

## A. Preliminaries.

### 1. Contraction Mapping.

$(\mathbb{X}, d)$  is a metric space.  $T: \mathbb{X} \rightarrow \mathbb{X}$  is a mapping. If  $\exists L > 0$  s.t.  $\forall x, y \in \mathbb{X}$ ,  $d(Tx, Ty) \leq L \cdot d(x, y)$ ,  $T$  is  $L$ -Lipschitz in  $\mathbb{X}$  w.r.t.  $d$ . Specifically, if  $L < 1$ , then  $T$  is a contraction mapping.

### 2. Fixed point.

$(\mathbb{X}, d)$  is a metric space.  $T: \mathbb{X} \rightarrow \mathbb{X}$  is a mapping. If  $Tx = x$ , then  $x$  is a fixed point of  $T$ .

### 3. Contraction Mapping Theorem.

①  $\mathbb{X}$  is complete.

②  $T$  contraction mapping.

Then there exists and exists only one fixed point of  $T$ .

(OBO)

4. The optimal Bellman operator  $H$  is a contraction mapping w.r.t. maximum norm.

Proof.

Recall the OBO  $H$  satisfies

$$H\varphi(s, a) = \sum_{s' \sim P(s', | s, a)} p(s', a, s') \cdot [r(s, a, s') + \gamma \cdot \max_{a' \in A} \varphi(s', a')],$$

for some  $Q$ -function  $\varphi$ .

For any two Q-functions  $q_1$  and  $q_2$ ,

$$\| Hq_1 - Hq_2 \|_\infty$$

$$= \max_{(S,a) \in S \times A} \left| \sum_{S' \sim P(S,a)} P(S'|S,a) \left[ \gamma \cdot \max_{a' \in A} q_1(S',a') - \max_{a' \in A} q_2(S',a') \right] \right|$$

$$\leq \max_{(S,a) \in S \times A} \gamma \sum_{S' \sim P(S,a)} P(S'|S,a) \left| \max_{a' \in A} q_1(S',a') - \max_{a' \in A} q_2(S',a') \right|. \quad (1)$$

$$\leq \max_{(S,a) \in S \times A} \gamma \sum_{S' \sim P(S,a)} P(S'|S,a) \left| \max_{a' \in A} q_1(S',a') - q_2(S',a') \right|$$

$$\leq \max_{(S,a) \in S \times A} \gamma \sum_{S' \sim P(S,a)} P(S'|S,a) \left| \max_{(S'',a'') \in S \times A} q_1(S'',a'') - q_2(S'',a'') \right|$$

$$= \gamma \cdot \| q_1 - q_2 \|_\infty,$$

where Eq (1) follows from Jensen's inequality.  $\square$

### 5. Q-Learning

$$\begin{aligned} Q_{t+1}^{\pi}(S_t, a_t) &= Q_t^{\pi}(S_t, a_t) - \alpha \cdot (Q_t^{\pi}(S_t, a_t) - \text{sample}_t), \\ &= (1-\alpha) \cdot Q_t^{\pi}(S_t, a_t) + \alpha \cdot \text{sample}_t, \end{aligned}$$

where  $\text{sample}_t = R_t(S_t, a_t, S_{t+1}) + \gamma \max_{a' \in A} Q(S_{t+1}, a')$ , and  $\{S_t\}$  is a sequence of states obtained by following policy  $\pi$  which satisfies

$$P_\pi(A_t = a | S_t = s) > 0$$

for any state action pairs  $(S, a) \in S \times A$ .

### B. Content.

#### Theorem 1. Convergence of Q-Learning

Consider a finite MDP  $(S, A, P, R, \gamma)$ , where  $S$  and  $A$  are finite.

and  $R$  is bounded and deterministic. If we use  $\alpha$ -learning algorithm to solve it with learning rate  $\alpha_t$ , which satisfies

$$\textcircled{1} \quad 0 < \alpha_t \leq 1$$

$$\textcircled{2} \quad \sum_t \alpha_t = +\infty$$

$$\textcircled{3} \quad \sum_t \alpha_t^2 < +\infty.$$

Then,  $Q_{t+1}$  converges to  $Q^*$  w.p. 1.

