

$$\text{定义 1: } Q_2 \leftarrow Q_2$$

$$U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{n-t} R_n$$

$$= \sum_{k=t}^n \gamma^{k-t} R_k$$

$$U_{t+1} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-t-1} R_n$$

$$= \sum_{k=t+1}^n \gamma^{k-t-1} R_k$$

$$\Rightarrow U_t = R_t + \gamma (R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-t-1} R_n)$$

$$= R_t + \gamma U_{t+1}$$

$$\text{Denote } S_{t+1} = \{S_{t+1}, S_{t+2}, \dots\} \quad A_{t+1} = \{A_{t+1}, A_{t+2}, \dots\}$$

$$Q_2(s_t, a_t) = \mathbb{E}_{S_{t+1}, A_{t+1}} [U_t | S_t = s_t, A_t = a_t]$$

$$\text{替代 } U_t = \mathbb{E}_{S_{t+1}, A_{t+1}} [R_t + \gamma U_{t+1} | S_t = s_t, A_t = a_t]$$

$$= \mathbb{E}_{S_{t+1}, A_{t+1}} [R_t | S_t = s_t, A_t = a_t]$$

$$+ \gamma \mathbb{E}_{S_{t+1}, A_{t+1}} [U_{t+1} | S_t = s_t, A_t = a_t]$$

R_t 是 S_t, A_t 和 S_{t+1} 的函数 (状态转移函数有随机性, 因此 S_{t+1} 也是重要的)

$$\mathbb{E}_{S_{t+1}, A_{t+1}} [R_t | S_t = s_t, A_t = a_t] = \mathbb{E}_{S_{t+1}} [R_t | S_t = s_t, A_t = a_t]$$

$$\text{而 } U_{t+1} = \sum_{k=t+1}^n \gamma^{k-t-1} R_k \text{ 是 } S_{t+1} = \{S_{t+1}, S_{t+2}, \dots, S_n\}$$

$$A_{t+1} = \{A_{t+1}, A_{t+2}, \dots, A_n\} \text{ 的函数}$$

$$\mathbb{E}_{S_{t+1}, A_{t+1}} [U_{t+1} | S_t = s_t, A_t = a_t]$$

$$= \mathbb{E}_{S_{t+1}, A_{t+1}} [\mathbb{E}_{S_{t+2}, A_{t+2}} [U_{t+1} | S_{t+1}, A_{t+1}] | S_t = s_t, A_t = a_t]$$

$$= \mathbb{E}_{S_{t+1}, A_{t+1}} [Q_2(S_{t+1}, A_{t+1}) | S_t = s_t, A_t = a_t]$$

$$\therefore Q_2(s_t, a_t) = R_t + \gamma \mathbb{E}_{S_{t+1}, A_{t+1}} [Q_2(S_{t+1}, A_{t+1}) | S_t = s_t, A_t = a_t]$$

$$= \mathbb{E}_{S_{t+1}, A_{t+1}} [R_t + \gamma Q_2(S_{t+1}, A_{t+1}) | S_t = s_t, A_t = a_t]$$

形式 2: $Q_\pi \leftarrow V_\pi$

$$\text{由 } V_\pi(s_t) = \sum_{a \in A} \pi(a|s_t) \cdot Q_\pi(s_t, a) = \mathbb{E}_{A_t} [Q(s_t, A_t)]$$

$$\text{同理 } V_\pi(s_{t+1}) = \mathbb{E}_{A_{t+1}} [Q(s_t, A_t)]$$

$$\begin{aligned} \text{对于 } \mathbb{E}_{S_{t+1}, A_{t+1}} [Q_\pi(s_{t+1}, A_{t+1}) | S_t = s_t, A_t = a_t] \\ &= \mathbb{E}_{S_{t+1}} [\mathbb{E}_{A_{t+1}} [Q_\pi(s_{t+1}, A_{t+1})] | S_t = s_t, A_t = a_t] \\ &= \mathbb{E}_{S_{t+1}} [V_\pi(s_{t+1}) | S_t = s_t, A_t = a_t] \end{aligned}$$

$$\therefore Q_\pi(s_t, a_t) = \mathbb{E}_{S_{t+1}, A_{t+1}} [R_t + \gamma V_\pi(s_{t+1}) | S_t = s_t, A_t = a_t]$$

形式 3: $Q_\pi \leftarrow V_\pi$

$$\text{由 } V_\pi(s_t) = \mathbb{E}_{A_t} [Q(s_t, A_t)]$$

$$\begin{aligned} V_\pi(s_t) &= \mathbb{E}_{A_t} [\mathbb{E}_{S_{t+1}, A_{t+1}} [R_t + \gamma V_\pi(s_{t+1}) | S_t = s_t, \underline{A_t = a_t}]] \\ &= \mathbb{E}_{S_{t+1}, A_t} [R_t + \gamma V_\pi(s_{t+1}) | S_t = s_t] \end{aligned}$$

最优贝尔曼方程.

最优策略函数. $\pi^* = \arg \max_{\pi} Q_\pi(s, a) \quad \forall s, a$

↓ 就是最优动作价值函数

$$Q_{\pi^*}(s_t, a_t) = \mathbb{E}_{S_{t+1}, A_{t+1}} [R_t + \gamma Q_{\pi^*}(s_{t+1}, A_{t+1}) | S_t = s_t, A_t = a_t]$$

没什么不同, 只是现在指定了策略.

原先 $Q(s_t, a_t)$ 是关于 $\pi \in \Pi$ 的函数, 现固定 $\pi = \pi^*$

$$\text{由 } \pi^* \text{ 性质, } Q_{\pi^*}(s_{t+1}, A_{t+1}) = \max_{A \in A} Q_{\pi^*}(s_{t+1}, A)$$

$$Q_{\pi^*}(s_t, a_t) = \mathbb{E}_{S_{t+1}} [R_t + \gamma \max_{A \in A} Q_{\pi^*}(s_{t+1}, A) | S_t = s_t, A_t = a_t]$$