

1. a. for all s $r_s = -1$. In this way, to max the profit of the agent, he will choose the path that goes through the min amount of states, i.e the shortest path.

The optimal value of each state would be the number of states to get to g + s .

b. $-(V_{old}^{\pi_g} - s) + \gamma$:

12	11	10	9
-3	10	9	8
10	9	8	7
11	12	13	14

$$\begin{aligned} c. V_{new}^{\pi} &= E\left[\sum_{l=0}^{\infty} \gamma^l r_{t+l} | S_t = s\right] = \\ &= E\left[\sum_{l=0}^{\infty} \gamma^l r_{t+l} | S_t = s\right] + E\left[\sum_{l=0}^{\infty} \gamma^l c | S_t = s\right] = \\ &V_{old}^{\pi} + \sum_{l=0}^{\infty} \gamma^l c \end{aligned}$$

d. Now $r_s = 1$, thus the optimal policy will be the longest path, thus the value will be ∞ and we will get stuck at the edges.

e. We will still wonder around, until γ^+ will be small enough and then we'll go to g .

f. $r_s \leq -5$

a. for $t \in [H]$:

$$V_t^{\pi_1}(x_t) - V_t^{\pi_2}(x_t) = Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) =$$

$$Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) + \underbrace{Q_t^{\pi_1}(x_t, \pi_2(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t))}_{\otimes} =$$

$$\begin{aligned}
 \otimes \quad Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) &\stackrel{p(x_t, \pi_2(x_t, t)) \stackrel{\text{def}}{=} \Omega}{=} \\
 & r_t(x_t, \pi_2(x_t, t)) - \mathbb{E}_{s' \sim \Omega} V_{t+1}^{\pi_1}(s') - r_t(x_t, \pi_2(x_t, t)) - \mathbb{E}_{s' \sim \Omega} V_{t+1}^{\pi_2}(s') \\
 &= \mathbb{E}_{s' \sim \Omega} V_{t+1}^{\pi_1}(s') - \mathbb{E}_{s' \sim \Omega} V_{t+1}^{\pi_2}(s') = \mathbb{E}_{s' \sim \Omega} [V_{t+1}^{\pi_1}(s') - V_{t+1}^{\pi_2}(s')]
 \end{aligned}$$

$$\Rightarrow V_t^{\pi_1}(x_t) - V_t^{\pi_2}(x_t) = Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) + \mathbb{E}_{s' \sim \Omega} [V_{t+1}^{\pi_1}(s') - V_{t+1}^{\pi_2}(s')]$$

Now, we will take the expectation of $\mathbb{E}_{x_t \sim \pi_2}$ on both sides:

$$\mathbb{E}_{x_t \sim \pi_2} [V_t^{\pi_1}(x_t) - V_t^{\pi_2}(x_t)] = \mathbb{E}_{x_t \sim \pi_2} [Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) + \mathbb{E}_{s' \sim \Omega} [V_{t+1}^{\pi_1}(s') - V_{t+1}^{\pi_2}(s')]]$$

$$\Rightarrow V_t^{\pi_1}(x_t) - V_t^{\pi_2}(x_t) = \mathbb{E}_{x_t \sim \pi_2} \{ Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) + \mathbb{E}_{s' \sim \Omega} [V_{t+1}^{\pi_1}(s') - V_{t+1}^{\pi_2}(s')] \} =$$

$$= \mathbb{E}_{x_t \sim \pi_2} \{ Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_2}(x_t, \pi_2(x_t, t)) \} +$$

$$\mathbb{E}_{x_{t+1} \sim \pi_2} \{ V_{t+1}^{\pi_1}(x_{t+1}) - V_{t+1}^{\pi_2}(x_{t+1}) \} =$$

$$= \sum_{T=t}^H \mathbb{E}_{x_T \sim \pi_2} \{ Q_T^{\pi_1}(x_T, \pi_1(x_T, T)) - Q_T^{\pi_2}(x_T, \pi_2(x_T, T)) \}$$