# 1.3a Computational Complexity

Notations:
① MDP: $M = (S, A, P, r, \gamma)$
② $L(P, r, \gamma)$ : total bit-size to specify M
   <span style="color:blue">↓</span>
   <span style="color:blue">specified with rational entries</span>

Def [strongly polynomial]

   An algorithm is said to be <span style="color:blue">strongly polynomial</span> if it return
   an optimal policy with a polynomial runtime independent with
   $$L(P, r, \gamma)$$

   First we consider classical iterative algorithms that compute $Q^*$.

1. Value iteration

   ① start at some $Q$
   ② iteratively apply $T$: $Q \leftarrow TQ$

   <span style="color:blue">Note: This is called as Q-value iteration.</span>

   To show this algorithm converges to $Q^*$, we have following statement.

---

Lem 1.10. [contraction] $T$ is a contraction mapping on $(\mathbb{R}^{|S||A|}, \|\cdot\|_\infty)$

---

   Pf: we need to show : for any two vectors $Q, Q' \in \mathbb{R}^{|S||A|}$,
   $$\|TQ - TQ'\|_\infty \leq \gamma \|Q - Q'\|_\infty .$$

   Assume $V_Q(s) > V_{Q'}(s)$ without loss of generality
   and let $a$ be the action s.t. $Q(s,a) = \max_{a' \in A} Q(s, a')$

   $$|V_Q(s) - V_{Q'}(s)| = Q(s,a) - \max_{a' \in A} Q'(s, a')$$

$$\leq Q(s,a) - Q'(s,a)$$

$$\leq \max_{a \in A} |Q(s,a) - Q'(s,a)|$$

$$\Rightarrow \|TQ - TQ'\|_\infty = \gamma \|PV_Q - PV_{Q'}\|_\infty$$

$$= \gamma \|P(V_Q - V_{Q'})\|_\infty$$

$$\leq \gamma \|V_Q - V_{Q'}\|_\infty$$

$$\leq \gamma \max_s |V_Q(s) - V_{Q'}(s)|$$

$$\leq \gamma \max_s \max_a |Q(s,a) - Q'(s,a)|$$

$$= \gamma \|Q - Q'\|_\infty$$

Since $\gamma \in [0,1)$, $T$ is a contraction mapping $\square$

---

Corollary: since $(\mathbb{R}^{|S||A|}, \|\cdot\|_\infty)$ is complete, using Lemma 1.10 we have that $T$ has a unique fixed point. According to Theorem 1.8, the fixed point can only be $Q^*$.

---

Lemma 1.11. (Q-Error Amplification) For any vector $Q \in \mathbb{R}^{|S||A|}$

$$V^{\pi_Q} \geq V^* - \frac{2\|Q - Q^*\|_\infty}{1 - \gamma} \mathbb{1}$$

where $\mathbb{1}$ denotes the vector of all ones.

Pf: Fix state $s$ and let $a = \pi_Q(s)$.

$$V^*(s) - V^{\pi_Q}(s) = Q^*(s, \pi^*(s)) - Q^{\pi_Q}(s, a)$$

$$= Q^*(s, \pi^*(s)) - Q^*(s, a)$$

$$+ Q^*(s, a) - Q^{\pi_\theta}(s, a)$$

$$= Q^*(s, \pi^*(s)) - Q^*(s, a) + \gamma E_{s' \sim p(s, a)}[V^*(s') - V^{\pi_\theta}(s')]$$

$$= Q^*(s, \pi^*(s)) - Q(s, \pi^*(s))$$

$$+ Q(s, \pi^*(s)) - Q^*(s, a)$$

$$+ \gamma E_{s' \sim p(s, a)}[V^*(s') - V^{\pi_Q}(s')]$$

$$\leq Q^*(s, \pi^*(s)) - Q(s, \pi^*(s))$$

$$+ Q(s, a) - Q^*(s, a)$$

$$+ \gamma E_{s' \sim p(s, a)}[V^*(s') - V^{\pi_Q}(s')]$$

$$\leq 2\|Q - Q^*\|_\infty + \gamma \|V^* - V^{\pi_Q}\|_\infty$$

$$\Rightarrow (1 - \gamma)\|V^* - V^{\pi_\theta}\|_\infty \leq 2\|Q - Q^*\|_\infty$$

$$\Rightarrow V^{\pi_Q} \geq V^* - \frac{2\|Q - Q^*\|_\infty}{1 - \gamma} \mathbb{1} \qquad \square$$

---

Theorem 1.12 [Q-value iteration convergence]

Set $Q^{(0)} = 0$. For $k = 0, 1, \cdots$. suppose:

$$Q^{(k+1)} = T Q^{(k)}$$

Let $\pi^{(k)} = \pi_{Q^{(k)}}$, For $k \geq \dfrac{\log \frac{2}{(1-\gamma)^2 \varepsilon}}{1 - \gamma}$,

$$V^{\pi^{(k)}} \geq V^* - \varepsilon \mathbb{1}.$$

Pf: since $Q^\pi(s, a) = E[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | S_0 = s, a_0 = a, \pi] \leq \frac{1}{1-\gamma} \quad \forall \pi$

$$\Rightarrow \|Q^*\|_\infty \leq \frac{1}{1 - \gamma}$$

$$\|Q^{(k)} - Q^*\|_\infty = \|T^k Q^{(0)} - T^k Q^*\|_\infty$$
$$\leq \gamma^k \|Q^{(0)} - Q^*\|_\infty \quad (\text{Lemma } 1.10)$$
$$= (1 - (1-\gamma))^k \|Q^*\|_\infty$$
$$\leq \frac{1}{1-\gamma} e^{-(1-\gamma)k} \quad (1 + x \leq e^x)$$

