

1. (1) 状态 $s = \{s_0, \dots, s_7\}$ 表示子及双板纳副能晋字等8个状态, 行动为 $A = \{A_0, A_1\}$
 A_0 为向上, A_1 为向下.

(2) 此时大象向上和向下的概率相同. 转移矩阵为

$$P = \begin{matrix} & s_0 & s_1 & s_2 & s_3 & s_4 & s_5 & s_6 & s_7 \\ \begin{matrix} s_0 \\ s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \\ s_6 \\ s_7 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix}$$

(3) 贝尔曼方程为 $V_s^* = r_s + \gamma P_{s,s'} V_{s'}^*$ $r_s = \{0, 0, 0, 0, 0, \frac{0+700}{2} = 350, 0\}$ 边界处 V_s^* 为 350. 其余为 0
 $V_0^* = V_0^*$ $V_1^* = 0 + \frac{1}{2} V_2^*$ $V_2^* = \frac{V_1^* + V_3^*}{2}$ $V_3^* = \frac{V_2^* + V_4^*}{2}$ $V_4^* = \frac{V_3^* + V_5^*}{2}$ $V_5^* = \frac{V_4^* + V_6^*}{2}$
 $V_6^* = \frac{V_5^* + V_7^*}{2} + 350$ $V_7^* = V_7^*$, 解得 6 个有效方程. 设 $V_0^* = 0$, $V_7^* = 0$. 共 8 个方程
解得 $V_0^* = 0$ $V_1^* = 100$ $V_2^* = 200$ $V_3^* = 300$ $V_4^* = 400$ $V_5^* = 500$ $V_6^* = 600$ $V_7^* = 0$

(4) $q_\pi(s, a_0) = r_s + \gamma \sum_{s'} P_{s,s'}^a V_{s'}^*(s)$ 向上行动价值为 $q_\pi(s, A_0) = r_s + \gamma V_{s+1}$ $r_s = \begin{cases} 700, & s=s_7 \\ 0, & s \neq s_7 \end{cases}$
解得 $q_\pi(s, A_0) = \{0, 200, 300, 400, 500, 600, 700\}$

同理. 向下行动价值 $q_\pi(s, A_1) = r_s + \gamma V_{s-1}$ $r_s = 0$

得 $q_\pi(s, A_1) = \{0, 0, 100, 200, 300, 400, 500, 0\}$

(5) 蒙特卡洛模拟方法. 估计状态价值. 之后可通过直接求解线性方程解出最优策略.

(6) 利用深度学习模型. 如 Deep Q Network. 优化策略模型. 以获得较好的策略或逼近最优策略

3. $V_s^* = r_s + \gamma P_{s,s'} V_{s'}^*$ 有 $V_0 = (\frac{1}{4} \times 3 + \frac{3}{4} \times -3) + \frac{1}{4} V_1 + \frac{3}{4} V_2 = \frac{1}{4} V_1 + \frac{3}{4} V_2 - \frac{3}{2}$ 同理 $V_2 = \frac{1}{2} V_1 + \frac{1}{2}$
解得 $V_1 = -1$, $V_2 = 1$

(2) 首次访问: $V(s_A) = \frac{1}{n} \sum G_t$ 第一个片段 $G_{10} = 3 + 2 - 4 + 4 - 3 = 2$ 第二个 $G_{11} = 3 - 3 = 0$
 $V(s_A) = \frac{1}{2} (2 + 0) = 1$.

同理 第一个片段 $G_{20} = 1 - 4 + 4 - 3 = -2$ 第二个片段 $G_{21} = -2 + 3 - 3 = -2$
 $V(s_B) = \frac{1}{2} (-2 - 2) = -2$

每次访问 片段一有 3 次访问 A. $G_{A0} = 3 + 2 - 4 + 4 - 3 = 2$, $G_{A1} = 2 - 4 + 4 - 3 = -1$, $G_{A2} = 4 - 3 = 1$
片段二有 1 次访问 A $G_{A3} = 3 - 3 = 0$ $V(s_A) = \frac{1}{4} (2 - 1 + 1 + 0) = \frac{1}{2}$

片段一有 2 次 B: $G_{B0} = -4 + 4 - 3 = -3$ $G_{B1} = -3$ 片段二有 2 次: $G_{B2} = -2 + 3 - 3 = -2$ $G_{B3} = -3$
 $V(s_B) = \frac{1}{4} (-3 - 3 - 2 - 3) = -\frac{11}{4}$



2. 同步价值迭代 $V_{n+1}(s) = \max_{a \in A} \{ r_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_n(s') \}$ V_0 初始值为 0

第 1 步: $V_1(1) = \max(-8 + 0.5 \times 0) = -8$

$V_1(2) = \max(2 + 0.5 \times 0, -3 + 0.5 \times 0) = 2$

$V_1(3) = \max(0.25 \times 4 + 0.75 \times 0, 8 + 0.5 \times 0) = 8$

第 2 步: $V_2(1) = \max(-8 + 0.5 \times 2) = -7$

$V_2(2) = \max(2 + 0.5 \times (-8), -3 + 0.5 \times 8) = 1$

$V_2(3) = \max(0.25 \times 4 + 0.75 \times 0 + 0.5 \times (0.25 \times (-8) + 0.75 \times 8), 8 + 0.5 \times 2) = 9$

A 策略只能为 ab, $q(A, ab) = (-8 + 0.5 \times 1) = -7.5$ $\pi(a|A)$ 为 (ab)

B 策略为 ba, bc $q(B, ba) = 2 + 0.5 \times 7 = 3.5$ $\pi(a|B)$ 为 (bc)

$q(B, bc) = -3 + 0.5 \times 9 = 1.5$

C 策略为 ca, cb $q(C, ca) = 0.25 \times 9 + 0.75 \times 0 + 0.5 \times (0.25 \times -7 + 0.75 \times 9) = 3.5$

$q(C, cb) = 8 + 0.5 \times 1 = 8.5$ $\pi(a|C)$ 为 (cb)

异步价值迭代: $V_1(1) = \max(-8 + 0.5 \times 0) = -8$

$V_1(2) = \max(2 + 0.5 \times -8, -3 + 0.5 \times 0) = -2$

$V_1(3) = \max(1 + 0.5 \times (0.25 \times -8 + 0.75 \times 0), 8 + 0.5 \times -2) = 7$

$V_2(1) = \max(-8 + 0.5 \times -2) = -9$

$V_2(2) = \max(2 + 0.5 \times -9, -3 + 0.5 \times 7) = 0.5$

$V_2(3) = \max(1 + 0.5 \times (0.25 \times -9 + 0.75 \times 1), 8 + 0.5 \times 0.5) = 8.25$

$q(A, ab) = -8 + 0.5 \times 0.5 = -7.75$, $\pi(a|A) = (ab)$

$q(B, ba) = 2 + 0.5 \times -9 = -2.5$

$q(B, bc) = -3 + 0.5 \times 8.25 = 1.125$

$q(C, ca) = 1 + 0.5 \times (0.25 \times -9 + 0.75 \times 8.25) = 2.46875$

$q(C, cb) = 8 + 0.5 \times 0.5 = 8.25$ $\pi(a|C) = (cb)$

这里 $\delta(x) = \begin{cases} 1 & a \neq x \\ 0 & a = x \end{cases}$

