

MC 和 TD 对比

MC 是高方差，零偏差：

需要完整的事件序列计算出回报后学习

好的收敛性(即使是使用逼近器也能保证收敛性)；

与价值函数初始值无关；

原理简单，使用方便；

要求事件达到终止状态或序列足够长

能够利用马尔可夫性，因此在马尔可夫环境下更有效

TD 是低方差，有偏差：

可以在智能体运行过程中的每一步在线学习

可以从不完整的事件序列中学习

通常情况是比 MC 更有效

TD(0) 能够收敛到 $v_\pi(s)$

(但是与逼近器结合后并不一定收敛)

受价值函数初始值影响

不能利用马尔可夫性，因此在非马尔可夫环境下更有效

偏差/方差的权衡：

回报 $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t} R_t = v_\pi(S_t)$ 的无偏估计

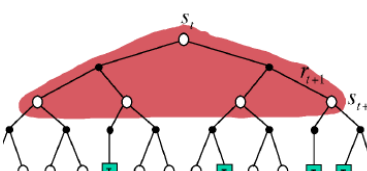
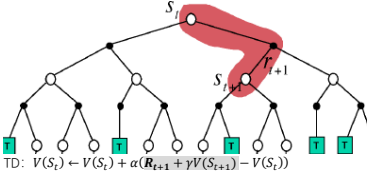
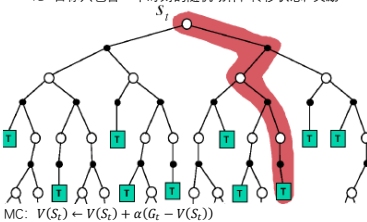
真实的 TD 目标 $R_{t+1} + \gamma v_\pi(S_{t+1}) = v_\pi(S_t)$ 的无偏估计

TD 目标 $R_{t+1} + \gamma V(S_{t+1}) = v_\pi(S_t)$ 的有偏估计

TD 目标与回报的方差差小很多

回报的计算包含整个序列中的随机动作、转移状态、奖励

TD 目标只包含一个时刻的随机动作、转移状态、奖励



动态规划更新: $V(S_t) \leftarrow E_t[R_{t+1} + \gamma V(S_{t+1})]$

自举法: 更新时包含一个预测量

MC 不使用自举 DP 使用自举 TD 使用自举

采样法: 使用采样的数据计算期望

MC 使用采样 DP 不使用采样 TD 使用采样

等价模型:

MC 收敛结果对应最小二乘误差，最佳匹配观测的回报

$$\sum_{k=1}^N \sum_{t=1}^{T_k} (G_t^k - V(S_t^k))^2$$

TD(0) 收敛结果对应最小似然马尔可夫模型 $(S, A, \hat{P}, \hat{R}, \gamma)$

$$\hat{P}_{s,s'} = \frac{1}{N(s,a)} \sum_{k=1}^N \sum_{t=1}^{T_k} 1(s_t^k, a_t^k, s_{t+1}^k = s, a, s')$$
$$\hat{R}_{s,s'} = \frac{1}{N(s,a)} \sum_{k=1}^N \sum_{t=1}^{T_k} 1(s_t^k, a_t^k = s, a) r_t^k$$

$n=1$ (TD) $G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$

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

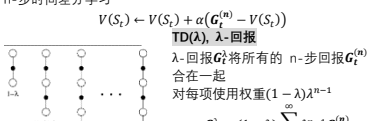
\vdots

$n=\infty$ (MC) $G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{t+T-1} R_T$

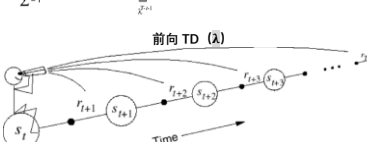
n-步的回报:

