

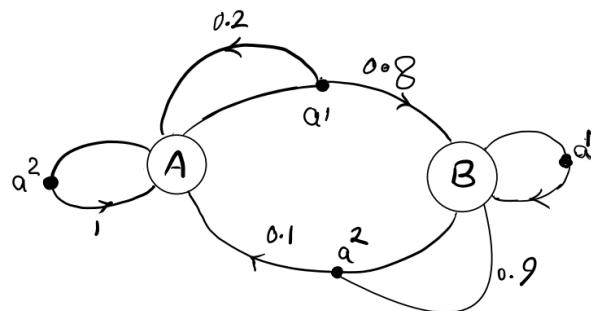


Problem 1.

Consider the following system with the state space $S = \{A, B\}$, and action space $A = \{a^1, a^2\}$. The state transition diagram is shown below, where $P(s' = B | s = A, a = a^1) = 0.8$, $P(s' = A | s = A, a = a^1) = 0.2$.

The reward is as follows:

$+2$	<i>moving to state B</i>
0	<i>moving to state A</i>
-1.5	<i>taking action a^1</i>
-1	<i>taking action a^2</i>



a) Construct transition matrices $M(a^1)$, $M(a^2)$ and compute $R_s^{a^1}$, $R_s^{a^2}$.

b) Perform matrix-form Policy Iteration method with initial policy $\pi^*(A) = a^2$, $\pi^*(B) = a^1$ and $\gamma = 0.9$ to compute π^* .

$$a) M(a^1) = \begin{bmatrix} A & B \\ A & B \end{bmatrix} \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix}$$

$$M(a^2) = \begin{bmatrix} A & B \\ A & B \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix}$$

$$R_{SS'}^{a^1} = \begin{bmatrix} -1.5 & 0.5 \\ -1.5 & 0.5 \end{bmatrix} \quad R_{SS'}^{a^2} = \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$RS^{a^1} = \begin{bmatrix} R(A, a^1) \\ R(B, a^1) \end{bmatrix} = (M(a^1) \odot R_{SS'}^{a^1}) \mathbf{1}_{2 \times 1}$$

$$= \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix}$$

$$RS^{a^2} = \begin{bmatrix} R(A, a^2) \\ R(B, a^2) \end{bmatrix} = (M(a^2) \odot R_{SS'}^{a^2}) \mathbf{1}_{2 \times 1}$$

$$= \begin{bmatrix} -1 \\ 0.2 \end{bmatrix}$$

b) $M(\lambda') = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad RS^{x^0} = \begin{bmatrix} -1 \\ 0.5 \end{bmatrix}$

$$V^{x^1} = (I - 8M(\lambda'))^{-1} RS^{x^1}$$

$$= \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} -1 \\ 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} -10 \\ 5 \end{bmatrix}$$

$$\pi^2 = \underset{a \in A}{\operatorname{argmax}} R_S^{a^1} + \gamma M(a^1) V^{\lambda^1}$$

$$= \operatorname{argmax} \left\{ R_S^{a^1} + \gamma M(a^1) V^{\lambda^1}, R_S^{a^2} + \gamma M(a^2) V^{\lambda^2} \right\}$$

$$= \operatorname{argmax} \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -10 \\ 5 \end{bmatrix}, \right.$$

$$\left. \begin{bmatrix} -0.1 \\ 0.8 \end{bmatrix} + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} -10 \\ 5 \end{bmatrix} \right\}$$

$$= \left\{ \begin{array}{l} a^1 \\ a^2 \end{array} \right\}$$

$\pi^2 \neq \lambda^1$, so PE continue.

