

多智能体强化学习实训

Multi-agent Reinforcement Learning

Lecture 2: Value Estimation

教师：张寅

zhangyin98@zju.edu.cn

助教：邓悦

devindeng@zju.edu.cn

王子瑞

ziseoiwong@zju.edu.cn

李奕澄

yichengli@zju.edu.cn

浙江大学计算机学院

课程大纲

- 蒙特卡洛方法
- 时序差分方法
- 资格迹方法
- 表格型时序差分方法：
 - Q-learning
 - SARSA

- 在现实问题中，通常不能假设对环境有完全的了解，即无法明确地给出状态转移和奖励函数因此需要一种直接从经验数据中学习价值和策略，无需构建马尔可夫决策过程模型的方法

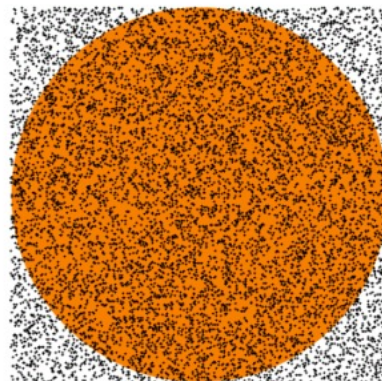
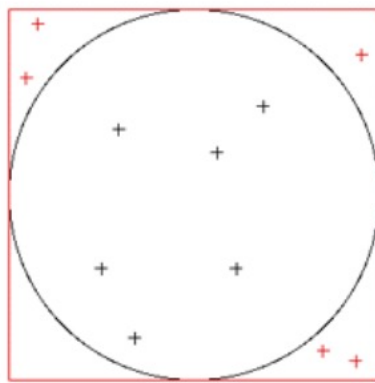
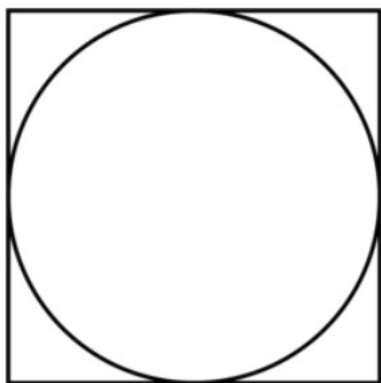


状态空间： $19^2! = 361! \approx 10^{768}$

- 蒙特卡洛方法是一种基于随机采样和统计的方法，得名于摩纳哥的蒙特卡洛赌场，该方法的特点是大量的随机模拟。蒙特卡洛方法在强化学习中的基本思想是通过多次采样来估计状态或动作的值函数，随后利用值函数进行策略改进

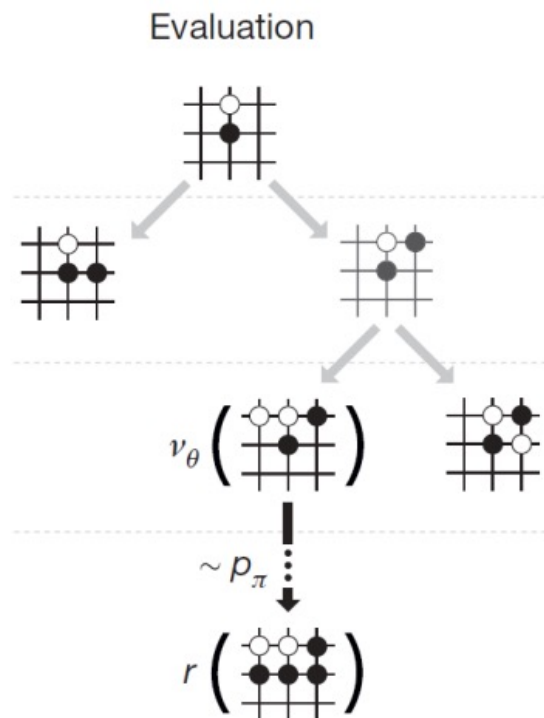
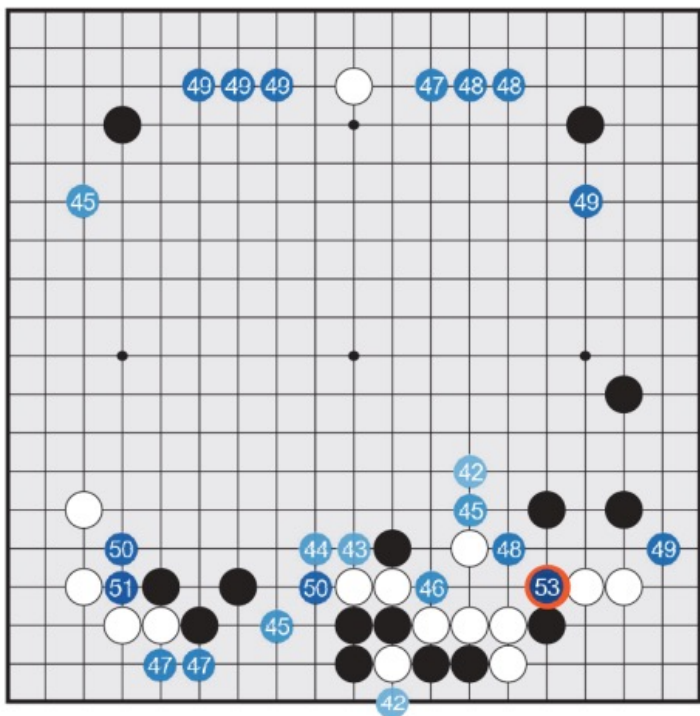


- 蒙特卡洛方法（Monte-Carlo methods，简称 MC）是一类广泛的计算算法蒙特卡洛方法依赖重复随机抽样来获得数值结果
- 例如，计算圆的面积



$$\text{圆的面积} \approx \text{正方形面积} \times \frac{\text{圆内点的数量}}{\text{圆外点的数量}}$$

- 例如，AlphaGo 中 AI 对围棋落子胜率的估计



当前棋局的胜率 $\approx \frac{\text{对弈胜利局数}}{\text{对弈总局数}}$ (从当前棋局出发)

- 目标：从策略 π 采样的历史经验中估计 V^π
- 回顾设定：累计奖励（return）是总折扣奖励

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T$$

- 回顾设定：值函数（value function）是期望累计奖励

$$V^\pi(s) = \mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots | s_0 = s, \pi]$$

$$= \mathbb{E}[G_t | s_0 = s, \pi] \approx \frac{1}{N} \sum_{i=1}^N G_t^i$$

- 目标：从策略 π 采样的历史经验中估计 V^π
- 回顾设定：累计奖励（return）是总折扣奖励

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T$$

- 回顾设定：值函数（value function）是期望累计奖励

$$V^\pi(s) = \mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots | s_0 = s, \pi]$$

$$= \mathbb{E}[G_t | s_0 = s, \pi] \approx \frac{1}{N} \sum_{i=1}^N G_t^i$$

蒙特卡洛方法使用策略 π 从状态 s 采样 N 个样本，并使用经验均值累计奖励近似期望累计奖励

- 具体实现：使用策略 π 采样时间步数量为 T 的回合

$$s_0^i \rightarrow \underset{a_0^i}{R_1^i} \rightarrow s_1^i \rightarrow \underset{a_1^i}{R_2^i} \rightarrow \dots \rightarrow s_T^i \sim \pi$$

- 对于每一个回合中的时间步 t 中的状态 s
 - 更新访问次数 $N(s) \leftarrow N(s) + 1$
 - 更新累计奖励return 的总和 $S(s) \leftarrow S(s) + G_t$
 - 估计累计奖励return 的均值 $V(s) = S(s)/N(s)$
 - 由大数定律，当 $N(s) \rightarrow \infty$ 有 $V(s) \rightarrow V^\pi(s)$

- 将公式整理为增量更新的形式，可得

$$N(s_t) \leftarrow N(s_t) + 1$$

$$V(s_t) \leftarrow V(s_t) + (G_t - V(s_t))/N(s_t)$$

- 对于非平稳的环境（即环境的动态会随时间发生变化），蒙特卡洛方法可以跟踪一个滑动窗口内的平均值

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

Sliding window -->



Slide one element forward



- 直接从经验回合进行学习，不需要模拟/搜索
- 模型无关（model-free），无需环境信息
- 核心思想简单直白：value = mean return
- 使用完整回合进行更新：只能应用于有限长度的马尔可夫决策过程，即所有的回合都应有终止状态

- 示例：二十一点 (Blackjack)
 - 状态空间
 - 当前点数之和 (12~21)
 - Dealer展示手牌 (Ace~10)
 - 玩家是否有Ace (Yes or No)



- 示例：二十一点 (Blackjack)
 - 动作空间：
 - “停止”：停止接收卡（回合终止）
 - “继续”：再拿一张卡（无放回）



蒙特卡洛方法



腾讯开悟

- 示例：二十一点 (Blackjack)
 - 奖励设定 (“停止” 动作) :
 - +1: 卡牌点数总和>Dealer
 - 0: 卡牌点数总和=Dealer
 - -1: 卡牌点数总和<Dealer



- 示例：二十一点（Blackjack）
 - 奖励设定（“继续”动作）：
 - -1: 卡牌点数总和 >21 （游戏终止）
 - 0: 其他情况
 - 如果卡牌点数 <12 自动执行“继续”动作

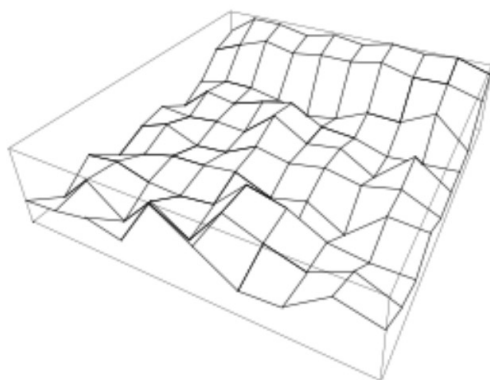


- 给定策略：若手牌点数 ≥ 20 ，执行“停止”动作，否则执行“继续”动作

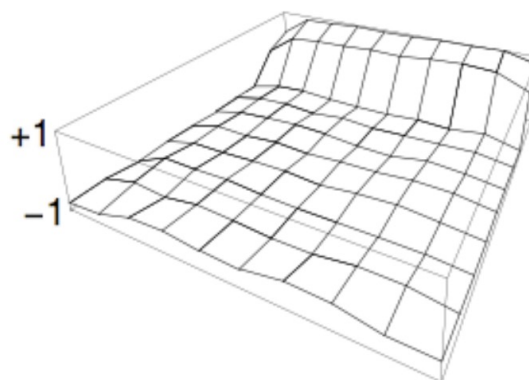
After 10,000 episodes

After 500,000 episodes

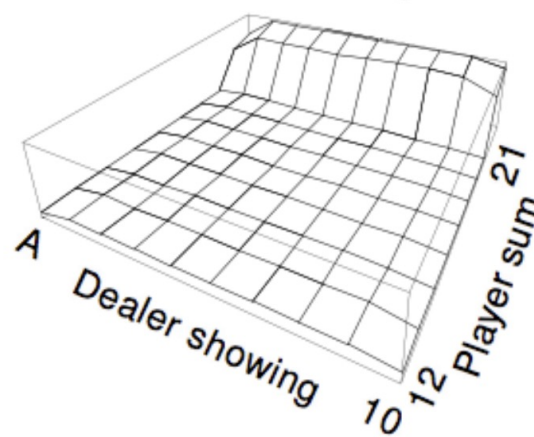
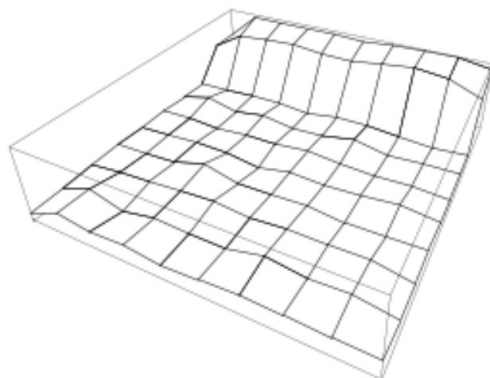
Usable
ace



