

① a)

First 2 iterations of the Value Iteration alg.:

Step 0:

$$V_0(s_1) = 10.0$$

$$V_0(s_2) = 1.0$$

$$V_0(s_3) = 0.0$$

Step 1:

$$q_1(s_1, a_1) = R(s_1, a_1) + \gamma (P(s_1, a_1, s_1) V_0(s_1) + P(s_1, a_1, s_2) V_0(s_2))$$

$$= 8.0 + 1.0 (0.2 \cdot 10.0 + 0.6 \cdot 1.0)$$

$$= 10.6$$

$$q_1(s_1, a_2) = R(s_1, a_2) + \gamma (P(s_1, a_2, s_1) V_0(s_1) + P(s_1, a_2, s_2) V_0(s_2))$$

$$= 10.0 + 1.0 (0.1 \cdot 10.0 + 0.2 \cdot 1.0)$$

$$= 11.2$$

$$q_1(s_2, a_1) = R(s_2, a_1) + \gamma (P(s_2, a_1, s_1) V_0(s_1) + P(s_2, a_1, s_2) V_0(s_2))$$

$$= 1.0 + 1.0 (0.3 \cdot 10.0 + 0.3 \cdot 1.0)$$

$$= 4.3$$

$$\begin{aligned}
 q_1(s_2, a_2) &= R(s_2, a_2) + \gamma (P(s_2, a_2, s_1) V_0(s_1) \\
 &\quad + P(s_2, a_2, s_2) V_0(s_2)) \\
 &= -1.0 + 1.0 (0.5 \cdot 1.0 + 0.3 \cdot 1.0) \\
 &= 4.3
 \end{aligned}$$

$$\begin{aligned}
 V_1(s_1) &= \max_{a \in A} \{ q_1(s_1, a) \} \\
 &= 11.2
 \end{aligned}$$

$$\begin{aligned}
 V_1(s_2) &= \max_{a \in A} \{ q_1(s_2, a) \} \\
 &= 4.3
 \end{aligned}$$

$$\pi_1(s_1) = a_2$$

$$\pi_1(s_2) = a_1 \quad (a_1 \text{ and } a_2 \text{ produce an equivalent greedy improvement})$$

Step 2:

$$\begin{aligned}
 q_2(s_1, a_1) &= R(s_1, a_1) + \gamma (P(s_1, a_1, s_1) V_1(s_1) \\
 &\quad + P(s_1, a_1, s_2) V_1(s_2)) \\
 &= 8.0 + 1.0 (0.2 \cdot 11.2 + 0.6 \cdot 4.3) \\
 &= 12.82
 \end{aligned}$$

$$\begin{aligned}
 q_2(s_1, a_2) &= R(s_1, a_2) + \gamma (P(s_1, a_2, s_1) V_1(s_1) \\
 &\quad + P(s_1, a_2, s_2) V_1(s_2)) \\
 &= 10.0 + 1.0 (0.1 \cdot 11.2 + 0.2 \cdot 4.3) \\
 &= 11.98
 \end{aligned}$$

$$\begin{aligned}
 q_2(s_2, a_1) &= R(s_2, a_1) + \gamma (P(s_2, a_1, s_1) V_1(s_1) \\
 &\quad + P(s_2, a_1, s_2) V_1(s_2)) \\
 &= 1.0 + 1.0 (0.3 \cdot 11.2 + 0.3 \cdot 4.3) \\
 &= 5.65
 \end{aligned}$$

$$\begin{aligned}
 q_2(s_2, a_2) &= R(s_2, a_2) + \gamma (P(s_2, a_2, s_1) V_1(s_1) \\
 &\quad + P(s_2, a_2, s_2) V_1(s_2)) \\
 &= -1 + 1 (0.5 \cdot 11.2 + 0.3 \cdot 4.3) \\
 &= 5.89
 \end{aligned}$$

$$\begin{aligned}
 V_2(s_1) &= \max_{a \in A} \{ q_2(s_1, a) \} \\
 &= 12.82
 \end{aligned}$$