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

n-步时间差分学习

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t))$$


前向更新的 TD(λ) 形式

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t))$$


更新价值函数向 λ -回报逼近，面向未来时刻计算 $G_t^{(n)}$ ，与 MC 一样

要求完整的事件序列

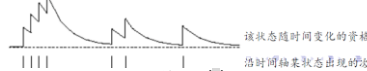
资格迹 Eligibility Traces

可信度：警铃和闪光是否和高电压直接相关？

受频率启发：对发生次数最多的状态给予最高的可信度

受近因启发：对时间上最近发生的状态给予最高的可信度

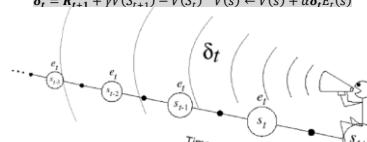
资格迹将两种方式结合

$$E_t(s) = 0 \quad E_t(s) = \gamma E_{t-1}(s) + 1(S_t = s)$$


后向 TD (λ)

为每个状态 s 记录它的资格迹，更新 s 的价值 $V(s)$

更新量与 TD 误差 δ_t 和资格迹 $E_t(s)$ 成正比

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \quad V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$


TD(A) 和 TD(0)

当 $\lambda=0$ 时，只有当前状态被更新

$$E_t(s) = 1(S_t = s) \quad V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

等同于 TD(0) 方法 $V(s) \leftarrow V(s) + \alpha \delta_t$

TD(A) 和 MC

当 $\lambda=1$ 时，可信度被推迟到事件结束，考虑每个事件都到达终止状态的环境，使用离线更新，对每个事件序列，TD(1) 的所有更新和 MC 的所有更新是一致的

当使用离线更新时，前向和后向 TD(A) 所有的更新是相同的

$$\sum_{t=1}^T \alpha \delta_t E_t(s) = \sum_{t=1}^T \alpha (G_t^{(n)} - V(S_t)) 1(S_t = s)$$

MC 和 TD(1)

考虑在一次事件中 s 在 k 时刻被唯一访问一次

该状态的 TD(1) 资格迹会在被访问之后随时间衰减

$$E_t(s) = \gamma E_{t-1}(s) + 1(S_t = s) = \begin{cases} 0 & \text{if } t < k \\ (\gamma)^{t-k} & \text{if } t \geq k \end{cases}$$

TD(1) 在更新时会在在线累计误差

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^{T-1} (\gamma)^{t-k} \delta_t = (G_k - V(S_k))$$

在事件结束时累计误差等于

$$\delta_k + \gamma \delta_{k+1} + \gamma^2 \delta_{k+2} + \dots + \gamma^{T-1-k} \delta_{T-1}$$

当 $\lambda=1$ 时，TD 误差可以转换成 MC 误差

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$
$$\delta_k + \gamma \delta_{k+1} + \gamma^2 \delta_{k+2} + \dots + \gamma^{T-1-k} \delta_{T-1} = G_t - V(S_t)$$

TD(A) 和 TD(1)

TD(1) 大致等价于每次经过的 MC 方法，误差会随在线运行逐步累积，如果价值函数只在事件结束时离线更新，那么所有的更新量等同于 MC 方法

前向和后向 TD(A)

考虑在一次事件中 s 在 k 时刻被唯一访问一次

该状态的 TD(A) 资格迹会在被访问之后随时间衰减

$$E_t(s) = \gamma \lambda E_{t-1}(s) + 1(S_t = s) = \begin{cases} 0 & \text{if } t < k \\ (\lambda \gamma)^{t-k} & \text{if } t \geq k \end{cases}$$

后向 TD(A) 在更新时会在在线累计误差

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^{T-1} (\lambda \gamma)^{t-k} \delta_t = (G_k^{(\lambda)} - V(S_k))$$

在事件结束时总共累计误差等于 λ -回报

如果 s 被多次访问， $E_t(s)$ 也会累积多次的误差

前向提供理论原理，后向给出算法实现：在线学习；每一时刻更新；可以适用事件序列中的一小段，非完整序列

目标：获得最优策略，基于模型：动态规划、价值迭代、策略迭代

目标：获得策略的价值函数，基于样本：迭代策略评估、蒙特卡洛方法、时间差分学习 TD(A)