+1
-1



No
usable
ace



- 蒙特卡洛方法
- **时序差分方法**
- 资格迹方法
- 表格型时序差分方法：
 - Q-learning
 - SARSA

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T = R_{t+1} + \gamma V(s_{t+1})$$

$$V(s_t) \leftarrow V(s_t) + \alpha (\underset{\substack{\uparrow \\ \text{观测值}}}{R_{t+1}} + \gamma \underset{\substack{\uparrow \\ \text{对未来的猜测}}}{V(s_{t+1})} - V(s_t))$$

- 时序差分方法 (Temporal Difference methods, 简称 TD)
能够直接使用经验回合学习, 同样也是模型无关的

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T = R_{t+1} + \gamma V(s_{t+1})$$

$$V(s_t) \leftarrow V(s_t) + \alpha(\underset{\substack{\uparrow \\ \text{观测值}}}{R_{t+1}} + \gamma \underset{\substack{\uparrow \\ \text{对未来的猜测}}}{V(s_{t+1})} - V(s_t))$$

- 时序差分方法 (Temporal Difference methods, 简称 TD) 能够直接使用经验回合学习, 同样也是模型无关的
- 与蒙特卡洛方法不同, 时序差分方法结合了自举 (bootstrapping), 能从不完整的回合中学习

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T = R_{t+1} + \gamma V(s_{t+1})$$

$$V(s_t) \leftarrow V(s_t) + \alpha \left(\underset{\substack{\uparrow \\ \text{观测值}}}{R_{t+1}} + \gamma \underset{\substack{\uparrow \\ \text{对未来的猜测}}}{V(s_{t+1})} - V(s_t) \right)$$

- 时序差分方法（Temporal Difference methods, 简称 TD）能够直接使用经验回合学习，同样也是模型无关的
- 与蒙特卡洛方法不同，时序差分方法结合了自举（bootstrapping），能从不完整的回合中学习
- 时序差分通过更新当前预测值，使之接近估计出来的累计奖励，而非真实累计奖励



- 目标：从策略 π 采样的历史经验中估计 V^π

- 蒙特卡洛方法更新值函数 $V(s_t)$ 使其接近 G_t

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

- 时序差分方法更新值函数 $V(s_t)$ 使其接近 $R_{t+1} + \gamma V(s_{t+1})$

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

- $R_{t+1} + \gamma V(s_{t+1})$ 为时序差分目标

- $\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$ 为时序差分误差

时序差分方法 vs. 蒙特卡洛方法



腾讯开悟

状态	经过的时间 (分钟)	预计所剩时间	预计总时间
离开公司	0	30	30
开始驾车, 下雨	5	35	40
离开高速公路	20	15	35
卡车后跟车	30	10	40
到达家所在街道	40	3	43
直奔家门	43	0	43

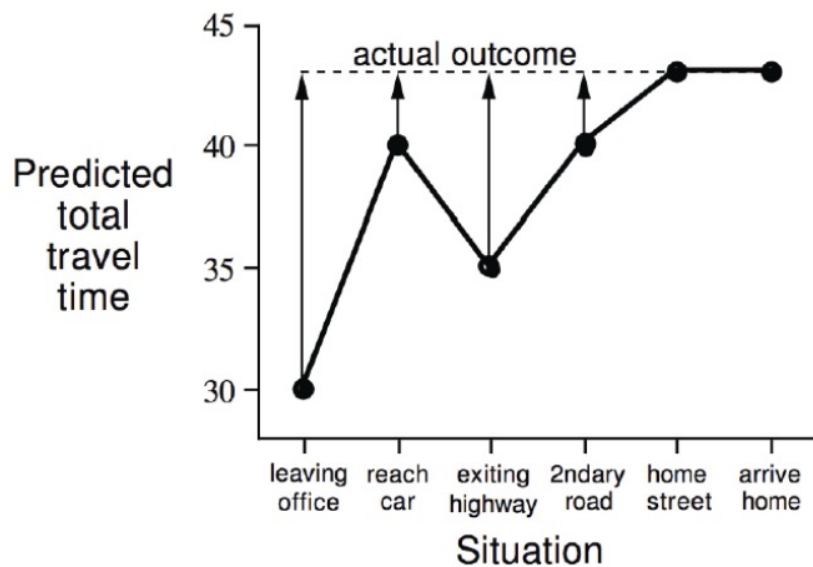


时序差分vs. 蒙特卡洛

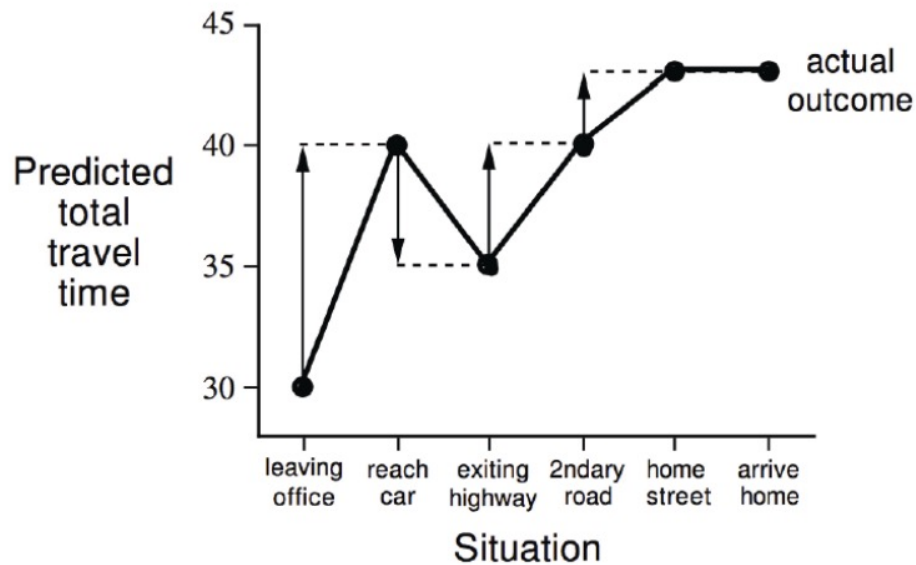


腾讯开悟

蒙特卡洛方法



时序差分方法



- 时序差分方法能够在每一步之后进行在线学习，蒙特卡洛方法必须等待回合终止，直到累计奖励已知

- 时序差分方法能够在每一步之后进行在线学习，蒙特卡洛方法必须等待回合终止，直到累计奖励已知
- 时序差分方法能够从不完整的序列中学习，蒙特卡洛方法只能从完整序列中学习

- 时序差分方法能够在每一步之后进行在线学习，蒙特卡洛方法必须等待回合终止，直到累计奖励已知
- 时序差分方法能够从不完整的序列中学习，蒙特卡洛方法只能从完整序列中学习
- 时序差分方法能够应用于无限长度的马尔可夫决策过程，蒙特卡洛方法只适用于有限长度

时序差分vs. 蒙特卡洛



- 累计奖励 $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$ 是 V^π 的无偏估计
 - 时序差分的真实目标 $R_{t+1} + \gamma V^\pi(s_{t+1})$ 也是 V^π 的无偏估计
 - 时序差分目标 $R_{t+1} + \gamma V(s_{t+1})$ 是 V^π 的有偏估计
- ↑
当前估计

