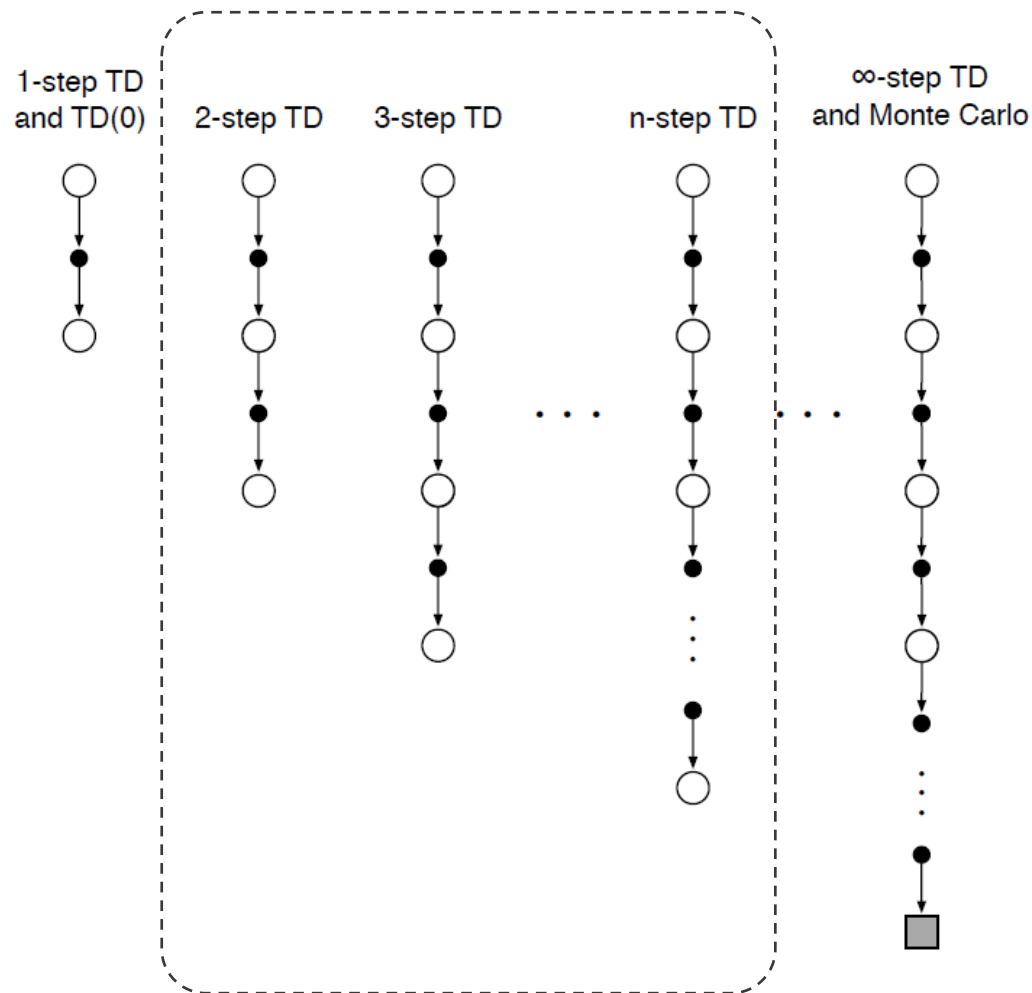
$$V(s_t) \leftarrow V(s_t) + \alpha(r_t + \gamma V(s_{t+1}) - V(s_t))$$

有没有方法介于两者之间？

# 多步时序差分预测



# 多步时序差分预测



# $n$ 步累计奖励

- 考虑下列  $n$  步累计奖励,  $n = 1, 2, \dots, \infty$

$$n = 1 \quad (\text{TD}) \quad g_t^{(1)} = r_t + \gamma V(r_{t+1})$$

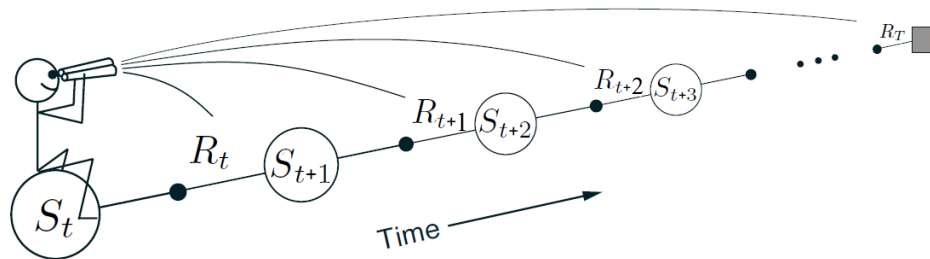
$$n = 2 \quad g_t^{(2)} = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2})$$

$\vdots$   $\vdots$

$$n = \infty \quad (\text{MC}) \quad g_t^{(\infty)} = r_t + \gamma r_{t+1} + \dots + \gamma^{T-t} r_T$$

