

强化学习与有监督学习的差别在于强化学习的loss函数不可导

## 16.7

假设  $T = 3$ ，则有随机变量  $a_1, a_2, a_3$  表示每一步采取的动作， $x_1, x_2, x_3$  表示采取对应动作进入的状态。 $a'_1, a'_2, a'_3, x'_1, x'_2, x'_3$  表示对应的具体取值。为了统一标记，把原公式里的  $x$  记为  $x'_0$

$$\begin{aligned}
 V_3^\pi(x'_0) &= \mathbb{E}_\pi \left[ \frac{1}{3} \sum_{t=1}^3 r_t \mid x_0 = x'_0 \right] \\
 &= \sum_{\substack{a'_1, a'_2, a'_3 \in A, \\ x'_1, x'_2, x'_3 \in X}} [p(a'_1, a'_2, a'_3, x'_1, x'_2, x'_3) \frac{1}{3} \sum_{t=1}^3 R_{x'_{t-1} \rightarrow x'_t}^{a'_t}] \\
 &= \sum_{\substack{a'_1, a'_2, a'_3 \in A, \\ x'_1, x'_2, x'_3 \in X}} [\pi(x'_0, a'_1) P_{x'_0 \rightarrow x'_1}^{a'_1} \pi(x'_1, a'_2) P_{x'_1 \rightarrow x'_2}^{a'_2} \pi(x'_2, a'_3) P_{x'_2 \rightarrow x'_3}^{a'_3} \frac{1}{3} (R_{x'_0 \rightarrow x'_1}^{a'_1} + R_{x'_1 \rightarrow x'_2}^{a'_2} + R_{x'_2 \rightarrow x'_3}^{a'_3})] \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{\substack{a'_2, a'_3 \in A, \\ x'_1, x'_2, x'_3 \in X}} [P_{x'_0 \rightarrow x'_1}^{a'_1} \pi(x'_1, a'_2) P_{x'_1 \rightarrow x'_2}^{a'_2} \pi(x'_2, a'_3) P_{x'_2 \rightarrow x'_3}^{a'_3} \frac{1}{3} (R_{x'_0 \rightarrow x'_1}^{a'_1} + R_{x'_1 \rightarrow x'_2}^{a'_2} + R_{x'_2 \rightarrow x'_3}^{a'_3})] \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} (R_{x'_0 \rightarrow x'_1}^{a'_1} + R_{x'_1 \rightarrow x'_2}^{a'_2} + R_{x'_2 \rightarrow x'_3}^{a'_3})] \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_0 \rightarrow x'_1}^{a'_1}] + \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_1 \rightarrow x'_2}^{a'_2}] + \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_2 \rightarrow x'_3}^{a'_3}] \right\} \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \frac{1}{3} R_{x'_0 \rightarrow x'_1}^{a'_1} + \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_1 \rightarrow x'_2}^{a'_2}] + \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_2 \rightarrow x'_3}^{a'_3}] \right\} \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \frac{1}{3} R_{x'_0 \rightarrow x'_1}^{a'_1} + \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_1 \rightarrow x'_2}^{a'_2}] + \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{3} R_{x'_2 \rightarrow x'_3}^{a'_3}] \right\} \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \frac{1}{3} R_{x'_0 \rightarrow x'_1}^{a'_1} + \frac{2}{3} \sum_{\substack{a'_2, a'_3 \in A, \\ x'_2, x'_3 \in X}} [p(a'_2, a'_3, x'_2, x'_3) \frac{1}{2} (R_{x'_1 \rightarrow x'_2}^{a'_2} + R_{x'_2 \rightarrow x'_3}^{a'_3})] \right\} \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left( \frac{1}{3} R_{x'_0 \rightarrow x'_1}^{a'_1} + \frac{2}{3} \mathbb{E}_\pi \left[ \frac{1}{2} \sum_{t=1}^2 r_t \mid x_0 = x'_1 \right] \right) \\
 &= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left( \frac{1}{3} R_{x'_0 \rightarrow x'_1}^{a'_1} + \frac{2}{3} V_2^\pi(x'_1) \right)
 \end{aligned}$$

