

# \* Assignment-1 \*

Nelakuditi Rahul Naga  
AI20BTECH11029

1) a) The states in the given Markov Reward process are  $s, 1, 3, 5, 6, 7, 8, w$ . This is because the states 2, 4 and 9 are equivalent to states 7, 8 and 3 respectively. The transition matrix  $P$  is given by :-

$$P = \begin{matrix} & \begin{matrix} s & 1 & 3 & 5 & 6 & 7 & 8 & w \end{matrix} \\ \begin{matrix} s \\ 1 \\ 3 \\ 5 \\ 6 \\ 7 \\ 8 \\ w \end{matrix} & \begin{bmatrix} 0 & 0.25 & 0.25 & 0 & 0 & 0.25 & 0.25 & 0 \\ 0 & 0 & 0.25 & 0.25 & 0 & 0.25 & 0.25 & 0 \\ 0 & 0 & 0 & 0.25 & 0.25 & 0.25 & 0.25 & 0 \\ 0 & 0 & 0.25 & 0 & 0.25 & 0.25 & 0.25 & 0 \\ 0 & 0 & 0.25 & 0 & 0 & 0.25 & 0.25 & 0.25 \\ 0 & 0 & 0.25 & 0 & 0 & 0.25 & 0.25 & 0.25 \\ 0 & 0 & 0.25 & 0 & 0 & 0 & 0.5 & 0.25 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \end{matrix}$$

b) The goal of this question is to calculate the expected number of die throws required to reach the goal state  $w$  from any other state. For this problem, the suitable reward function and discount factor  $\gamma$  are as follows:-

$$R(s) = \begin{cases} -1, & s \neq w \\ 0, & s = w \end{cases}$$

$$\gamma = 1$$

→ The Bellman evaluation equation for a Markov Reward process is :-

$$v = (I - \gamma P)^{-1} R$$

⇒ The vector " $v$ " would give the required average number of moves to reach goal state  $w$  from a state  $s$ .

→ As the number of moves are required from states other than  $w$ , we will

consider only the submatrix that constitutes the first 7 rows & columns of matrix  $P$  as the transition matrix in the Bellman equation. Subsequently,  $R$  becomes a column vector with length 7 consisting of all -1's. Upon substituting the matrices in the Bellman equation we obtain:

$$v(s) : \begin{bmatrix} -7.083 \\ -7 \\ -6.67 \\ -6.67 \\ -5.33 \\ -5.33 \\ -5.33 \end{bmatrix}$$

→ But as we require the expected no. of moves, we take modulus of above values and the expected moves are thus given by:

$$v(s) : \begin{bmatrix} 7.083 \\ 7 \\ 6.67 \\ 6.67 \\ 5.33 \\ 5.33 \\ 5.33 \end{bmatrix}$$

2) a) We are required to formulate the given problem as an MDP. We know that an MDP is given by the tuple  $\langle S, A, P, R, r \rangle$ . For the given problem, they are given by:

i.) State space 'S' is the enumeration of all possible number of working machines on a particular day. Therefore it is given by  $\{0, 1, 2, \dots, N-1, N\}$ .

ii.) Action space 'A' is given by  $\{\text{call the Repair Man, Don't call the Repair Man}\}$ .

iii.) Our goal in this problem is to maximize the profits earned. Therefore the

Reward function can be formulated as follows :-

$$R(s_t, a_t) = \begin{cases} s_t - \frac{N}{2}, & \text{if } a_t : \text{Call the Repair Man} \\ s_t, & \text{if } a_t : \text{Don't call the Repair Man} \end{cases}$$

iv.) As there are 2 possible actions - call the Repair Man (CR) and Don't call the Repair Man (DCR) we will have 2 state transition matrices, one for each action. They are given by :-

$$P^{CR} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & \dots & N-1 & N \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ \vdots \\ N-1 \\ N \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & \dots & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 & 1 \end{bmatrix} \end{matrix}$$

$$P^{DCR} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & \dots & N-1 & N \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ \vdots \\ N-1 \\ N \end{matrix} & \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & \dots & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{1}{N} & \frac{1}{N} & \frac{1}{N} & \dots & \frac{1}{N} & 0 \\ \frac{1}{N+1} & \frac{1}{N+1} & \frac{1}{N+1} & \dots & \frac{1}{N+1} & \frac{1}{N+1} \end{bmatrix} \end{matrix}$$

b) The problem deals with maximizing long-time profit, the horizon is infinite. Hence, for the expected rewards/profit to converge it is better to use discounted setting for the above MDP formulation.