时序差分方法

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

- 低方差，有偏差
 - 更高效
 - 最终收敛到 V^π
 - 对初始值更敏感

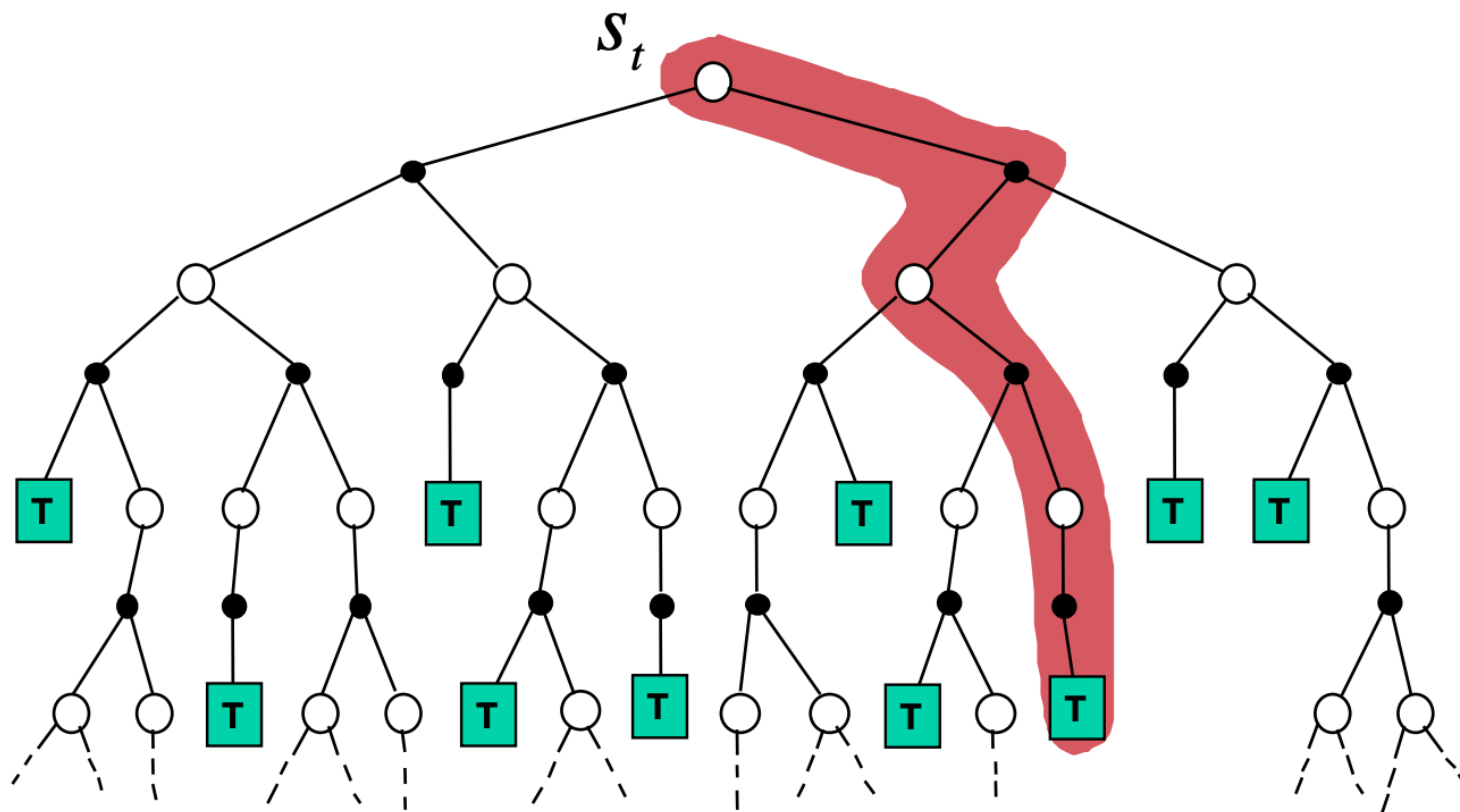
蒙特卡洛方法

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

- 高方差，无偏差
 - 良好收敛性
 - 对初始值不敏感
 - 易于理解和使用

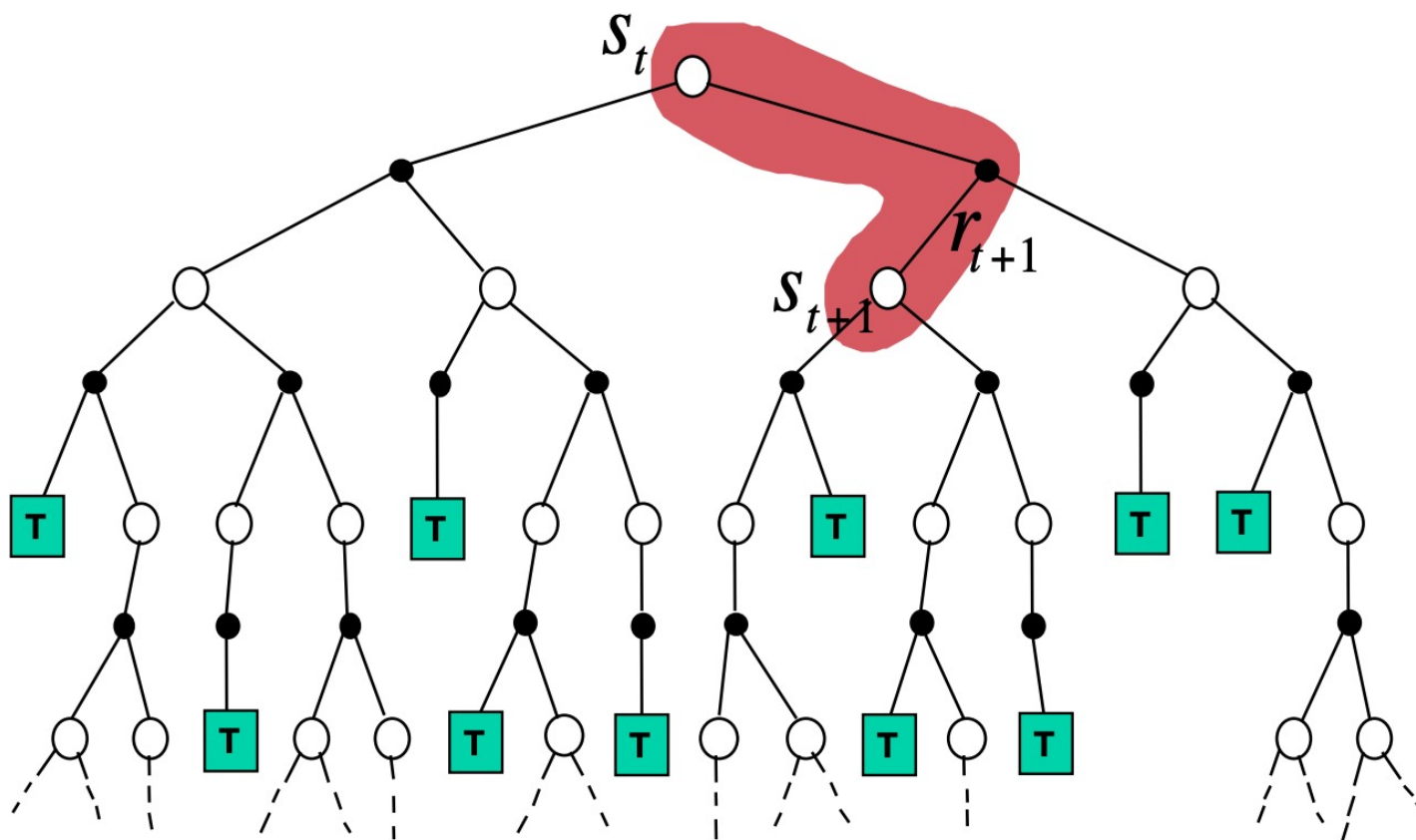
- 蒙特卡洛方法

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



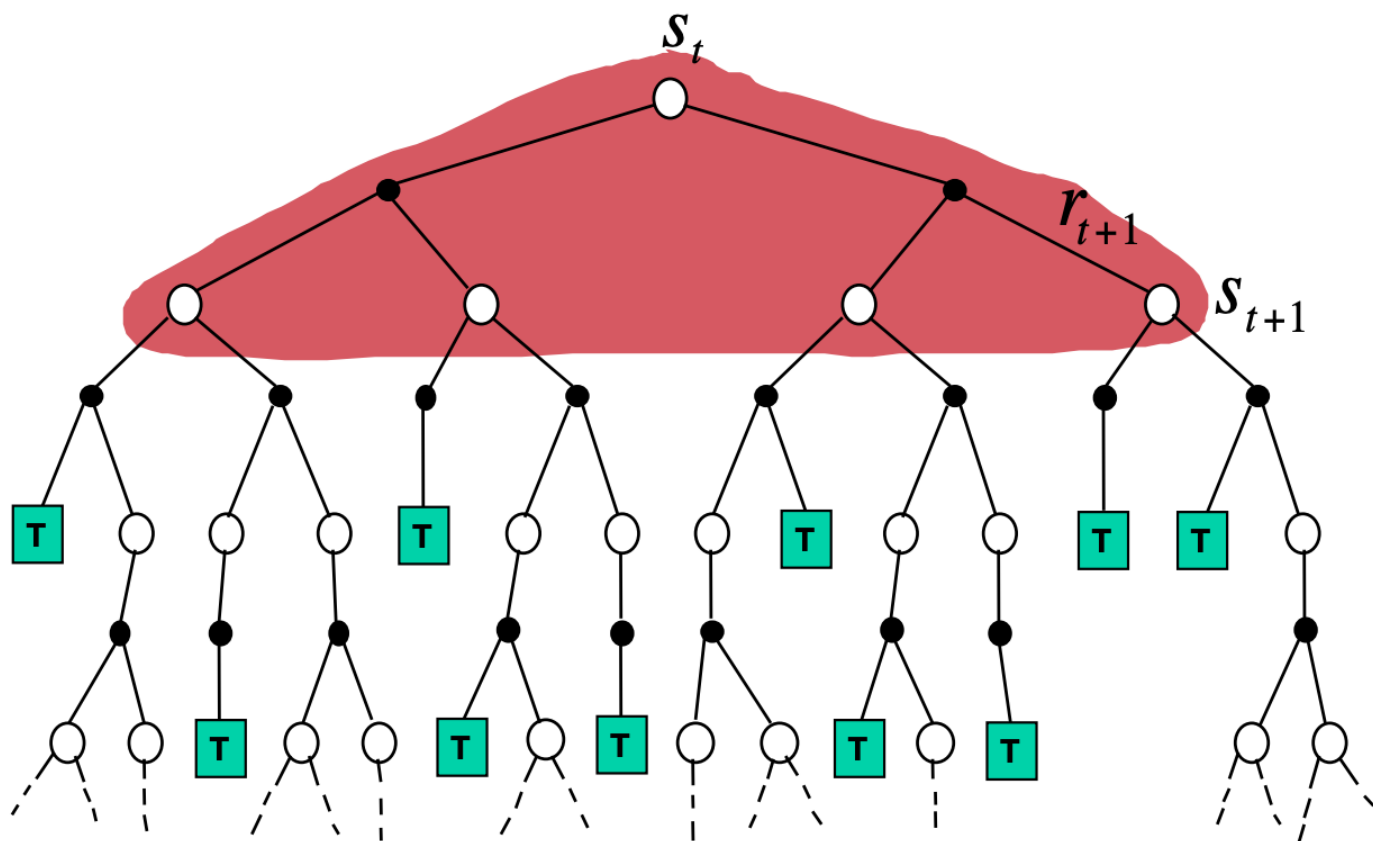
- 时序差分方法

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



- 动态规划

