




Figure 1: MDP problem with (action, reward) pairs.

Problem 1a [2 points]

Find the optimal policy at the initial state S_A with discount factor $\gamma = 0.001$. Justify your answer.

episode의 length가 3으로 제한되어 있으므로 모든 policy에 대한 utility를 직접 계산한다.

$$\begin{aligned}
 & \pi(S_A) = - \quad \text{reward} \\
 & \begin{aligned}
 & \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \quad 5 + r \cdot 5 + r^2 \cdot 5 \\
 & \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \quad 5 + r \cdot 5 \\
 & \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \xrightarrow{-} \textcircled{S_A} \quad 5 \\
 & \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \xrightarrow{+} \textcircled{S_C} \quad 5
 \end{aligned} \\
 & \text{total utility of } \pi(S_A) = - : 20 + 10\gamma + 5\gamma^2 \\
 & \text{value} : \frac{20 + 10\gamma + 5\gamma^2}{4}
 \end{aligned}$$

$$\begin{aligned}
 & \pi(S_A) = + \\
 & \begin{aligned}
 & \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \xrightarrow{+} \textcircled{S_C} \xrightarrow{+} \textcircled{S_D} \quad 0 + r \cdot 0 + r^2 \cdot 16 \\
 & \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \xrightarrow{+} \textcircled{S_C} \xrightarrow{-} \textcircled{S_B} \quad 0 + r \cdot 0 + r^2 \cdot 0 \\
 & \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \xrightarrow{-} \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \quad 0 + r \cdot 0 + r^2 \cdot 0 \\
 & \textcircled{S_A} \xrightarrow{+} \textcircled{S_B} \xrightarrow{-} \textcircled{S_A} \xrightarrow{-} \textcircled{S_A} \quad 0 + r \cdot 0 + r^2 \cdot 5
 \end{aligned} \\
 & \text{total utility of } \pi(S_A) = + : 21\gamma^2 \\
 & \text{value} : \frac{21\gamma^2}{4}
 \end{aligned}$$

$$\text{when } \gamma = 0.001, \quad \begin{array}{l} \pi(S_A) = - \Rightarrow V\pi(S_A) \approx \frac{10}{4} = 2.5 \\ \pi(S_A) = + \Rightarrow V\pi(S_B) \approx 0 \end{array}$$

$$\therefore \pi_{\text{opt}}(S_A) = - \quad \text{when } \gamma = 0.001$$

Problem 1b [2 points]

Find the optimal policy at the initial state S_A with discount factor $\gamma = 0.999$. Justify your answer.

$$\text{When } \gamma = 0.999, \quad \begin{array}{l} \pi(S_A) = - \Rightarrow V\pi(S_A) \approx 8.75 \\ \pi(S_A) = + \Rightarrow V\pi(S_B) \approx 5.24 \end{array}$$

$$\therefore \pi_{\text{opt}}(S_A) = - \quad \text{when } \gamma = 0.999$$

Problem 1c [2 points]

What is the optimal policy at the initial state S_B ? Explain your answer in terms of discount factor $\gamma \in (0, 1)$.

$$\pi(S_B) = - \text{인정}$$

reward

$$S_B \rightarrow S_A \rightarrow S_A \rightarrow S_A \quad 0 + r \cdot 5 + r^2 \cdot 5$$

$$S_B \rightarrow S_A \rightarrow S_A \xrightarrow{+} S_B \quad 0 + r \cdot 5 + r^2 \cdot 0$$

$$S_B \rightarrow S_A \xrightarrow{+} S_B \rightarrow S_A \quad 0 + r \cdot 0 + r^2 \cdot 0$$

$$S_B \rightarrow S_A \xrightarrow{+} S_B \xrightarrow{+} S_C \quad 0 + r \cdot 0 + r^2 \cdot 0$$

$$\begin{aligned} \text{utility} &= 10r + 5r^2 \\ \text{value} &= \frac{10r + 5r^2}{1} \end{aligned}$$

$$\pi(S_B) = + \text{인정}$$

reward

$$S_B \xrightarrow{+} S_C \xrightarrow{+} S_D \xrightarrow{+} S_D \quad 0 + r \cdot 16 + \gamma^2 \cdot 0$$

$$S_B \xrightarrow{+} S_C \xrightarrow{+} S_D \xrightarrow{-} S_C \quad 0 + r \cdot 16 + \gamma^2 \cdot 0$$

$$S_B \xrightarrow{+} S_C \xrightarrow{-} S_B \xrightarrow{+} S_C \quad 0 + r \cdot 0 + \gamma^2 \cdot 0$$

$$S_B \xrightarrow{+} S_C \xrightarrow{-} S_B \xrightarrow{-} S_A \quad 0 + r \cdot 0 + \gamma^2 \cdot 0$$

$$\text{utility} = 32r$$

$$\text{value} = 8r$$

$$0 < \gamma < 1, \quad 8r > \frac{10r + 5r^2}{4}$$

$$\therefore \pi_{\text{opt}}(S_B) = +$$