$$M(\lambda^2) = \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \quad R_S^{\lambda^2} = \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix}$$

$$V^{\lambda^2} = (I - \gamma M(\lambda^2))^{-1} R_S^{\lambda^2}$$

$$= \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \right)^{-1} \cdot \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$$

$$\pi^3 = \operatorname{argmax} \{ R_S a + \gamma M(a) V^{2^2} \}$$

$$= \{ R_S a^1 + \gamma M(a^1) V^{2^2}, R_S a^2 + \gamma M(a^2) V^{2^2} \}$$

$$= \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 4.5 \\ 5 \end{bmatrix}, \right.$$

$$\left. \begin{bmatrix} -1 \\ 0.8 \end{bmatrix} + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 4.5 \\ 5 \end{bmatrix} \right\}$$

$$= \left\{ \begin{array}{l} a^1 \\ a^2 \end{array} \right\}$$

$\because \pi^3 \neq \pi^2$, so PE continue

$$M(\pi^3) = \begin{bmatrix} 0.2 & 0.8 \\ 0.1 & 0.9 \end{bmatrix} \quad R_C \pi^3 = \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix}$$

$$V^{2^3} = [I - \gamma M(\pi^3)]^{-1} R_S \pi^3$$

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0.1 & 0.9 \end{bmatrix}^{-1} \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix}$$

$$= \begin{bmatrix} 6.54 \\ 7.31 \end{bmatrix}$$

$$\pi^4 = \arg\max \{ R_s^a + \gamma V(a) \cup \pi^3 \}$$

$$= \arg\max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 6.54 \\ 7.31 \end{bmatrix}, \right.$$

$$\left. \begin{bmatrix} -1 \\ 0.8 \end{bmatrix} + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 6.54 \\ 7.31 \end{bmatrix} \right\}$$

$$= \begin{Bmatrix} a^1 \\ a^2 \end{Bmatrix}$$

$$\text{if } \pi^4 = \pi^3 \quad \therefore \pi^* = \pi^4 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}$$

Problem 2.

For the system defined in Problem 2, perform matrix-form Value Iteration method with $V_0(s)=0$, $\gamma=0.9$ and $\theta=0.5$ to compute V^* and π^* .

$$V_1 = \max R_S^a + S M(a) V_0$$

$$= \max \{ R_S^{a^1} + S M(a^1) V_0, R_S^{a^2} + S M(a^2) V_0 \}$$

$$= \max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0.8 \end{bmatrix} + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}$$

$$= \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix}$$

$$\left| \left| V_1 - V_0 \right| \right|_{\max} = \left| \left| \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right| \right| \quad (0.5 \text{ is not the max})$$

∴ Continue

$$V_2 = \max R_S^a + S M(a) V_1$$

$$= \max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.8 \end{bmatrix} + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix} \right\}$$

$$= \begin{bmatrix} 0.69 \\ 1.46 \end{bmatrix}$$

$$\left| \left| V_2 - V_1 \right| \right|_{\max} = \left| \left| \begin{bmatrix} 0.69 \\ 1.46 \end{bmatrix} - \begin{bmatrix} 0.1 \\ 0.8 \end{bmatrix} \right| \right| \quad (0.5 \text{ is not the max})$$

∴ Continue

∴ Continue

$$V_3 = \max R_S^A + \text{SM}(a)V_2$$

$$= \max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0.69 \\ 1.46 \end{bmatrix}, \begin{bmatrix} -1 \\ 0.2 \end{bmatrix} \right. \\ \left. + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 0.69 \\ 1.46 \end{bmatrix} \right\}$$

$$= \boxed{\begin{bmatrix} 1.28 \\ 2.05 \end{bmatrix}}$$

$$\frac{|V_4 - V_3|}{\max} = \frac{\left| \left(\begin{bmatrix} 1.28 \\ 2.05 \end{bmatrix} - \begin{bmatrix} 0.69 \\ 1.46 \end{bmatrix} \right) \right|}{\max} < 0.015$$

not true

∴ continue

$$V_4 = \max R_S^A + \text{SM}(a)V_3$$

$$= \max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1.28 \\ 2.05 \end{bmatrix}, \begin{bmatrix} -1 \\ 0.2 \end{bmatrix} \right. \\ \left. + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 1.28 \\ 2.05 \end{bmatrix} \right\}$$

$$= \boxed{\begin{bmatrix} 1.28 \\ 2.05 \end{bmatrix}}$$

$$\|V_4 - V_3\| = \left\| \begin{bmatrix} 1.81 \\ 2.02 \end{bmatrix} - \begin{bmatrix} 1.28 \\ 2.00 \end{bmatrix} \right\|_{\max} \text{ is } \max$$

