

AAI Assignment 4

Name: 吉辰卿

SID: 12332152

1. solution

proof: To prove this problem, we should reveal that the optimal policy remains unchanged for this modified MDP.

Firstly, we denote the state-value function for a policy π in the original MDP is $V_{\pi}(s)$ and in the modified MDP is $V'_{\pi}(s)$.

$$\text{so: } V_{\pi}(s) = E_{\pi} [G_{\pi} | \pi, St=s] = E_{\pi} [R_{t+1} + \gamma R_{t+2} + \dots | \pi, St=s] \quad \text{--- (1)}$$

Then, because for the modified MDP, the new reward function will become.

therefore, the state-value function in the modified MDP under the policy π will become $V'_{\pi}(s)$: $R'(s) = \alpha R(s) + \beta$
($\alpha > 0$)

$$\begin{aligned} V'_{\pi}(s) &= E_{\pi} [G'_{\pi} | \pi, St=s] = E_{\pi} [\cancel{\alpha R_{t+1} + \beta} + \gamma R'_{t+1} + \gamma R'_{t+2} + \dots | \pi, St=s] \\ &= E_{\pi} [\alpha R_{t+1} + \beta + \gamma \alpha R_{t+2} + \beta + \dots | \pi, St=s] \quad \text{--- (2)} \end{aligned}$$

let's observe the structure of (1) and (2) and we can find β is only considered as a constant term and the α is only considered as a coefficient of the expected reward. so, we can get:

$$V'_{\pi}(s) = E_{\pi} [G_{\pi} | \pi, St=s] \quad \text{and} \quad V'_{\pi}(s) = E_{\pi} [G'_{\pi} | \pi, St=s]$$

So: G'_{π} in the modified MDP is only a linear combination of the G_{π} in the original MDP. That is: $G'_{\pi} = \alpha G_{\pi} + \beta$

Since the relationship between G'_{π} and G_{π} is linear, which and $\alpha > 0$, which is linear positive correlation. Therefore, the modified reward function $R'(s)$ will not affect the formulation of the optimal policy.

So, the modified MDP has the same optimal policy as the original MDP.

2. solution

ii) when MDP stops after two steps:

$$G_t = R_{t+1} + \gamma R_{t+2}$$

so: the state-value function is: $V_{\pi}(s) = E[G_t | S_t=s, \pi]$

$$= E[R_{t+1} + \gamma V_{\pi}(s') | S_t=s, \pi]$$

In this problem, we denote H as the state high, L as the state low, W as the action wait, S as the action search, R as the action recharge.

And for stochastic $\pi(a|s)$, ① will become:

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s, a, s') (R(s, a, s') + \gamma V_{\pi}(s'))$$

And: $\pi(S|H) = 1$, $\pi(S|L) = 0.4$ and $\pi(R|L) = 0.6$, and others are 0.

$$\text{So: } v_{\pi}(H) = \pi(S|H) \times [P(H, S, H) (R(H, S, H) + \gamma v'_{\pi}(H)) + P(H, S, L) (R(H, S, L) + \gamma v'_{\pi}(L))] \\ + \pi(W|H) \times [P(H, W, H) (R(H, W, H) + \gamma v'_{\pi}(H)) + P(H, W, L) (R(H, W, L) + \gamma v'_{\pi}(L))] \\ + \pi(R|H) \times [P(H, R, H) (R(H, R, H) + \gamma v'_{\pi}(H)) + P(H, R, L) (R(H, R, L) + \gamma v'_{\pi}(L))]$$

$$v_{\pi}(L) = \pi(S|L) \times [P(L, S, H) (R(L, S, H) + \gamma v'_{\pi}(H)) + P(L, S, L) (R(L, S, L) + \gamma v'_{\pi}(L))] \\ + \pi(R|L) \times [P(L, R, H) (R(L, R, H) + \gamma v'_{\pi}(H)) + P(L, R, L) (R(L, R, L) + \gamma v'_{\pi}(L))] \\ + \pi(W|L) \times [P(L, W, H) (R(L, W, H) + \gamma v'_{\pi}(H)) + P(L, W, L) (R(L, W, L) + \gamma v'_{\pi}(L))]$$

$$v'_{\pi}(H) = \pi(S|H) \times [P(H, S, H) \times R(H, S, H) + P(H, S, L) \times R(H, S, L)]$$

$$+ \pi(W|H) \times [P(H, W, H) \times R(H, W, H) + P(H, W, L) \times R(H, W, L)]$$

$$+ \pi(R|H) \times [P(H, R, H) \times R(H, R, H) + P(H, R, L) \times R(H, R, L)]$$

$$v'_{\pi}(L) = \pi(S|L) \times [P(L, S, H) \times R(L, S, H) + P(L, S, L) \times R(L, S, L)]$$

$$+ \pi(R|L) \times [P(L, R, H) \times R(L, R, H) + P(L, R, L) \times R(L, R, L)]$$

$$+ \pi(W|L) \times [P(L, W, H) \times R(L, W, H) + P(L, W, L) \times R(L, W, L)]$$