Proof of Theorem 1:

Theorem 2.

The random process  $\{\Delta_t\}$ , taking values in  $\mathbb{R}^n$  and defined as

$$\Delta_{t+1}(x) = (1 - \alpha_t) \cdot \Delta_t(x) + \alpha_t \cdot F_t(x).$$

If the following assumptions hold :

$$\textcircled{1} \quad 0 < \alpha_t \leq 1, \quad \sum_t \alpha_t = +\infty, \quad \sum_t \alpha_t^2 < +\infty.$$

$$\textcircled{2} \quad \|E[F_t(x) | \mathcal{F}_t]\|_\infty \leq \gamma \cdot \|\Delta_t\|_\infty, \quad \text{with } \gamma < 1.$$

$$\textcircled{3} \quad V[F_t(x) | \mathcal{F}_t] \leq C \cdot (1 + \|\Delta_t\|_\infty)^2, \quad \text{for } C > 0,$$

where  $\mathcal{F}_t = \{\Delta_t, \Delta_{t-1}, \dots, \Delta_1, F_{t-1}, F_{t-2}, \dots, F_1\}$  stands for the past at time  $t$ . Then  $\{\Delta_t\}$  converges to 0 w.p. 1.

Recall the  $\alpha$ -learning update:

$$Q_{t+1}(S_t, a_t) = (1 - \alpha_t) \cdot Q_t(S_t, a_t) + \alpha_t \cdot (R(S_t, a_t, S') + \gamma \cdot \max_{a' \in A} Q_t(S', a')).$$

where  $S' \sim P_C \cdot |S_t, a_t\rangle$ .

Substracting from both sides the quantity  $Q^*(S_t, a_t)$  shows that

$$Q_{t+1}(S_t, a_t) - Q^*(S_t, a_t) = (1 - \alpha_t) \cdot [Q_t(S_t, a_t) - Q^*(S_t, a_t)]$$

$$+ \alpha_t (R(S_t, a_t, S') + \gamma \cdot \max_{a' \in A} (Q_t(S', a') - Q^*(S', a'))) \quad (2)$$

Let

$$\Delta_t(s, a) = Q_t(s, a) - Q^*(s, a).$$

and

$$F_t(s, a) = R(s, a, s') + \gamma \max_{a' \in A} Q_t(s', a') - Q^*(s, a),$$

with  $s' \sim P(\cdot | s, a)$ . Then Eq.(2) could be reformulated as

$$\Delta_{t+1}(s_t, a_t) = (1 - \delta_t) \Delta_t(s_t, a_t) + \delta_t F_t(s_t, a_t).$$

a. Assumption (1) in Theorem 2. is satisfied.

b.  $E[F_t(s_t, a_t) | \mathcal{F}_t]$

$$= E_{s' \sim P(\cdot | s_t, a_t)} [R(s_t, a_t, s') + \gamma \max_{a' \in A} Q_t(s', a') - Q^*(s_t, a_t)].$$

$$= (H Q_t)(s_t, a_t) - Q^*(s_t, a_t).$$

$$= (H Q_t)(s_t, a_t) - (H Q^*)(s_t, a_t), \quad (3)$$

where Eq.(3) follows from that  $Q^*$  is the fixed point of  $H$ .

Then

$$\|E[F_t | \mathcal{F}_t]\|_\infty = \|H Q_t - H Q^*\|_\infty \leq \gamma \|Q_t - Q^*\|_\infty = \gamma \|\Delta_t\|_\infty.$$

Assumption (2) in Theorem 2. is satisfied.