c) The policy given is to never call the repair man. We are also given that  $N=5$ . The Bellman Evaluation equation in matrix form is given by :-

$$V^\pi = (I - \gamma P^\pi)^{-1} R^\pi$$

→ We will take  $\gamma = 0.5$ .

→ Also, we know that :-

$$P^\pi(s'|s) = \sum_{a \in A} \pi(a|s) P_{ss'}^a \quad \text{and} \quad R^\pi(s) = \sum_{a \in A} \pi(a|s) \sum_{s'} P_{ss'}^a R_{ss'}^a$$



Upon calculating we have: (taking  $N=5$ )

$$P^\pi : \begin{matrix} & 0 & 1 & 2 & 3 & 4 & 5 \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 \\ \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & 0 \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \end{bmatrix} \end{matrix}$$

$$R^\pi : [0, 1, 2, 3, 4, 5]^T$$

⇒ Upon substituting in the Bellman Evaluation equation we obtain:

$$V^\pi : [0, \frac{4}{3}, \frac{8}{3}, 4, \frac{16}{3}, \frac{20}{3}]^T$$

d) The policy iteration algorithm involves initialization of policy and then repeated evaluation and improvement of policy steps.

→ The initial policy in this problem is the no-repair policy.

→ We have already evaluated this policy in part (c) of this question.

→ Now to get an improved policy we will find the greedy policy  $\pi' = \text{greedy}(V^\pi)$  where  $V^\pi$  is the value function calculated above for policy  $\pi$  (no repair policy). We know that:

$$\pi'(s) = \text{Greedy}(V^\pi) : \begin{cases} 1, & \text{if } a = \arg \max_{a \in A} \left[ \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s')) \right] \\ 0, & \text{otherwise} \end{cases}$$

→ Let us denote the actions Don't call Repair Man and call the Repair Man as DCR and CR respectively. Let us also denote  $\sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s'))$  by  $Q_a^s$ .

\* Upon calculation, the results are as follows:

i) for  $s=0$ ;  $V^{CR} = \frac{5}{6}$ ,  $V^{DCR} = 0 \Rightarrow$  optimal action is CR

ii) for  $s=1$ ;  $V^{CR} = \frac{11}{6}$ ,  $V^{DCR} = \frac{4}{3} \Rightarrow$  optimal action is CR

iii) for  $s=2$ ;  $V^{CR} = \frac{17}{6}$ ,  $V^{DCR} = \frac{8}{3} \Rightarrow$  optimal action is CR

iv) for  $s=3$  ;  $V^{CR} = \frac{23}{6}$  ,  $V^{DCR} = 4 \Rightarrow$  Optimal action is DCR

v.) for  $s=4$  ;  $V^{CR} = \frac{29}{6}$  ,  $V^{DCR} = \frac{16}{3} \Rightarrow$  Optimal action is DCR

vi.) for  $s=5$  ;  $V^{CR} = \frac{35}{6}$  ,  $V^{DCR} = \frac{20}{3} \Rightarrow$  Optimal action is DCR

$\therefore$  The improved policy  $\pi'$  is as follows:

$$\pi'(s) : \begin{cases} \text{Call the Repair Man, if } s = 0, 1, 2 \\ \text{Don't call the Repair Man, if } s = 3, 4, 5 \end{cases}$$

3) a) The Bellman evaluation equation is given by :

$$V^{\pi}(s) : \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^{\pi}(s')]$$

We will choose  $\gamma = 1$  for the given problem. The state transition matrix corresponding to policy  $\pi_1$  is given by :

$$\begin{array}{c} \begin{array}{ccccc} & A & B & C & D \\ \begin{array}{c} A \\ B \\ C \\ D \end{array} & \begin{bmatrix} 0 & 0.9 & 0.1 & 0 \\ 0.1 & 0 & 0 & 0.9 \\ 0.9 & 0 & 0 & 0.1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \end{array}$$

