

~~$$V_n(s) = E_x(b_t | S_t = s) = E_x(R_{t+1} +$$~~

~~$$(1) \text{ 同理可得: } P(S_{t+1} | S_t = s), \forall t \in \mathbb{Z}$$~~

$$V_n(s) = E_x(b_t | S_t = s) = E_x(R_{t+1} + \gamma V_n(S_{t+1}) | S_t = s)$$

$$= \sum_{a \in A} \pi(a|s) E_x(R_{t+1} + \gamma V_n(S_{t+1}) | S_t = s, A_t = a)$$

$$= \sum_{a \in A} \pi(a|s) \sum_{j=1}^{\infty} \gamma^{j-1} E_x(R_{t+j} | S_t = s, A_t = a)$$

~~$$= \sum_{a \in A} \pi(a|s) P(S_{t+1} = s')$$~~

~~$$E_x(R_{t+1} | S_t = s) \quad \forall t \in \mathbb{Z}$$~~

~~$$\textcircled{1} E_x(R_t | S_t) \quad \forall t \in \mathbb{Z}$$~~

~~$$\textcircled{2} \text{ 证明 } E_x(R_{t+1} | S_t = s) = E_x(R_{t+1} | S_{t+1} = s')$$~~

~~$$E_x(R_{t+1} | S_t = s) = E_x(R_{t+1} | S_{t+1} = s')$$~~

~~$$E_x = E_x(R_{t+1} | S_{t+1} = s') \quad (j \leq i)$$~~

~~$$E_x(R_{t+1} | S_t = s) = \sum_{a \in A} \pi(a|s) E_x(R_{t+1} | S_t = s, A_t = a)$$~~

~~$$= \sum_{s' \in \mathcal{S}} P(S_{t+1} = s' | S_t = s) E_x(R_{t+1} | S_{t+1} = s')$$~~

$$R_{t+1} | S_{t+1} = s' = \sum_{s' \in \mathcal{S}} P(S_{t+1} = s' | S_{t+1} = s') \cdot E_x(R_{t+1} | S_{t+1} = s')$$

∴ 证毕, 且 $j=0$

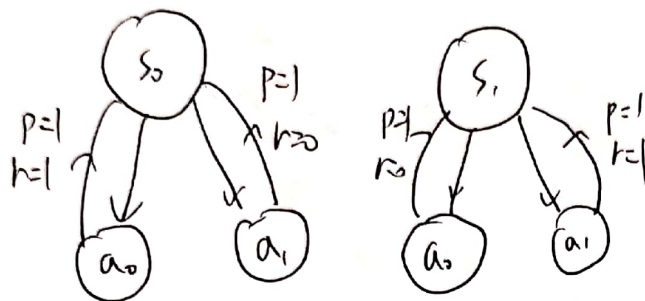
$$R_{t+1} | S_t = s$$

$$\text{且 } V_n(s) = E_x(b_t | S_t = s) \quad \forall t \in \mathbb{Z}$$



扫描全能王 创建

(2)



$$\begin{aligned} \pi_0: \pi_0(a_0|s_0) &= 1 \\ \pi_0(a_1|s_0) &= 0 \\ \pi_0(a_0|s_1) &= 1 \\ \pi_0(a_1|s_1) &= 0 \end{aligned}$$

$$\begin{aligned} \pi_1: \pi_1(a_0|s_0) &= 0 \\ \pi_1(a_1|s_0) &= 1 \\ \pi_1(a_0|s_1) &= 0 \\ \pi_1(a_1|s_1) &= 1 \end{aligned}$$

$$\begin{aligned} V_{\pi_0}(s_0) &= 1 + r + \dots + r^n \\ &= \lim_{n \rightarrow \infty} 1 \cdot \frac{1-r^{n+1}}{1-r} = \frac{1}{1-r} \end{aligned}$$

$$\begin{aligned} V_{\pi_1}(s_0) &= 0 \\ V_{\pi_1}(s_1) &= 1 + r + \dots + r^n \\ &= \lim_{n \rightarrow \infty} 1 \cdot \frac{1-r^{n+1}}{1-r} = \frac{1}{1-r} \end{aligned}$$

$$V_{\pi_0}(s_1) = 0$$

$$\begin{aligned} V_{\pi_0}(s_0) &> V_{\pi_1}(s_0) \\ V_{\pi_0}(s_1) &< V_{\pi_1}(s_1) \end{aligned}$$

不相等

(3)

证明:

$$V_{\pi}(s)$$

$$= \max_a \left(r_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s') \right)$$

由贝尔曼方程



扫描全能王 创建

$$(3) \textcircled{1} \forall s \in S, \mathbb{E}_{a \sim \pi(s)} [Q_\pi(s, a)] \geq \mathbb{E}_{a \sim \pi'(s)} [Q_\pi(s, a)] \Rightarrow (\forall s \in S, V_\pi(s) \geq V_{\pi'}(s))$$

$$V_\pi(s) = \mathbb{E}_{a \sim \pi(s)} [Q_\pi(s, a)] \leq \mathbb{E}_{a \sim \pi'(s)} [Q_\pi(s, a)]$$

$$= \mathbb{E}_{a \sim \pi'(s)} [r(s, a) + \gamma V_\pi(s_1) | s=s]$$

$$= \mathbb{E}_{a \sim \pi'(s)} [r(s, a) + \gamma \mathbb{E}_{a_1 \sim \pi(s_1)} [Q_\pi(s_1, a_1)] | s=s]$$

$$\leq \mathbb{E}_{a \sim \pi'(s)} [r(s, a)] + \gamma \mathbb{E}_{a_1 \sim \pi'(s_1)} [Q_\pi(s_1, a_1) | s=s]$$

$$= \mathbb{E}_{a \sim \pi'(s)} [r(s, a) + \gamma r(s, a_1) + \gamma^2 V_\pi(s_2) | s=s]$$

⊗ 归纳法, 对 $\forall T \geq 1$,

$$V_\pi(s) \leq \mathbb{E}_{a \sim \pi'(s)} \left[\sum_{t=0}^{T-1} \gamma^t \mathbb{E}_{a_t \sim \pi(s_t)} [r(s_t, a_t)] + \gamma^T V_\pi(s_{T+1}) \right]$$

$$= V_{\pi'}(s)$$

② 证明 π^* 为 optimal:

π^* : 对 $\forall (s, a) \in S \times A$, 有 $a \in \arg \max_{a' \in A} Q_\pi(s, a')$

证明: 若 $\exists (s_0, a_0) \notin \arg \max_{a' \in A} Q_\pi(s_0, a')$, 则 π 不是最优

构造 π' :

$$\begin{cases} \pi'(s) = \pi(s) & s \neq s_0 \\ \pi'(s_0) \in \arg \max_{a' \in A} Q_\pi(s_0, a') \end{cases}$$

由 $\pi'(s_0) \in \arg \max_{a' \in A} Q_\pi(s_0, a')$ 可知 π' 比 π 更优, 即 π 不是最优

$$\text{故对 } \forall s, \mathbb{E}_{a \sim \pi'(s)} [Q_\pi(s, a)] \geq \mathbb{E}_{a \sim \pi(s)} [Q_\pi(s, a)],$$

$$\text{故对 } \forall s \in S, V_{\pi'}(s) \geq V_\pi(s)$$

π' 为最优

证毕

同时, 对于一个非最优策略 π , 可以用上述构造性 π' 使 $V_{\pi'}(s) > V_\pi(s)$

证 (2) $\forall s \in S, \pi^*(s) = \arg \max_{a \in A} Q^*(s, a)$ 为最优