我们想获得最优策略，但是 MDP 的模型未知，或者问题空间太大只能考虑和智能体最相关的状态：无模型方法；使用智能体的经验样本：最优价值函数

上节课：无模型预测，估计一个模型未知 MDP 问题的价值函数

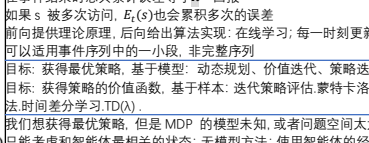
这节课：无模型控制，最优价值函数 MDP 问题的价值函数

在策略(on-policy) 和离策略(off-policy) 学习

在策略学习：根据策略 π 产生的样本来学习关于 π 的相关知识

离策略学习：根据另一策略 μ 产生的样本来学习关于 π 的相关知识

策略迭代：



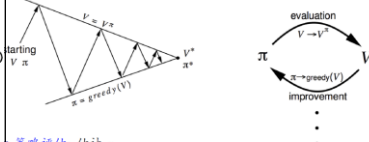
策略评估：估计 v_π

例如迭代策略评估

策略提升：生成新的 $\pi' \geq \pi$

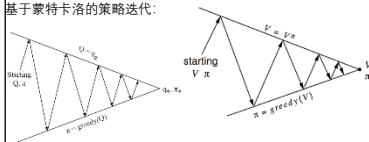
例如贪心策略提升

基于蒙特卡洛的策略迭代：



策略评估：MC 策略评估， $Q = q_\pi$

策略提升：贪心策略提升



策略评估：MC 策略评估， $Q = q_\pi$

策略提升：贪心策略提升

动作-价值函数实现无模型的策略迭代：

基于 $V(s)$ 进行贪心策略提升需要 MDP 的模型信息

$$\pi^*(s) = \arg \max_{a \in A} (R(s,a) + \gamma P_{s,a}^{\pi^*} V(s))$$

基于 $Q(s,a)$ 进行贪心策略提升是无模型的

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

e-贪心探索

最简单实现连续探索的方法，所有的 m 个动作都有非零概率被选中执行，以 $1-\epsilon$ 概率选择贪心动作，以 ϵ 概率随机选择一个动作

$$\pi(a|s) = \begin{cases} \epsilon/m + 1-\epsilon & \text{if } a = \arg \max_{a \in A} Q(s,a) \\ \epsilon/m & \text{otherwise} \end{cases}$$

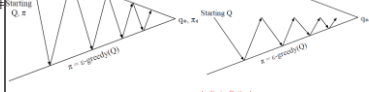
给定任意 e-贪心策略 π ，根据 q_π 构造出新的 e-贪心策略 π' 具有更好的性能，即 $v_{\pi'}(s) \geq v_\pi(s)$

$$q_\pi(s, \pi'(s)) = \sum_{a \in A} \pi'(a|s) q_\pi(s,a)$$

$\geq \epsilon/m \sum_{a \in A} q_\pi(s,a) + (1-\epsilon) \arg \max_{a \in A} q_\pi(s,a)$

$$\geq \epsilon/m \sum_{a \in A} q_\pi(s,a) + (1-\epsilon) \sum_{a \in A} \frac{\pi(a|s) - \epsilon/m}{1-\epsilon} q_\pi(s,a)$$
$$= \sum_{a \in A} \pi(a|s) q_\pi(s,a) = v_\pi(s) \rightarrow v_{\pi'}(s) \geq v_\pi(s)$$

MC 策略迭代：MC 控制



策略评估：MC 策略评估， $Q = q_\pi$

策略提升：贪心策略提升

无限探索，无穷时收敛为贪心策略 (Greedy in the Limit with Infinite Exploration, GLIE)