Clearly, as D is the terminal state with reward 100, we have  $V^{\pi_1}(D) = 100$ . Now using Bellman evaluation equation we have :

$$V^{\pi_1}(A) : 0.9[V^{\pi_1}(B) - 10] + 0.1[V^{\pi_1}(C) - 10]$$

$$V^{\pi_1}(B) : 0.9[V^{\pi_1}(D) - 10] + 0.1[V^{\pi_1}(A) - 10]$$

$$V^{\pi_1}(C) : 0.9[V^{\pi_1}(A) - 10] + 0.1[V^{\pi_1}(D) - 10]$$

Upon solving we obtain :  $V^{\pi_1} : [75.609, 87.56, 68.048, 100]^T$

→ The state transition matrix for policy  $\pi_2$  is :

$$\begin{array}{c} \begin{array}{ccccc} & A & B & C & D \\ \begin{array}{c} A \\ B \\ C \\ D \end{array} & \begin{bmatrix} 0 & 0.1 & 0.9 & 0 \\ 0.9 & 0 & 0 & 0.1 \\ 0.1 & 0 & 0 & 0.9 \\ 0 & 0 & 0 & 1 \end{bmatrix} \end{array}$$

We now have :

$$V^{\pi_2}(D) : 100$$

$$V^{\pi_2}(A) : 0.9[V^{\pi_2}(C) - 10] + 0.1[V^{\pi_2}(B) - 10]$$

$$V^{\pi_2}(B) : 0.9[V^{\pi_2}(A) - 10] + 0.1[V^{\pi_2}(D) - 10]$$

$$V^{\pi_2}(C) : 0.9[V^{\pi_2}(D) - 10] + 0.1[V^{\pi_2}(A) - 10]$$

Upon solving we obtain :  $V^{\pi_2} : [75.609, 68.048, 87.56, 100]^T$



→ The state transition matrix for policy  $\pi_3$  is :-

$$\begin{array}{c|cccc} & A & B & C & D \\ \hline A & 0 & 0.42 & 0.58 & 0 \\ B & 0.1 & 0 & 0 & 0.9 \\ C & 0.1 & 0 & 0 & 0.9 \\ D & 0 & 0 & 0 & 1 \end{array}$$

We trivially have  $V^{\pi_3}(D) = 100$ . We also have :-

$$V^{\pi_3}(A) = 0.42[V^{\pi_3}(B) - 10] + 0.58[V^{\pi_3}(C) - 10]$$

$$V^{\pi_3}(B) = 0.9[V^{\pi_3}(D) - 10] + 0.1[V^{\pi_3}(A) - 10]$$

$$V^{\pi_3}(C) = 0.1[V^{\pi_3}(A) - 10] + 0.9[V^{\pi_3}(D) - 10]$$

upon solving we obtain :-  $V^{\pi_3} = [77.77, 87.77, 87.77, 100]^T$

b) It is easy to observe that  $V^{\pi_3}(s) \geq V^{\pi_1}(s)$  and  $V^{\pi_3}(s) \geq V^{\pi_2}(s)$  for all  $s \in \{A, B, C, D\}$ . Hence, we conclude that  $\pi_3$  is the best policy among the given policies.

c) It is easy to see that the policies  $\pi_1$  and  $\pi_2$  are not comparable. This is because  $V^{\pi_1}(B) > V^{\pi_2}(B)$  but  $V^{\pi_1}(C) < V^{\pi_2}(C)$ . Hence, we conclude that all policies are not comparable.

d) We are given two deterministic policies  $\pi_1$  &  $\pi_2$  of an MDP  $M$ . A new policy  $\pi$  that is better than policies  $\pi_1$  and  $\pi_2$  can be constructed as follows:-

