

强化学习 HW1

18300290007 加兴华

1. (5 pts) Show that the infinite horizon discounted state value  $v_\pi(s)$  has the following alternative expression:

$$v_\pi(s) = \mathbb{E}_{N \sim \text{Geo}(1-\gamma)} \left[ \mathbb{E} \left[ \sum_{t=0}^{N-1} \gamma^t r(s_t, a_t, s_{t+1}) \mid s_0 = s \right] \right],$$

where  $\text{Geo}(1-\gamma)$  denotes the geometric distribution with parameter  $1-\gamma$ . In word, we can rewrite  $v_\pi(s)$  into an undiscounted form where the length of trajectory obeys the geometric distribution. In addition, compute  $\mathbb{E}[N]$  which is referred to as planning horizon.

$$\text{即 } r_t = r(s_t, a_t, s_{t+1})$$

$$V_\pi(s) = \mathbb{E}_{\substack{a_t \sim \pi(\cdot | s_t) \\ s_{t+1} \sim p(\cdot | s_t, a_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

$$= \mathbb{E} \left[ \sum_{t=0}^{\infty} \left( \frac{\gamma^t - \gamma^\infty}{1-\gamma} \right) (1-\gamma) r_t \mid s_0 = s \right]$$

$$= \mathbb{E} \left[ \sum_{t=0}^{\infty} \left( \sum_{k=t}^{\infty} \gamma^k \right) (1-\gamma) r_t \mid s_0 = s \right]$$

$$= \sum_{t=0}^{\infty} \sum_{k=t}^{\infty} \mathbb{E} \left[ \gamma^k (1-\gamma) r_t \mid s_0 = s \right]$$

$$= \sum_{k=0}^{\infty} \sum_{t=0}^k \mathbb{E} \left[ \gamma^k (1-\gamma) r_t \mid s_0 = s \right]$$

$$= \sum_{k=0}^{\infty} \gamma^k (1-\gamma) \mathbb{E} \left[ \sum_{t=0}^k r_t \mid s_0 = s \right] =: \mathbb{E}_{N \sim \text{Geo}(1-\gamma)} \mathbb{E} \left[ \sum_{t=0}^{N-1} r_t \mid s_0 = s \right]$$

$$\mathbb{E}N = \sum_{k=1}^{\infty} k \cdot \gamma^{k-1} (1-\gamma)$$

$$= \frac{1}{\gamma} \sum_{k=1}^{\infty} (k+1) \cdot \gamma^k (1-\gamma) - \sum_{k=1}^{\infty} \gamma^{k-1} (1-\gamma)$$

$$= \frac{1}{\gamma} (\mathbb{E}N - (1-\gamma)) - 1$$

$$= \frac{1}{\gamma} \mathbb{E}N - \frac{1}{\gamma}$$

$$= \frac{1}{1-\gamma}$$

2. (10 pts) Prove the Bellman equation for state value and action value functions (i.e., Theorem 1 in the latest version of Lecture1.pdf).

$$V_\pi(s) = \mathbb{E}_{\substack{a_t \sim \pi(\cdot | s_t) \\ s_{t+1} \sim p(\cdot | s_t, a_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

$$= \mathbb{E} \left[ r_0 + \gamma \cdot \sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_0 = s \right]$$

$$= E \left[ r_0 + \gamma \cdot E \left[ E \left( \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s' \right) \right] \mid s_0 = s \right]$$

$$= E \left[ r_0 + \gamma \cdot E V_{\pi}(s') \mid s_0 = s \right]$$

$$= E \left[ r_0 + \gamma \cdot V_{\pi}(s') \mid s_0 = s \right] \quad \text{由于 } V_{\pi} \text{ 为常数}$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) (r(s, a, s') + \gamma \cdot V_{\pi}(s'))$$

$$q_{\pi}(s, a) = E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]$$

$$= E \left[ r_0 + \gamma E \left[ E \left( \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s' \right) \right] \mid s_0 = s, a_0 = a \right]$$

$$= E \left[ r_0 + \gamma \cdot V_{\pi}(s') \mid s_0 = s, a_0 = a \right]$$

$$= \sum_{s'} P(s'|s, a) (r(s, a, s') + \gamma \cdot V_{\pi}(s'))$$

3. (5 pts) Show that the Bellman operator

$$T_{\pi} v = r_{\pi} + \gamma P^{\pi} v$$

is a contraction with respect to infinity norm.

$$\| T_{\pi} v_1 - T_{\pi} v_2 \|_{\infty} = \gamma \| P^{\pi} (v_1 - v_2) \|_{\infty}, \quad \forall v_1, v_2.$$

考虑  $\| P^{\pi} x \|_{\infty}$ :

$\because P^{\pi}$  为转移阵且行和 = 1.

$$\therefore \| P^{\pi} x \|_{\infty} = \max_i \left| \sum_j P_{ij}^{\pi} x_j \right|$$

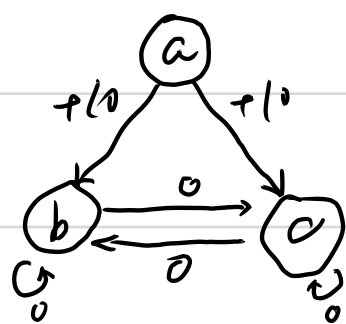
$$\leq \max_i \left| \sum_j P_{ij}^{\pi} |x_j| \right|$$

$$= \max_i \sum_j P_{ij}^{\pi} |x_j|$$

$$\leq \max_j |x_j| = \|x\|_{\infty} \quad \text{由于 } \max_j |x_j| = \left( \sum_j P_{ij}^{\pi} \right) \cdot \max_j |x_j| \geq \sum_j P_{ij}^{\pi} |x_j|$$

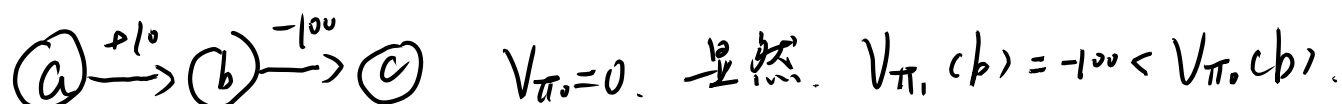
$\therefore T_{\pi}$  为收缩算子

4. (5 pts) Whether the optimal policy is unique? Prove or disprove by a counter example.



如图. 对 a, b, c 而言, 任何策略的一步回报都相同.  
因此任何策略都最优, 不唯一.

5. (optional) Let  $\pi_k$  be the policy extracted from the  $k$ -th iteration of the value iteration. Is it always that  $v_{\pi_{k+1}}(s) \geq v_{\pi_k}(s)$ ,  $\forall s$ ? Prove or disprove by a counter example.



6. (5 pts) Prove the policy improvement result (i.e., Theorem 2 in the latest version of Lecture2.pdf).

$$\begin{aligned}
 q_{\pi}(s, \pi'(s)) &= \sum_{s'} P(s'|s, \pi'(s)) (r(s, \pi'(s), s') + \gamma V_{\pi}(s')) \\
 &= \max_a \sum_{s'} P(s'|s, a) (r(s, a, s') + \gamma V_{\pi}(s')) \\
 &=: \max_a f(a) =: f(a^*) \\
 &= \left( \sum_a \pi(a|s) \right) \cdot f(a^*) \quad \text{由于 } \sum \pi(a|s) = 1 \\
 &\geq \sum_a \pi(a|s) \cdot f(a) \quad \text{由于 } f(a) \geq 0, \forall a \\
 &= V_{\pi}(s).
 \end{aligned}$$

得证.