

---

# A Note on Double Q-Learning

---

**Author Po-Wei Huang**  
Department of Computer Science  
National Yang Ming Chiao Tung University  
a0716084.cs07@nycu.edu.tw

## 1 Introduction

The main research challenges tackled by this paper is Q-learning's poor performance in stochastic MDPs due to large overestimations of the action values. This paper shows that why overestimation in Q-learning occurs and propose another algorithm called Double Q-learning to avoid this. Instead of using single estimator like Q-learning, double Q-learning uses an alternative method, double estimator, to find the maximum expected value. This paper also shows that double Q-learning converges to the optimal value and can handle the overestimation problem, but may have underestimation problem, so the double Q-learning method is not always better than the Q-learning method.

## 2 Problem Formulation

### 2.1 Notations

- MDP:  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ 
  - $\mathcal{S}$ : state space
  - $\mathcal{A}$ : action space
  - $P$ : fixed state transition distribution,  $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$
  - $R$ : fixed reward distribution,  $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{R}$
  - $\gamma \in [0, 1)$ : discount factor
- Q-learning:
  - $Q_t(s, a)$ : the value of the action  $a$  in state  $s$  at time  $t$ .
  - $r_t$ : the reward drawn from  $R$ , where  $E[r_t | (s, a, s') = (s_t, a_t, s_{t+1})] = R_{sa}^{s'}$
  - $\alpha_t(s, a) \in [0, 1]$ : learning rate
- Estimating the maximum expected value
  - $X = (X_1, \dots, X_M)$ : a set of  $M$  random variables
  - $\max_i E[X_i]$ : the maximum expected value of the variables
  - $S = \cup_{i=1}^M S_i$ : a set of iid samples,  $S_i$  is the subset containing samples for  $X_i$
  - $\mu_i$ : an estimator for  $X_i$
  - $f_i$ : the probability density function of  $X_i$
  - $F_i$ : the cumulative distribution function of  $f_i$ ,  $F_i = \int_{-\infty}^x f_i(x)dx$

It is impossible to exactly determine  $\max_i E[X_i]$  for  $X = (X_1, \dots, X_M)$  in most cases. How can we estimate the maximum expected value? This section shows two kinds of estimator: single estimator and double estimator.

The sample average from  $S_i$  for each  $X_i$  is an unbiased estimator for the expected values:

$$E[X_i] = E[\mu_i] \approx \mu_i(S) \stackrel{def}{=} \frac{1}{|S_i|} \sum_{s \in S_i} s \quad (1)$$

The maximum expected value:

$$\max_i E[X_i] = \max_i \int_{-\infty}^x x f_i(x) dx \quad (2)$$

## 2.2 The Single Estimator

Q-learning approximates the value  $\max_i E[X_i] = \max_i E[\mu_i] \approx \max_i \mu_i(S)$  by maximizing over the estimated action values in that state.  $\max_i \mu_i$  is distributed according to some PDF  $f_{\max}^\mu$  that is dependent on the PDFs of the estimators  $f_i^\mu$ . To determine this PDF, consider the CDF  $F_{\max}^\mu(x)$ :

$$F_{\max}^\mu \stackrel{\text{def}}{=} P(\max_i \mu_i \leq x) = \Pi_{i=1}^M P(\mu_i \leq x) \stackrel{\text{def}}{=} \Pi_{i=1}^M F_i^\mu(x) \quad (3)$$

$\max_i \mu_i(S)$  is an unbiased estimate for:

$$E[\max_j \mu_j] = \int_{-\infty}^{\infty} x f_{\max}^\mu(x) dx = \int_{-\infty}^{\infty} x \frac{d}{dx} \Pi_{i=1}^M F_i^\mu(x) dx = \sum_j^M \int_{-\infty}^{\infty} x f_j^\mu(s) \Pi_{i \neq j}^M F_i^\mu(x) dx \quad (4)$$

However, it is not an unbiased estimate for  $\max_j E[\mu_j]$ .