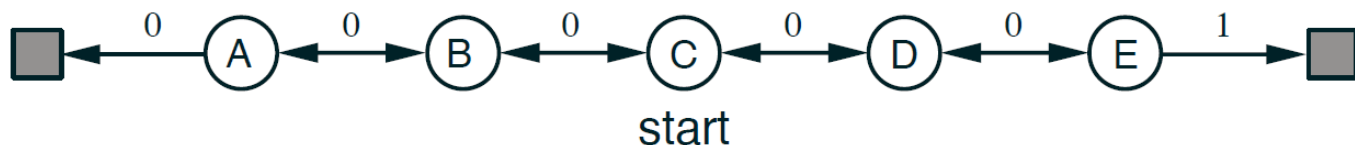
- 定义  $n$  步累计奖  $g_t^{(n)} = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n V(s_{t+n})$

- $n$  步时序差分学习  $V(s_t) \leftarrow V(s_t) + \alpha (g_t^{(n)} - V(s_t))$



# 随机游走的例子里使用 $n$ 步时序差分

## 随机游走

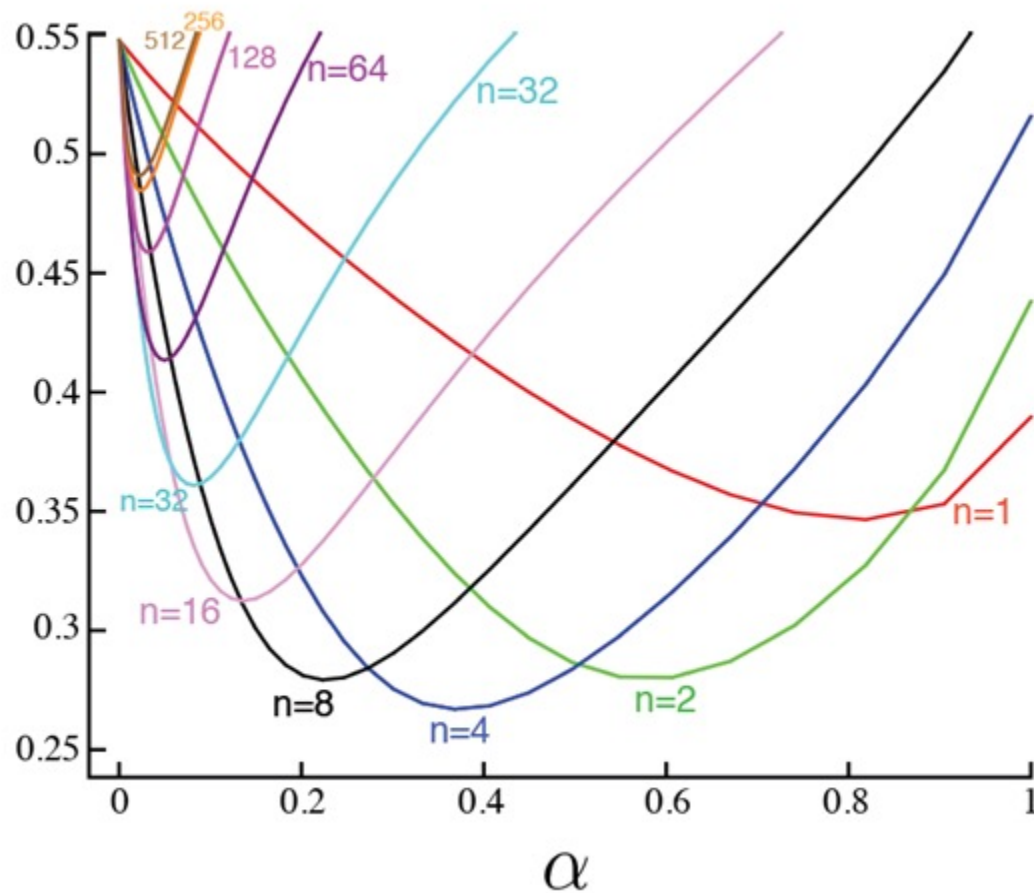


- 每个片段都从中间状态  $C$  开始。
- 每一时刻都有均等的概率向左走或者向右走。
- 到最左侧或者最右侧时，片段结束。
- 如果走到最右侧，得到的奖励为1；如果走到最左侧，得到的奖励为0。