$$V(S_t) \leftarrow \mathbb{E}_{\pi} [R_{t+1} + \gamma V(S_{t+1})]$$

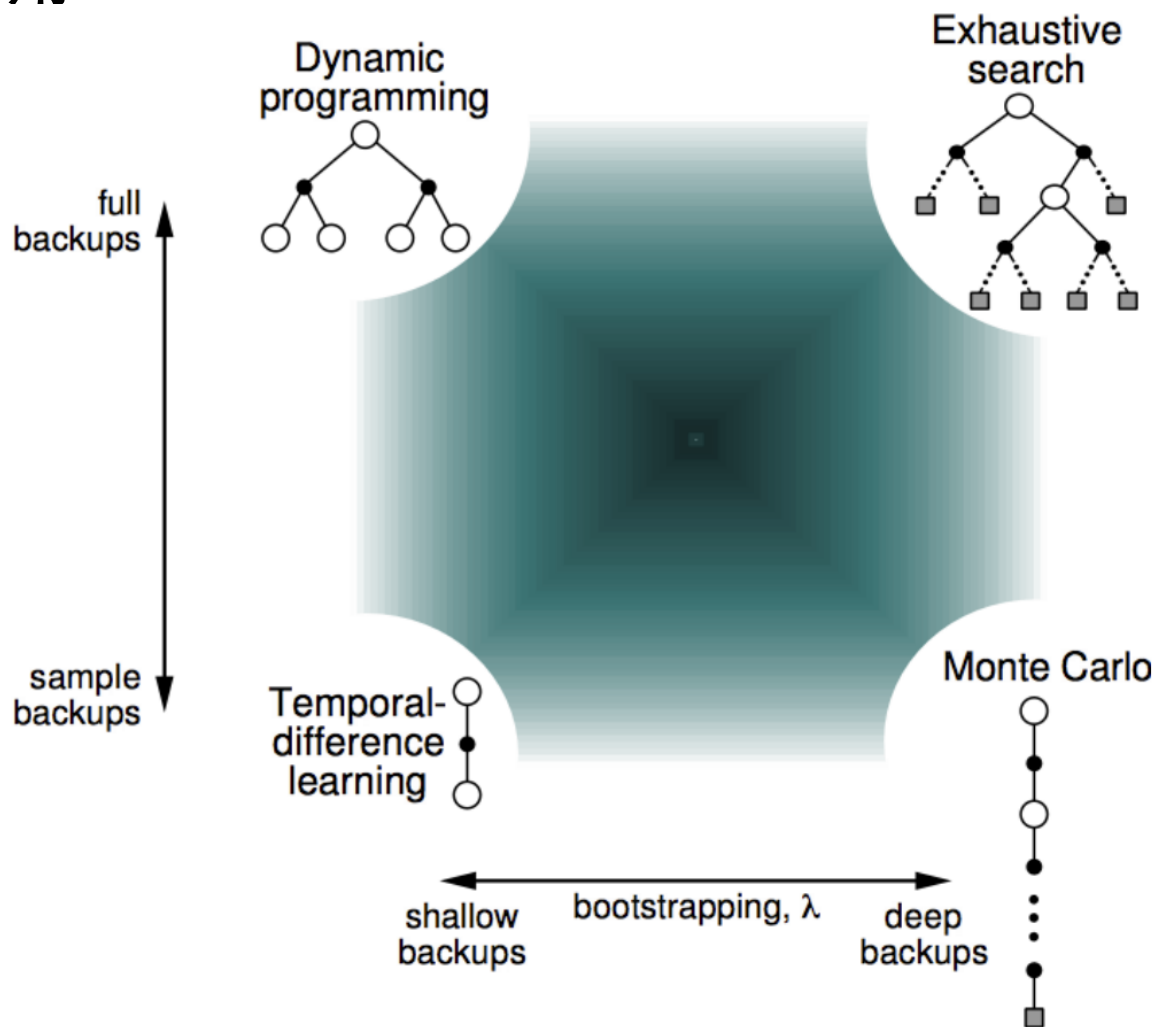


时序差分vs. 蒙特卡洛



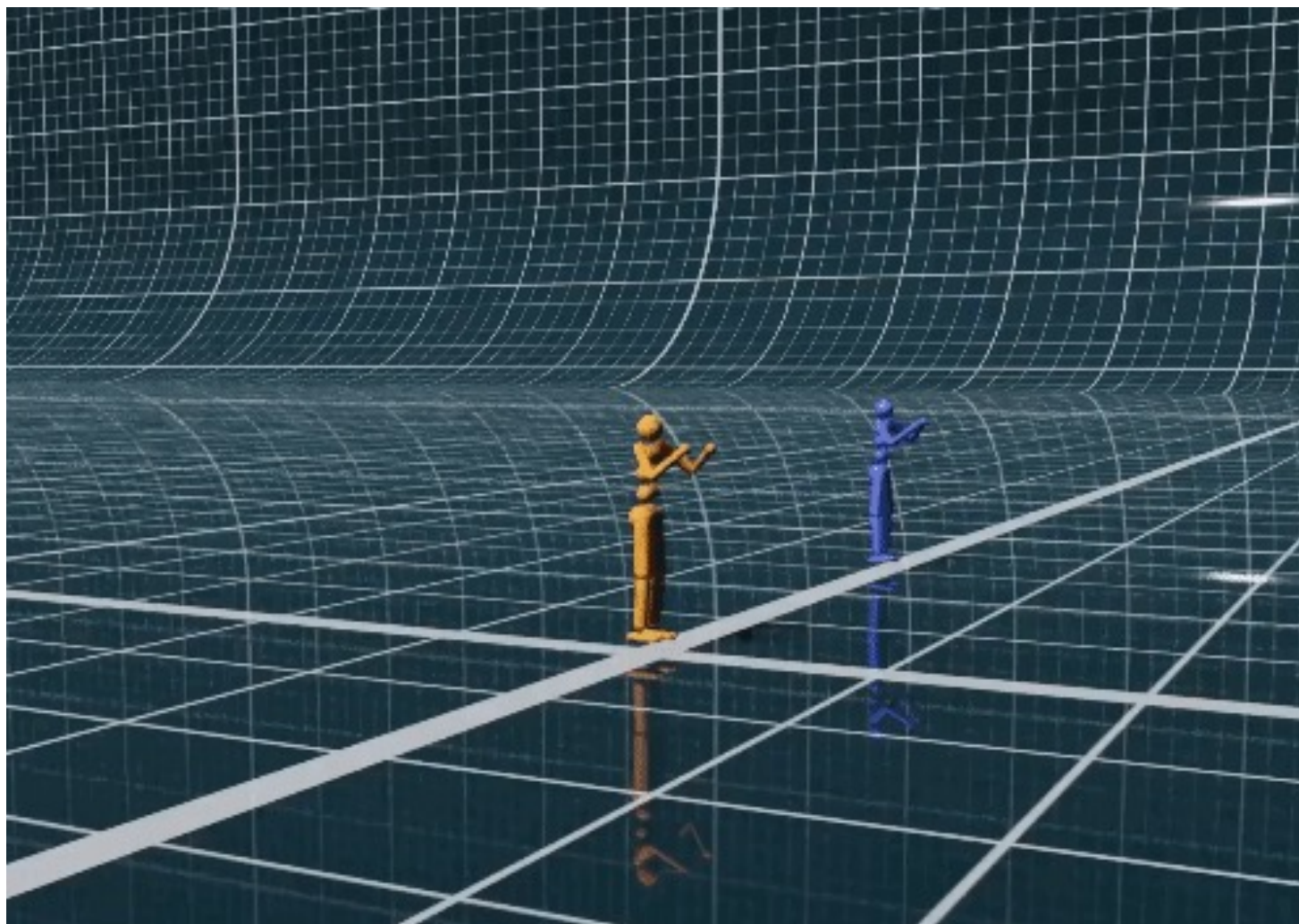
腾讯开悟

- 统一视角



- 蒙特卡洛方法
- 时序差分方法
- **资格迹方法**
- 表格型时序差分方法：
 - Q-learning
 - SARSA

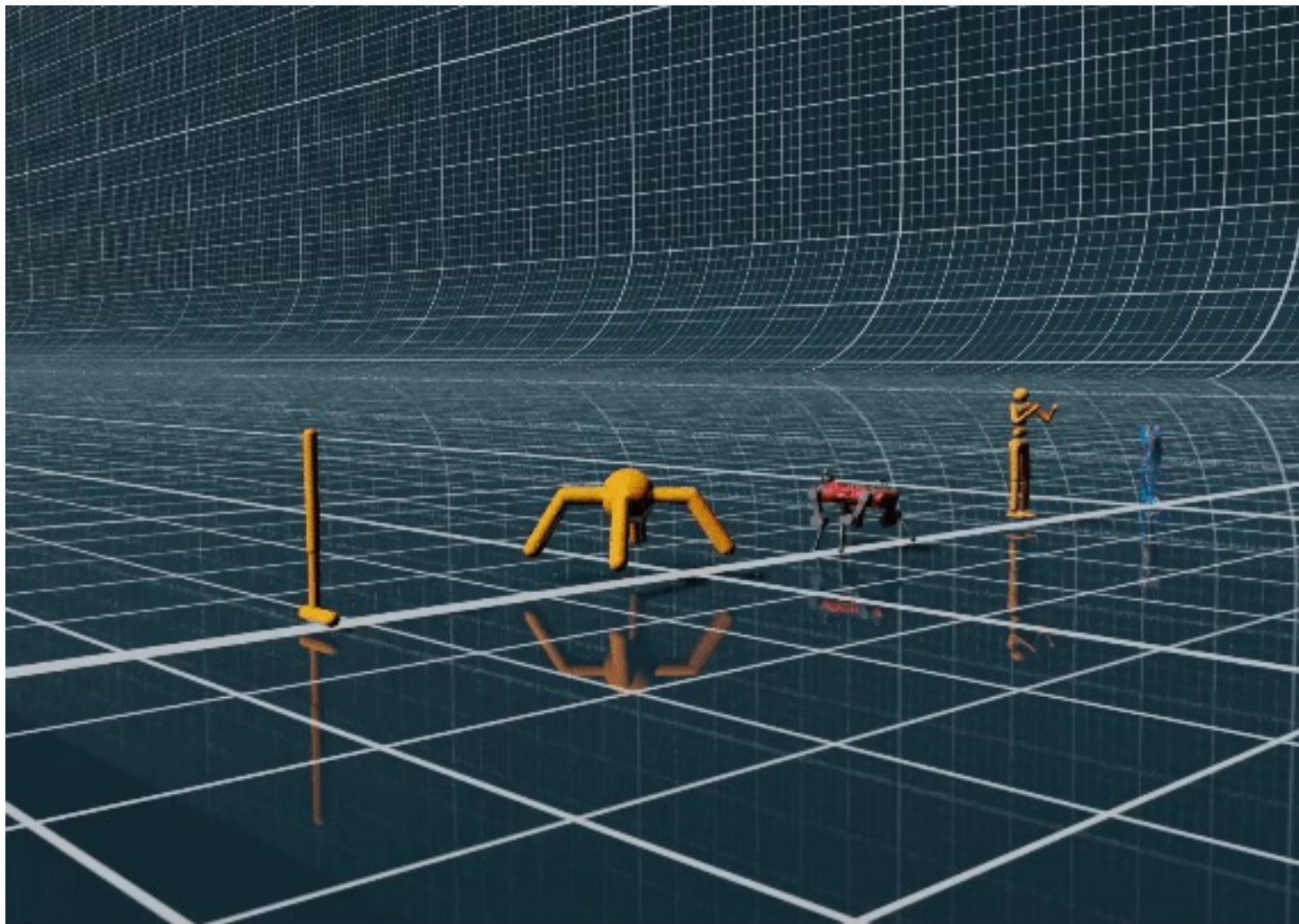
资格迹方法



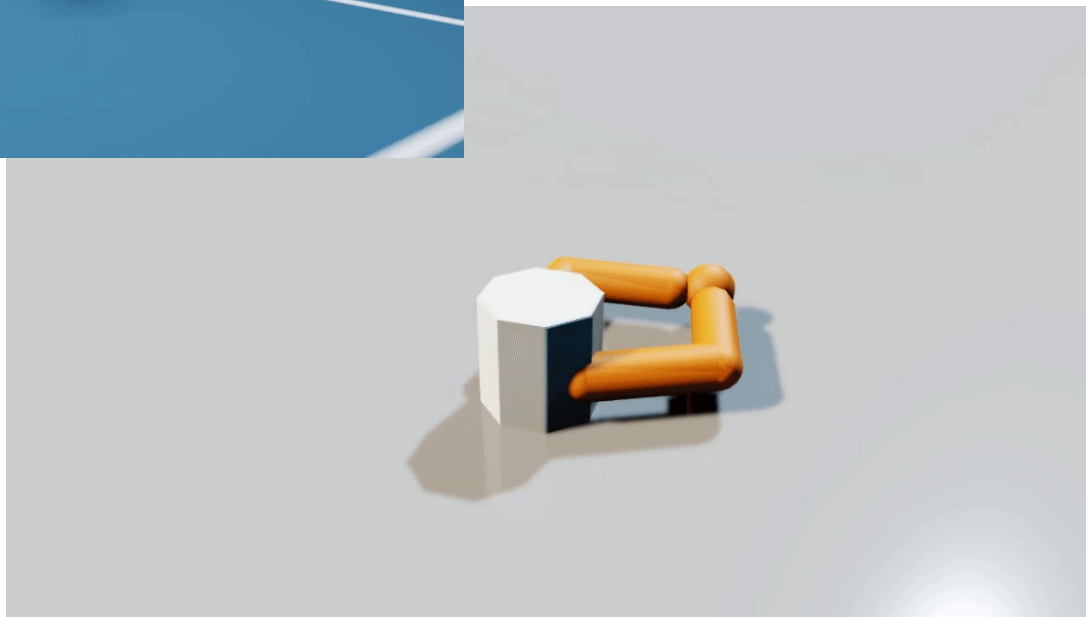
资格迹方法



蒙特卡洛方法



资格迹方法



- 蒙特卡洛方法（高方差，无偏差）

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$

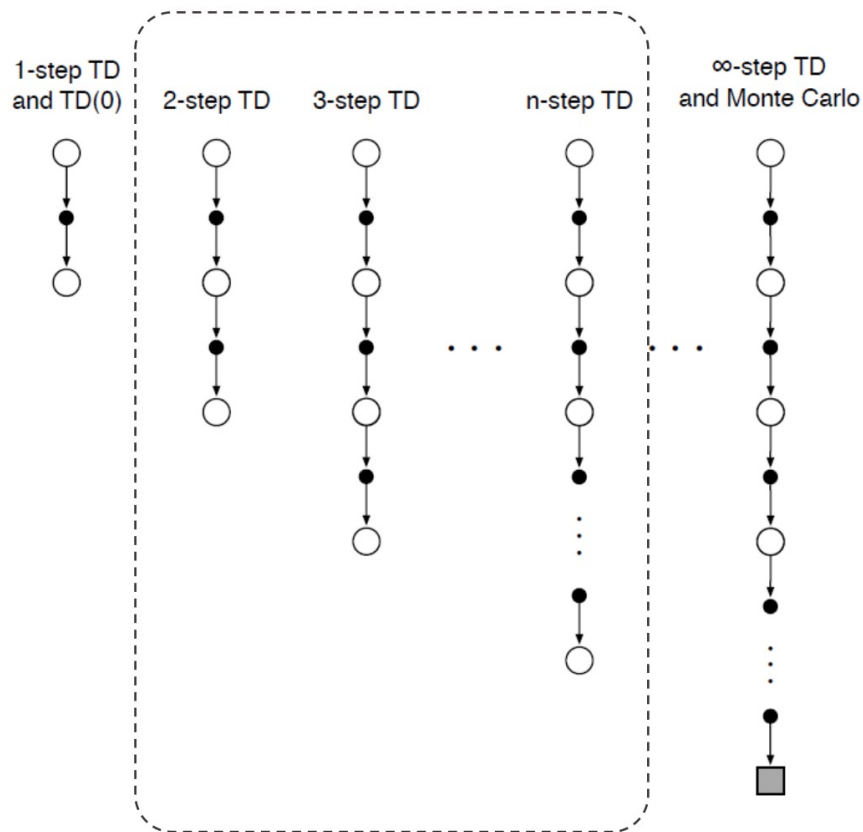
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