## 2.3 The Double Estimator

An alternative method to approximate  $\max_i E[X_i]$ , uses two sets of estimators:  $\mu^A = \{\mu_1^A, \dots, \mu_M^A\}$  and  $\mu^B = \{\mu_1^B, \dots, \mu_M^B\}$ . Both sets of estimators are updated with a subset of the samples we draw,  $S = S^A \cup S^B$ ,  $S^A \cap S^B = \phi$ ,  $\mu_i^A(S) = \frac{1}{|S^A|} \sum_{s \in S^A} s$ , and  $\mu_i^B(S) = \frac{1}{|S^B|} \sum_{s \in S^B} s$ , assume that the samples are split in a proper manner, both  $\mu^A$  and  $\mu^B$  are unbiased. Let the set of maximal estimates in  $\mu^A(S)$ :  $\text{Max}^A(S) \stackrel{\text{def}}{=} \{j | \mu_j^A(S) = \max_i \mu_i^A(S)\}$ . Since  $\mu^B$  is an independent, unbiased set of estimators.  $E[\mu_j^B] = E[X_j] \forall j$ , including all  $j \in \text{Max}^A$ . Pick an estimator  $a^*$  that maximizes  $\mu^A$ :  $\mu_{a^*}^A(S) \stackrel{\text{def}}{=} \max_i \mu_i^A(S)$ , then we can use  $\mu_{a^*}^B(S)$  as an estimate for  $\max_i E\{\mu_i^B\}$ , then obtain the approximation:  $\max_i E\{X_i\} = \max_i E\{\mu_i^B\} \approx \mu_{a^*}^B$ . Assume that the underlying PDFs are continuous,  $f_j^A$  and  $F_j^A$  are the PDF and CDF of  $\mu_i^A$ :

$$P(j = a^*) = \int_{-\infty}^{\infty} P(\mu_j^A = x) \Pi_{i \neq j}^M P(\mu_i^A < x) dx \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f_j^A(x) \Pi_{i \neq j}^M F_i^A(x) dx \quad (5)$$

The expected value of the double estimator's approximation is:

$$\sum_j^M P(j = a^*) E[\mu_j^B] = \sum_j^M E[\mu_j^B] \int_{-\infty}^{\infty} f_j^A(x) \Pi_{i \neq j}^M F_i^A(x) dx \quad (6)$$

For discrete PDFs, use sums to replace the integrals.

Comparing the double estimator  $\sum_j^M E[\mu_j^B] \int_{-\infty}^{\infty} f_j^A(x) \Pi_{i \neq j}^M F_i^A(x) dx$  to single estimator  $\sum_j^M \int_{-\infty}^{\infty} x f_j^\mu(s) \Pi_{i \neq j}^M F_i^\mu(x) dx$ , the double estimator uses  $E[\mu_j^B]$  in place of  $x$  within integral.

- The single estimator overestimates because  $x$  within integral and the CDF product  $\Pi_{i \neq j}^M F_i^\mu(x)$  is monotonically increasing.
- The double estimator underestimates because  $\sum_j^M P(j = a^*) = 1$  and the approximation  $\sum_j^M P(j = a^*) E[\mu_j^B]$  can be viewed as a weighted estimate of unbiased expected values.

**Lemma 1** Let  $X = \{X_1, \dots, X_M\}$  be a set of random variables and let  $\mu^A = \{\mu_1^A, \dots, \mu_M^A\}$  and  $\mu^B = \{\mu_1^B, \dots, \mu_M^B\}$  be two sets of unbiased estimators such that  $E[\mu_i^A] = E[\mu_i^B] = E[X_i]$ , for all  $i$ . Let  $\mathcal{M} \stackrel{\text{def}}{=} \{j | E[X_j] = \max_i E[X_i]\}$  be the set of elements that maximize the expected values. Let  $a^*$  be an element that maximizes  $\mu^A$ :  $\mu_{a^*}^A = \max_i \mu_i^A$ . Then  $E[\mu_{a^*}^B] = E[X_{a^*}] \leq \max_i E[X_i]$ . Furthermore, the inequality is strict if and only if  $P(a^* \notin \mathcal{M}) > 0$ .

