

► Soft action value $q_{\pi}^i(s, a)$: $[a_i]$ is chosen, then entropy equal to $\tau \cdot \pi(a_i | s) = 1 \cdot \pi(a_i = a_i)$

$$q_{\pi}^i(s, a) = E_{\pi} \left[r(s_0, a_0, s_1) + \sum_{t=1}^{\infty} \gamma^t r_t(s_t, a_t, s_{t+1}) | s_0 = s, a_0 = a \right]$$

For any $v \in \mathbb{R}^S$, the soft Bellman optimality operator T^{π} is defined by

$$[T^{\pi} v](s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(r(s, a, s') + \gamma \sum_{s'} P(s'|s, a) [r(s, a, s') + \tau v(s')] \right)$$

where the maximum value is attained iff $\pi(s) \propto \exp(q_{\pi}(s, \cdot) / \tau)$

► T^{π} is γ -contraction with respect to ℓ_{∞} -norm.

$$V_{\pi}^i(s) = E_{\pi} [\sum_t \gamma^t r_{t,i} | s_0 = s] \quad q_{\pi}^i(s, a) = E_{\pi} [r_0 + \sum_t \gamma^t r_{t,i} | s_0 = s, a_0 = a]$$

$$V_{\pi}^i(s) = E_{\pi} [q_{\pi}^i(s, a) - \tau \log \pi(a|s)]$$

$$= E_{\pi} [q_{\pi}^i(s, a)] - H(\pi(\cdot | s))$$

< 每遇到一个随机变量 a 就惩罚一次 >

$$l^i: [T^i v](s) = \tau \log(\|\exp(q(s, \cdot) / \tau)\|_1)$$

$$\arg \max (\dots) = \frac{1}{\tau} \cdot \exp(q(s, \cdot) / \tau)$$

① $\sum \pi(a|s) [q(s, a) - \tau \log \pi(a|s)]$ 对 $\pi(a|s)$ 求导:

$$\Rightarrow q(s, a) - \tau \log \pi(a|s) - \pi(a|s) \cdot \tau \cdot \frac{1}{\tau \pi(a|s)} = 0$$

$$\Rightarrow \frac{q(s, a)}{\tau} = \log \pi(a|s) + 1 \quad ? \quad \checkmark \quad \pi(a|s) = \exp\left[\frac{q(s, a)}{\tau}\right] \cdot e$$

target $\pi(\cdot | s) = \frac{1}{Z} \cdot \exp\left[\frac{q(s, \cdot)}{\tau}\right]$ & $\sum \pi(a|s) = 1$ 因式

$$[T^i v](s) = \sum_a \frac{1}{Z} \exp\left[\frac{q(s, a)}{\tau}\right] [q(s, a) - q(s, a) + \tau \log C]$$

$$= \frac{1}{Z} \cdot \sum_a \exp\left[\frac{q(s, a)}{\tau}\right] \cdot \tau \log C$$

$$Z(s) \leftarrow C = \sum_a \exp\left[\frac{q(s, a)}{\tau}\right]$$

正则化

$$\text{或者: } E_{\pi} [q(s, a) - \tau \log \pi(a|s)] = \sum_a \pi(a|s) [\tau \log \pi(a|s) - \log \exp\left[\frac{q(s, a)}{\tau}\right]]$$

$$\downarrow$$
$$\max \tau = \max \left[\tau \log \pi(a|s) - \log \exp\left[\frac{q(s, a)}{\tau}\right] \right]$$
$$\Leftrightarrow \min KL\left(\pi(a|s) \parallel \frac{\exp\left[\frac{q(s, a)}{\tau}\right]}{Z(s)}\right) \quad \log \frac{\exp\left[\frac{q(s, a)}{\tau}\right]}{Z(s)} + \log Z(s)$$

2: T^i 是压缩算子.

先 T_{π}^i 是压缩算子也. $r(s, a, s') - \tau \log \pi(a|s)$ 由 $JH(z) = E_{\pi} [\log \pi(a|s)]$

$$[T_{\pi}^i v](s) = E_{\pi} [r(s, a, s') + \gamma V_{\pi}(s')]$$

$$[T_{\pi}^i v_1 - T_{\pi}^i v_2](s) = E_{\pi} E_{s'} [\gamma V_{\pi}(s') - \gamma V_{\pi}(s')]$$