# 多步时序差分法的表现

拥有19个状态的随机游走游戏中，在10个片段(episode)结束时的RMS误差

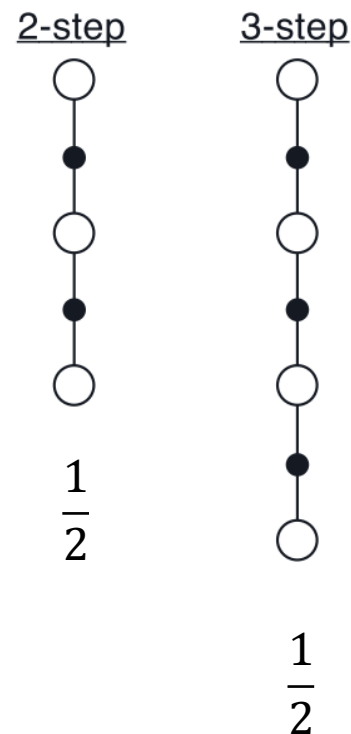


# 平均 $n$ 步累计奖励

- 我们可以进一步对不同  $n$  下的  $n$  步累计奖励求平均值
- 例如，求 2 步和 3 步时的平均累计奖励

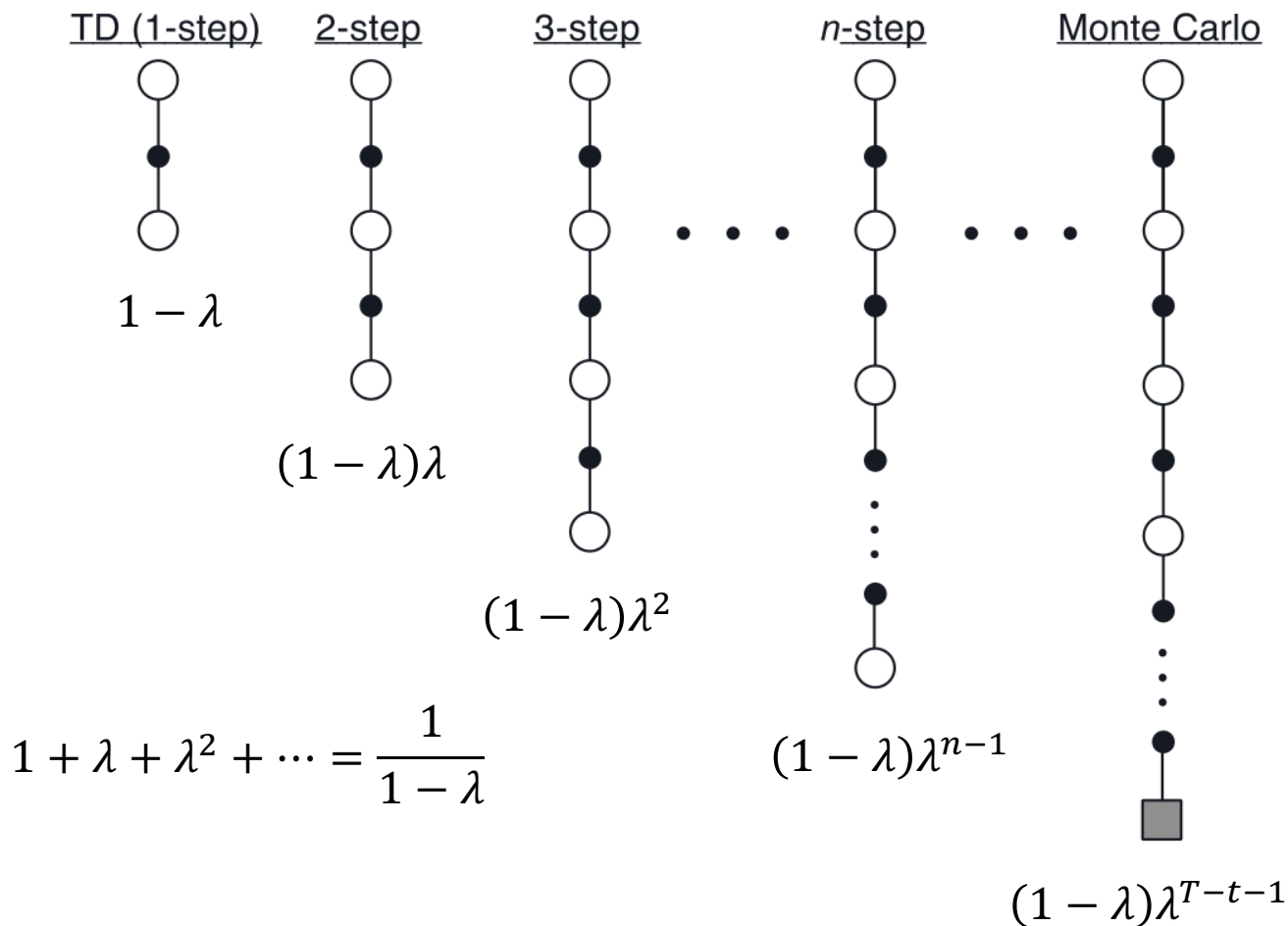
$$\frac{1}{2}g^{(2)} + \frac{1}{2}g^{(3)}$$

