

DP

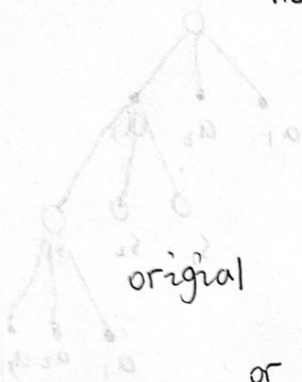
策略改进定理

$$\pi, \pi', \text{ if } \forall s \in S, q_{\pi}(s, \pi') \geq V_{\pi}(s) \Rightarrow V_{\pi'}(s) \geq V_{\pi}(s)$$

Proof:

$$q_{\pi}(s, a) = \sum_{s', r} P(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

$$= E_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = a]$$



original $q_{\pi}(s, a) = E_{\pi} [G_t | S_t = s, A_t = a]$

or

$$q_{\pi}(s, a) = \sum_{s' \in S} P(s, a, s') \left[R(s, a) + \gamma \sum_{a'} \pi(s', a') q_{\pi}(s', a') \right]$$

tips:

$$V_{\pi}(s')$$

$$q_{\pi}(s, a)$$

one parameter

about s and a.

$$G_t = R_{t+1} + \gamma G_{t+1}$$

Bellman equation

$$\Rightarrow E_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = a]$$

substitute π'

eg:

$$\Rightarrow E [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi'(s)]$$

$$\Rightarrow E_{\pi'} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

R_{t+1} has been controlled by π'

now

$$\begin{aligned}
 V_{\pi}(s) &= E_{\pi} [G_t | s_t = s] \quad \boxed{MC} \\
 &\downarrow \\
 &= E_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \\
 &= E_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \quad \boxed{TD} \\
 &= \sum_a \pi(a|s) \sum_{s', r} \underbrace{p(s', r | s, a)}_{\downarrow V_{\pi}} (r + \gamma V_{\pi}(s')) \quad \boxed{DP}
 \end{aligned}$$

DP

$\{V_k\} \rightarrow V_{\pi}$ 自举

$$V_k = \begin{pmatrix} V_k(s_1) \\ V_k(s_2) \\ \vdots \\ V_k(s_n) \end{pmatrix}$$

$$V_1 = \begin{pmatrix} V_1(s_1) \\ \vdots \\ V_1(s_n) \end{pmatrix}$$

必须要

MC 采样

$$V_{\pi}(s) = E_{\pi} [G_t | s_t = s] \quad G_t$$

$$\boxed{s_0, a_0, R_{t+1}, s_{t+1}, a_{t+1}, R_{t+2}, \dots}$$

eg:

必须走完才能得到 G_t

$$V_1(s_1) = V_1(s_1) + \alpha (G_t - V_1(s_1))$$

处于 s_1 的状态下,
更新 $V(s_1)$ (假设在 V_1) 真实 当前估计
↓
第一轮.

TD

mouse

克服更新延迟的问题.

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

TD-0

| | | | | | | | | |
|-------|-------|-----------|-----------|-----------|-----------|-----|-------|-------|
| s_t | A_t | R_{t+1} | s_{t+1} | A_{t+1} | R_{t+2} | ... | R_T | s_T |
|-------|-------|-----------|-----------|-----------|-----------|-----|-------|-------|

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

先走一步, 得到 R_{t+1} . 不走了, 直接去表里找 $V(s_{t+1})$

$$V(s_3) \leftarrow V(s_3) + \alpha (R_4 + \gamma V(s_4) - V(s_3))$$

比如 s_4 的 $V(s_4)$

其实从表里拿.

| |
|----------|
| $V(s_1)$ |
| $V(s_2)$ |
| $V(s_3)$ |
| $V(s_4)$ |

注意, 这里的 $V(s_4)$ 并不是一定的

因为状态 (下一个) 是根据 π 抽样得到的

可能是 s_1, s_2, s_3, s_4 中的任何一个.

控制. 单纯的 $V(s)$ 是不能够修改 π 的. (greedy)

因此, 先求 $Q(s, a)$

SARSA

进入这个状态, 抽样得到 a_6

$$Q(s_t, A_t)$$

| | | | | | |
|-------|-------|-----------|-----------|-----------|-----------|
| s_3 | a_3 | 2 | s_5 | a_6 | |
| s_t | A_t | R_{t+1} | s_{t+1} | A_{t+1} | R_{t+2} |

$$= Q(s_t, A_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(s_t, A_t)]$$

(SARSA)

这是根据 π 随机抽样得到的

(同轨迹更新)

高维执行:

07

当我们进入 S_t 的时候, 不按策略函数来采样.

我有 Q-table, 我计算 $Q(S_t, a_1)$, 我自己行!

$Q(S_t, a_2)$

找最大值的 action.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \beta \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

Q-learning

同轨 SARSA (听教练的)

高维 Q-learning (不听)

期望 SARSA

教练 (b) or π

| S_{t+1} | A_{t+1} | R_t |
|-----------|-----------|-------|
| S_t | a_1 | 0.7 |
| S_t | a_2 | 0.3 |

Q-table

| S | A | $Q(S, a)$ |
|-------|-------|-----------|
| S_t | a_1 | 100 |
| S_t | a_2 | 99 |

用期望来做为 \leftarrow (根据教练, π 来到 Q-table 相对应的 Q)

$$100 \times 0.7 + 99 \times 0.3$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \beta (E_b [Q(S_{t+1}, a) | S_{t+1}]) - Q(S_t, A_t)]$$

还是使用了 π , 还是同轨的

~~图 5.11~~

(I) 默认情况下 $\pi = b$, 同轨

(II) 若 π 是贪心策略, 也即 $\pi(s_{t+1}) = \arg \max_a Q(s_{t+1}, a) \triangleq a^*$

则期望为高轨策略

$$E_{\pi} [Q(s_{t+1}, a) | s_{t+1}] = \max_a Q(s_{t+1}, a)$$

because

$$\Rightarrow \pi(a | s_{t+1}) = \begin{cases} 1 & a = \arg \max_a Q(s_{t+1}, a) \\ 0 & \text{else} \end{cases}$$

Q-learning 是 期望 SARSA 的特例.