智能体能无限次数地探索所有的状态-动作对： $\lim_{k \rightarrow \infty} N_k(s,a) = \infty$

策略在无穷时收敛到贪心策略

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = 1(a = \arg \max_{a \in A} Q_k(s,a))$$

e-贪心策略的探索率 ϵ 能随时间衰减到 $\epsilon_k = \frac{1}{k}$ ，则它就满足 GLIE

GLIE 蒙特卡洛控制

使用策略 μ 采第 k 次事件 $\{S_1, A_1, R_1, \dots, S_T\} \sim \mu$

对事件中的每个状态 s 和动作 A_t

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

基于新得到的动作-价值函数对策略进行提升

$$\epsilon \leftarrow \frac{1}{k} \quad \pi \leftarrow \text{greedy}(Q)$$

GLIE 蒙特卡洛控制能够收敛到最优动作-价值函数

$$Q(s,a) \leftarrow q_*(s,a)$$

MC vs TD 控制

TD 学习和 MC 方法相比有如下优势：

低方差，在线学习，无完整事件序列

在控制方法中使用 TD 替代 MC：

基于 TD 方法训练 $Q(S,A)$ ，使用 e-贪心方法进行策略提升，每个时刻都状态更新

基于 Sarsa 方法更新动作-价值函数

动作-价值函数的贝尔曼期望方程：

$$q_\pi(s,a) = R_s^a + \gamma \sum_{s',a'} P_{ss'}^{aa'} \pi(a'|s') q_\pi(s',a')$$

在某一时刻智能体观察到的一段事件 (S,A,R,S',A') 基于样本的 TD 更新：

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

基于 Sarsa 的在策略控制：

在每一时刻：

策略评估：

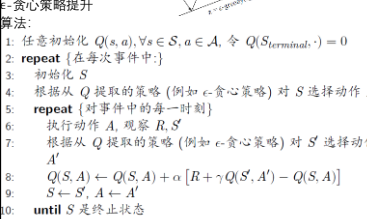
使用 Sarsa 方法， $Q \approx q_\pi$

策略提升：

e-贪心策略提升

算法：

1. 任意初始化 $Q(s,a)$, $\forall s \in S, a \in A$, 令 $Q(S_{\text{terminal}}, \cdot) = 0$
2. repeat {在每次事件中:}
3. 初始化 S
4. 根据从 Q 提取的策略 (例如 e-贪心策略) 对 S 选择动作 A
5. repeat {对事件中的每一时刻:}
6. 执行动作 A , 观察 R, S'
7. 根据从 Q 提取的策略 (例如 e-贪心策略) 对 S' 选择动作 A'
8. $Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma Q(S',A') - Q(S,A)]$
9. $S \leftarrow S', A \leftarrow A'$
10. until S 是终止状态
11. until $Sarsa$ 的收敛性:



Sarsa 的收敛性：

如下条件满足时 Sarsa 算法是收敛到最优动作-价值函数：

策略序列 $\pi_t(s)$ 是 GLIE 的

更新步长 α_t 满足 Robbins-Monro 序列要求

$$\sum_{t=1}^{\infty} \alpha_t = \infty \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Sarsa vs MC

注意：蒙特卡洛方法并不能直接应用在这个问题上

并不是所有的策略都能到达终止状态，如果学到的某一策略使智能体停留在原地不动，MC 在基于该策略对事件采样时永远都无法结束

Sarsa 算法没有这样的问题

每一步都在学习，很快就会发现这样的策略性能很差，然后转向其它策略

n-步 Sarsa

考虑 n-步的回报， $n=1, 2, \dots$ ：

$n=1$ (Sarsa) $q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$

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

$n=\infty$ (MC) $q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{t+T-1} R_T$

定义 n-步的 Q-回报

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

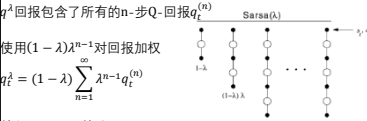
n-步 Sarsa 算法更新 $Q(s,a)$ ，向 n-步 Q-回报逼近

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

前向 Sarsa(X)