Not the \therefore continue

$$V_5 = \max R_S^a + S M(a) V_4$$

$$= \max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 1.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1.81 \\ 2.02 \end{bmatrix}, \begin{bmatrix} -1 \\ 0.8 \end{bmatrix} \right.$$

$$\left. + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 1.81 \\ 2.02 \end{bmatrix} \right\}$$

$$\|V_5 - V_4\|_{\max} = 0.47 < 0.5 \quad \text{if } V^* = V_5 = \begin{bmatrix} 2.28 \\ 3.05 \end{bmatrix}$$

$$V^* = \operatorname{argmax}_{a \in A} \left\{ R_S^{a'} + S M(a') V^*, R_S^{a''} + S M(a'') V^* \right\}$$

$$= \arg \max \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2.28 \\ 3.05 \end{bmatrix}, \right. \\ \left. \begin{bmatrix} -1 \\ 0.8 \end{bmatrix} + 0.9 \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 2.28 \\ 3.05 \end{bmatrix} \right\}$$

$$\Rightarrow \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

$$\therefore \pi^* = \begin{bmatrix} a_1' \\ a_2' \end{bmatrix}$$

Problem 3.

Consider an MDP with two states $\{A, B\}$ and two actions $\{a^1, a^2\}$. The system state transitions are governed through the following transition matrices:

$$M(a^1) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, M(a^2) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

The reward is as follows

$\begin{cases} +5 \\ 0 \\ -1 \\ 0 \end{cases}$	$\begin{array}{l} \text{moving to state } B \\ \text{moving to state } A \\ \text{taking action } a^2 \\ \text{taking action } a^1 \end{array}$
--	---

Consider an initial Policy $\pi^0 = \begin{bmatrix} \pi^0(A) \\ \pi^0(B) \end{bmatrix} = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}$, $\gamma=0.9$ and episode length 5. Perform

Monte Carlo Policy Iteration method to obtain the best policy.

* You need to show all trajectories, the approximation of Q-values and Policy Improvement till the time that Policies in two consecutive iterations stays the same.

episode 1

$A, \pi(A)=a^1, A, r=0$

episode 2

$B, \pi(B)=a^2, A, r=-1$

$A, \pi(A)=a^1, A, r=0$

$A, \pi(A)=a^1, A, r=0$

$A, \pi(A)=a^1, A, r=0$

$A, \pi(A)=a^1, A, r=0$

episode 3

$$A, \lambda(A) = a^2, B, r = 4$$

$$B, \lambda(B) = a^2, A, r = -1$$

$$A, \lambda(A) = a^1, B, r = 0$$

$$A, \lambda(A) = a^1, B, r = 0$$

$$A, \lambda(A) = a^1, B, r = 0$$

$$G^1 A, a^1 = 0 \quad G^2 A, a^1 = 0 \quad G^3 A, a^1 = 0 \quad G^4 A, a^1 = 0$$

$$Q^{x^0}(A, a^1) = 0$$

$$G^2 B, a^2 = -1 \quad G^3 B, a^2 = -1 \quad G^4 B, a^2 = -1$$

$$Q^{x^0}(B, a^2) = -1$$

$$G^3 A, a^2 = 4 + 5(-1) = 3 \cdot 1$$

$$Q^{x^0}(A, a^2) = 3 \cdot 1$$

$$G^4 B, a^1 = 5 + 6(-1) = 4 \cdot 1$$

episode 4

$$B, \lambda(B) = a^1; B, r = 5$$

$$B, \lambda(B) = a^2; A, r = -1$$

$$A, \lambda(A) = a^1; A, r = 0$$

$$A, \lambda(A) = a^1; A, r = 0$$

$$A, \lambda(A) = a^1; A, r = 0$$

$$G^1 A, a^1 = 0 \quad G^2 A, a^1 = 0 \quad G^3 A, a^1 = 0 \quad G^4 A, a^1 = 0$$