→ Suppose the states in MDP are  $s_1, s_2, \dots, s_n$  where  $s_n$  is the terminal state without loss of generality. Now policy  $\pi$  is as follows:-

$$\pi(s_i) = \begin{cases} \pi_1(s_i), & \text{if } V^{\pi_1}(s_i) > V^{\pi_2}(s_i) \\ \pi_2(s_i), & \text{if } V^{\pi_2}(s_i) > V^{\pi_1}(s_i) \\ \pi_1(s_i) \text{ with probability } p, \pi_2(s_i) \text{ with probability } 1-p, & \text{o/w} \end{cases}$$

for  $i = 1, 2, \dots, n$  and  $p \in [0, 1]$ . If  $p = 0$  (or)  $1$ ,  $\pi$  becomes a deterministic

Policy else if  $p \in (0, 1)$ , the policy  $\pi$  is stochastic. We are constructing an action suggested by the new policy by greedily choosing the deterministic policy which <sup>is</sup> ~~suggests~~ an action that maximizes the value function.

→ When both policies suggest an action that is equally optimal, we choose from either of them with appropriate probabilities as mentioned above.

→ Hence, the new policy  $\pi$  is better than both  $\pi_1$  and  $\pi_2$ .

4) a) We will 4 possible pairs of  $\gamma$  and  $\eta$  :-  $(0.9, 0)$ ,  $(0.9, 0.5)$ ,  $(0.1, 0)$  and  $(0.1, 0.5)$ . We will analyze what is the optimal chosen in each of the four scenarios.

i)  $\gamma = 0.9$  and  $\eta = 0$  i.e; high discount factor and zero noise environment.

→ As  $\gamma$  is high, the agent focuses on rewards that can be obtained in faraway future. Also, as there is no noise in the environment, the agent can go along the dashed path i.e; by risking the cliff. This is because in zero noise environment, the agent will successfully avoid the negative payoff terminal states while going along the dashed path.

∴ for high  $\gamma$  and low  $\eta$ , agent prefers distant exit but risks the cliff.

ii)  $\gamma = 0.9$  and  $\eta = 0.5$  i.e; high discount factor and noisy environment.

→ As  $\gamma$  is high, the agent focuses on rewards that can be obtained in faraway future. Also, as the environment is noisy, the agent chooses to go along the solid path i.e; by avoiding the cliff. This is because in a noisy environment, there is high probability that the agent will fall off into the negative payoff terminal states if it chooses to go along the dashed path.

∴ for high  $\gamma$  and high  $\eta$ , agent prefers distant exit and avoids the cliff.



iii.)  $\gamma = 0.1$  and  $\eta = 0$  i.e; low discount factor and zero noise environment.

→ As  $\gamma$  is low, the agent focuses on rewards that can be obtained in immediate future. Also, as there is no noise in the environment, the agent chooses to go along the dashed path i.e; by risking the cliff. This is because in zero noise environment, the agent will successfully avoid the negative pay-off terminal states while going along the dashed path.

∴ For low  $\gamma$  and low  $\eta$ , agent prefers close exit but risks the cliff.

iv.)  $\gamma = 0.1$  and  $\eta = 0.5$  i.e; low discount factor and noisy environment.

→ As  $\gamma$  is low, the agent focuses on rewards that can be obtained in immediate future. Also, as the environment is noisy, the agent chooses to go along the solid path i.e; by avoiding the cliff. This is because in a noisy environment, there is high probability that the agent will fall off into the negative payoff terminal states if it chooses to go along the dashed path.

∴ For low  $\gamma$  and high  $\eta$ , agent prefers close exit and avoids the cliff.

5) a) We are given action value functions  $Q_1^\pi(s)$  and  $Q_2^\pi(s)$  corresponding to MDP's  $M_1$  and  $M_2$ . Now,  $M_3$  is a composite MDP of  $M_1$  &  $M_2$  and  $Q_3^\pi(s)$  is its corresponding action-value function. The expressions for the action value functions are as follows:

$$Q_1^\pi(s, a) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1} \mid s_t = s, a_t = a \right); \delta \sim R_1(s, a)$$

$$Q_2^\pi(s, a) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1} \mid s_t = s, a_t = a \right); \delta \sim R_2(s, a)$$

$$Q_3^\pi(s, a) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1} \mid s_t = s, a_t = a \right); \delta \sim (R_1 + R_2)(s, a)$$

→ From above equations, we can conclude that it is possible to combine the action-value functions  $Q_1^\pi(s)$  and  $Q_2^\pi(s)$  linearly (since the expectation operator is linear) to calculate  $Q_3^\pi(s)$ .