$$\left\{ \begin{array}{l} \|T_{\pi}^i v_1 - T_{\pi}^i v_2\|_{\infty} = \gamma \max E_{\pi} E_{s'} [V_{\pi}(s') - V_{\pi}(s')] \leq \gamma \max E_{\pi} E_{s'} |V_{\pi}(s') - V_{\pi}(s')| \quad \textcircled{A} \\ \gamma \|v_1 - v_2\|_{\infty} = \gamma \max |v_1(s) - v_2(s)| \quad \textcircled{B} \end{array} \right.$$

$$\text{假设 } s_k = \arg \max_s |v_1(s) - v_2(s)|$$

$$\textcircled{A}: \gamma \cdot |v_1(s_k) - v_2(s_k)|$$

$$\textcircled{B}: E_{s'} |v_1(s') - v_2(s')| < |v_1(s_k) - v_2(s_k)|$$

$$E_{\pi} [\dots] < E_{\pi} |v_1(s_k) - v_2(s_k)| =$$

$$\max_s |v_1(s) - v_2(s)| = \text{Done.}$$

$$\text{再解 } T^i: [T^i v](s) = \max_a E_{\pi} E_{s'} [r_0 + \gamma V(s')]$$

$$[T^i v_1 - T^i v_2](s) \leq \max_a E_{\pi} E_{s'} [\gamma V_1(s') - \gamma V_2(s')] \xrightarrow{\max(A_1) - \max(A_2) \leq \max(A_1 - A_2)}$$

同上可证 V .

- Soft Policy Evaluation ①
- Soft Policy Improvement ②
- Soft Policy Iteration ③

$$\textcircled{1} \quad q_{\pi}^i = T_{\pi}^i q_{\pi}^i \quad (\text{证明有解}) \quad T_{\pi}^i \text{ 压缩算子}$$

$$(网络站) \quad q_{\pi}^i(s, a) = E_{\pi} [r(s, a, s') + \gamma V_{\pi}(s')]$$

$$= E_{\pi} [r(s, a, s') + \gamma E_{\pi} [q_{\pi}^i(s', a') - \tau \log \pi(a'|s')]]$$

$$= E_{\pi} [r(s, a, s') + \gamma E_{\pi} [q_{\pi}^i(s', a') + \tau H(\pi)]]$$

常数下常数

\therefore 在 (Bellman Backup) 仍然收敛 \Rightarrow Q 收敛

经典RL中如何收敛? 其实和 r 收敛 r 不收敛的压缩算子证明

$$\textcircled{2} \quad \pi^i = \exp\left[\frac{q_{\pi}^i(s, \cdot)}{\tau}\right] / Z(s), \quad \text{证明有提升} \quad \text{st. } q_{\pi}^i(s, a) \leq q_{\pi}^i(s, a')$$

$$V_{\pi}^i(s) = E_{\pi} [q_{\pi}^i(s, a) - \log \pi(a|s)]$$

$$= -KL(\pi(\cdot | s) \parallel \exp\left[\frac{q_{\pi}^i(s, \cdot)}{\tau}\right] / Z) + \log(Z)$$

$$\leq E_{\pi} [q_{\pi}^i(s, a) - \log \pi^i(a|s)] = \log(Z)$$

$$= H(\pi^i) + E_{\pi} [q_{\pi}^i(s, a)]$$

$$q_{\pi}^i(s, a) = E_{\pi} [r_0 + \gamma V_{\pi}^i(s')]$$

$$\leq E_{\pi} [r_0 + \gamma H(\pi^i) + \gamma E_{\pi} [q_{\pi}^i(s, a)]]$$

$$= E_{\pi} [r_0 + \gamma H(\pi^i) + \gamma \cdot \tau + E_{\pi} [V_{\pi}^i(s')]]$$

...

$$\leq E_{\pi} [r_0 + \gamma H(\pi^i) + \gamma \tau + \gamma H(\pi^i)] = q_{\pi}^i(s, a)$$

$$\textcircled{3} \quad \begin{array}{l} \text{Policy Evaluation} \\ \text{Policy Improvement} \end{array} \quad \text{收敛到最优} \quad \text{Policy Iteration}$$

$$V_{\pi}^i(s) = \max_a V_{\pi}^i(s) = \tau \log(\|\exp\left[\frac{q_{\pi}^i(s, \cdot)}{\tau}\right]\|_1)$$

$$[T^i v](s) = \max [T_{\pi}^i v](s)$$

$$\left\{ \begin{array}{l} T_{\pi}^i v = v \quad \text{Policy Evaluation} \\ T^i v = v \quad \text{Optimal} \end{array} \right.$$

[经典RL]

[正则RL]
相同