- 上式结合两种不同时间步长的信息
- 我们是否能够结合所有不同时间步长的信息呢？



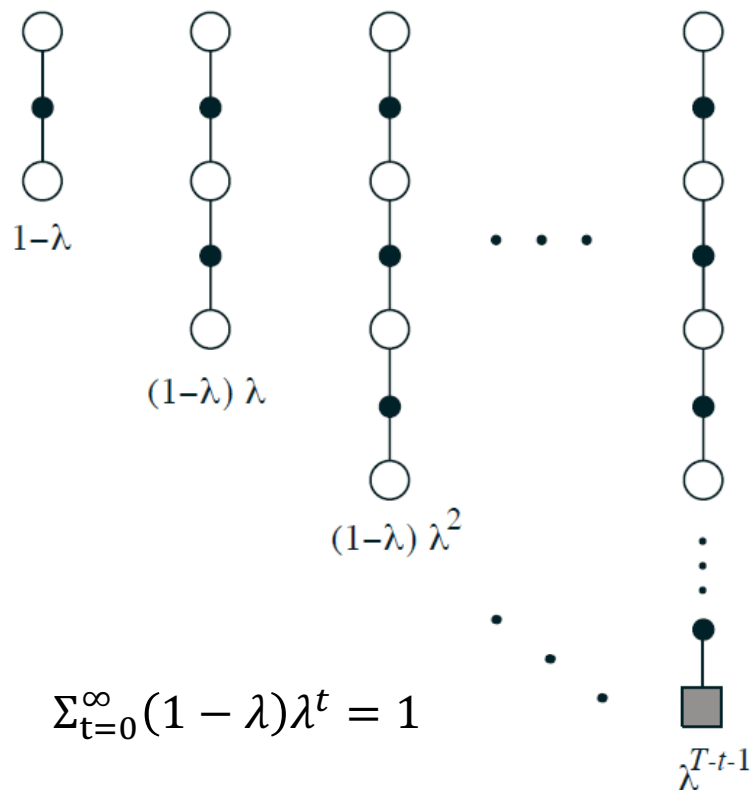
# 使用平均 $n$ 步累计奖励的 TD( $\lambda$ ) 算法

TD( $\lambda$ ),  $\lambda$  - 累计奖励



# 使用平均 $n$ 步累计奖励的 TD( $\lambda$ ) 算法

TD( $\lambda$ ),  $\lambda$  - 累计奖励

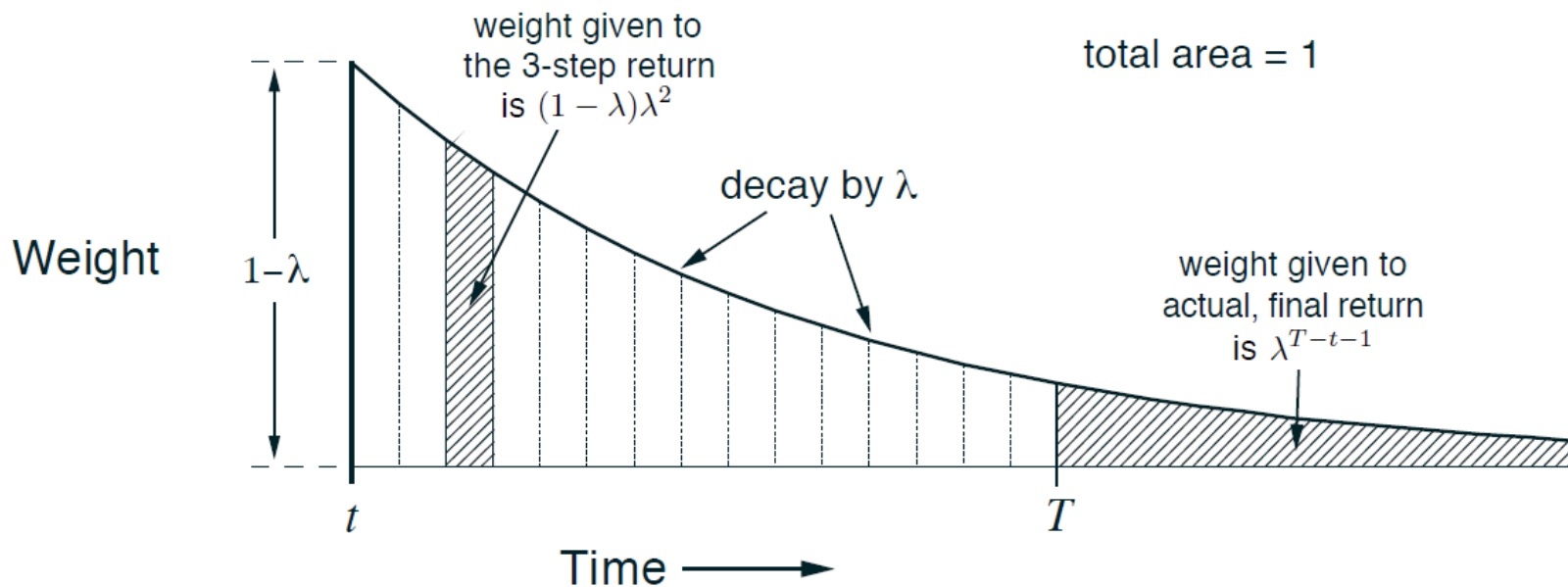


- $\lambda$  - 累计奖励  $g_t^\lambda$  结合了所有  $n$  步  
累计奖励  $g_t^{(n)}$
- 使用权重  $(1-\lambda)\lambda^{n-1}$