**Proof.**

$$1. E[\mu_{a^*}^B] \leq \max_i E[X_i]:$$

if  $a^* \in M$ , then  $E[X_{a^*}] = \max_i E[X_i]$ , otherwise  $E[X_{a^*}] < \max_i E[X_i]$   
so  $E[X_{a^*}] \leq \max_i E[X_i]$

$$2. E[\mu_{a^*}^B] < \max_i E[X_i] \Leftrightarrow P(a^* \notin M) > 0:$$

$$\begin{aligned} E[\mu_{a^*}^B] &= P(a^* \in M)E[X_{a^*}|a^* \in M] + P(a^* \notin M)E[X_{a^*}|a^* \notin M] \\ E[X_{a^*}|a^* \in M] &= \max_i E[X_i] \\ (\Rightarrow) E[\mu_{a^*}^B] < \max_i E[X_i] &\Rightarrow P(a^* \notin M) < 0: \end{aligned}$$

if  $P(a^* \notin M) = 0$

then

$$\begin{aligned} E[\mu_{a^*}^B] &= P(a^* \in M)E[X_{a^*}|a^* \in M] + P(a^* \notin M)E[X_{a^*}|a^* \notin M] \\ &= 1 \cdot E[X_{a^*}|a^* \in M] + 0 \cdot E[X_{a^*}|a^* \notin M] \\ &= \max_i E[X_i] \\ &\text{contradiction with } E[\mu_{a^*}^B] < \max_i E[X_i], \text{ so } P(a^* \notin M) > 0 \end{aligned} \tag{7}$$

$$(\Leftrightarrow) P(a^* \notin M) < 0 \Rightarrow E[\mu_{a^*}^B] < \max_i E[X_i]:$$

$$\begin{aligned} \max_i E[X_i] - E[\mu_{a^*}^B] &= E[X_{a^*}|a^* \in M] - (P(a^* \in M)E[X_{a^*}|a^* \in M] + P(a^* \notin M)E[X_{a^*}|a^* \notin M]) \\ &= P(a^* \notin M)(E[X_{a^*}|a^* \in M] - E[X_{a^*}|a^* \notin M]) > 0 \\ (\because P(a^* \notin M) > 0 \text{ and } E[X_{a^*}|a^* \in M] > E[X_{a^*}|a^* \notin M]) \\ \text{so } E[\mu_{a^*}^B] &< \max_i E[X_i] \end{aligned} \tag{8}$$

$P(a^* \notin \mathcal{M})$  happens when the variables have different expected values but their distributions overlap. The double estimator is unbiased when the variables are iid since then all expected values are equal and  $P(a^* \in \mathcal{M}) = 1$ .

## Double Q-learning

Q-learning can be interpreted as using the single estimator to estimate the value of the next state:

$$\max_a Q_t(s_{t+1}, a) \approx E[\max_a Q_t(s_{t+1}, a)] \approx \max_a E[Q_t(s_{t+1}, a)] \tag{9}$$

$\max_a Q_t(s_{t+1}, a)$  is an unbiased sample, drawn from an iid distribution with mean  $E[\max_a Q_t(s_{t+1}, a)]$ . Empirically, Q-learning can suffer from large overestimations. The next section will present Double Q-learning Algorithm to avoid these overestimation issues and will show that it can converge to the optimal policy.

### 3 Theoretical Analysis

#### 3.1 Double Q-learning Algorithm

---

**Algorithm 1** Double Q-learning

---

```

1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end

```

---

Double Q-learning stores two Q functions:  $Q^A$  and  $Q^B$ . Each Q function is updated with a value from the other Q function for the next state.  $a^*$  is the maximal valued action in state  $s'$  according to  $Q^A$ . However, instead of using  $Q^A(s', a^*)$  to update  $Q^A$  like Q-learning, this algorithm uses  $Q^B(s', a^*)$ . Since  $Q^A$  and  $Q^B$  have the same update method but with a different set of experience samples, this can be considered an unbiased estimate for the value of this action.

Like the double estimator, double Q-learning may underestimate because action  $a^*$  may not be the action that maximizes the expected Q function  $\max_a E[Q^A(s', a)]$ . In general,  $\max_a E[Q^B(s', a^*)] \leq \max_a E[Q^A(s', a^*)]$ .