For $k \geq \dfrac{\log \frac{2}{(1-\gamma)^2 \varepsilon}}{1-\gamma} \implies -(1-\gamma)k \leq \log \dfrac{(1-\gamma)^2 \varepsilon}{2}$

By Lemma 1.11, 
$$V^{\pi(k)} \geq V^* - \frac{2\|Q^{(k)} - Q^*\|_\infty}{1-\gamma} \mathbb{1}$$
$$\geq V^* - \frac{2}{(1-\gamma)^2} e^{-(1-\gamma)k} \mathbb{1}$$
$$\geq V^* - \varepsilon \mathbb{1} \qquad \square$$

---

## 2. Policy Iteration

① start from an arbitrary policy $\pi_0$

② repeat following : for $k = 0, 1, 2, \cdots$

(i) Compute $Q^{\pi_k}$

(ii) Update policy $\pi_{k+1} = \pi_{Q^{\pi_k}}$

Lemma 1.13. ① $Q^{\pi_{k+1}} \geq T Q^{\pi_k} \geq Q^{\pi_k}$

② $\|Q^{\pi_{k+1}} - Q^*\|_\infty \leq \gamma \|Q^{\pi_k} - Q^*\|_\infty$

Pf: ① (i) first we show $TQ^{\pi_k} \geq Q^{\pi_k}$

note that $\pi_k$ is alway deterministic.

$$TQ^{\pi_k}(s,a) = r(s,a) + \gamma E_{s' \sim p(s,a)}\left[\max_{a' \in A} Q^{\pi_k}(s',a')\right]$$

$$\geq r(s,a) + \gamma E_{s' \sim p(s,a)}\left[Q^{\pi_k}(s', \pi_k(a'))\right]$$

$$= Q^{\pi_k}(s,a)$$

(ii) now we prove $Q^{\pi_{k+1}} \geq TQ^{\pi_k}$

$$Q^{\pi_k} = r + \gamma P^{\pi_k} Q^{\pi_k}$$

$$\Rightarrow Q^{\pi_k} \leq r + \gamma P^{\pi_{k+1}} Q^{\pi_k} \quad [\text{By def of } \pi_{k+1}]$$

By iterate the ineq., we have

$$Q^{\pi_k} \leq r + \gamma P^{\pi_{k+1}} (r + \gamma P^{\pi_{k+1}} Q^{\pi_k})$$

$$\leq \cdots$$

$$\leq \sum_{t=0}^{\infty} \gamma^t (P^{\pi_{k+1}})^t r \quad \cdots \text{ since } \gamma^t P^{\pi_{k+1}} Q^{\pi_k} \to 0$$

$$= Q^{\pi_{k+1}} \quad \cdots \text{ since } Q^{\pi_{k+1}} = r + \gamma P^{\pi_{k+1}} Q^{\pi_{k+1}}$$

$$\Rightarrow Q^{\pi_{k+1}}(s,a) = r(s,a) + \gamma E_{s' \sim p(s,a)}\left[Q^{\pi_{k+1}}(s', \pi_{k+1}(s'))\right]$$

$$\geq r(s,a) + \gamma E_{s' \sim p(s,a)}\left[Q^{\pi}(s', \pi_{k+1}(s'))\right]$$

$$= r(s,a) + \gamma E_{s' \sim p(s,a)}\left[\max_{a' \in A} Q^{\pi}(s',a')\right]$$

$$= TQ^{\pi},$$

which completes the proof of ①.

②

$$\|Q^* - Q^{\pi_{k+1}}\|_\infty \leq \|TQ^* - TQ^\pi\|_\infty \leq \gamma \|Q^* - Q^\pi\|_\infty \quad \square$$

$\underset{\text{Lem 10}}{\Uparrow}$

---

**Theorem 1.14 [Policy iteration convergence]**

Let $\pi_0$ be an initial policy. For $k \geq \dfrac{\log \frac{1}{(1-\gamma)\varepsilon}}{1-\gamma}$, the policy

iteration has its bound:

$$Q^{\pi_k} \geq Q^* - \varepsilon \mathbb{1}.$$

---

Pf:
$$Q^* - Q^{\pi_k} \leq \|Q^* - Q^{\pi_k}\|_\infty \mathbb{1}$$

$$\leq \gamma \|Q^* - Q^{\pi_{k-1}}\|_\infty \mathbb{1}$$

$$\leq \cdots$$

$$\leq \gamma^k \|Q^* - Q^{\pi_0}\|_\infty \mathbb{1}$$

$$\leq \gamma^k \|Q^*\|_\infty \mathbb{1} \qquad \text{since } 0 \leq Q^{\pi_0} \leq Q^*$$

$$\leq [1 - (1-\gamma)]^k \cdot \frac{1}{1-\gamma} \mathbb{1}$$

$$\leq e^{-(1-\gamma)k} \cdot \frac{1}{1-\gamma} \cdot \mathbb{1}$$

$$\leq \varepsilon \mathbb{1} \qquad \text{since } k \geq \frac{\log \frac{1}{(1-\gamma)\varepsilon}}{1-\gamma} \quad \square$$