

。时序差分算法

(1) MC Method:

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

其中 G_t 要完成一次采样才能知道。

(2) 时序差分:

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

$$V_{\pi}(S) = E_{\pi}[G_t | S_t = S]$$

$$= E_{\pi}[\sum \gamma^k R_{t+k} | S_t = S]$$

$$= E_{\pi}[R_t + \sum_{k=1}^{\infty} \gamma^k R_{t+k} | S_t = S]$$

$$= E_{\pi}[R_t + \gamma V_{\pi}(S_{t+1}) | S_t = S]$$

*: 单步更新, 完成了策略评估。

(3) 策略提升: 估计 $Q(s, a)$, 然后取 \max .

Sarsa Algorithm: $S - a - r - S' - a'$

for $e = 1 \rightarrow E$:

初始状态 S .

ϵ -greedy \rightarrow 选择 a .

for $t = 1 \rightarrow T$:

得到 r' 与 S' (环境反馈)

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

$$S \leftarrow S', a \leftarrow a'$$

*: 主要维护了一个 Q -table.

(4) 多步时序差分:

$$G_t = r_t + \gamma r_{t+1} + \dots + \gamma^n Q(S_{t+n}, a_{t+n})$$

(5) Q-Learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

实际上, Q-learning 直接估计 Q^* (off-policy)

Sarsa 在估计当前 ϵ -greedy 下的 $Q(s, a)$ (on-policy)

(b) 采样策略: π_{sample}

目标策略: π_{target}

on-policy: $\pi_{\text{sample}} = \pi_{\text{target}}$

off-policy: $\pi_{\text{sample}} \neq \pi_{\text{target}}$