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

- TD( $\lambda$ )  
 $V(s_t) \leftarrow V(s_t) + \alpha(g_t^\lambda - V(s_t))$

# 使用平均 $n$ 步累计奖励的 TD( $\lambda$ ) 算法

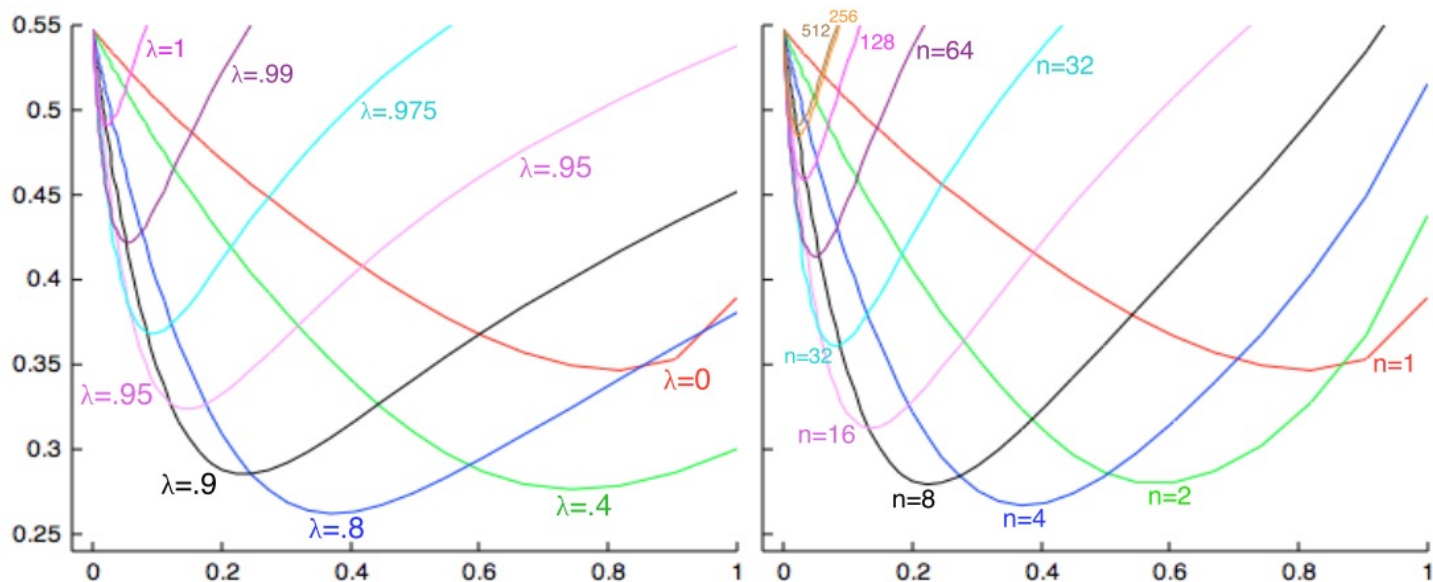


$$g_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} g_t^{(n)} + \lambda^{T-t-1} g_t$$

- 当  $\lambda = 1$  时,  $g_t^\lambda = g_t$ , 相当于蒙特卡洛方法
- 当  $\lambda = 0$  时,  $g_t^\lambda = g_t^{(1)}$ , 相当于单步时序差分方法