Therefore:

$$v'_{\pi}(H) = 1 \times [\alpha \times \gamma_{\text{search}} + (1-\alpha) \times \gamma_{\text{search}}]$$

$$= 0.5 \times 3 + 0.5 \times 3 = 3$$

$$v'_{\pi}(L) = 0.4 \times [(1-\beta) \times (-3) + \beta \times \gamma_{\text{search}}] + 0.6 \times [1 \times 0]$$

$$= 0.4 \times [0.7 \times (-3) + 0.3 \times 3] = -0.48$$

$$v_{\pi}(H) = 1 \times [\alpha (\gamma_{\text{search}} + \gamma v'_{\pi}(H)) + (1-\alpha) (\gamma_{\text{search}} + \gamma v'_{\pi}(L))]$$

$$= 0.5 \times (3 + 0.8 \times 3) + 0.5 \times (3 + 0.8 \times (-0.48)) = 4.008$$

$$v_{\pi}(L) = 0.4 \times [(1-\beta) (-3 + \gamma v'_{\pi}(H)) + \beta (\gamma_{\text{search}} + \gamma v'_{\pi}(L))] + 0.6 \times [1 \times (0 + \gamma v'_{\pi}(H))]$$

$$= 0.4 \times [0.7 \times (-3 + 0.8 \times 3) + 0.3 \times (3 - 0.8 \times 0.48)] + 0.6 \times 0.8 \times 3$$

$$= 1.58592$$

so, the state-value function is: $V_{\pi}(H) = 1.008$, $V_{\pi}(L) = 1.58572$.

(2) The expression of the action-value function in the case the MDP stops after a single step is:

$$Q_{\pi}(s, a) = E[G_{t+1} | \pi, s, a] = E[R_{t+1} | \pi, s, a]$$

So, we use the same notation as (1).

$$Q_{\pi}(H, W) = P(H, W, H) \times R(H, W, H) + P(H, W, L) \times R(H, W, L)$$

$$Q_{\pi}(H, S) = P(H, S, H) \times R(H, S, H) + P(H, S, L) \times R(H, S, L)$$

$$Q_{\pi}(H, R) = P(H, R, H) \times R(H, R, H) + P(H, R, L) \times R(H, R, L)$$

$$Q_{\pi}(L, W) = P(L, W, H) \times R(L, W, H) + P(L, W, L) \times R(L, W, L)$$

$$Q_{\pi}(L, S) = P(L, S, H) \times R(L, S, H) + P(L, S, L) \times R(L, S, L)$$

$$Q_{\pi}(L, R) = P(L, R, H) \times R(L, R, H) + P(L, R, L) \times R(L, R, L)$$

Therefore:

$$Q_{\pi}(H, W) = 1 \times r_{wait} = 1 \times 0 = 0$$

$$Q_{\pi}(H, S) = 2 \times r_{search} + (1-2) \times r_{search} = 0.5 \times 3 + 0.5 \times 3 = 3$$

$$Q_{\pi}(L, W) = 1 \times r_{wait} = 1 \times 0 = 0$$

$$Q_{\pi}(L, R) = 1 \times 0 = 0$$

$$Q_{\pi}(L, S) = (1-\beta) \times (-3) + \beta \times r_{search} = -3 \times 0.7 + 0.3 \times 3 = -1.2$$

Therefore, the action-value function for each value pair

is: $Q_{\pi}(H, W) = 0$, $Q_{\pi}(H, S) = 3$, $Q_{\pi}(L, W) = 0$, $Q_{\pi}(L, S) = -1.2$ and $Q_{\pi}(L, R) = 0$.

3. solution

- (1) For a policy that always takes the action 'produce' in the 'Good' state, the conditional probability distribution over actions given states can be expressed by: (assume 'G' is good state, 'B' is broken state, 'P' is produce action, 'I' is inactive action and 'R' is repair action)

$$\pi(P|G) = 1 \quad \pi(I|G) = 0 \quad \pi(R|B) = 1 \quad \pi(R|G) = 0$$

$$\text{so: } v_{\pi}(s) = E[G_t | S_t = s, \pi] = E[R_{t+1} + \gamma v_{\pi}(s') | \pi, S_t = s]$$

$$\begin{aligned} \text{Therefore: } v_{\pi}(G) &= \pi(P|G) \times [P(G, P, G) \times (R(G, P, G) + \gamma v_{\pi}(G)) \\ &\quad + P(G, P, B) \times (R(G, P, B) + \gamma v_{\pi}(B))] \\ &\quad + \pi(I|G) \times [P(G, I, G) \times (R(G, I, G) + \gamma v_{\pi}(G)) + P(G, I, B) \\ &\quad \times (R(G, I, B) + \gamma v_{\pi}(B))] \\ &= 1 \times [0.8 \times (2 + \gamma v_{\pi}(G)) + 0.2 \times (-5 + \gamma v_{\pi}(B))] \\ &= 1.6 + 0.8\gamma v_{\pi}(G) + 0.2\gamma v_{\pi}(B) \end{aligned}$$

it's a recursive process.

