

# 清华大学实验报告

系别 \_\_\_\_\_ 班号 \_\_\_\_\_ 姓名 \_\_\_\_\_ (同组姓名: \_\_\_\_\_)

作实验日期 2021年 12 月 18 日 教师评定: \_\_\_\_\_

2021.12.18

第8次作业.

2. (1) 首先计算更新后的状态价值函数  $V_2(s)$ .

$$V_2(A) = r + \gamma \sum_{s' \in S} P_{A|A}^s V_1(s') = -8 + 0.5 \times 1.0 \times V_1(B) = -8 + 0.5 \times 1.0 \times 2 = -7$$

$$V_2(B) = \max \{ 2 + 0.5 \times 1.0 \times V_1(A), -2 + 0.5 \times 1.0 \times V_1(C) \} = \max \{ 3, -1 \} = 3$$

$$V_2(C) = \max \{ 8 + 0.5 \times 1.0 \times V_1(B), 0 \times 0.75 + 4 \times 0.25 + 0.5 \times 0.25 \times V_1(A) + 0.5 \times 0.75 \times V_1(C) \} = \max \{ 9, 2 \} = 9$$

$$\text{则 } V_2(s) = (-7, 3, 9)^T$$

由确定性策略:  $\pi_2(a=ab|s=A)=1$ .

$$\pi_2(a=ba|s=B)=1, \pi_2(a=bc|s=B)=0$$

$$\pi_2(a=cb|s=C)=1, \pi_2(a=ca|s=C)=0.$$

(2) 同样地, 计算  $V'_2(s)$

$$V'_2(A) = V_2(A) = -7$$

$$V'_2(B) = \max \{ 2 + 0.5 \times 1.0 \times V'_2(A), -2 + 0.5 \times 1.0 \times V_1(C) \} = \max \{ -1.5, -1 \} = -1$$

$$V'_2(C) = \max \{ 8 + 0.5 \times 1.0 \times V'_2(B), 0 \times 0.75 + 4 \times 0.25 + 0.5 \times 0.25 \times V'_2(A) + 0.5 \times 0.75 \times V_1(C) \} = \max \{ 7.5, 0.875 \} = 7.5$$

$$\text{则 } V'_2(s) = (-7, -1, 7.5)^T$$

由确定性策略:  $\pi'_2(a=ab|s=A)=1$ .

$$\pi'_2(a=ba|s=B)=0, \pi'_2(a=bc|s=B)=1$$

$$\pi'_2(a=cb|s=C)=1, \pi'_2(a=ca|s=C)=0.$$

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

$$4 \rightarrow 5 \quad V(4) = V(4) + 0.5(-1 + 1 \times V(5) - V(4)) = 0 + 0.5 \times (-1) = -0.5$$

$$5 \rightarrow 4 \quad V'(5) = V(5) + 0.5(-1 + 1 \times V'(4) - V(5)) = 0 + 0.5 \times (-1 + 1 \times -0.5) = -0.75$$

$$4 \rightarrow 3 \quad V''(4) = V'(4) + 0.5(-1 + 1 \times V(3) - V'(4)) = -0.5 + 0.5 \times (-1 + 0.5) = -0.75$$

$$3 \rightarrow \text{terminate} \quad V'(3) = V(3) + 0.5(-1 + 1 \times 0 - V(3)) = -0.5$$

更新后的  $V$  值:

0	0	0
-0.5	-0.75	-0.75
0	0	0

(2)

上	-4	-3	-1	-3	-4	-2	-4
右	-3	-3	-2	-4	-2	-3	-3
下	-4	-3	-4	-2	-2	-3	-4
左	-3	-2	-3	-3	-4	-3	-2
	1	2	3	4	5	6	7

SARSA算法:

① 4 上 → 7 右

$$Q(4, 7) = Q(4, 7) + \alpha [-1 + \gamma Q(7, 右) - Q(4, 7)]$$

$$= -2 + 1 \times [-1 + 1 \times (-2) - (-2)] = -3$$

② 7 右 → 6 上

$$Q(7, 右) = Q(7, 右) + 1 \times [-1 + \gamma Q(6, 上) - Q(7, 右)]$$

$$= -2 + 1 \times [-1 + 1 \times (-2) - (-2)] = -3$$

③ 6 上 → 3 上

$$Q(6, 上) = Q(6, 上) + 1 \times [-1 + \gamma Q(3, 上) - Q(6, 上)]$$

$$= -2 + 1 \times [-1 + 1 \times (-1) - (-2)] = -2$$

④ 3 上 → 终止

$$Q(3, 上) = Q(3, 上) + 1 \times [-1 + \gamma Q(\text{terminate}, 上) - Q(3, 上)]$$

$$= -1 + 1 \times [-1 + 1 \times 0 - (-1)] = -1$$

故更新后的Q表为

上	-4	-3	-1	-3	-4	-2	-4
右	-3	-3	-2	-4	-2	-3	-3
下	-4	-3	-4	-3	-2	-3	-4
左	-3	-2	-3	-3	-4	-3	-3
	1	2	3	4	5	6	7

Q-learning算法:

由于两种算法均使用贪心方式选择动作, 故得到的Q表应完全相同. 更新后的Q表如上表所示.