#### 3.2 Convergence in the Limit

**Theorem 1** *In the limit Double Q-learning converges to the optimal policy.*

**Theorem 1.** Assume the conditions below are fulfilled. Then, in a given ergodic MDP, both  $Q^A$  and  $Q^B$  as updated by Double Q-learning as described in Algorithm 1 will converge to the optimal value function  $Q^*$  as given in the Bellman optimality equation (2) with probability one if an infinite number of experiences in the form of rewards and state transitions for each state action pair are given by a proper learning policy. The additional conditions are: 1) The MDP is finite, i.e.  $|S \times A| < \infty$ . 2)  $\gamma \in [0, 1)$ . 3) The Q values are stored in a lookup table. 4) Both  $Q^A$  and  $Q^B$  receive an infinite number of updates. 5)  $\alpha_t(s, a) \in [0, 1]$ ,  $\sum_t \alpha_t(s, a) = \infty$ ,  $\sum_t (\alpha_t(s, a))^2 < \infty$  w.p.1, and  $\forall (s, a) \neq (s_t, a_t) : \alpha_t(s, a) = 0$ . 6)  $\forall s, a, s' : \text{Var}\{R_{sa}^{s'}\} < \infty$ .

**Conditions**

1. The MDP is finite, i.e.  $|S \times A| < \infty$
2.  $\gamma \in [0, 1)$
3. The Q values are stored in a lookup table.
4. Both  $Q^A$  and  $Q^B$  receive an infinite number of updates.
5.  $\alpha_t(s, a) \in [0, 1]$ ,  $\sum_t \alpha_t(s, a) = \infty$ ,  $\sum_t \alpha_t(s, a)^2 < \infty$  w.p.1 and  $\forall (s, a) \neq (s_t, a_t) : \alpha_t(s, a) = 0$
6.  $\forall s, a, s' : \text{Var}[R_{sa}^{s'}] < \infty$

**Proof** This paper proves Theorem 1 by making use of lemma 2, which was also used to prove convergence of Sarsa:

## Lemma 2

**Lemma 2.** Consider a stochastic process  $(\zeta_t, \Delta_t, F_t)$ ,  $t \geq 0$ , where  $\zeta_t, \Delta_t, F_t : X \rightarrow \mathbb{R}$  satisfy the equations:

$$\Delta_{t+1}(x_t) = (1 - \zeta_t(x_t))\Delta_t(x_t) + \zeta_t(x_t)F_t(x_t) , \quad (8)$$

where  $x_t \in X$  and  $t = 0, 1, 2, \dots$ . Let  $P_t$  be a sequence of increasing  $\sigma$ -fields such that  $\zeta_0$  and  $\Delta_0$  are  $P_0$ -measurable and  $\zeta_t, \Delta_t$  and  $F_{t-1}$  are  $P_t$ -measurable,  $t = 1, 2, \dots$ . Assume that the following hold: 1) The set  $X$  is finite. 2)  $\zeta_t(x_t) \in [0, 1]$ ,  $\sum_t \zeta_t(x_t) = \infty$ ,  $\sum_t (\zeta_t(x_t))^2 < \infty$  w.p.1 and  $\forall x \neq x_t : \zeta_t(x) = 0$ . 3)  $\|E\{F_t|P_t\}\| \leq \kappa\|\Delta_t\| + c_t$ , where  $\kappa \in [0, 1]$  and  $c_t$  converges to zero w.p. 1. 4)  $\text{Var}\{F_t(x_t)|P_t\} \leq K(1 + \kappa\|\Delta_t\|)^2$ , where  $K$  is some constant. Here  $\|\cdot\|$  denotes a maximum norm. Then  $\Delta_t$  converges to zero with probability one.

## Conditions

1. The set  $X$  is finite.
2.  $\zeta_t(x_t) \in [0, 1]$ ,  $\sum_t \zeta_t(x_t) = \infty$ ,  $\sum_t \zeta_t(x_t)^2 < \infty$  w.p.1,  $\forall x \neq x_t, \zeta_t(x) = 0$
3.  $\|E[F_t|P_t]\| \leq \kappa\|\Delta_t\| + c_t$ , where  $\kappa \in [0, 1]$  and  $c_t \rightarrow 0$  w.p.1
4.  $\text{Var}[F_t(x_t)|P_t] \leq K(1 + \kappa\|\Delta_t\|)^2$ , where  $K$  is some constant.

