# A4-QA

## Q1

- G0 = 4
- G1 = 6
- G2 = 8
- G3 = 4
- G4 = 2
- G5 = 0

## Q2

### 1

$$V_\pi(s) = E_\pi \left[ r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \ldots \mid s_t = s \right]$$
$$= E_{a \sim \pi(s, \cdot)} \left[ E_\pi \left[ r_{t+1} + \gamma r_{t+2} + \ldots \mid s_t = s, A_t = a \right] \right]$$
$$= \sum_{a \in A} \pi(s, a) \times \theta_\pi(s, a)$$
$$= \sum_{a \in A} \pi(s, a) \, \theta_\pi(s, a)$$

$$\theta_\pi(s, a) = E_\pi \left[ r_{t+1} + \gamma r_{t+2} + \ldots \mid s_t = s, A_t = a \right]$$
$$= E_{s' \sim P(\cdot \mid s, a)} \left[ r(s, a, s') + \gamma E_\pi \left[ r_{t+2} + r_{t+3} + \ldots \mid s_{t+1} = s' \right] \right]$$
$$= \sum_{s \in S} P(s' \mid s, a) \left[ r(s, a, s') + \gamma \times V_\pi(s') \right]$$

### 2

$$V_\pi(s) = E \left[ r_{t+1} + \gamma V_\pi(s') \mid s_t = s \right]$$
$$= \sum_{a \in A} \pi(s, a) \sum_{s' \in S} P(s' \mid s, a) \left[ r(s, a, s') + \gamma V_\pi(s') \right]$$
$$= E_{a \sim \pi(s, \cdot)} E_{s' \sim P(\cdot \mid s, a)} \left[ r(s, a, s') + \gamma V_\pi(s') \right]$$

$$\theta_\pi(s) = E_\pi \left( r_{t+1} + \gamma \theta_\pi(s', a') \mid s_t = s, A_t = a \right)$$
$$= \sum_{s' \in S} P(s' \mid s, a) \left[ r(s, a, s') + \gamma \sum_{a' \in A} \pi(s', a') \theta_\pi(s', a') \right]$$
$$= E_{s' \sim P(\cdot \mid s, a)} \left[ r(s, a, s') + \gamma E_{a' \sim \pi(s', \cdot)} \left[ \theta_\pi(s', a') \right] \right]$$

# Q3

## 1

$$V_1 = \frac{1}{4}(0 + rV_3) + \frac{1}{4}(0 + rV_2) + \frac{1}{4}(-1 + 0) + \frac{1}{4}(-1 + 0) = \frac{r}{4}(V_2 + V_3) - \frac{1}{2}$$
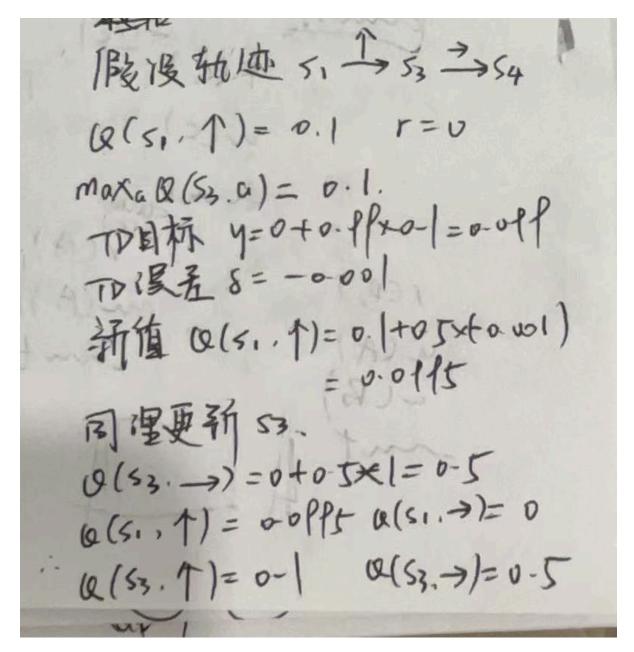
$$V_2 = \frac{1}{4}(1 + 0) + \frac{1}{4}(-1 + 0) + \frac{1}{4}(-1 + 0) + \frac{1}{4}(0 + V_1 r) = -\frac{1}{4} + \frac{r}{4}V_1$$

$$V_3 = -\frac{1}{4} + \frac{r}{4}V_1$$

解得 $\begin{cases} V_1 = 0.7108 \\ V_2 = V_3 = -0.4260 \end{cases}$

## 2

$$q_\pi(S_3, 向右) = 1 + r \cdot 0 = 1$$

$$q_\pi(S_3, 向下) = 0 + r \cdot V_1 \approx -0.7038$$

$$q_\pi(S_3, 向左/向上) = -1 + r \cdot 0 = -1$$

∴ 最优为向右. $\pi'(S_3) = 向右.$

假设轨迹 $s_1 \xrightarrow{\uparrow} s_3 \xrightarrow{\rightarrow} s_4$

$Q(s_1, \uparrow) = 0.1$     $r = 0$

$\max_a Q(s_3, a) = 0.1$.

TD目标 $y = 0 + 0.9 \times 0.1 = 0.09$

TD误差 $\delta = -0.001$

新值 $Q(s_1, \uparrow) = 0.1 + 0.5 \times (-0.001)$

$\qquad\qquad\qquad = 0.0995$

同理更新 $s_3$.

$Q(s_3, \rightarrow) = 0 + 0.5 \times 1 = 0.5$

$Q(s_1, \uparrow) = 0.0995$   $Q(s_1, \rightarrow) = 0$

$\therefore Q(s_3, \uparrow) = 0.1$     $Q(s_3, \rightarrow) = 0.5$

最后策略为 $s_1$ 时向上，$s_3$ 时向右