$$Q^{x^0}(A, a^1) = 0$$

$$G^2 B, a^2 = -1 \quad G^3 B, a^2 = -1 \quad G^4 B, a^2 = -1$$

$$Q^{x^0}(B, a^2) = -1$$

$$Q^{x^0}(B, a^1) = 4 \cdot 1$$

$$\pi'(s) = \arg\max_{a \in A} Q_{\pi^0}(s, a)$$

$$\pi' = \begin{bmatrix} a^2 \\ a^1 \end{bmatrix}$$

episode 1

episode 2

$$A, \pi(A) = a^1, A, t = 0 \quad B, \pi(B) = a^2, A, t = -1$$

$$A, \pi(A) = a^2, B, t = 4 \quad A, \pi(A) = a^2, B, t = 4$$

$$B, \pi(B) = a^1, B, t = 5 \quad B, \pi(B) = a^1, B, t = 5$$

$$B, \pi(B) = a^1, B, t = 5 \quad B, \pi(B) = a^1, B, t = 5$$

$$B, \pi(B) = a^1, B, t = 5 \quad B, \pi(B) = a^1, B, t = 5$$

episode 3

episode 4

$$\begin{array}{ll}
 A, \pi(A) = a^2, B, r = 4 & B, \pi(B) = a^1, B, r = 5 \\
 B, \pi(B) = a^1, B, r = 5 & B, \pi(B) = a^1, B, r = 5 \\
 B, \pi(B) = a^1, B, r = 5 & B, \pi(B) = a^1, B, r = 5 \\
 B, \pi(B) = a^1, B, r = 5 & B, \pi(B) = a^1, B, r = 5 \\
 B, \pi(B) = a^1, B, r = 5 & B, \pi(B) = a^1, B, r = 5
 \end{array}$$

$$G^1 A, a^1 = 0 + 5 \cdot 4 + 5^2 \cdot 5 + 5^3 \cdot 5 + 5^4 \cdot 5 = 14.58$$

$$Q(A, a^1) = 14.58$$

$$G^2 B, a^2 = -1 + 5 \cdot 4 + 5^2 \cdot 5 + 5^3 \cdot 5 + 5^4 \cdot 5 = 13.58$$

$$Q(B, a^2) = 13.58$$

$$G^1 A, a^2 = 4 + 5 \cdot 5 + 5^2 \cdot 5 + 5^3 \cdot 5 = 16.105$$

$$G^2 A, a^2 = 4 + 5 \cdot 5 + 5^2 \cdot 5 + 5^3 \cdot 5 = 16.195$$

$$G^3 A, a^2 = 4 + 5 \cdot 5 + 5^2 \cdot 5 + 5^3 \cdot 5 + 5^4 \cdot 5 = 19 \cdot 48$$

$$Q(A, a^2) = \frac{G^0 A, a^2 + G^1 A, a^2 + G^2 A, a^2}{3} = 17.29$$

$$G^1 B, a^1 = 5 + 5 \cdot 5 + 5^2 \cdot 5 = 13 \cdot 55$$

$$G^2 B, a^1 = 5 + 5 \cdot 5 + 5^2 \cdot 5 = 13 \cdot 55$$

$$G^3 B, a^1 = 5 + 5 \cdot 5 + 5^2 \cdot 5 + 5^3 \cdot 5 = 17 \cdot 195$$

$$G^4 B, a^1 = 5 + 5 \cdot 5 + 5^2 \cdot 5 + 5^3 \cdot 5 + 5^4 \cdot 5 = 20 \cdot 4755$$

$$Q(B, a^1) = \frac{G^1 B, a^1 + G^2 B, a^1 + G^3 B, a^1 + G^4 B, a^1}{4} = 16 \cdot 19$$

$$\pi_t^* = \arg \max Q(S, a)$$

$$\therefore \pi^2 = \begin{bmatrix} a^2 \\ a' \end{bmatrix}$$

$$\therefore \pi_0^2 = \pi^1$$

\therefore the optimal policy is

equal to $\pi^2 = \begin{bmatrix} a^2 \\ a' \end{bmatrix}$