- 时序差分方法（低方差，有偏差）

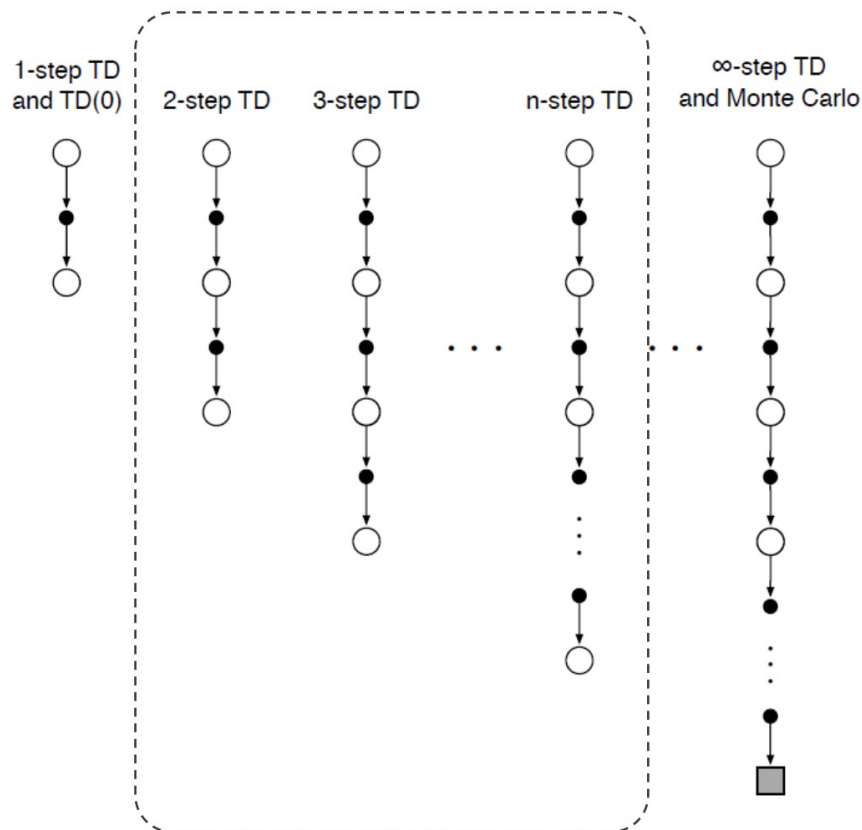
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

有没有方法介于两者之间？

- 多步时序差分方法



- 多步时序差分方法



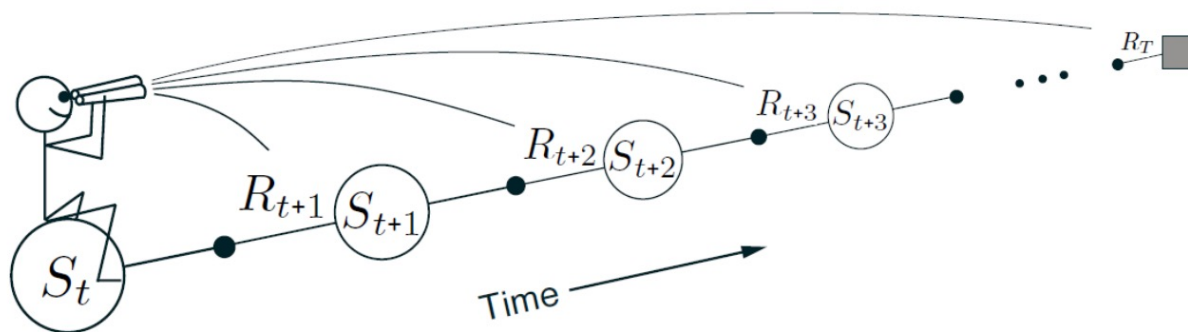
- 无限步时序差分方法等效于蒙特卡洛方法

- 为了实现方差与偏差的平衡，一种可行的方案是将蒙特卡洛方法与时序差分方法融合，实现多步时序差分
- 定义n步累计奖励

$$G_t^n = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(s_{t+n})$$

- 可得出n步时序差分学习

$$V(s_t) \leftarrow V(s_t) + \alpha(G_t^n - V(s_t))$$

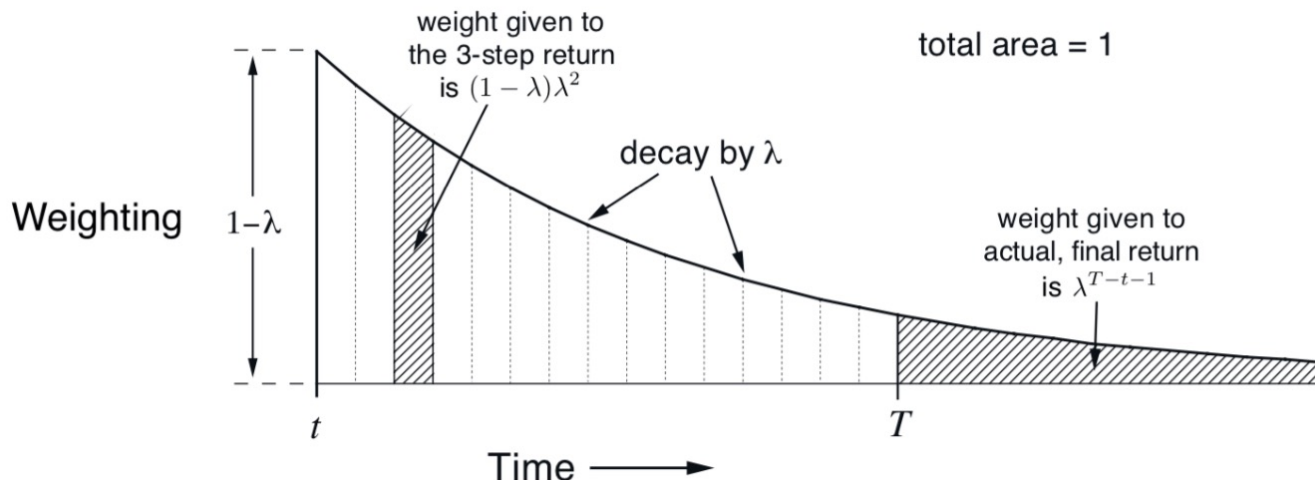


- 资格迹方法 (Eligibility Traces methods) 统一了时序差分和蒙特卡罗方法
- 资格迹方法通常使用一个超参数 $\lambda \in [0, 1]$ 控制值估计蒙特卡罗还是时序差分，通常来说，当 $\lambda = 1$ 时资格迹方法等价于蒙特卡罗方法，当 $\lambda = 0$ 等价于时序差分方法
- 介于两者之间的方法通常比任何一种极端方法都要好

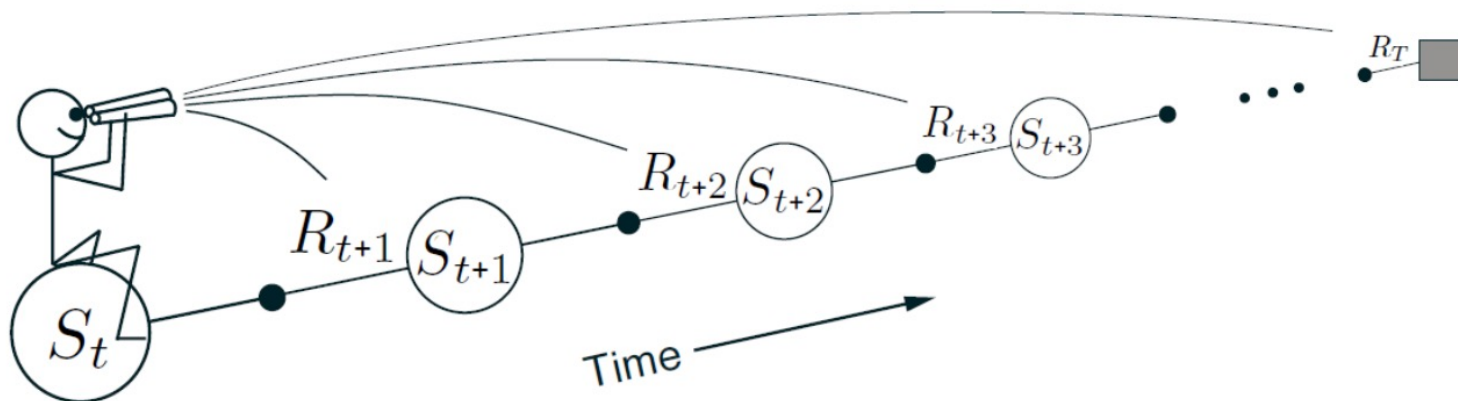
- TD- λ 是一种比较常见的资格迹方法，不同时间步差分估计值的权重随着其时间步的增加而衰减

$$G_t^n = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V(s_{t+n})$$

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

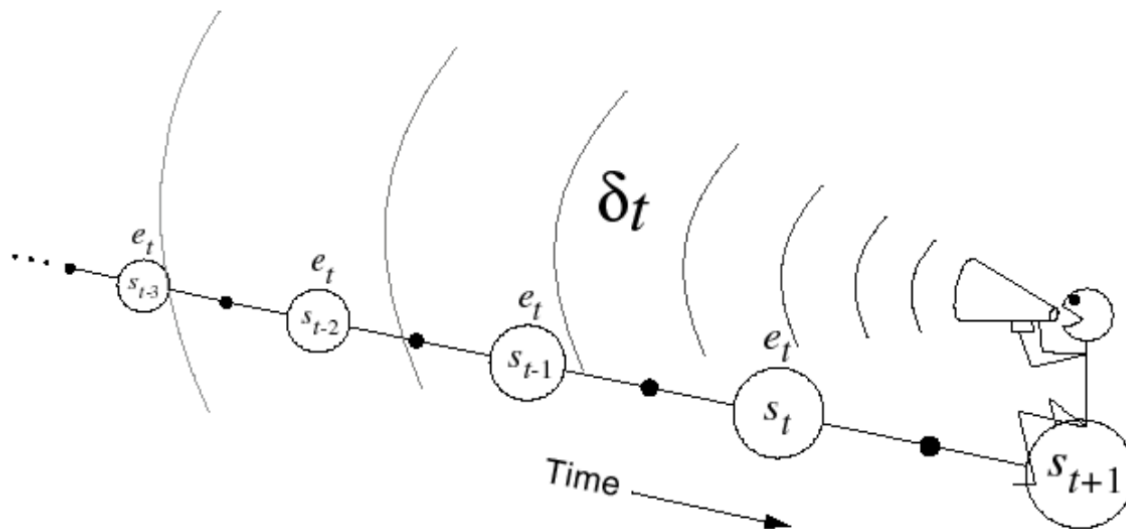


- TD- λ 的前向视角：



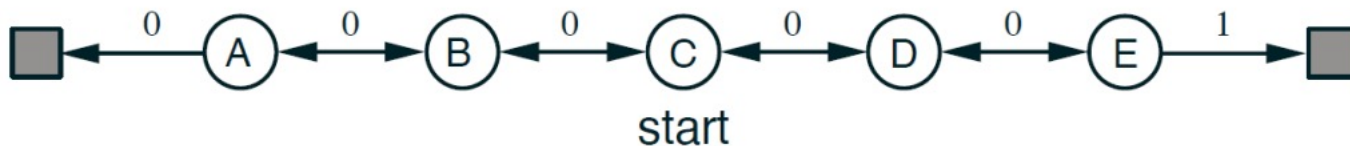
- 朝着 G_t^λ 的方向更新值函数
- 通过观测未来的数据计算 G_t^λ
- 与蒙特卡洛方法类似，只能计算完整的回合

- TD- λ 的后向视角:

