

No. 1. (a)

$$Q(3, left) = Q(3, left) + \alpha [r + \gamma^2 Q(2, right) - Q(3, left)]$$

$$\text{例有 } 7 + [-1 + 0.9 \times 8 - 7]$$

得  $Q(3, left) = 6.2$

(b) 状态 1  $\rightarrow$  right  $\rightarrow$  状态 2

$$Q(1, right) = Q(1, right) + \alpha [r + \gamma^2 Q(2, up) - Q(1, right)]$$

$$\text{例有 } 3 + 0.2 \times [-1 + 0.8 \times 6 - 3]$$

得  $Q(1, right)$  更新为 3.6

② 状态 2  $\rightarrow$  up  $\rightarrow$  状态 5

$$Q(2, up) = Q(2, up) + \alpha [r + \gamma^2 Q(5, right) - Q(2, up)]$$

$$\text{例有 } 6 + 0.2 \times [-1 + 0.8 \times 8 - 6]$$

得  $Q(2, up)$  更新为 5.88

③ 状态 5  $\rightarrow$  right  $\rightarrow$  状态 6

$$Q(5, right) = Q(5, right) + \alpha [r + \gamma^2 Q(6) - Q(5, right)]$$

$$\text{例有 } 8 + 0.2 \times [10 + 0 - 8]$$

得  $Q(5, right)$  更新为 8.4

Next, 设?处状态为 " $S_3$ "

$$Q(S_3, \text{study}) = R + \gamma V = 10 + 1 \times 0 = 10$$

$$\begin{aligned} Q(S_3, \text{Pub}) &= R + \gamma V \\ &= 1.02 + 0.2V_1 + 0.4V_2 + 0.4V_3 \\ &= \cancel{1.02} + 1.84 + 0.4V_3 \end{aligned}$$

$$\begin{aligned} V_3 &= \pi(\text{study}) \cdot Q(S_3, \text{study}) + \pi(\text{Pub}) \cdot Q(S_3, \text{Pub}) \\ &= 0.5 \times 10 + 0.5 \times (1.84 + 0.4V_3) \end{aligned}$$

解得  $V_3 = 4$

(b) 由 a)  $Q(S_3, \text{study}) = 10$

$$Q(S_3, \text{Pub}) = 1.84 + 0.4 \times 7.4 = 4.8$$

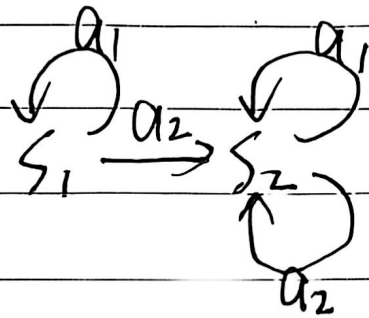
$$Q(S, a) = R(S, a) + \gamma V_S$$

$$S_1: Q(S_1, \text{Play}) = -1 - 2.3 = -3.3 \quad Q(S_1, \text{Study}) = -2 + 2.7 = 0.7$$

$$S_2: Q(S_2, \text{Sleep}) = 0 + 0 = 0 \quad Q(S_2, \text{Study}) = -2 + 4 = 2$$

$$S_4: Q(S_4, \text{Play}) = -1 - 2.3 = -3.3 \quad Q(S_4, \text{Quit}) = 0 + 1.3 = 1.3$$

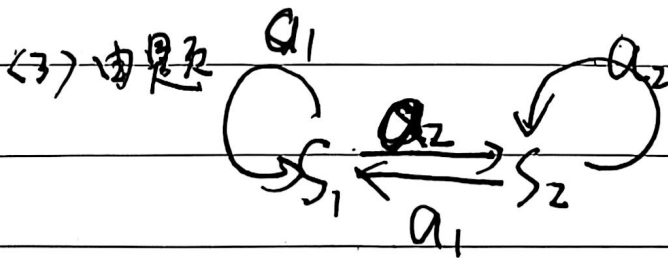
3. 分析 (1) 由题



$S_1 \xrightarrow{a_2} S_2$  时奖励 1

(2) 由题 转移函数同 (1)

只在  $S_2$  执行任意动作时奖励 1



$S_1$  执行  $a_1$ ,  $S_2$  执行  $a_2$  时奖励 1

综上 (a) MDP1 不存在无限大  $S_1$  值函数, 最大在  $S_1$  选择  $a_2$   $V(S_1)$

MDP2 存在无限大策略, 在  $S_1$  选择  $a_2$

MDP3 存在无限大策略, 在  $S_1$  选择  $a_1$  在  $S_2$  选择  $a_2$

(b) 设  $\gamma = 1/2$ , 值函数收敛,

$$\text{MDP1: } V(S_2) = 0 \quad Q(S_1, a_1) = \gamma^* V(S_1)$$

$$Q(S_1, a_2) = \gamma^* V(S_2)$$

$$V(S_1) = \max(\gamma^* V(S_1), 1) = 1$$

最优策略:  $S_1$  选择  $a_2$   $S_2$  任意

$$\text{MDP}_2: Q(s_1, a_1) = 0 + \gamma V(s_1) \quad Q(s_2, a_1) = 1 + \gamma V(s_2)$$

$$Q(s_1, a_2) = 0 + \gamma V(s_2) \quad Q(s_2, a_2) = 1 + \gamma V(s_2)$$

$$V(s_2) = 1 + \gamma V(s_2)$$

$$V(s_1) = \max(\gamma V(s_1), \frac{\gamma}{1-\gamma}) \quad V(s_2) = \frac{1}{1-\gamma}$$

$$= \frac{\gamma}{1-\gamma}$$

最优策略:  $s_1$  选择  $a_2$ ,  $s_2$  任意。

$$\text{MDP}_3: Q(s_1, a_1) = 1 + \gamma V(s_1) \quad Q(s_2, a_1) = \gamma V(s_1)$$

$$Q(s_1, a_2) = \gamma V(s_2) \quad Q(s_2, a_2) = 1 + \gamma V(s_2)$$

$$V(s_1) = \max(\gamma V(s_2), 1 + \gamma V(s_1))$$

$$V(s_2) = \max(\gamma V(s_1), 1 + \gamma V(s_2))$$

$$\text{假设 } V(s_1) = 1 + \gamma V(s_1) \text{ 则 } V(s_1) = V(s_2) = \frac{1}{1-\gamma}$$

$$Q(s_1, a_2) = \frac{\gamma}{1-\gamma} < Q(s_1, a_1) = \frac{1}{1-\gamma} \text{ 选择 } a_1$$

同理  $s_2$  选择  $a_2$

最优策略:  $s_1$  选择  $a_1$

$s_2$  选择  $a_2$