I will prove that in the algorithm,  $Q^A$  will converge to  $Q^*$ , the optimal value, and by the symmetry of A and B,  $Q^B$  will converge to  $Q^*$ , too.

Let  $P_t = Q_0^A, Q_0^B, s_0, a_0, r_1, s_1, \dots, s_t, a_t, X = S \times A, \Delta_t = Q_t^A - Q^*$

$$\begin{aligned} Q_{t+1}^A(s_t, a_t) &= Q_t^A(s_t, a_t) + \alpha(s_t, a_t)(r_t + \gamma Q_t^B(s_{t+1}, a^*) - Q^A(s_t, a_t)) \\ &\Rightarrow Q_{t+1}^A(s_t, a_t) = (1 - \alpha(s_t, a_t))Q_t^A(s_t, a_t) + \alpha(s_t, a_t)(r_t + \gamma Q_t^B(s_{t+1}, a^*)) \end{aligned}$$

subtract  $Q^*(s_t, a_t)$  from both sides of the equation, we get:

$$\Delta_{t+1}(s_t, a_t) = (1 - \alpha(s_t, a_t))\Delta_t(s_t, a_t) + \alpha(s_t, a_t)(r_t + \gamma Q_t^B(s_{t+1}, a^*) - Q^*(s_t, a_t))$$

Let  $\zeta_t = \alpha_t, F_t(s_t, a_t) = r_t + \gamma Q_t^B(s_{t+1}, a^*) - Q_t^*(s_t, a_t), a^* = \arg\max_a Q^A(s_{t+1}, a^*), \gamma = \kappa$   
Check the four conditions in lemma 2 is satisfied:

1. Satisfied by the condition 1 of Theorem 1.
2. Satisfied by the condition 5 of Theorem 1.
3. Let  $F_t^Q(s_t, a_t) = r_t + \gamma Q_t^A(s_{t+1}, a^*) - Q_t^*(s_t, a_t)$ , the value of  $F_t$  in Q-learning:

$$F_t(s_t, a_t) = F_t^Q(s_t, a_t) + \gamma(Q_t^B(s_{t+1}, a^*) - Q_t^A(s_{t+1}, a^*))$$

$$\text{Let } \Delta^{BA} = Q_t^B(s_{t+1}, a^*) - Q_t^A(s_{t+1}, a^*)$$

Show that  $\|E[F_t^Q(s_t, a_t)|P_t]\|_\infty \leq \gamma\|\Delta_t\|_\infty$  first:

The update for  $Q^*$  is:  $Q_{t+1}^*(s, a) = (1 - \alpha)Q_t^*(s, a) + \alpha(r + \gamma \cdot \max_a Q_t^*(s', a))$ ,  
 $\therefore Q_{t+1}^* = Q_t^* \therefore Q^*(s, a) = r + \gamma \cdot \max_a Q^*(s', a)$

$$\begin{aligned} \|E[F_t^Q(s_t, a_t)|P_t]\|_\infty &= \|E[r_t + \gamma Q_t^A(s_{t+1}, a^*) - Q^*(s_t, a_t)]\|_\infty \\ &= \max_{a_t} |E[r_t + \gamma \max_a Q_t^A(s_{t+1}, a) - Q^*(s_t, a_t)]| \\ &= \max_{a_t} |E[r_t + \gamma \max_{a_t} Q_t^A(s_{t+1}, a_t) - (r_t + \gamma \cdot \max_{a_t} Q^*(s_{t+1}, a_t))]| \\ &= \gamma \max_{a_t} |E[\max_{a_t} Q_t^A(s_{t+1}, a_t) - \max_{a_t} Q^*(s_{t+1}, a_t)]| \\ &\leq \gamma \max_{a_t} |E[Q_t^A(s_{t+1}, a_t) - Q^*(s_{t+1}, a_t)]| \\ &\leq \gamma \max_{a_t} E[|Q_t^A(s_{t+1}, a_t) - Q^*(s_{t+1}, a_t)|] \text{ (Jensen's inequality)} \\ &\leq \gamma \|Q_t^A - Q^*\|_\infty = \gamma \|\Delta_t\|_\infty \end{aligned}$$

