

$$V_{\pi}(s) = E_{\pi}[r_t | s_t = s] + \gamma \sum_{s' \in S} \pi(a|s) \sum_{a \in A} p(s'|s, a) V_{\pi}^a(s') \quad (1)$$

$$A = \{\text{stay, switch}\}, \quad S = \{H, O\}$$

$$\pi(a|s) = \frac{1}{2} \forall a, s \in A, s \in S \quad (\text{uniform policy})$$

$$V^{\pi}(H) = E_{\pi}[r | s=H] + \gamma [\pi(\text{stay}|H) [P(H|H, \text{stay}) V^{\pi}(H) + P(O|H, \text{stay}) V^{\pi}(O)] + \pi(\text{switch}|H) [P(H|H, \text{switch}) V^{\pi}(H) + P(O|H, \text{switch}) V^{\pi}(O)]]$$

$$\frac{0+1}{2} = \frac{1}{2} \cdot \left[\frac{1}{2} [1 \cdot V^{\pi}(H) + 0 \cdot V^{\pi}(O)] + \frac{1}{2} [0 \cdot V^{\pi}(H) + 1 \cdot V^{\pi}(O)] \right]$$

$$\frac{1}{2} [0.2 \cdot V^{\pi}(H) + 0.8 V^{\pi}(O)] =$$

$$\frac{1}{2} \cdot \frac{1}{2} (1.2 V^{\pi}(H) + 0.8 V^{\pi}(O)) = \frac{1}{2} \cdot 0.3 V^{\pi}(H) + 0.2 V^{\pi}(O) =$$

$$0. \Rightarrow V^{\pi}(H) = 0.2 V^{\pi}(O) + \frac{1}{2}$$

$$V^{\pi}(H) = \frac{2}{7} V^{\pi}(O) + \frac{5}{7}$$

$$V^{\pi}(O) = \frac{1}{3} \left(\frac{2}{7} V^{\pi}(O) + \frac{5}{7} \right) + \frac{4}{3} \Rightarrow 3V^{\pi}(O) = \frac{2}{7} V^{\pi}(O) + \frac{5}{7} + 4$$

$$\left(2 + \frac{5}{7}\right) V^{\pi}(O) = \frac{5}{7} + 4 \Rightarrow V^{\pi}(O) = 1.736$$

$$V^{\pi}(H) = \frac{2}{7} \cdot 1.736 + \frac{5}{7} = 1.21$$

$$V^{\pi}(s) = \begin{cases} 1.21 & s=H \\ 1.736 & s=O \end{cases}$$

(2) נתון π האומר אם הנגזרים יביאו תגובה או לא (האם הם יביאו תגובה או לא) $H \rightarrow$ תגובה (האם הם יביאו תגובה או לא) π \rightarrow האם הם יביאו תגובה או לא.

$$\pi^*(a|s) = \begin{cases} 1 & (s=H, a=switch) \text{ or } (s=O, a=stay) \\ 0 & \text{else} \end{cases}$$

(3) נבדוק האם π הוא אופטימלי π \rightarrow האם הוא אופטימלי.

$$V^{\pi^*}(s) = \max_{a' \in A} \sum_{s' \in S} E[r|s, a'] + \gamma \sum_{s' \in S} P(s'|s, a') V^{\pi^*}(s')$$

$$V^{\pi^*}(H) = \max \left\{ r(H, stay) + \gamma [P(H|H, stay) V^{\pi^*}(H) + P(O|H, stay) V^{\pi^*}(O)], r(H, switch) + \gamma [P(H|H, switch) V^{\pi^*}(H) + P(O|H, switch) V^{\pi^*}(O)] \right\}$$

$$V^{\pi^*}(O) = \max \left\{ r(O, stay) + \gamma [P(H|O, stay) V^{\pi^*}(H) + P(O|O, stay) V^{\pi^*}(O)], r(O, switch) + \gamma [P(H|O, switch) V^{\pi^*}(H) + P(O|O, switch) V^{\pi^*}(O)] \right\}$$

~~Value~~ $V^{\pi}(s)$ \rightarrow $V^{\pi}(s)$ \rightarrow $V^{\pi}(s)$

$$V^{\pi}(H) = 1 + \frac{1}{2} [0 \cdot [P(H|H, \text{stay}) V^{\pi}(H) + P(O|H, \text{stay}) V^{\pi}(O)] + 1 \cdot [0.2 V^{\pi}(H) + 0.8 V^{\pi}(O)]] = 1 + \frac{1}{2} (0.2 V^{\pi}(H) + 0.8 V^{\pi}(O))$$

$$V^{\pi}(O) = 2 + \frac{1}{2} [1 \cdot [0 V^{\pi}(H) + 1 \cdot V^{\pi}(O)] + 0 \cdot [P(H|O, \text{switch}) \cdot V^{\pi}(H) + P(O|O, \text{switch}) V^{\pi}(O)]] = 2 + \frac{1}{2} (V^{\pi}(O))$$

\Downarrow

$$1.3 V^{\pi}(O) = 2 \Rightarrow V^{\pi}(O) = \frac{4}{3}$$

$$V^{\pi}(H) = 1 + \frac{1}{2} (0.2 V^{\pi}(H) + 0.8 \cdot \frac{4}{3}) = 1 + 0.9 V^{\pi}(H) + 1.066$$

\Downarrow

$$0.9 V^{\pi}(H) = 2.066$$

\Downarrow

$$V^{\pi}(H) = 2.96$$

Value Iteration

$$V^{\pi}(H) = \max \left\{ 0 + \frac{1}{2} (1 \cdot 2.96), 1 + \frac{1}{2} (0.2 \cdot 2.96 + 0.8 \cdot \frac{4}{3}) \right\} =$$

$$\max \{ 1.48, 1.83 \} = 1.83 \Rightarrow \arg \max_a = \text{switch}$$

$$V^{\pi}(O) = \max \left\{ 2 + \frac{1}{2} (0 + 1 \cdot \frac{4}{3}), 0 + \frac{1}{2} (1 \cdot 2.96) + 0 \right\} =$$

$$\max \{ 2.66, 1.48 \} = 2.66 \Rightarrow \arg \max_a = \text{stay}$$