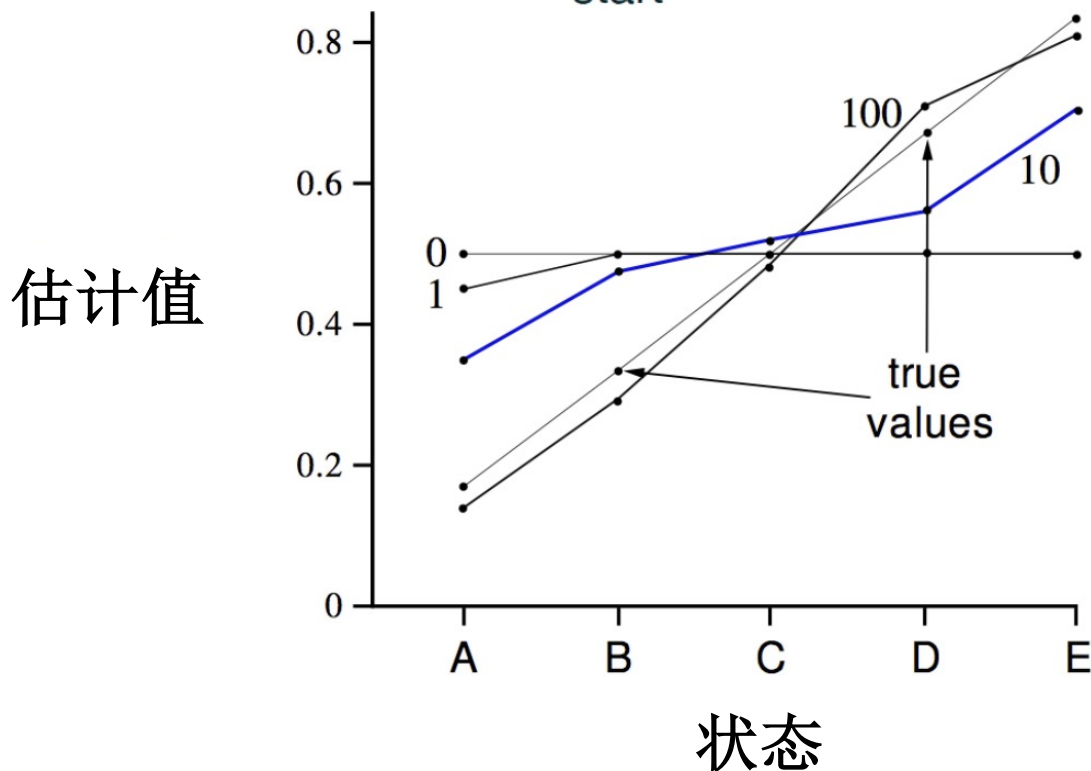
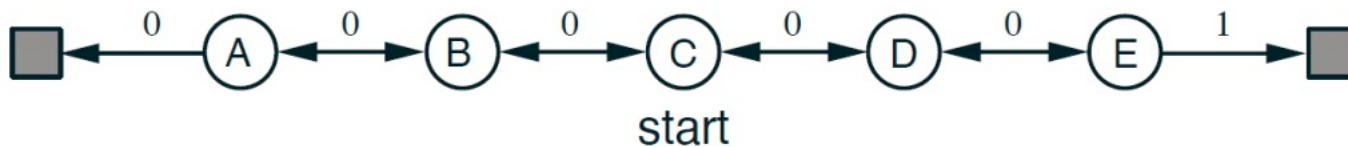


- 对每个状态 s 保持资格迹
- 更新每个状态 s 的价值函数
- 与TD-error和资格迹 $E_t(s)$ 呈比例关系

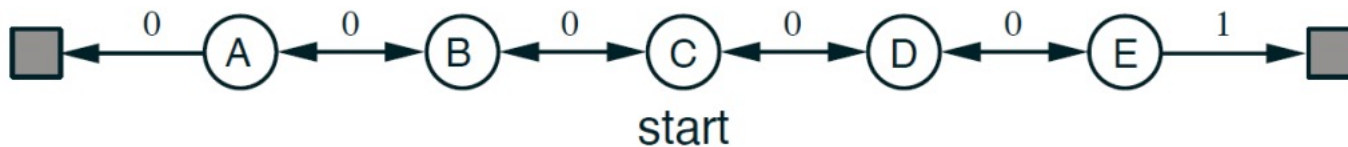
- 示例1：随机游走
 - 每个回合都从中间状态C开始
 - 每一时刻都有均等的概率向左走或者向右走
 - 到最左侧或者最右侧时，回合结束
 - 如果走到最右侧，得到的奖励为1
 - 如果走到最左侧，得到的奖励为0



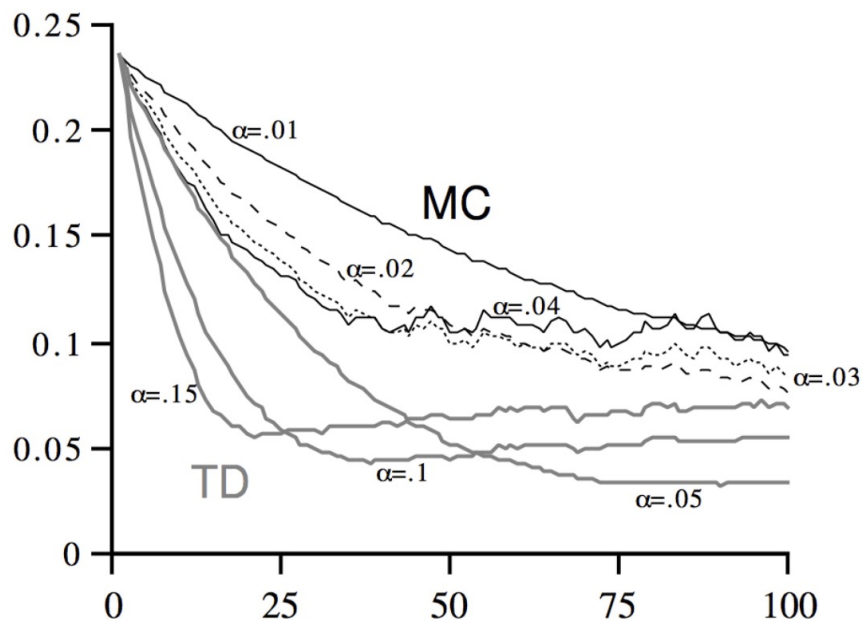
- 示例1：随机游走



- 示例1：随机游走



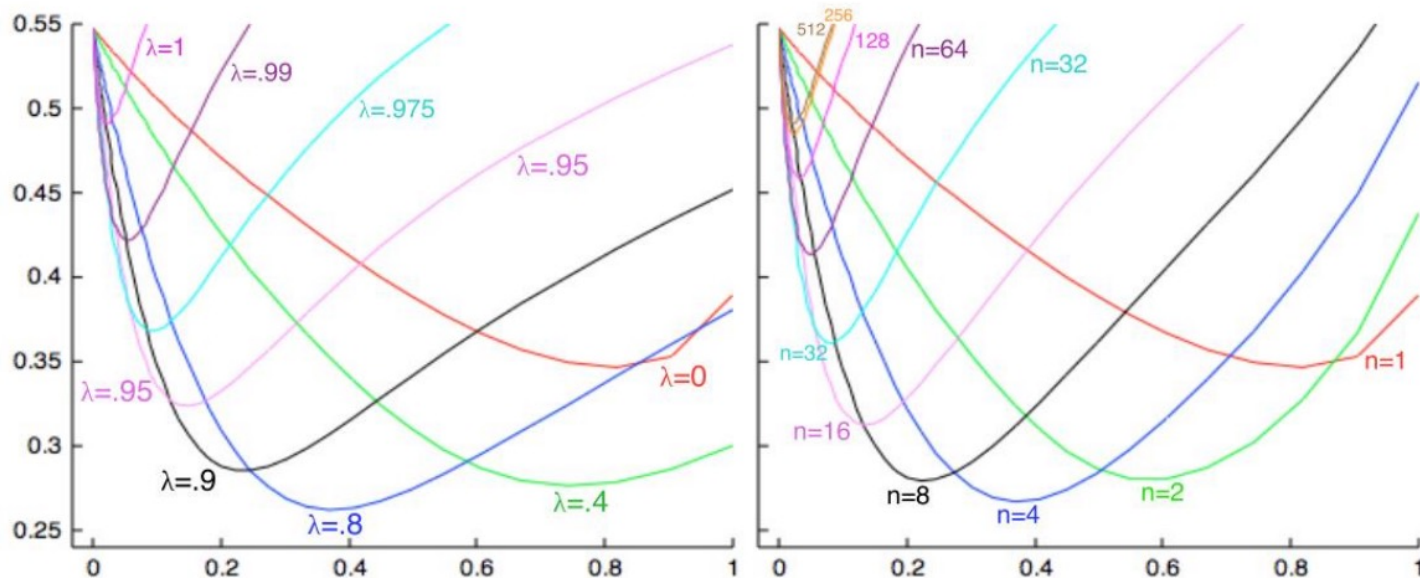
状态对应的
均方误差



回合数

- 示例1：随机游走

前10个片段
(episode)
结束时的
Offline
RMS误差

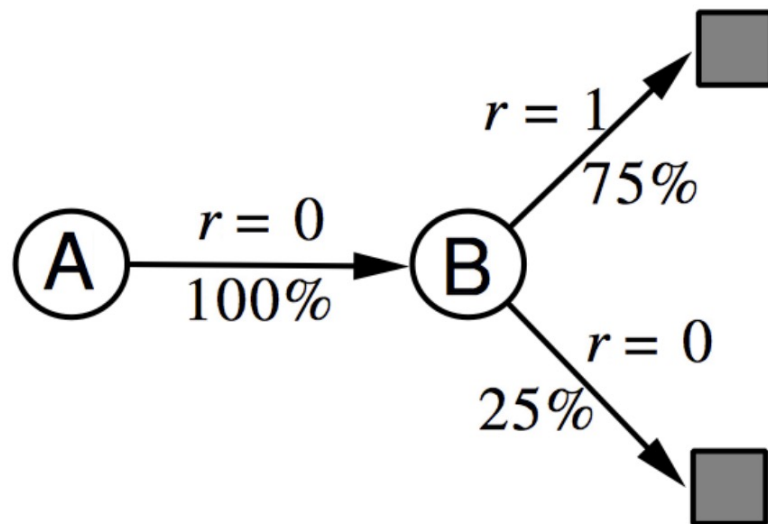


19个状态的随机游走实验结果

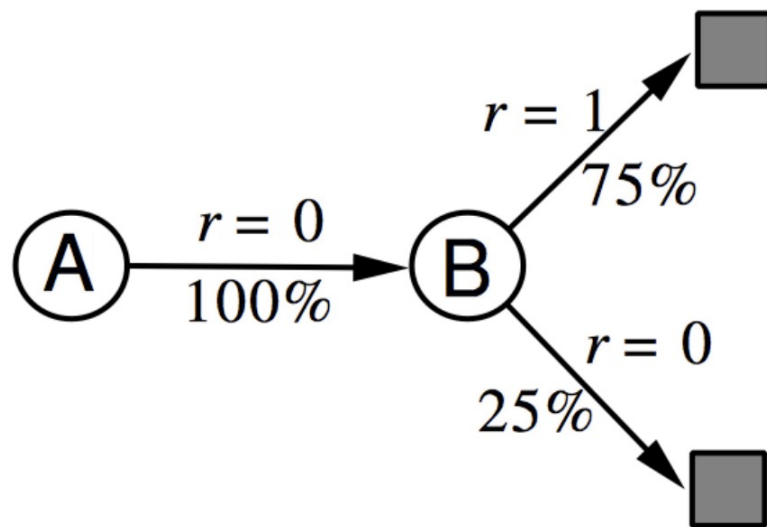
- TD- λ 与 n 步时序差分有类似的效果

- 示例2：AB环境
 - 无折扣因子
 - 一共8个回合
 - $A \rightarrow r=0 \rightarrow B \rightarrow r=0$
 - $B \rightarrow r=1$ (6次重复)
 - $B \rightarrow r=0$
 - 求状态A, B的价值?

- 示例2：AB环境
 - 无折扣因子
 - 一共8个回合的样本
 - $A \rightarrow r=0 \rightarrow B \rightarrow r=0$
 - $B \rightarrow r=1$ (6次重复)
 - $B \rightarrow r=0$
 - 求状态A的价值 V_A ?



