

Prove $V^*(s) = \max_{a \in A} Q^*(s, a)$

By definition, $\# \pi$

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) \cdot Q^\pi(s, a)$$

since $Q^*(s, a) \geq Q^\pi(s, a) \# \pi$.

(*)

$$V^*(s) \leq \sum_{a \in A} \pi(a|s) \cdot Q^*(s, a) \leq \max_{a \in A} Q^*(s, a).$$

Thus,

$$V^*(s) \leq \max_{a \in A} Q^*(s, a) \quad \text{--- } \textcircled{1}$$

Assume $V^*(s) < \max_{a \in A} Q^*(s, a)$, let $a^* = \arg \max_{a \in A} Q^*(s, a)$, then

$$Q^*(s, a^*) \geq \max_{a \in A} Q^*(s, a) > V^*(s)$$

since $V^*(s)$ is the optimallity for all state s .

$$V^*(s) \geq Q^*(s, a^*)$$

which cause a contradiction., thus $V^*(s) \neq \max_{a \in A} Q^*(s, a)$

hence $V^*(s) = \max_{a \in A} Q^*(s, a)$ QED.

$$V^*(s) = \max_{a \in A} \left[R_{s,a} + \gamma \sum_{s' \in S} P_{ss'}^a \cdot V^*(s') \right]$$

$$= \max_{a \in A} \left[R_{s,a} + \gamma \sum_{s' \in S} P_{ss'}^a \cdot \max_{a \in A} \right]$$

$$V^*(s) = \max_{a \in A} Q^*(s, a).$$

$$V^*(s) = \mathbb{E}[G(t)]$$

prove(a)

By definition, $\neq \pi$

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) \cdot Q^\pi(s, a).$$

Since $Q^*(s, a) \geq Q^\pi(s, a)$. $\neq \pi$,

$$V^\pi(s) \leq \sum_{a \in A} \pi(a|s) \cdot Q^*(s, a) \leq \max_{a \in A} Q^*(s, a).$$

Hence,

$$V^\pi(s) \leq \max_{a \in A} Q^*(s, a)$$

prove(b) =

Now we only need to prove that $V^*(s) \geq \max Q^*(s, a)$.

We prove it by proving $V^*(s) < \max Q^*(s, a)$ by contradiction.

Let $a^* = \arg\max_{a \in A} Q^*(s, a)$, then

$$Q^*(s, a^*) \geq \max Q^*(s, a).$$

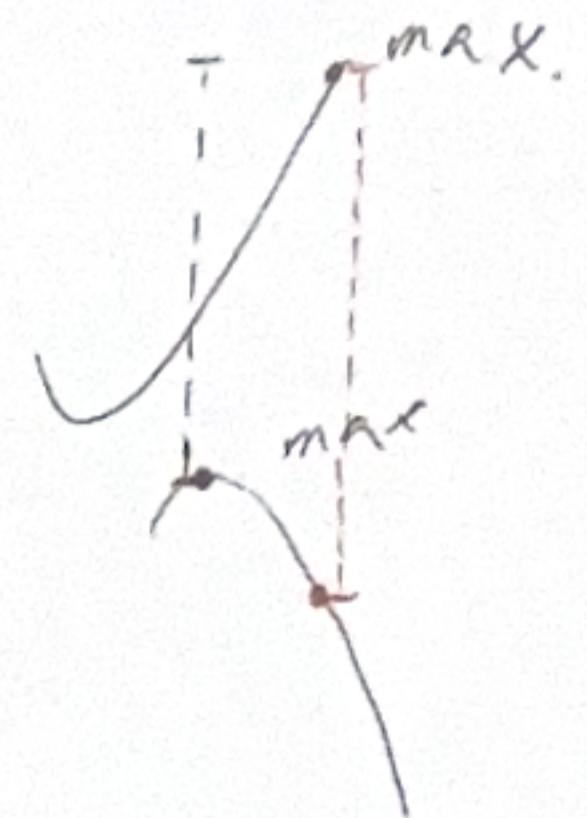
To show T^* is a γ -contraction operator, we simply prove.

$$\|T^*(Q) - T^*(Q')\|_\infty \leq \gamma \|Q - Q'\|_\infty, \forall Q, Q'$$

where

$$\cancel{\|Q - Q'\|} = \max$$

$$\|Q - Q'\|_\infty = \max_{s,a} \|Q(s,a) - Q'(s,a)\|$$



$$\cancel{\|T^*(Q) - T^*(Q')\|_\infty} =$$

$$\begin{aligned} [T^*(Q) - T^*(Q')] &= R_{s,a} + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q(s', a') - (\cancel{R_{s,a}}) \\ &\quad \cancel{R_{s,a}} + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q'(s', a') \\ &= \gamma \sum_{s'} P_{ss'}^a [\max_{a'} (Q(s', a')) - \max_{a'} (Q'(s', a'))] \\ &\leq \gamma \sum_{s'} P_{ss'}^a \max_{a'} (Q(s', a') - Q'(s', a')) \end{aligned}$$

Since for all (s', a') ,

$$Q(s', a') - Q'(s', a') \leq \|Q - Q'\|_\infty$$

Hence,

$$\begin{aligned} \gamma \sum_{s'} P_{ss'}^a \max_{a'} (Q(s', a') - Q'(s', a')) &\leq \gamma \|Q - Q'\|_\infty \sum_{s'} P_{ss'}^a \cancel{0} \\ &= \gamma \|Q - Q'\|_\infty \end{aligned}$$

Since for all $Q, Q' \ s, a$ pair.

$$[T^*(Q) - T^*(Q')] \leq \gamma \|Q - Q'\|_\infty$$

Hence

$$\max_{s,a} |T^*(Q) - T^*(Q')| = \|T^*(Q) - T^*(Q')\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

$$RHS = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_\theta}} \sum_a \pi_\theta(a|s) \cdot [f(s, a)]$$

$$= \frac{1}{1-\gamma} \sum_s d_{\mu}^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \cdot [f(s, a)]$$

$$= \frac{1}{1-\gamma} \sum_s (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | s_0 \sim \mu, \pi_\theta) \cdot \sum_a \pi_\theta(a|s) \cdot [f(s, a)]$$

$$= \cancel{\sum_{t=0}^{\infty} \sum}$$

$$= \sum_{t=0}^{\infty} \gamma^t \left[\sum_s \Pr(s_t = s | s_0 \sim \mu, \pi_\theta) \sum_a \pi_\theta(a|s) \cdot f(s, a) \right]$$

$$= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[f(s_t, a_t) | \text{being in state } s \text{ in time } t]$$

$$= \sum_{t=0}^{\infty} \gamma^t \left[\mathbb{E}_{s_t \sim p_{\mu}^{\pi_\theta}, a_t \sim \pi_\theta(\cdot|s_t)} [f(s_t, a_t)] \right]$$

$$= \mathbb{E}_{s \sim p_{\mu}^{\pi_\theta}} \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t)$$

\mathbb{E}_s

choosing
given state s , prob of a .

Not affect by t .

Probability of being in
state s in time t choosing
action a .