# TD( $\lambda$ ) vs. $n$ 步时序差分

前10个片段  
(episode)  
结束时的  
Offline  
RMS误差



19个状态的随机游走实验结果

- 在使用最佳的  $\alpha$  与  $\lambda$  值时，离线  $\lambda$ -累计奖励算法具有稍好一些的实验结果

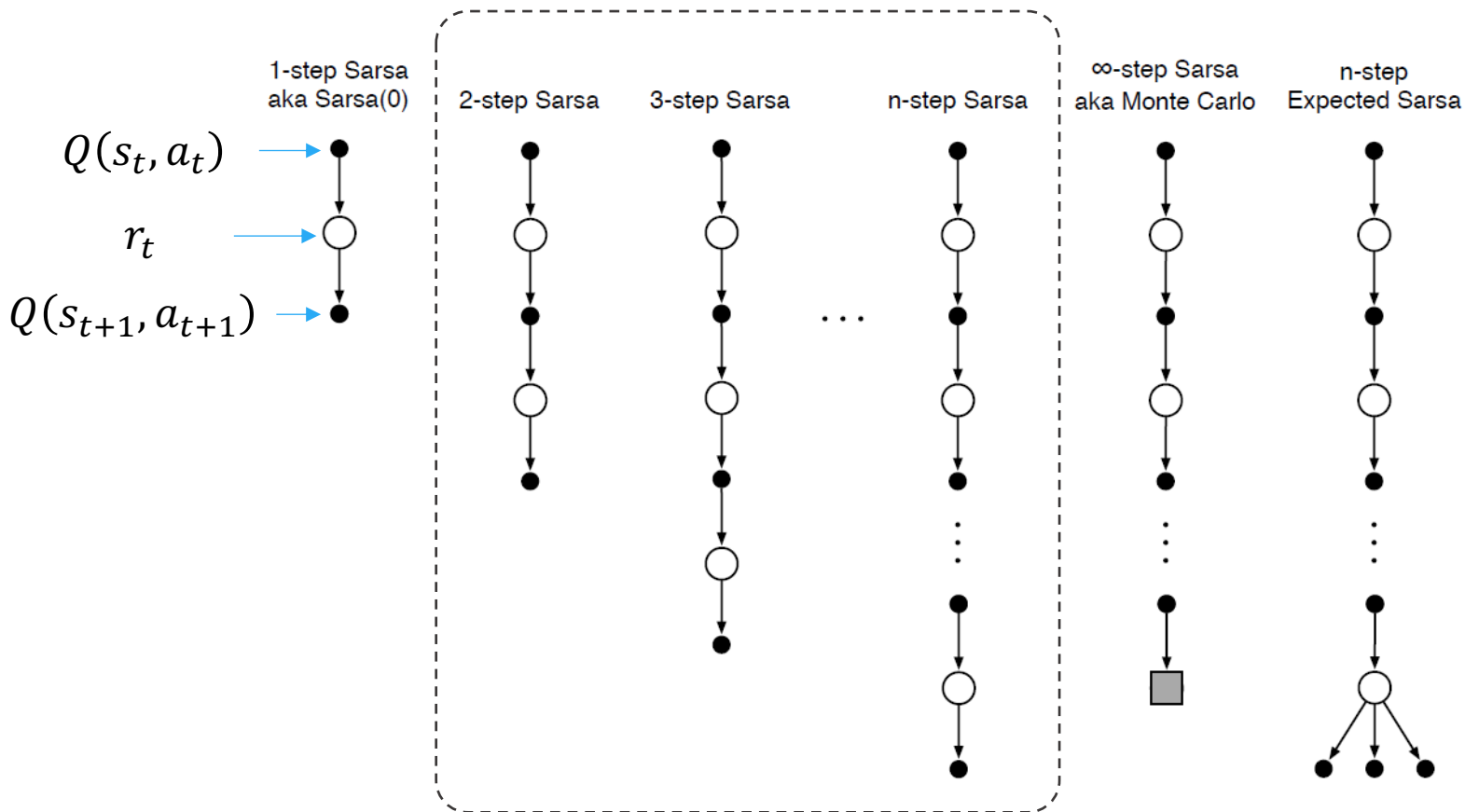
02

# 多步Sarsa

# 多步Sarsa

□  $n$  步的思想是否只能用在预测上？我们能不能把这种思想放在控制上？

将这种思想用在 Sarsa 上，我们得到了  $n$  步 Sarsa 算法





## 多步Sarsa

---

- 类似与  $n$  步TD , 我们有 :

$$g_{t:t+n} = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n Q(s_{t+n}, a_{t+n})$$

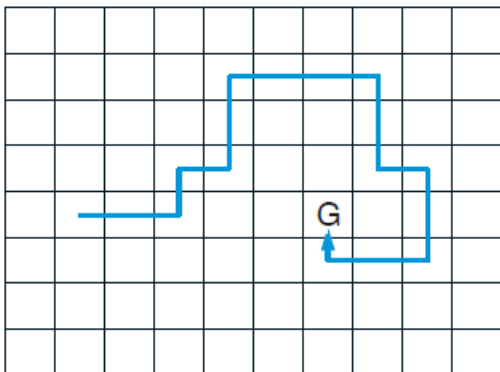
其中  $n \geq 1, 0 \leq t < T - n$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [g_{t:t+n} - Q(s_t, a_t)]$$

