

强化学习

值迭代

价值函数，给定当前状态 s ，(当前行动 a)，和策略 π

$$V^\pi(s) = E_{z \sim \pi} [R(z) | s]$$

$$Q^\pi(s, a) = E_{z \sim \pi} [R(z) | s, a]$$

Bellman 方程

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

$$= E_{s' \sim P(\cdot | s, a)} \max_{\pi} \left[r_0 + \sum_{t=1}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

给定 s' ，选 a' ，使

$r_0 + \sum_{t=1}^{\infty} \gamma^t r_t$ 最大

$$\Rightarrow a' = \underset{a}{\operatorname{argmax}} Q^*(s', a)$$

$$= E_{s' \sim P(\cdot | s, a)} \left[r_0(s', s, a) + \gamma \max_{a'} Q^*(s', a') \right]$$

深度值迭代

$$\nabla_{\theta} L(\theta) = \nabla_{\theta} E_{s, a \sim p(\cdot)} \left[(y - Q(s, a; \theta))^2 \right]$$

$$= E_{s, a \sim p(\cdot)} \left[-2(y - Q(s, a; \theta)) \nabla_{\theta} Q(s, a; \theta) \right]$$