

example 1: Policy iteration
 Policy iteration: Q value iteration

① V -value iteration: Q -value iteration

در حالت V -value iteration، ما به دنبال $V(s)$ هستیم که $R(s', a, s)$ و $P(s', a, s)$ را شامل می شود. $V(s)$ به عنوان $V_t(s)$ شروع می شود و به $V_{t+1}(s)$ به روز می شود.

$$V(s) = \max_{a \in A} \sum_{s' \in S} P(s' | a, s) (R(s', a, s) + \gamma V(s'))$$

برای $V_t(s)$ به $V_{t+1}(s)$ به روز می شود.

$$V_{t+1}(s) = \max_{a \in A} \sum_{s' \in S} P(s' | a, s) (R(s', a, s) + \gamma V_t(s'))$$

② در حالت Q -iteration، ما به دنبال $Q(s, a)$ هستیم که $R(s', a, s)$ و $P(s', a, s)$ را شامل می شود. $Q(s, a)$ به عنوان $Q_t(s, a)$ شروع می شود و به $Q_{t+1}(s, a)$ به روز می شود.

$$V(s) = \max_{a \in A} Q(s, a) \rightarrow \sum_{s' \in S} P(s' | a, s) (R(s', a, s) + \gamma \max_{a' \in A} Q(s', a'))$$

$Q(s, a) = \dots$ time component $\rightarrow V$ -iteration

$$Q_{t+1}(s, a) = \sum_{s' \in S} P(s' | a, s) (R(s', a, s) + \gamma \max_{a' \in A} Q_t(s', a'))$$

در حالت Q -iteration، ما به دنبال $Q(s, a)$ هستیم که $R(s', a, s)$ و $P(s', a, s)$ را شامل می شود. $Q(s, a)$ به عنوان $Q_t(s, a)$ شروع می شود و به $Q_{t+1}(s, a)$ به روز می شود.

$$Q(s_{t+1}, a_{t+1}) = (1 - \alpha) Q_{old}(s_t, a_t) + \alpha (R(s_{t+1}, a_{t+1}, s_t) + \gamma \max_{a' \in A} Q(s_{t+1}, a'))$$

$$Q_{t+1}(s_{t+1}, a_{t+1}) = (1 - \alpha) Q_{old}(s_t, a_t) + \alpha (R(s_{t+1}, a_{t+1}, s_t) + \gamma \max_{a' \in A} Q(s_{t+1}, a'))$$

$$= \sum_{s' \in S} P(s' | a_{t+1}, s_t) (R(s', a_{t+1}, s_t) + \gamma \max_{a' \in A} Q(s', a'))$$

سوال چرا greedy؟ مشکلات؟

greedy بهتر است زیرا به جست و جوی ادامه می دهد تا مسائل رسیدن به حالت بهینه را افزایش دهد

در هر مقدار E اولی که برابر $\frac{1}{2}$ شود $\frac{1}{2}$ تکرار می شود و این فرآیند بهینه است

اگر در این $\frac{1}{2}$ $\frac{1}{2}$ باشد و به جست و جوی بهینه می آید و این الگوریتم بهینه است

این امر بهینه برای مثال الگوریتم $\frac{1}{2}$ در هر مقدار اول مقدار واقعی هر $\frac{1}{2}$ را می دهد و در هر تکرار $\frac{1}{2}$ را می دهد

هر زمان که $\frac{1}{2}$ را می دهد و احتمال $\frac{1}{2}$ است و از $\frac{1}{2}$ است و این روش بهینه است تا به اندازه $\frac{1}{2}$ می رسد

می بینیم که بهینه می شود و $\frac{1}{2}$ را می دهد و $\frac{1}{2}$ را می دهد و از راه حل های آن است

در صورتی که از این الگوریتم استفاده کنیم می توانیم بهینه را پیدا کنیم و بهینه را پیدا کنیم و بهینه را پیدا کنیم

بهینه

$$\forall s: E_{a \sim \pi'} [Q^\pi(s, a)] \geq E_{a \sim \pi} [Q^\pi(s, a)] \Rightarrow$$

$$\forall s: V^{\pi'}(s) \geq V^\pi(s)$$

$$V^\pi = Q^\pi(s, \pi(s)) = E_{a \sim \pi} [Q^\pi(s, a)] \leq Q^\pi(s, a')$$

$$\pi' = \underset{a \in A}{\operatorname{argmax}} Q^\pi(s, a)$$

$$V^{\pi'} - V^\pi = \sum_{i=0}^{\infty} (V^{\pi_{i+1}} - V^{\pi_i})$$

$$= \sum_{i=0}^{\infty} \gamma^i \sum_{s' \in S} P(s_{i+1} = s' | s_i = s, \pi') (Q^\pi(s', \pi'(s')) - Q^\pi(s', \pi(s)))$$

$$= \sum_{i=0}^{\infty} \gamma^i \sum_{s' \in S} P(s_{i+1} = s' | s_i = s, \pi') \underbrace{A^\pi(s', \pi')}_{\gamma \rightarrow \gamma^i \rightarrow \gamma^i \rightarrow \gamma^i} \rightarrow \frac{1}{1-\gamma} A^{\pi'}(s')$$

$$V^{\pi'} - V^\pi \geq 0 \rightarrow V^{\pi'} \geq V^\pi$$