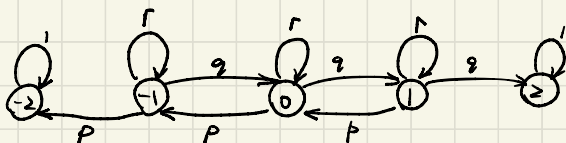


HW7

支-并 新雅-智1 2020010018

1. a) 状态空间 $S = \{-2, -1, 0, 1, 2\}$
状态转移矩阵

$$\begin{matrix} & -2 & -1 & 0 & 1 & 2 \\ \begin{matrix} -2 \\ -1 \\ 0 \\ 1 \\ 2 \end{matrix} & \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ p & r & q & 0 & 0 \\ 0 & 0 & p & r & q \\ 0 & 0 & p & r & q \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \end{matrix}$$



b) 乙执1分的情况下,再赛2局结束

① → ① → ② 只有这一种可能

$$P = r \cdot q$$

2. a) $v(A) = -1 + r[P_{AB}v(B) + P_{AC}v(C)] = 0.1v(B) + 0.1v(C) - 1$

$$v(B) = -1 + r[P_{BB}v(B) + P_{BD}v(D)] = 0.1v(B) + 0.1v(D) - 1$$

$$v(C) = -1 + r[P_{CA}v(A) + P_{CD}v(D)] = 0.1v(A) + 0.1v(D) - 1$$

$$v(D) = 0 + r[P_{DD}v(D)] = 0.2v(D)$$

有 $\begin{cases} v(A) = -1.2223 \\ v(B) = -1.09 \approx -1.1111 \\ v(C) = -1.1223 \\ v(D) = 0 \end{cases}$

	A	B	C	D
A	0	0.5	0.5	0
B	0	0.5	0	0.5
C	0.5	0	0	0.5
D	0	0	0	1

b) 1° 价值迭代

i 初始化 $v(s)$ 均为 0

ii 反复执行 $v(s) \leftarrow R(s) + r \sum_{S'} P_{ss'} v(s')$

iii 当 $\Delta v(s) < \text{threshold}$, 可认为已收敛, 停止迭代

2° 策略迭代

i 初始化任一策略 π

ii 反复执行 $v_{\pi}(s) = R(s) + r \sum_{S'} P_{ss'} v_{\pi}(s')$ (计算当前策略下 $v(s)$)

$\pi(s) \leftarrow \arg \max_a [R(s, a) + r \sum_{S'} P_{ss'} v_{\pi}(s')]$ (用新状态价值去更新策略)

iii 同上收敛.

$$3. \quad q_{\pi}(4, \text{left}) = R_t + v_{\pi}(3)$$

$$q_{\pi}(7, \text{right}) = R_t + v_{\pi}(0)$$

下用 V_i 代替 $v_{\pi}(i)$

$$V_0 = 0$$

$$V_1 = R_t + \frac{1}{4}(V_0 + V_1 + V_2 + V_4)$$

$$V_2 = R_t + \frac{1}{4}(V_1 + 2V_2 + V_5)$$

$$V_3 = R_t + \frac{1}{4}(V_0 + V_3 + V_4 + V_6)$$

$$V_4 = R_t + \frac{1}{4}(V_1 + V_2 + V_5 + V_7)$$

$$V_5 = R_t + \frac{1}{4}(V_2 + V_4 + V_5 + V_8)$$

$$V_6 = R_t + \frac{1}{4}(V_3 + 2V_6 + V_7)$$

$$V_7 = R_t + \frac{1}{4}(V_4 + V_6 + V_7 + V_8)$$

$$V_8 = 0$$

$$\text{解得 } v_{\pi}(3) = V_3 = -7$$

$$v_{\pi}(0) = V_0 = 0$$

$$\text{故有 } q_{\pi}(4, \text{left}) = R_t + v_{\pi}(3) = -8$$

$$q_{\pi}(7, \text{right}) = R_t + v_{\pi}(0) = -1$$