$q_t^{(n)}$ 回报包含了所有的 n-步 Q-回报 $q_t^{(n)}$

使用 $(1-\lambda)\lambda^{n-1}$ 对回报加权

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


前向 Sarsa(X) 算法

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

和 TD(X) 算法一样，我们将资格迹引入到在线的控制算法中

但是 Sarsa(X) 算法不同的地方在于为每个状态-动作对都记录资格迹-之前只是对状态记录

$$E_t(s,a) = 0 \quad E_t(s,a) = \gamma \lambda E_{t-1}(s,a) + 1(S_t = s, A_t = a)$$

在每一时刻会对每个状态-动作对更新 $Q(s,a)$ 更新

更新量与 TD 误差 δ_t 和资格迹 $E_t(s,a)$ 成正比

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

后向 Sarsa(X)

一步 Sarsa 只对最终导致高奖励的最后一步动作强化它的价值

资格迹方法能够对事件中的多个动作强化它们的价值

强化的幅度(箭头大小) 随着离奖励步数的增加而衰减

衰减率等于 λ

这里 $\lambda = 1, \lambda = 0.9$

高策略学习

想要对目标策略 $\mu(a|s)$ 进行评估，计算 $v_\pi(s)$ 或 $q_\pi(s,a)$ ，但是智能体实际的行为策略是 $\mu(a|s)$ ： $\{S_1, A_1, R_1, \dots, S_T\} \sim \mu$

为什么要考虑这种情况？

有时智能体需要观察人类或别的智能体的行为去学习

有时需要重复利用旧策略 $\pi_1, \pi_2, \dots, \pi_{t-1}$ 产生的经验去学习

执行探索性的策略去学习最优策略

执行单一的策略去学习多个策略

Q-学习

考虑基于动作-价值 $Q(s,a)$ 的离策略学习，不再使用重要性采样

智能体下一时刻执行的动作是由行为策略产生 $A_{t+1} \sim \mu(S_t)$

但是学习算法考虑的是由另一个目标策略产生的后继动作

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

更新 $Q(S_t, A_t)$ 向另一个后继动作的价值逼近

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

Q-学习实现离策略控制：

目标策略 μ 是由 $Q(s,a)$ 生成的贪心策略

$$\pi(S_{t+1}) = \arg \max_{a \in A} Q(S_{t+1}, a)$$

行为策略 μ 是由 $Q(s,a)$ 生成的例如 e-贪心策略

那么 Q-学习的目标就是

$$R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a) = R_{t+1} + \gamma \arg \max_{a \in A} Q(S_{t+1}, a)$$

Q-学习控制算法

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

Q-学习控制能够收敛到最优动作-价值函数。

Q-学习算法实现高策略控制：

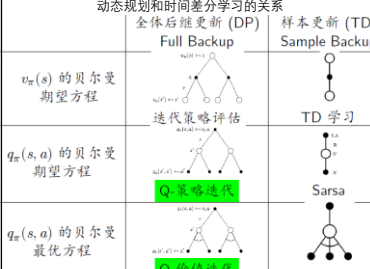
1. 任意初始化 $Q(s,a)$, $\forall s \in S, a \in A$, 令 $Q(S_{\text{terminal}}, \cdot) = 0$
2. repeat {在每次事件中:}
3. 初始化 S
4. repeat {对事件中的每一时刻:}
5. 根据从 Q 提取的策略 (例如 e-贪心策略) 对 S 选择动作 A
6. 执行动作 A , 观察 R, S'
7. $Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma \max_{a' \in A} Q(S', a') - Q(S,A)]$
8. $S \leftarrow S'$
9. until S 是终止状态
10. until

Sarsa vs Q-学习

Q-学习正确地学到了最优路径，即沿着悬崖的边缘行走，但是由于使用的 e-贪心策略会采取一定概率的随机动作，有时智能体跌入悬崖

Sarsa 学到的是安全路径，在学习过程中考虑了随机探索动作的影响，即使 Sarsa 学到的安全路径比 Q-学习的最优路径行走步数要多，但是每次获得的奖励和对比 Q-学习的高

动态规划和时间差分学习的关系



全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)

全局后维更新 (DP) 样本更新 (TD)