c.  $\mathbb{V}_{s' \sim P(\cdot | s_t, a_t)} [F_t(s_t, a_t) | \mathcal{F}_t]$

$$= E_{s' \sim P(\cdot | s_t, a_t)} [(F_t(s_t, a_t) - E_{s' \sim P(\cdot | s_t, a_t)} [F_t(s_t, a_t) | \mathcal{F}_t])^2 | \mathcal{F}_t]$$

$$= E_{s' \sim P(\cdot | s_t, a_t)} [R(s_t, a_t, s') + \gamma \max_{a' \in A} Q_t(s', a') - Q^*(s_t, a_t)]$$

$$- \sum_{s' \sim P(\cdot | s_t, a_t)} P(s_t, a_t, s') [R(s_t, a_t, s') + \gamma \max_{a' \in A} Q_t(s', a') - Q^*(s_t, a_t)]^2 | \mathcal{F}_t.$$

$$= \mathbb{E}_{S' \sim P(\cdot | S_t, a_t)} \left[ R(S_t, a_t, S') + \gamma \max_{a' \in A} Q_t(S', a') - Q^*(S_t, a_t) - (H(Q_t)(S_t, a_t)) \right. \\ \left. + Q^*(S_t, a_t) \right]^2 | \mathcal{F}_t]$$

$$= \mathbb{E}_{S' \sim P(\cdot | S_t, a_t)} \left[ R(S_t, a_t, S') + \gamma \max_{a' \in A} Q_t(S', a') - (H(Q_t)(S_t, a_t)) | \mathcal{F}_t \right]$$

$$= \mathbb{V}_{S' \sim P(\cdot | S_t, a_t)} \left[ R(S_t, a_t, S') + \gamma \max_{a' \in A} Q_t(S', a') | \mathcal{F}_t \right].$$

$$= \mathbb{V}_{S' \sim P(\cdot | S_t, a_t)} \left[ R(S_t, a_t, S') + \gamma \cdot \overline{Q^*(S', \arg\max_{a' \in A} Q_t(S', a'))} \right]$$

$$+ \underbrace{\left( \gamma \max_{a' \in A} Q_t(S', a') \right) - \overline{Q^*(S', \arg\max_{a' \in A} Q_t(S', a'))}}_{x_{t,2}}.$$

Let random variable  $x_{t,1}(S') = R(S_t, a_t, S') + \gamma \cdot \overline{Q^*(S', \arg\max_{a' \in A} Q_t(S', a'))}$

$$x_{t,2}(S') = \left( \gamma \max_{a' \in A} Q_t(S', a') \right) - \overline{Q^*(S', \arg\max_{a' \in A} Q_t(S', a'))}$$

$$\text{Then } \sup_{S' \in S} |x_{t,2}(S')| = \sup_{S' \in S} \left| \left( \gamma \max_{a' \in A} Q_t(S', a') \right) - \overline{Q^*(S', \arg\max_{a' \in A} Q_t(S', a'))} \right|$$

$$\leq \sup_{S' \in S} \gamma \left| \max_{a' \in A} Q_t(S', a') - Q^*(S', a') \right|$$

$$\leq \sup_{S' \in S} \gamma \left| \max_{\substack{(S'', a'') \in S \times A \\ (S'', a'') \in S \times A}} Q_t(S'', a'') - Q^*(S'', a'') \right|$$

$$= \gamma \| \Delta_t \|_\infty.$$

It follows that

$$\mathbb{V}_{s' \sim \mathbb{P}( \cdot | s_t, a_t)} [F_t(s_t, a_t) | \mathcal{H}_t]$$

$$= \mathbb{V}_{s' \sim \mathbb{P}( \cdot | s_t, a_t)} [x_{t+1} + x_{t+2} | \mathcal{H}_t]$$

$$\leq \sup_{s' \in S} |x_{t+1}(s') + x_{t+2}(s')|^2$$

$$\leq \sup_{s' \in S} (|x_{t+1}(s')| + |x_{t+2}(s')|)^2$$

$$\leq (C' + \gamma^2 \| \Delta t \|_\infty)^2$$

$$= C (1 + \| \Delta t \|_\infty)^2,$$

for some constants  $C'$  and  $C$ . Assumption ③ in Theorem 2

is satisfied and applying Theorem 2 immediately shows that

$\lim_{t \rightarrow +\infty} \Delta t = 0$ , which further implies that  $\lim_{t \rightarrow +\infty} Q_t = Q^*$  and

concludes the proof.

□