3. 首次访问的蒙特卡洛预测

(1).

$v(A) = \frac{1}{2}[G_1 + G_2]$.

$G_1 = 3+2-4+4-3 = 2$    $G_2 = 3-3 = 0$.

故 $v(A) = 1$

$v(B) = \frac{1}{2}[G_1 + G_2]$.

$G_1 = -4+4-3 = -3$    $G_2 = -2+3-3 = -2$

故 $v(B) = -2.5$

每次访问的蒙特卡洛预测:

$v(A) = \frac{1}{4}(G_1 + G_2 + G_3 + G_4)$

$G_1 = 3+2-4+4-3 = 2$.    $G_2 = 2-4+4-3 = -1$    $G_3 = 4-3 = 1$.    $G_4 = 3-3 = 0$

故 $v(A) = 0.5$

$v(B) = \frac{1}{4}(G_1 + G_2 + G_3 + G_4)$

$G_1 = -4+4+3 = 3$    $G_2 = -3$    $G_3 = -2+3-3 = -2$.    $G_4 = -3$.

故 $v(B) = -2.75$

(2). $R_A = \frac{1}{4}(3+2+4+3) = 3$.    $R_B = \frac{1}{4}(-4-3-2-3) = -3$.

由统计可估计:   $P_{AA} = 0.25$.    $P_{AB} = 0.75$.

$P_{BA} = 0.5$.    $P_{Bt} = 0.5$.

故马尔可夫回报过程图应为:



$p_{AA} = 0.25$    $p_{BA} = 0.5$    $p_{BB} = 0$

$R_A = 3$    $R_B = -3$

$p_{AB} = 0.75$

$p_{At} = 0$    $p_{Bt} = 0.5$

$R_t = 0$

(3). 由于:
$$\begin{bmatrix} V_A \\ V_B \\ V_t \end{bmatrix} = \begin{bmatrix} R_A \\ R_B \\ R_t \end{bmatrix} + r \cdot P \cdot \begin{bmatrix} V_A \\ V_B \\ V_t \end{bmatrix}$$

即:
$$\begin{bmatrix} V_A \\ V_B \\ V_t \end{bmatrix} = \begin{bmatrix} 3 \\ -3 \\ 0 \end{bmatrix} + \begin{bmatrix} \frac{1}{4} & \frac{3}{4} & 0 \\ \frac{1}{3} & 0 & \frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} V_A \\ V_B \\ V_t \end{bmatrix}.$$

由于 $V_t = R_t + V_t = V_t$

不妨设 $V(t) = 0$

则:
$$\begin{cases} V_A = 3 + \frac{1}{4} V_A + \frac{3}{4} V_B \\ V_B = -3 + \frac{1}{2} V_A + \frac{1}{2} V_t \end{cases} \Rightarrow \begin{cases} V_A = 2 \\ V_B = -2 \end{cases}$$

## 4.

(1). 考虑每步:  $4 \to 5$ : $V(4) = V(4) + \alpha(R_t + rV(5) - V(4))$
$$= 0.5(-1) = -0.5$$

$5 \to 4$ :  $V(5) = V(5) + \alpha(R_t + rV(4) - V(5))$
$$= 0.5(-1-0.5) = -0.75$$

$4 \to 3$ :  $V(4) = V(4) + \alpha(R_t + rV(3) - V(4))$
$$= -0.5 + 0.5(-1+0.5) = -0.75$$

$3 \to term$ :  $V(3) = V(3) + \alpha(R_t + rV(term) - V(3))$
$$= 0.5(-1) = -0.5$$

V 值:

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

(2). SARSA 算法与 Q-learning 算法在贪心方式下结果相同.

$S_1 = 4$.  $A_1 = $ 向下,
$S_2 = 7$   $A_2 = $ 向左

$Q(4,3) = Q(4,3) + \alpha(R_2 + rQ(7,4) - Q(4,3))$

$$= -2 - 1 - 2 + 2 = -3.$$

$S_3 = 6$, $A_3 = $ 向上.

$$Q(7,4) = Q(7,4) + \alpha \left( R_3 + r Q(6,1) - Q(7,4) \right)$$

$$= -2 - 1 - 2 + 2 = -3$$

$S_4 = 3$, $A_4 = $ 向上.

$$Q(6,1) = Q(6,1) + \alpha \left( R_4 + r Q(3,1) - Q(6,1) \right)$$

$$= -2 - 1 - 1 + 2$$

$$= -2$$

$S_5 = term.$

$$Q(3,1) = Q(3,1) + \alpha \left( R_5 + r Q(term, a) - Q(3,1) \right)$$

$$= -1 - 1 + 1$$

$$= -1$$

继2. 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 |