$$\text{so: } v_{\pi}(G) = 1.6 + 0.8\gamma v_{\pi}(G) + 0.2\gamma v_{\pi}(B) \quad (v_{\pi}(G) \text{ is similar to } v_{\pi}(G) \text{ and } v_{\pi}(B) \text{ and is a recursive process})$$

$$\begin{aligned} v_{\pi}(B) &= \pi(R|B) \times [P(B, R, G) \times (R(B, R, G) + \gamma v_{\pi}(G)) + P(B, R, B) \times \\ &\quad (R(B, R, B) + \gamma v_{\pi}(B))] \\ &= 1 \times [0.8 \times (-1 + \gamma v_{\pi}(G)) + 0.2 \times (-5 + \gamma v_{\pi}(B))] \\ &= -1 + 0.8\gamma v_{\pi}(G) + 0.2\gamma v_{\pi}(B) \end{aligned}$$

$$\text{so: } v_{\pi}(B) = -1 + 0.8\gamma v_{\pi}(G) + 0.2\gamma v_{\pi}(B) \quad (v_{\pi}(B) \text{ is similar to } v_{\pi}(B) \text{ and is a recursive process})$$

- (2) For the optimal value function in 'Good' state and 'Broken' state denoted by $V^*(G)$ and $V^*(B)$.

$$V^*(G) = \max_{a \in A} \sum_{s' \in S} P(G, a, s') (R(G, a, s') + \gamma V^*(s'))$$

$$V^*(B) = \max_{a \in A} \sum_{s' \in S} P(B, a, s') (R(B, a, s') + \gamma V^*(s'))$$

$$\text{Therefore, } V^*(G) = \max_{a \in A} \sum_{s' \in S} P(G, a, s') (R(G, a, s') + \gamma V^*(s'))$$

$$= \max \left\{ \left[P(G, P, G) \cdot (R(G, P, G) + \gamma V'_\pi(G)) + P(G, P, B) \cdot (R(G, P, B) + \gamma V'_\pi(B)) \right], \right. \\ \left. P(G, I, G) \cdot (R(G, I, G) + \gamma V'_\pi(G)) \right\}$$

$$= \max \left\{ [0.8 \times (2 + \gamma V'_\pi(G)) + 0.2 \cdot \gamma V'_\pi(B)], \gamma V'_\pi(G) \right\}$$

So: ~~is~~

$$= \max \left\{ [1.6 + 0.8\gamma V'_\pi(G) + 0.2\gamma V'_\pi(B)], \gamma V'_\pi(G) \right\}$$

$$= \max \left\{ [1.6 + 0.8\gamma V^*_\pi(G) + 0.2\gamma V^*_\pi(B)], \gamma V^*_\pi(G) \right\}$$

$$= \max \left\{ (1.6 + 0.8\gamma V^*_\pi(G) + 0.2\gamma V^*_\pi(B)), \gamma V^*_\pi(G) \right\}$$

$$\text{So: } V^*(G) = \max \left\{ 1.6 + 0.8\gamma V^*(G) + 0.2\gamma V^*(B), \gamma V^*(G) \right\}$$

$$\text{Then: } V^*(B) = \max_{a \in A} \sum_{s' \in S} P(B, a, s') (R(B, a, s') + \gamma V^*(s'))$$

$$= P(B, R, G) (R(B, R, G) + \gamma V^*(G)) + P(B, R, B) (R(B, R, B) + \gamma V^*(B))$$

$$= 0.8 \cdot \gamma V^*(G) + 0.2 \cdot (-5 + \gamma V^*(B))$$

$$= -1 + 0.8\gamma V^*(G) + 0.2\gamma V^*(B)$$

$$\text{So: } V^*(B) = -1 + 0.8\gamma V^*(G) + 0.2\gamma V^*(B)$$

Therefore:

$$\left\{ \begin{array}{l} V^*(G) = \max \left\{ (1.6 + 0.8\gamma V^*(G) + 0.2\gamma V^*(B)), \gamma V^*(G) \right\} \\ V^*(B) = -1 + 0.8\gamma V^*(G) + 0.2\gamma V^*(B) \end{array} \right.$$

(And $V^*(G)$, $V^*(B)$ is the ~~max~~ optimal value function by the next state related to "G" and "B").