- 示例2：AB环境
 - 根据蒙特卡洛方法，必须使用完整的回合进行计算，因此状态A的价值 $V_A=0$
 - 单步的时序差分方法融合了状态B的价值 V_B ，因此状态A的价值 $V_A=0.75$



- 蒙特卡洛方法
- 时序差分方法
- 资格迹方法
- **表格型时序差分方法：**
 - **Q-learning**
 - **SARSA**

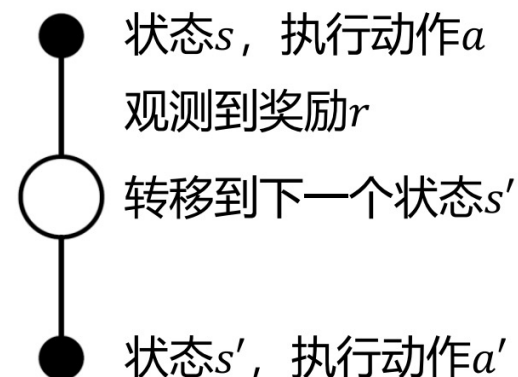
- 强化学习由策略评估和策略改进组成
- 策略评估的目标是知道什么决策是好的，即估计 $V^\pi(s_t)$
- 策略提升的目标是根据值函数选择好的行动

- 基于V函数，
$$\pi(s) = \arg \max_{a \in A} \sum_{s' \in S} P(s'|s, a) V(s')$$

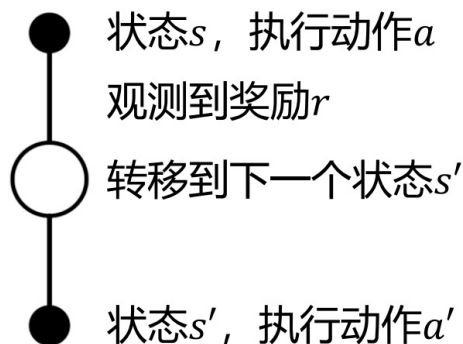
↑
需要知道环境模型

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

↑
直接作用于决策

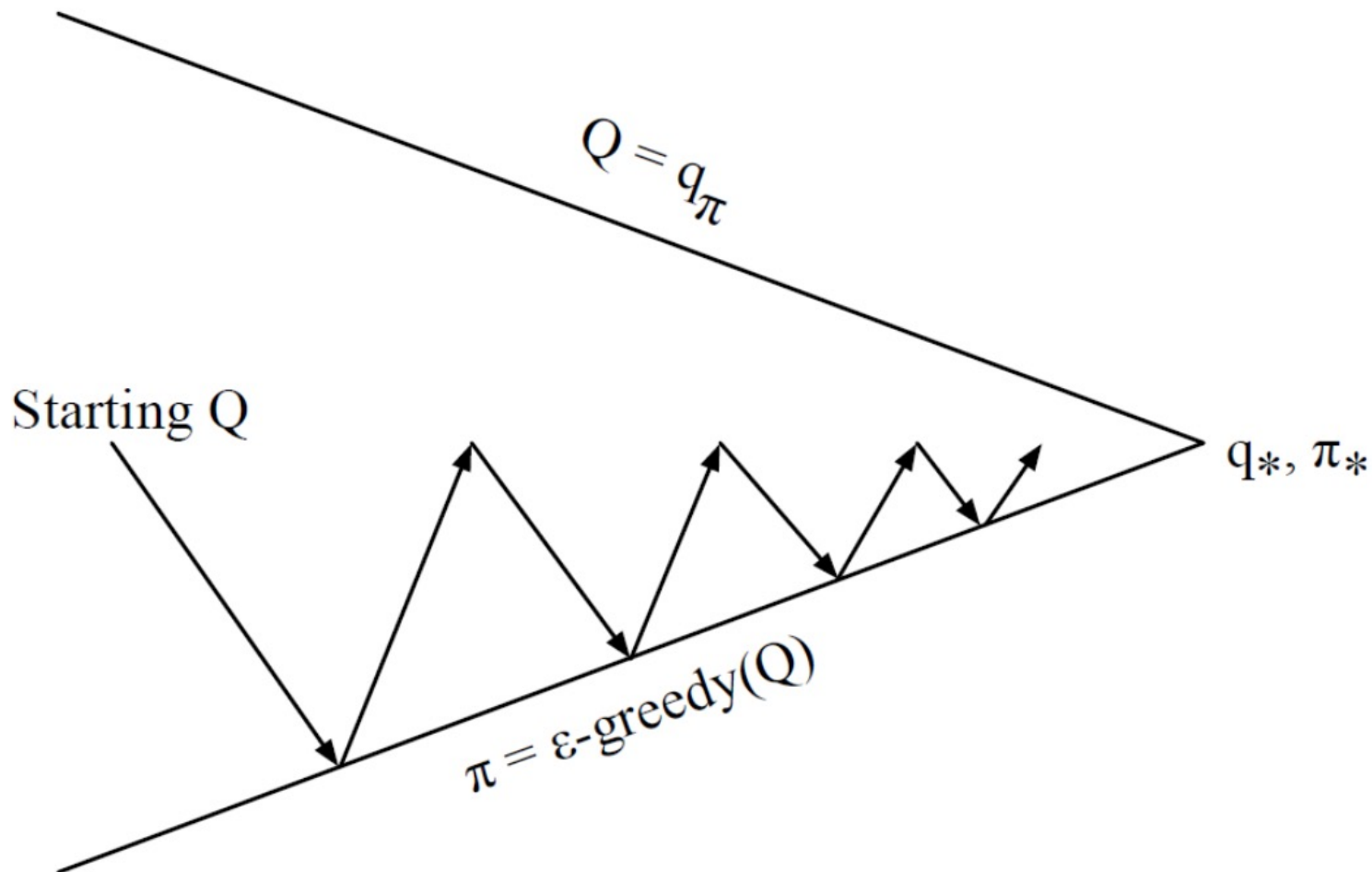


- SARSA是一种针对表格环境中的时序差分方法
- 对策略执行的每个（状态-动作-奖励-状态-动作）元组



- SARSA的策略评估为更新状态-动作值函数
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$
- SARSA的策略改进为 ϵ -greedy

表格型时序差分：SARSA



- 迭代更新Q函数的估计值和策略

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., ϵ -greedy)

 Repeat (for each step of episode):

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ϵ -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

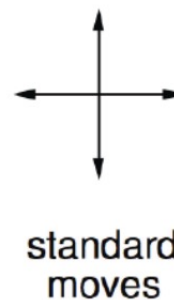
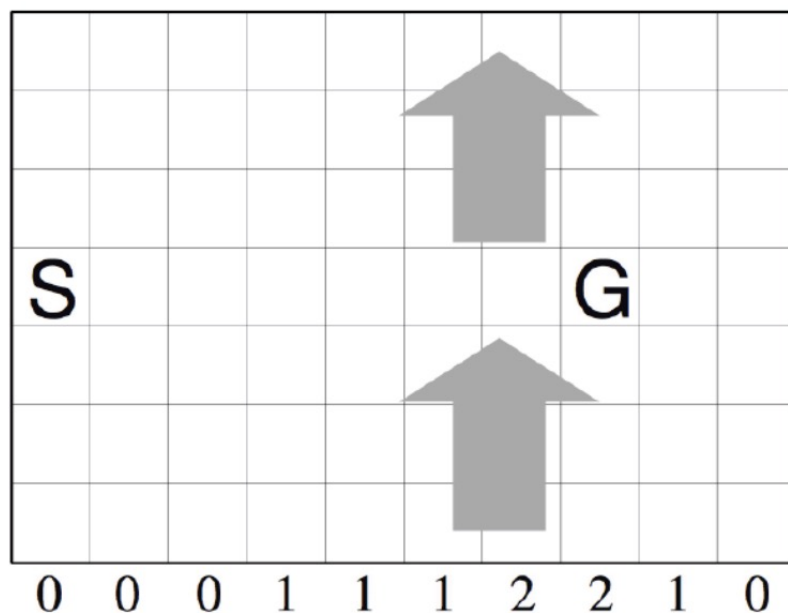
- 在线策略时序差分控制（on-policy TD control）使用当前策略进行动作采样，即SARSA算法中的两个动作“A”都是由当前策略选择的

表格型时序差分：SARSA



腾讯开悟

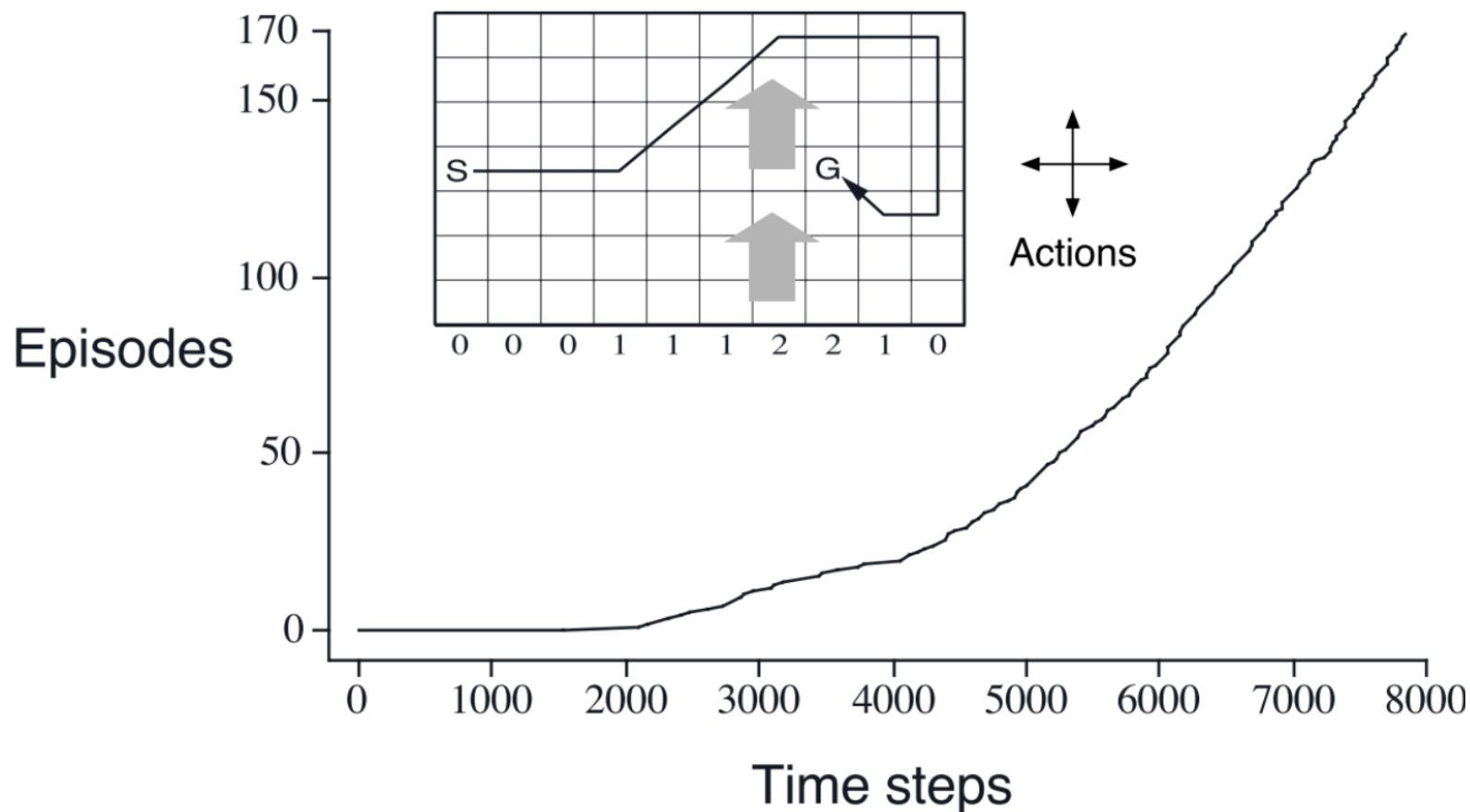
- SARSA示例：Windy Grid World
 - 每步的奖励为-1，直到智能体抵达目标网格
 - 折扣因子 $\gamma = 1$



表格型时序差分：SARSA



腾讯开悟



- 随着训练的进行，SARSA策略越来越快速地抵达目标

- Q-learning学习状态动作值函数 $Q(s, a) \in \mathbb{R}$ ，是一种离线策略（off-policy）方法

给定策略（可以不同于当前策略）

$$Q(s_t, a_t) = \sum_{t=0}^T \gamma^t \underset{\substack{\uparrow \\ \text{奖励函数}}}{R(s_t, a_t)}, a_t \sim \underset{\substack{\downarrow \\ \text{给定策略}}}{\mu(s_t)}$$