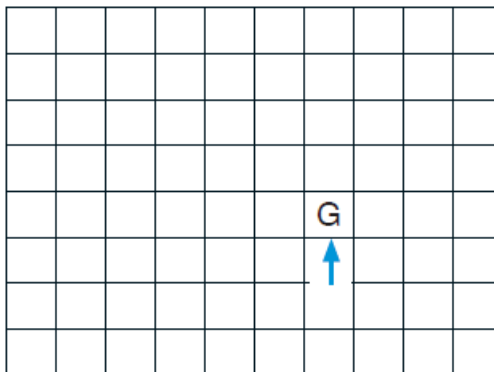
其中  $0 \leq t < T - n$

# 多步Sarsa的例子

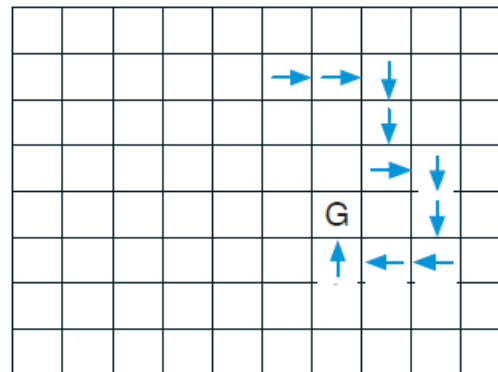
Path taken



Action values increased by one-step Sarsa



Action values increased by 10-step Sarsa



03

**使用重要性采样的  
多步离线学习法**

## 回顾：采用重要性采样的离线学习

- 在离线学习的场景下，往往会有两个策略：
  - $\pi$ : 通常是对于目前的动作值估计函数的一个贪婪策略（也有可能是一个随机策略）
  - $b$ : 用来收集数据的，更有探索性的一个策略（常用  $\epsilon$ -greedy 来实现）
- 由于收集数据的策略和我们要实际优化的策略并不是同一个策略，我们需要加上一项修正（以值函数为例）：

$$\Delta V' = \frac{\pi(a_t|s_t)}{b(a_t|s_t)} \Delta V$$

## 回顾：采用重要性采样的离线学习

- 在重要性采样的角度对值函数的推导：

$$\mathbb{E}_{s \sim \pi}[f(s)] = \int p(s|\pi)f(s)ds = \int \frac{p(s|\pi)}{p(s|b)}p(s|b)f(s)ds = \mathbb{E}_{s \sim b} \left[ \frac{p(s|\pi)}{p(s|b)} f(s) \right]$$

$$P(a_t, s_{t+1}, a_{t+1}, \dots, s_h | s_t, a_{t:h} \sim b) = \prod_{k=t}^{h-1} b(a_k | s_k) p(s_{k+1} | s_k, a_k)$$

$$P(a_t, s_{t+1}, a_{t+1}, \dots, s_h | s_t, a_{t:h} \sim \pi) = \prod_{k=t}^{h-1} \pi(a_k | s_k) p(s_{k+1} | s_k, a_k)$$

$$g_t^{(n)} = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n \underline{V(s_{t+n})}$$

➔ 
$$\rho_{t:h} = \prod_{k=t}^{\min(h, T-1)} \frac{\pi(a_k | s_k)}{b(a_k | s_k)}$$

$$V(s_t) \leftarrow V(s_t) + \alpha \rho_{t:t+n-1} [g_t^{(n)} - V(s_t)]$$

# 使用重要性采样的多步离线学习法

$$\rho_{t:h} = \prod_{k=t}^{\min(h, T-1)} \frac{\pi(a_k | s_k)}{b(a_k | s_k)}$$

- 对于状态值函数  $V$  :

$$V(s_t) \leftarrow V(s_t) + \alpha \rho_{t:t+n-1} [g_t^{(n)} - V(s_t)]$$

- 对于动作值函数  $Q$  :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \rho_{t:t+n-1} [g_t^{(n)} - Q(s_t, a_t)]$$

# 使用重要性采样的多步离线学习法

## Algorithm 1: 离线 $n$ 步Sarsa算法

**Input:** 任意探索性的策略 $b$ ,  $Q$ 函数, 目标策略 $\pi$ , 常数 $n$

**Output:** 估计的 $Q$ 函数,  $Q \approx q_\pi$

```
1 初始化;
2 for 对每个 episode do
3   for  $t = 0, 1, 2, \dots$  do
4     if  $t < T$  then
5       Take  $A_t$ , get  $(R_{t+1}, S_{t+1})$ ;
6       if  $S_{t+1}$  是终止状态 then
7          $T \leftarrow t + 1$ ;
8       else
9         根据策略  $b(\cdot | S_{t+1})$  选择  $A_{t+1}$  并保存;
10       $\tau \leftarrow t - n + 1$  if  $\tau \geq 0$  then
11         $\rho \leftarrow \prod_{i=\tau+1}^{\min(\tau+n-1, T-1)} \frac{\pi(A_i | S_i)}{b(A_i | S_i)}$ ;
12         $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ ;
13        if  $\tau + n < T$  then
14           $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ ;
15         $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha \rho [G - Q(S_\tau, A_\tau)]$ ;
16        更新  $\pi$ , 确保其对  $Q$  具有贪婪性质;
```