b) We are given optimal policies  $\pi_1^*$  and  $\pi_2^*$  corresponding to MDP's  $M_1$  and  $M_2$  respectively. The optimal policies for  $M_1$  and  $M_2$  can be obtained by solving :-

$$\pi_1^* = \arg \max_{\pi} \left[ E_{\pi} \left( \sum_{t=0}^{\infty} \gamma^t \cdot \delta_{t+1} \right) \right] ; \delta \sim R_1$$

$$\pi_2^* = \arg \max_{\pi} \left[ E_{\pi} \left( \sum_{t=0}^{\infty} \gamma^t \cdot \delta_{t+1} \right) \right] ; \delta \sim R_2$$

→ Now the optimal policy  $\pi_3^*$  for MDP  $M_3$  is given by :-

$$\pi_3^* = \arg \max_{\pi} \left[ E_{\pi} \left( \sum_{t=0}^{\infty} \gamma^t \cdot \delta_{t+1} \right) \right] ; \delta \sim R_1 + R_2$$

→ Due to non-linearity of the "max" operator, we can conclude that it is not possible to simply combine optimal policies  $\pi_1^*$  and  $\pi_2^*$  in order to obtain  $\pi_3^*$ .

c) Consider an arbitrary policy  $\pi$ . We have :-

$$v_1^{\pi}(s) = E_{\pi} \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1} \mid s_t = s \right) ; \delta \sim R_1(s)$$

$$v_2^{\pi}(s) = E_{\pi} \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1} \mid s_t = s \right) ; \delta \sim R_2(s)$$

$$v_3^{\pi}(s) = E_{\pi} \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1} \mid s_t = s \right) ; \delta \sim (R_1 + R_2)(s)$$

It is easy to see that  $v_3^{\pi}(s) = v_1^{\pi}(s) + v_2^{\pi}(s)$ . Now it is given that  $\pi^*$  is an optimal policy for MDP's  $M_1$  &  $M_2$ . We will prove that this is also an optimal policy of  $M_3$  using proof by contradiction as follows:-

→ Assume that there exists a <sup>such</sup> ~~better~~ policy  $\pi'$  that  $\pi^*$  performs ~~worser~~ than  $\pi'$  on a state  $s \in S$  i.e;  $v_3^{\pi'}(s) > v_3^{\pi^*}(s)$ . Hence we have :-

$$v_1^{\pi'}(s) + v_2^{\pi'}(s) > v_1^{\pi^*}(s) + v_2^{\pi^*}(s)$$

which implies that  $v_1^{\pi'}(s) > v_1^{\pi^*}(s)$  (or)  $v_2^{\pi'}(s) > v_2^{\pi^*}(s)$  which is a contradiction since  $\pi^*$  is an optimal policy for both  $M_1$  and  $M_2$ . Hence,

$\pi^*$  is an optimal policy for  $M_3$ .

d) We know that value functions under policy  $\pi$  for MDP's  $M_1$  and  $M_2$  are given as follows:

$$v_1^\pi(s) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1}^1 \mid s_t = s \right) ; \delta^1 \sim R_1(s)$$

$$v_2^\pi(s) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot \delta_{t+k+1}^2 \mid s_t = s \right) ; \delta^2 \sim R_2(s)$$

$\Rightarrow$  As it is given that  $R_1(s, a, s') - R_2(s, a, s') = \epsilon$ , we have:-

$$v_1^\pi(s) - v_2^\pi(s) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot (\delta_{t+k+1}^1 - \delta_{t+k+1}^2) \mid s_t = s \right)$$

$$v_1^\pi(s) - v_2^\pi(s) = E_\pi \left( \sum_{k=0}^{\infty} \gamma^k \cdot \epsilon \mid s_t = s \right)$$

$$v_1^\pi(s) - v_2^\pi(s) = \epsilon \cdot \sum_{k=0}^{\infty} \gamma^k = \frac{\epsilon}{1-\gamma}$$

$$\therefore v_1^\pi(s) = v_2^\pi(s) + \frac{\epsilon}{1-\gamma}$$