- 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
 - 目标策略 $\pi(a_t | s_t)$ 进行值函数评（ $V^\pi(s_t)$ 或 $Q^\pi(s_t, a_t)$ ）
 - 行为策略 $\mu(a_t | s_t)$ 收集数据
- 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
 - 遵循一个策略时学习多个策略

- 具体实现
- 使用行为策略 $\mu(\cdot | s_t)$ 选择动作 a_t
- 使用当前策略 $\pi(\cdot | s_{t+1})$ 选择后续动作 a'_{t+1} ，计算目标

$$Q'(s_t, a_t) = R_t + \gamma Q(s_{t+1}, a'_{t+1})$$

- 使用时序差分更新 $Q(s_t, a_t)$ 的值以拟合目标状态-动作值

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

↑
动作来自当前策略 π

- 目标策略 π 是关于 $Q(s, a)$ 的贪心策略

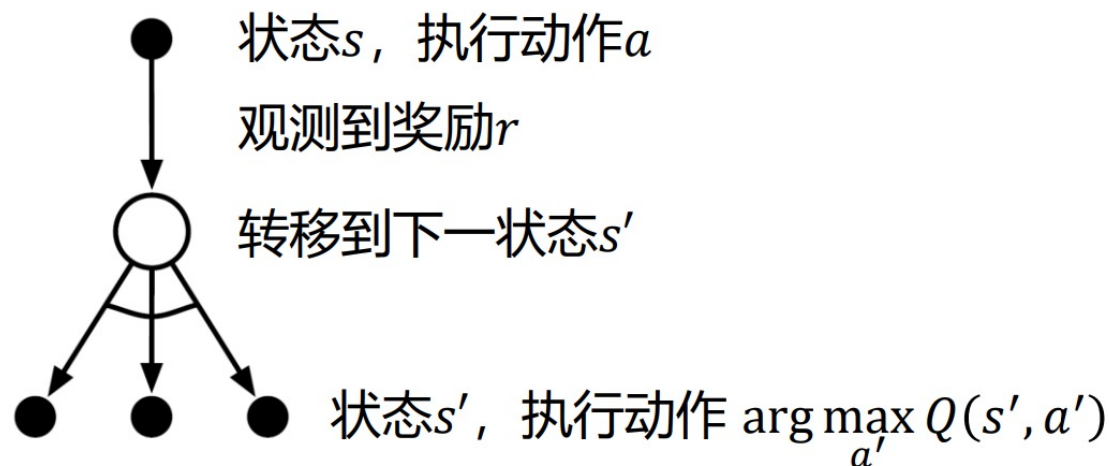
$$\pi(s_{t+1}) = \arg \max_{a'} Q(s_{t+1}, a'_{t+1})$$

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

$$Q'(s_t, a_t) = R_t + \gamma \max_{a'} Q(s_{t+1}, a'_{t+1})$$

- Q-learning学习更新方式

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



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

- 定理：Q-learning能够收敛到最优状态-动作值函数

$$Q(s_t, a_t) \rightarrow Q^*(s_t, a_t)$$

- 回顾：收缩算子（contraction operator）
- 定义H算子

$$HQ(s_t, a_t) = R_t + \gamma \mathbb{E}_{s_{t+1} \sim p(\cdot | s, a)} [\max_{a'} Q(s_{t+1}, a_{t+1})]$$

- 最优值函数 Q^* 是H的不动点，意味着： $Q^* = HQ^*$

- 直接从Q函数证明

$$\begin{aligned} & \| \mathbf{H}q_1 - \mathbf{H}q_2 \|_\infty \quad (\mathbf{H}q)(x, a) = \sum_{y \in \mathcal{X}} P_a(x, y) [r(x, a, y) + \gamma \max_{b \in \mathcal{A}} q(y, b)] \\ &= \max_{x, a} \left| \sum_{y \in \mathcal{X}} P_a(x, y) [r(x, a, y) + \gamma \max_{b \in \mathcal{A}} q_1(y, b) - r(x, a, y) + \gamma \max_{b \in \mathcal{A}} q_2(y, b)] \right| \\ &= \max_{x, a} \gamma \left| \sum_{y \in \mathcal{X}} P_a(x, y) [\max_{b \in \mathcal{A}} q_1(y, b) - \max_{b \in \mathcal{A}} q_2(y, b)] \right| \\ &\leq \max_{x, a} \gamma \sum_{y \in \mathcal{X}} P_a(x, y) \left| \max_{b \in \mathcal{A}} q_1(y, b) - \max_{b \in \mathcal{A}} q_2(y, b) \right| \\ &\leq \max_{x, a} \gamma \sum_{y \in \mathcal{X}} P_a(x, y) \max_{z, b} |q_1(z, b) - q_2(z, b)| \\ &= \max_{x, a} \gamma \sum_{y \in \mathcal{X}} P_a(x, y) \|q_1 - q_2\|_\infty \\ &= \gamma \|q_1 - q_2\|_\infty. \end{aligned}$$

因此 $\| \mathbf{H}q_1 - \mathbf{H}q_2 \|_\infty \leq \gamma \|q_1 - q_2\|_\infty$

- 数列 $\{x_n\}$ 收敛的充分必要条件是：对于任意给定的正数 ϵ ，总存在正整数 N ，使得当 $n > N$ ， $m > N$ 时有 $|x_n - x_m| < \epsilon$ ，把满足该条件的数列 $\{x_n\}$ 称为柯西数列
- 数列 $\{x_n\}$ 收敛的充分必要条件是，该数列中的元素随着序数的增加而愈发靠近，即足够靠近后的任意两项都无限接近

$$\|\mathbf{H}q_1 - \mathbf{H}q_2\|_{\infty} \leq \gamma \|q_1 - q_2\|_{\infty}$$

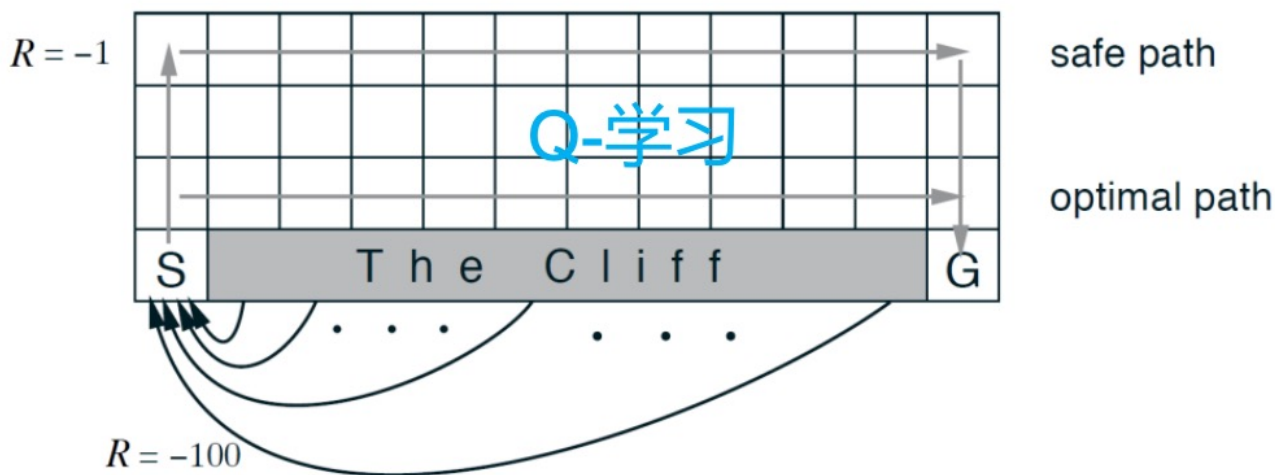
表格型时序差分：Q-learning



腾讯开悟

- Q-learning示例：悬崖漫步（Cliff-walking）
 - 折扣因子 $\gamma = 1$
 - 每一步移动奖励为-1
 - 踏入悬崖会产生-100奖励，并将智能体送回开始处

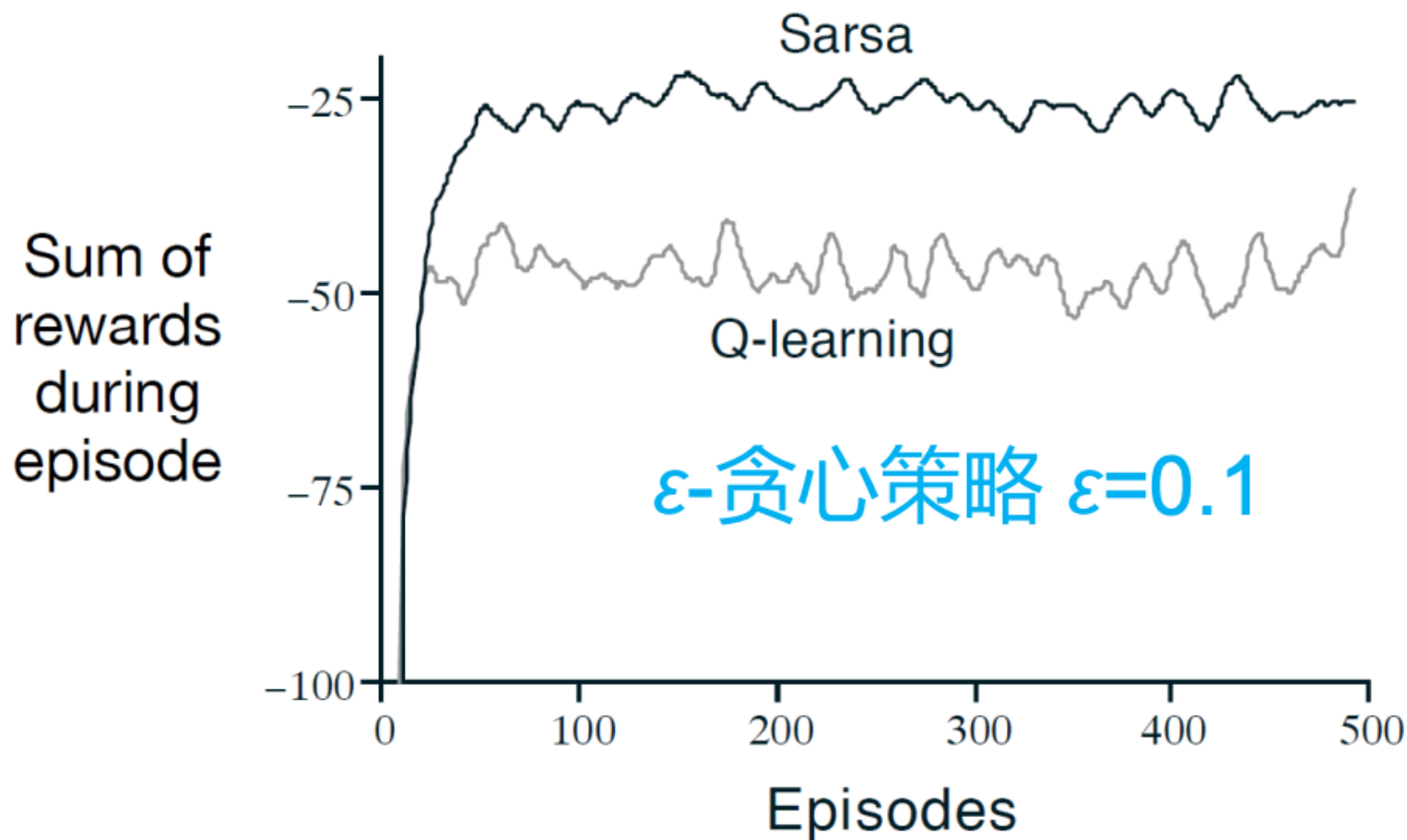
SARSA



表格型时序差分：Q-learning



腾讯开悟



提问环节



腾讯开悟

需要完成：

- 运行 gym 中的 Cliff Walking 环境
- 基于 Cliff Walking 实现表格型 Q-learning 和 SARSA 算法
- 测试 Q-learning 和 SARSA 算法的表现与区别
- 推荐依赖：
 - numpy（大于等于1.21.0，小于2.0.0）
 - gym（大于等于0.24.0）
 - matplotlib（不需要指定版本）
 - tqdm（不需要指定版本）