Sarsa 算法框架

计算重要性采样系数

计算  $n$  步目标

更新  $Q$  函数

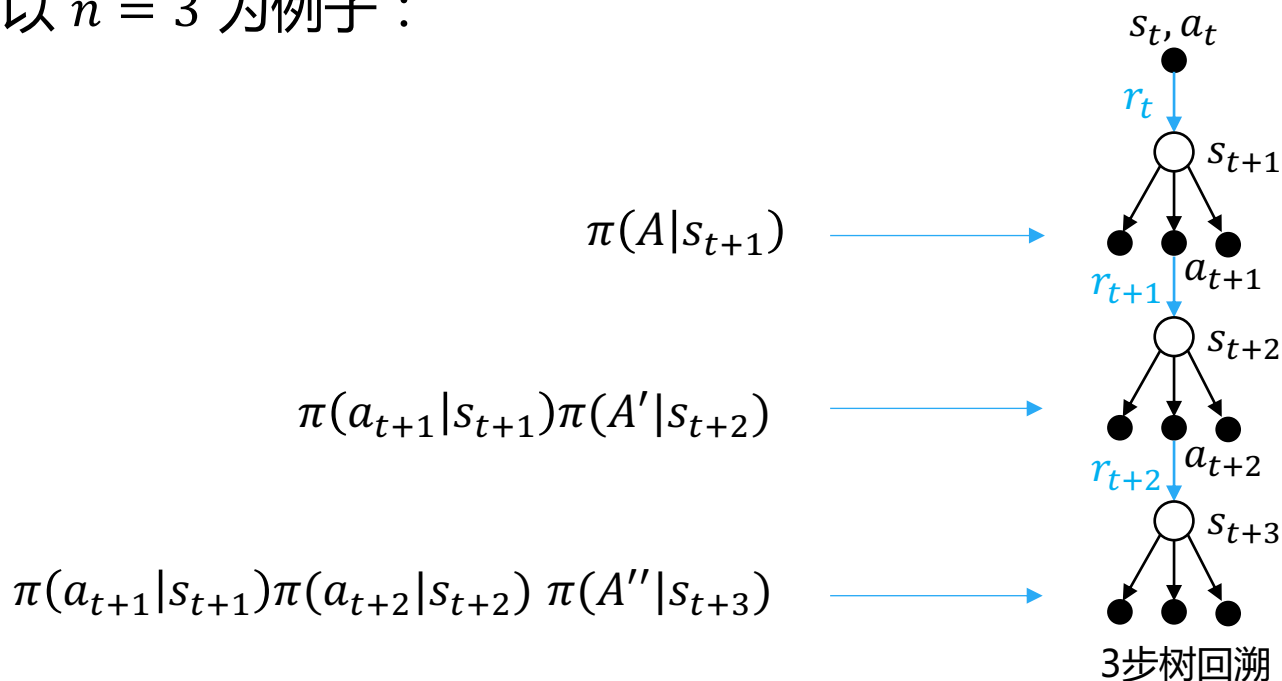
04

# 多步树回溯算法



# 多步树回溯算法 (N-step Tree Backup Algorithm)

- 是否存在一种不需使用重要性采样的离线算法？
  - $Q$ 学习和期望 Sarsa (expected Sarsa)
- 我们可以将这种形式扩展到  $n$  步，称之为多步树回溯算法。
- 以  $n = 3$  为例子：



# 多步树回溯算法推导

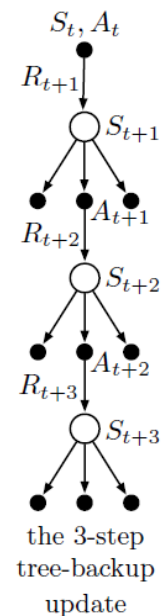
□ 对于  $n = 1$  : 
$$g_{t:t+1} = r_t + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a)$$

□ 对于  $n = 2$  :

$$\begin{aligned} g_{t:t+2} &= r_t + \gamma \sum_{a \neq a_{t+1}} \pi(a|s_{t+1}) Q(s_{t+1}, a) \\ &\quad + \gamma \pi(a_{t+1}|s_{t+1}) \left( r_{t+1} + \gamma \sum_a \pi(a|s_{t+2}) Q(s_{t+2}, a) \right) \\ &= r_t + \gamma \sum_{a \neq a_{t+1}} \pi(a|s_{t+1}) Q(s_{t+1}, a) + \gamma \pi(a_{t+1}|s_{t+1}) g_{t+1:t+2} \end{aligned}$$

□ 推广到  $n$  步 :

$$\begin{aligned} g_{t:t+n} &= r_t + \gamma \sum_{a \neq a_{t+1}} \pi(a|s_{t+1}) Q(s_{t+1}, a) + \gamma \pi(a_{t+1}|s_{t+1}) g_{t+1:t+n} \\ Q(s_t, a_t) &\leftarrow Q(s_t, a_t) + \alpha (g_{t:t+n} - Q(s_t, a_t)) \end{aligned}$$



# 无模型强化学习总结

---

- 无模型强化学习在黑盒环境下使用，往往具有更广的应用范围
- 蒙特卡洛方法通过采样直接估计价值函数
- 时序差分学习通过下一步的价值估计来更新当前一步的价值估计
- 无模型强化学习的on-policy和off-policy算法区别
  - SARSA、Q学习
- 多步自助法可看作介于蒙特卡洛方法和单步时序差分方法之间的一种方法
- 讨论：无模型强化学习和基于模型的强化学习方法之间的对比

THANK YOU