

# 人工智能基础第八次作业

3.

11) 首次访问的蒙特卡洛预测:

$$v(A) = \frac{1}{2} \times [(+3+2-4+4-3) + (+3-3)] = 1$$

$$v(B) = \frac{1}{2} \times [(-4+4-3) + (-2+3-3)] = -2.5$$

每次访问的蒙特卡洛预测:

$$v(A) = \frac{1}{4} \times [(+3+2-4+4-3) + (+2-4+4-3) + (+4-3) + (+3-3)] = 0.5$$

$$v(B) = \frac{1}{4} \times [(-4+4-3) + (-3) + (-2+3-3) + (-3)] = -2.75$$

12) 共8次状态转移过程

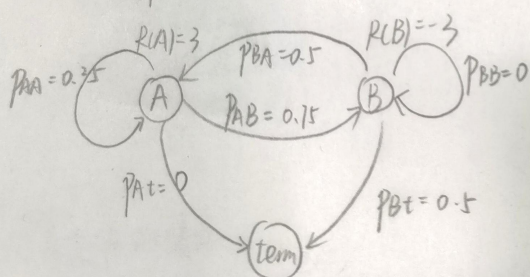
其中  $A \rightarrow A$  1次  $A \rightarrow B$  3次  $B \rightarrow A$  2次  $B \rightarrow \text{terminate}$  2次

$$\text{故 } p_{AA} = \frac{1}{1+3} = 0.25, p_{AB} = \frac{3}{1+3} = 0.75, p_{At} = 0$$

$$p_{BA} = \frac{2}{2+2} = 0.5, p_{Bt} = \frac{2}{2+2} = 0.5, p_{Bb} = 0$$

$$R(A) = \frac{1}{4} \times (3+2+4+3) = 3$$

$$R(B) = \frac{1}{4} \times (-4-3-2-3) = -3$$



13) 状态价值不变期望方程  $V_{\pi}(S) = \sum_{a \in A} \pi(a|S) (R_S^a + \gamma \sum_{S' \in S} P_{SS'}^a V_{\pi}(S'))$

其中  $S = A, B, \text{term}$

$$\begin{bmatrix} V(A) \\ V(B) \\ V(t) \end{bmatrix} = \begin{bmatrix} R(A) \\ R(B) \\ R(t) \end{bmatrix} + \gamma \times \begin{bmatrix} p_{AA} & p_{AB} & p_{At} \\ p_{BA} & p_{BB} & p_{Bt} \\ p_{tA} & p_{tB} & p_{tt} \end{bmatrix} \begin{bmatrix} V(A) \\ V(B) \\ V(t) \end{bmatrix}$$



班级:

姓名:

编号:

科目:

$$\begin{bmatrix} V(A) \\ V(B) \\ V(t) \end{bmatrix} = \begin{bmatrix} -3 \\ -3 \\ 0 \end{bmatrix} + 1 \times \begin{bmatrix} 0.25 & 0.75 & 0 \\ 0.5 & 0 & 0.5 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} V(A) \\ V(B) \\ V(t) \end{bmatrix}$$

$$\text{得} \begin{cases} V(A) = 2 \\ V(B) = -2 \end{cases}$$

4.

1) 4→5.

$$\begin{aligned} V(4) &= V(4) + 0.5 \times (-1) + 1 \times V(5) - V(4) \\ &= 0 + 0.5 \times (-1) + 0 - 0 = -0.5 \end{aligned}$$

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

$$4 \rightarrow 3: V(4) = V(4) + 0.5 \times (-1) + 1 \times V(3) - V(4) = -0.5 + 0.5 \times (-1) + 0 - 0.5 = -0.75$$

$$3 \rightarrow t: V(3) = V(3) + 0.5 \times (-1) + 1 \times V(t) - V(3) = 0 + 0.5 \times (-1) + 0 - 0 = -0.5$$

更新后V值为

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

2)

① SARSA算法:

$$\text{状态4: 向下运动} \quad Q(4, \downarrow) = Q(4, \downarrow) + 1 \times (-1) + 1 \times Q(7, \leftarrow) - Q(4, \downarrow)$$

$$(\text{下-状态: 7, 向左}) \quad = -2 + 1 \times (-1) + 1 \times (-2) - (-2) = -3$$

$$\text{状态7: 向左运动} \quad Q(7, \leftarrow) = Q(7, \leftarrow) + 1 \times (-1) + 1 \times Q(6, \uparrow) - Q(7, \leftarrow)$$

$$(\text{下-状态: 6, 向上}) \quad = -2 + 1 \times (-1) + 1 \times (-2) - (-2) = -3$$

$$\text{状态6: 向上运动} \quad Q(6, \uparrow) = Q(6, \uparrow) + 1 \times (-1) + 1 \times Q(3, \uparrow) - Q(6, \uparrow)$$

$$(\text{下-状态: 3, 向上}) \quad = -2 + 1 \times (-1) + 1 \times (-1) - (-2) = -2$$



状态3: 向上运动  
(下状态: 状态4)

$$Q(3, \uparrow) = Q(3, \uparrow) + \gamma \times (-1 + \gamma \times Q(1, \uparrow) + Q(3, \uparrow))$$

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

$$= -1$$

更新后的Q表

-4	-3	-1	-3	-4	-2	-4
-3	-3	-2	-4	-2	-3	-3
-4	-3	-4	-3	-2	-3	-4
-3	-2	-3	-3	-4	-3	-3

② Q-learning 算法:

因为也采用贪心方式选择动作, 故更新过程和更新后的Q表与SARSA算法相同

0	0	0
0	0	0
0	0	0