1. $q^H(s,a) \overset{def}{=\!=} \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma \max_{a'} q(s',a')] \Leftrightarrow q^H = H(q)$.

To proof $H$ is a contradiction mapping. just proof as following:

$\forall q_1 \cdot q_2 \in X$. $d(q_1^H, q_2^H) \leq d(q_1 \cdot q_2)$

$d(q_1, q_2) = \max_{s,a} |q_1(s,a) - q_2(s,a)|$. ---①

$d(q_1^H, q_2^H) = \max_{s,a} |q_1^H(s,a) - q_2^H(s,a)|$ ---②

$= \max_{s,a} |\sum_{s'} T(s,a,s')[R(s,a,s') + \gamma \max_{a'} q_1(s',a')]$ ---③

$\quad - \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma \max_{a'} q_2(s',a')]|$ ---④

$= \max_{s,a} |\sum_{s'} T(s,a,s') \gamma [\max_{a'} q_1(s',a') - \max_{a'} q_2(s',a')]|$ ---⑤

$< \max_{s,a} |\sum_{s'} T(s,a,s')[\max_{a'} q_1(s',a') - \max_{a'} q_2(s',a')]|$ ---⑥

$\leq \max_{s,a} \sum_{s'} T(s,a,s') |\max_{a'} q_1(s',a') - \max_{a'} q_2(s',a')|$ ---⑦

Since $\forall s',a'$ $|\max_{a'} q_1(s',a') - \max_{a'} q_2(s',a')|$

$\leq \max_{a'} |q_1(s',a') - q_2(s',a')|$ ---- Lemma $*$

⑦ $\leq \max_{s,a} \sum_{s'} T(s,a,s') \cdot \max_{a'} |q_1(s',a') - q_2(s',a')|$

$\leq \max_{s,a} \sum_{s'} T(s,a,s') \cdot \max_{s',a'} |q_1(s',a') - q_2(s',a')|$ ---⑧

$= \max_{s,a} \max_{s',a'} |q_1(s',a') - q_2(s',a')| = \max_{s,a} |q_1(s,a) - q_2(s,a)|$ ---⑨

Proof of lemma *

$\forall s', a'$  $\left| \max_{a'} q_1(s', a') - \max_{a'} q_2(s', a') \right| \leq \max_{a'} \left| q_1(s', a') - q_2(s', a') \right|$

Without loss of generality, we assume that $\max_{a'} q_1(s', a') > \max_{a'}(s', a')$

And we assume that, $\begin{cases} a_1 = \arg\max q_1(s', a') \\ a_2 = \arg\max q_2(s', a') \end{cases}$

Then we need to proof: $q_1(s', a_1) - q_2(s', a_2) \leq \max_{a'} \left| q_1(s', a') - q_2(s', a') \right|$

$q_1(s', a_1) - q_2(s', a_2) = q_1(s', a_1) - q_2(s', a_1) + \underbrace{[q_2(s', a_1) - q_2(s', a_2)]}_{\leq 0}$

$\leq q_1(s', a_1) - q_2(s', a_1)$

$\leq \max_{a'} \left| q_1(s', a') - q_2(s', a') \right|$

Q.E.D

2.

To proof $Q$ converges to $Q^*$, just proof $\lim_{t \to \infty} \left( Q_t(s, a) - Q^*(s, a) \right) = 0$.

$\Delta_t(s, a) \overset{def}{=} Q_t(s, a) - Q^*(s, a)$

According to the update rool of $Q_t$, we have:

$\Delta_{t+1}(s, a) = Q_{t+1}(s, a) - Q^*(s, a)$

$= (1 - \alpha_t) \cdot Q_t(s, a) + \alpha_t \cdot sample_t - Q^*(s, a)$

$= (1 - \alpha_t) \cdot \Delta_t(s, a) + \alpha_t \cdot \underbrace{\left[ (R(s, a, s') + \gamma \max_{a'} Q_t(s', a')) - Q^*(s, a) \right]}_{\overset{def}{=} F_t(s, a)} \quad \cdots (*)$

Then we proof $(*)$. Satisfy Lemma 1.

① $\alpha_t \in (0.1)$, $\sum \alpha_t = \infty$. $\sum \alpha_t^2 < \infty$.  satisfied.

② To proof: $\|E(F_t | \mathcal{F}_t)\|_\infty \leq \gamma \|\Delta_t\|_\infty$

$\Leftrightarrow \max\limits_{s,a} |E(F_t(s,a) | \mathcal{F}_t)| \leq \gamma \max\limits_{s,a} |\Delta_t(s,a)|$

$\Leftrightarrow \max\limits_{s,a} \left| E[R(s,a,s') + \gamma \max\limits_{a'} Q_t(s',a') - Q^*(s,a)] \right| \leq \gamma \max\limits_{s,a} |\Delta_t(s,a)|$

Since $Q^*(s,a) = E[R(s,a,s') + \gamma \max\limits_{a'} Q^*(s',a')]$  Then

$LHS = \max\limits_{s,a} \left| E[R(s,a,s') + \gamma \max\limits_{a'} Q_t(s',a')] - Q^*(s,a) \right|$

$= \max\limits_{s,a} \left| E[R(s,a,s') + \gamma \max\limits_{a'} Q_t(s',a')] - E[R(s,a,s') + \gamma \max\limits_{a'} Q^*(s',a')] \right|$

$= \max\limits_{s,a} \left| E[\gamma \max\limits_{a'} Q_t(s',a') - \gamma \max\limits_{a'} Q^*(s',a')] \right|$

$\left.\begin{array}{c}\\ \\ \end{array}\right\}$ according to lemma * in Q1.

$\leq \max\limits_{s,a} E[\gamma \max\limits_{a'} (Q_t(s',a') - Q^*(s',a'))]$

$= \max\limits_{s,a} E[\gamma \max\limits_{a'} \Delta_t(s',a')] \leq \gamma \max\limits_{s,a} \Delta_t.$

Q.E.D. satisfied.

③ To proof: $V[F_t(s,a) | \mathcal{F}_t] \leq C(1 + \|\Delta_t\|_\infty)^2$.

$\Leftrightarrow V[F_t(s,a) | \mathcal{F}_t] \leq C[1 + \max\limits_{s,a} \Delta_t(s,a)]^2$.

To simplify the representation, I will ignore the $\mathcal{F}_t$.

From ②, we get that:

$E(F_t(s,a)) = E[\gamma \max\limits_{a'} Q_t(s',a') - \gamma \max\limits_{a'} Q^*(s',a')]$

$= \gamma E[\max\limits_{a'} Q_t(s',a') - \max\limits_{a'} Q^*(s',a')]$

$F_t(s,a) = Q_{t+1}(s,a) - Q^*(s,a)$

$LHS = E[F_t(s,a) - E(F_t(s,a))]^2$

$= E\{Q_{t+1}(s,a) - Q^*(s,a) - E[Q_{t+1}(s,a) - Q^*(s,a)]\}^2$

$$= E\left\{Q_{t+1}(s,a) - E[Q_{t+1}(s,a)]\right\}^2$$

$$= V[Q_{t+1}(s,a)]$$

$$= V[R(s,a,s') + \gamma \max_{a'} Q_t(s',a')]$$

Since $R(s,a,s')$ and $Q_t(s',a')$ are all bounded. thus. verify.

$$V[R(s,a,s') + \gamma \max_{a'} Q_t(s',a')] < C \cdot 1 < C(1+\|\Delta_t\|_\infty)^2$$