推广一下就成为书里的公式：

$$\begin{aligned}
V_T^\pi(x'_0) &= \mathbb{E}_\pi \left[ \frac{1}{T} \sum_{t=1}^T r_t \mid x_0 = x'_0 \right] \\
&= \sum_{\substack{a'_1, \dots, a'_T \in A, \\ x'_1, \dots, x'_T \in X}} [p(a'_1, \dots, a'_T, x'_1, \dots, x'_T) \frac{1}{T} \sum_{t=1}^T R_{x'_{t-1} \rightarrow x'_t}^{a'_t}] \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \sum_{\substack{a'_2, \dots, a'_T \in A, \\ x'_2, \dots, x'_T \in X}} [p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) \frac{1}{T} \sum_{t=1}^T R_{x'_{t-1} \rightarrow x'_t}^{a'_t}] \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \sum_{\substack{a'_2, \dots, a'_T \in A, \\ x'_2, \dots, x'_T \in X}} [p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) \frac{1}{T} R_{x'_0 \rightarrow x'_1}^{a'_1} + p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) \frac{1}{T} \sum_{t=2}^T R_{x'_{t-1} \rightarrow x'_t}^{a'_t}] \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \frac{1}{T} R_{x'_0 \rightarrow x'_1}^{a'_1} + \sum_{\substack{a'_2, \dots, a'_T \in A, \\ x'_2, \dots, x'_T \in X}} [p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) \frac{1}{T} \sum_{t=2}^T R_{x'_{t-1} \rightarrow x'_t}^{a'_t}] \right\} \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \frac{1}{T} R_{x'_0 \rightarrow x'_1}^{a'_1} + \frac{T-1}{T} \sum_{\substack{a'_2, \dots, a'_T \in A, \\ x'_2, \dots, x'_T \in X}} [p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) \frac{1}{T-1} \sum_{t=2}^T R_{x'_{t-1} \rightarrow x'_t}^{a'_t}] \right\} \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \left\{ \frac{1}{T} R_{x'_0 \rightarrow x'_1}^{a'_1} + \frac{T-1}{T} V_{T-1}^\pi(x'_1) \right\}
\end{aligned}$$

## 16.8

$$\begin{aligned}
V_\gamma^\pi(x'_0) &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid x_0 = x'_0 \right] \\
&= \mathbb{E}_\pi \left[ r_1 + \sum_{t=1}^{\infty} \gamma^t r_{t+1} \mid x_0 = x'_0 \right] \\
&= \mathbb{E}_\pi \left[ r_1 + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} r_{t+1} \mid x_0 = x'_0 \right] \\
&= \sum_{\substack{a'_1, \dots, a'_T \in A, \\ x'_1, \dots, x'_T \in X}} [p(a'_1, \dots, a'_T, x'_1, \dots, x'_T) (R_{x'_0 \rightarrow x'_1}^{a'_1} + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R_{x'_t \rightarrow x'_{t+1}}^{a'_t})] \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} \sum_{\substack{a'_2, \dots, a'_T \in A, \\ x'_2, \dots, x'_T \in X}} [p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) (R_{x'_0 \rightarrow x'_1}^{a'_1} + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R_{x'_t \rightarrow x'_{t+1}}^{a'_t})] \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} (R_{x'_0 \rightarrow x'_1}^{a'_1} + \sum_{\substack{a'_2, \dots, a'_T \in A, \\ x'_2, \dots, x'_T \in X}} [p(a'_2, \dots, a'_T, x'_2, \dots, x'_T) (\gamma \sum_{t=1}^{\infty} \gamma^{t-1} R_{x'_t \rightarrow x'_{t+1}}^{a'_t})]) \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} (R_{x'_0 \rightarrow x'_1}^{a'_1} + \gamma \mathbb{E}_\pi \left[ \sum_{t=1}^{\infty} \gamma^t r_{t+1} \mid x_0 = x'_1 \right]) \\
&= \sum_{a'_1 \in A} \pi(x'_0, a'_1) \sum_{x'_1 \in X} P_{x'_0 \rightarrow x'_1}^{a'_1} (R_{x'_0 \rightarrow x'_1}^{a'_1} + \gamma V_\gamma^\pi(x'_1))
\end{aligned}$$

## 16.29

把  $\frac{1}{t+1}$  替代为  $\alpha$

$$\begin{aligned}
Q_{t+1}^\pi(x, a) &= Q_t^\pi(x, a) + \alpha(r_{t+1} - Q_t^\pi(x, a)) \\
&= (1 - \alpha)Q_t^\pi(x, a) + \alpha r_{t+1} \\
&= (1 - \alpha)((1 - \alpha)Q_{t-1}^\pi(x, a) + \alpha r_t) + \alpha r_{t+1} \\
&= (1 - \alpha)^2 Q_{t-1}^\pi(x, a) + \alpha((1 - \alpha)r_t + r_{t+1}) \\
&= \alpha(r_{t+1} + (1 - \alpha)r_t + (1 - \alpha)^2 r_{t-1} + (1 - \alpha)^3 r_{t-2} + \dots + (1 - \alpha)^t r_1)
\end{aligned}$$

每一步奖赏的系数之和：

$$S = \alpha \frac{1 - (1 - \alpha)^{t+1}}{1 - (1 - \alpha)} = \alpha \frac{1 - (1 - \alpha)^{t+1}}{\alpha} = 1 - (1 - \alpha)^{t+1} \rightarrow 1$$