Then show that  $\Delta^{BA} = Q_t^B(s_{t+1}, a^*) - Q_t^A(s_{t+1}, a^*)$  converges to 0:

$$\text{Let } F_t^B(s_t, a_t) = r_t + \gamma Q_t^A(s_{t+1}, a^*) - Q_t^B(s_t, a_t),$$

$$F_t^A(s_t, a_t) = r_t + \gamma Q_t^B(s_{t+1}, a^*) - Q_t^A(s_t, a_t)$$

For Update A or Update B,

$$\begin{aligned}\Delta_{t+1}^{BA}(s_t, a_t) &= \Delta_t^{BA}(s_t, a_t) + \alpha(s_t, a_t)F_t^B(s_t, a_t) \text{ or} \\ \Delta_{t+1}^{BA}(s_t, a_t) &= \Delta_t^{BA}(s_t, a_t) - \alpha(s_t, a_t)F_t^A(s_t, a_t)\end{aligned}$$

$$E\Delta_{t+1}^{BA}(s_t, a_t)P_t = \Delta_t^{BA}(s_t, a_t) + E[\alpha(s_t, a_t)(F_t^B(s_t, a_t) - F_t^A(s_t, a_t))|P_t]$$

Let  $\zeta^{BA} = 1/2\alpha_t$ ,  $F^{BA}(s_t, a_t) = \gamma(Q_t^A(s_t, b^*) - Q_t^B(s_t, a^*)) \geq 0$  ( $\because$  the symmetry of A and B)

$$\begin{aligned}&= \Delta_t^{BA}(s_t, a_t) + \zeta^{BA}(s_t, a_t)E[F^{BA}(s_t, a_t) - \Delta_t^{BA}(s_t, a_t)|P_t] \\ &= (1 - \zeta^{BA}(s_t, a_t))\Delta_t^{BA}(s_t, a_t) + \zeta^{BA}(s_t, a_t)E[F^{BA}(s_t, a_t)|P_t] \\ &\because E[F^{BA}(s_t, a_t)|P_t] = \gamma EQ_t^A(s_{t+1}, b^*) - Q_t^B(s_{t+1}, a^*)P_t \leq \gamma EQ_t^A(s_{t+1}, a^*) - Q_t^B(s_{t+1}, a^*)P_t \leq \gamma \|\Delta_t^{BA}\| \\ &\therefore \text{By lemma 2, } \Delta_t^{BA} \rightarrow 0\end{aligned}$$

Because  $F_t(s_t, a_t) = F_t^Q(s_t, a_t) + \gamma\Delta_t^{BA}$ ,  
 $\|E[F_t(s_t, a_t)P_t]\|_\infty = \|E[F_t^Q(s_t, a_t) + \gamma\Delta_t^{BA}]\|_\infty \leq \gamma\|\Delta_t\|_\infty$   
This condition is satisfied.

4. From the condition 6 of theorem 1,  $Var(R_{sa}^{s'})$  is bounded, so  $Var(r_t)$ ,  $Var(Q_t^B)$ ,  $Var(Q_t^*)$  are bounded, too.  
 $F_t(s_t, a_t) = r_t + \gamma Q_t^B(s_{t+1}, a^*) - Q_t^*(s_t, a_t)$   
 $\Rightarrow Var(F_t(s_t, a_t)) = Var(r_t) + \gamma Var(Q_t^B) - Var(Q_t^*) < \infty$ , this condition is satisfied.  
By lemma 2,  $Q^A$  will converge to  $Q^*$   
Because of the symmetry of A and B,  $Q^B$  will converge to  $Q^*$ .  
By theorem 1, both  $Q^A$  and  $Q^B$  will converge to the optimal value.

## 4 Conclusion

Q-learning has overestimation bias, the double Q-learning algorithm in this paper uses a double estimator approach to determine the value of the next state, and it sometimes underestimates the action values, but does not suffer from the overestimation bias.

Since double Q-learning is negatively biased and Q-learning is positively biased, it may be possible to construct an unbiased off-policy reinforcement-learning algorithm in the future.