$$\begin{aligned}
 V_2(s_2) &= \max_{a \in A} \{ q_2(s_2, a) \} \\
 &= 5.89
 \end{aligned}$$

$$\pi_2(s_1) = a_1$$

$$\pi_2(s_2) = a_2$$

b) for $k=3, 4, \dots$ we have:

$$\begin{aligned}
 & q_k(s_1, a_1) - q_k(s_1, a_2) \\
 &= \underbrace{R(s_1, a_1) - R(s_1, a_2)}_{=-2} + \gamma \underbrace{\left(\sum_{s' \in N} \underbrace{P(s_1, a_1, s') - P(s_1, a_2, s')}_{>0} V_{k-1}(s') \right)}_{>0} \\
 &= -2 + (0.1 V_{k-1}(s_1) + 0.4 V_{k-1}(s_2))
 \end{aligned}$$

at $k=3$, $V_{k-1}(s_1) = 12.82$, $V_{k-1}(s_2) = 5.89$.

$$\Rightarrow 0.1 \cdot 12.82 + 0.4 \cdot 5.89 = 3.638 > 2$$

$$\Rightarrow q_{k=3}(s_1, a_1) - q_{k=3}(s_1, a_2) > 0$$

$$\Rightarrow \pi_{k=3}(s_1) = a_1$$

Since $V_{k+1}(s) \geq V_k(s)$ for all k , all s
 by Banach's Thm, then for all $k=3, 4, \dots$,

$$q_k(s_1, a_1) - q_k(s_1, a_2) > 0$$

$$\Rightarrow \pi_k(s_1) = a_1$$



Similar argument for $\pi_k(s_2)$.

#2

• state at time t : $s \in S = \{ ("employed", j_1), \dots, ("employed", j_n), ("unemployed", j_1), \dots, ("unemployed", j_n) \}$

Which represent the state of being "employed" at job i , $i=1, \dots, n$, or "unemployed" and offered job i , $i=1, \dots, n$, at the beginning of day t .

• action at time t : $a_t \in A = \{ "nothing", "accept", "decline" \}$

The set of possible actions depends on the state.

We express that as $A(s)$, where $s \in S$.

$$A("unemployed", j_i) = \{ "accept", "decline" \}, \quad i=1, \dots, n$$

$$A("employed", j_i) = \{ "nothing" \}, \quad i=1, \dots, n$$

$$P(s, a, s'): S \times A \times S \rightarrow [0, 1]$$

$$P(s, a, s') = \begin{cases} (p_h), \text{ if } s = (\text{"unemployed"}, j_i), a = \text{"decline"} \\ \quad s' = (\text{"unemployed"}, j_h), i=1, \dots, n; h=1, \dots, n \\ (x \cdot p_h), \text{ if } s = (\text{"unemployed"}, j_i), a = \text{"accept"} \\ \quad s' = (\text{"unemployed"}, j_h), i=1, \dots, n; h=1, \dots, n \\ (1-x), \text{ if } s = (\text{"unemployed"}, j_i), a = \text{"accept"} \\ \quad s' = (\text{"employed"}, j_i), i=1, \dots, n \\ (x \cdot p_h), \text{ if } s = (\text{"employed"}, j_i), a = \text{"nothing"} \\ \quad s' = (\text{"unemployed"}, j_h), i=1, \dots, n; h=1, \dots, n \\ (1-x), \text{ if } s = (\text{"employed"}, j_i), a = \text{"nothing"} \\ \quad s' = (\text{"employed"}, j_i) \\ \text{else} \end{cases}$$

$$R_T(s, a, s'): S \times A \times S \rightarrow \mathbb{R}^+$$

$$R_T(s, a, s') = \begin{cases} w_0, \text{ if } s = (\text{"unemployed"}, j_i), \\ \quad a = \text{"decline"}, i=1, \dots, n \\ w_i, \text{ if } s = (\text{"employed"}, j_i), i=1, \dots, n \\ w_i, \text{ if } s = (\text{"unemployed"}, j_i), i=1, \dots, n, \\ \quad a = \text{"accept"} \end{cases}$$

* Bellman Optimality Equation:

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