
A Note on Maxmin Q-Learning

Shang-Hsuan Yang

Department of Computer Science
National Chiao Tung University
sandy861003.eo05@nctu.edu.tw

1 Introduction

This paper proposed a new variant of Q-Learning, called *Maxmin Q-Learning*, in order to address the overestimation bias issues in Q-Learning. The key idea is to maintain N estimates of the action values, and use the minimum of these estimates in the Q-learning target. [1]

The overestimation bias may not always be detrimental. In some cases, overestimation bias can help encourage exploration for overestimated actions, which might be beneficial in highly stochastic areas if they correspond to high-value regions. On the other hand, if highly stochastic areas also have low values, overestimation bias might cause an agent to over-explore, and an underestimation bias may help prevent the situation.

Therefore, instead of fully moving towards underestimation, the method proposed in this paper tries to figure out how strongly we should correct for overestimation bias, and how to determine—or control—the level of bias .

The main contributions of this paper are: 1) It provides an easily-implemented method to determine overestimation or underestimation by simply tuning the parameters (i.e. the number of action-value estimates). 2) It introduces a new Generalized Q-learning framework to prove the convergence of *Maxmin Q-learning* as well as other Q-learning methods that use same number of action-value estimates as Maxmin Q-learning does.

Overall, the paper is interesting and the idea is novel. The theoretical analyses seem to be solid as well. And many other methods are listed to compare with Maxmin Q-learning, which makes the insight of this paper much clearer. I'm looking forward to a future work based on this paper which doesn't rely on parameter-tuning technique.

2 Problem Formulation

The paper formalizes the problem as a *Markov Decision Process (MDP)*, $(\mathcal{S}, \mathcal{A}, \mathbf{P}, r, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathbf{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the transition probabilities, $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward mapping, and $\gamma \in [0, 1]$ is the discount factor. At each time step t , the agent observes a state $S_t \in \mathcal{S}$ and takes an action $A_t \in \mathcal{A}$ and then transitions to a new state $S_{t+1} \in \mathcal{S}$ according to the transition probabilities \mathbf{P} and receives a scalar reward $R_{t+1} = r(S_t, A_t, S_{t+1}) \in \mathbb{R}$. The goal of the agent is to find a policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ that maximizes the expected return starting from some initial state.

Q-learning is an off-policy algorithm which attempts to learn the state-action values $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ for the optimal policy. It tries to solve for

$$Q^*(s, a) = \mathbb{E} \left[R_{t+1} + \max_{a' \in \mathcal{A}} Q^*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

The optimal policy is to act greedily with respect to these action values: from each s select a from

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

The update rule for an approximation Q for a sampled transition $(s_t, a_t, r_{t+1}, s_{t+1})$ is:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left(Y_t^Q - Q(s_t, a_t) \right) \quad \text{for } Y_t^Q \stackrel{\text{def}}{=} r_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q(s_{t+1}, a')$$

where α is the step-size.

The overestimations result from a positive bias that is introduced because Q-learning uses the maximum action value $\max_{a' \in \mathcal{A}} Q(s_{t+1}, a')$ as an approximation for the maximum expected action value [2]. For compactness, we write Q_{sa}^i instead of $Q^i(s, a)$. Each Q_{sa}^i has random approximation error e_{sa}^i .

$$Q_{sa}^i = Q_{sa}^* + e_{sa}^i$$

Assume that e_{sa}^i is a uniform random variable $U(-\tau, \tau)$ for some $\tau > 0$. The uniform random assumption was used by [3]. If e_{sa}^i is a constant instead of a random variable, then it has no influence on the optimal policy. Moreover, due to stochasticity, it must be a random variable. Thus, the uniform assumption is required. However, the tuning of τ is still remained to be concerned.

3 Theoretical Analysis

Maxmin Q-learning keeps N estimates of the action values, Q^i , and use the minimum of these estimates to update Q-learning target: $\max_{a'} \min_{i \in \{1, \dots, N\}} Q^i(s', a')$. The goal is to control overestimation by tuning N . It wants the overestimation decreases when N increases, and for some $N > 1$, this estimator switches from an overestimate, in expectation, to an underestimate. For $N = 1$, the update is simply Q-learning, and so likely has overestimation bias.

The full algorithm is summarized in Algorithm 1.

Algorithm 1: Maxmin Q-learning

Input : step size α , exploration parameter $\epsilon > 0$, number of action-value functions N

Initialize N action-value functions $\{Q^1, \dots, Q^N\}$ randomly

Initialize empty replay buffer D

Observe initial state s

while Agent is interacting with the Environment **do**

$Q^{\min}(s, a) \leftarrow \min_{k \in \{1, \dots, N\}} Q^k(s, a), \forall a \in \mathcal{A}$

Choose action a by ϵ -greedy based on Q^{\min}

Take action a , observe r, s'

Store transition (s, a, r, s') in D

Select a subset S from $\{1, \dots, N\}$ (e.e., randomly select on i to update)

for $i \in S$ **do**

Sample random mini-batch of transitions (s_D, a_D, r_D, s'_D) from D

Get update target: $Y^{MQ} \leftarrow r_D + \gamma \max_{a' \in \mathcal{A}} Q^{\min}(s'_D, a')$

Update action-value Q^i : $Q^i(s_D, a_D) \leftarrow Q^i(s_D, a_D) + \alpha [Y^{MQ} - Q^i(s_D, a_D)]$

end

$s \leftarrow s'$

end

Let M denotes the number of applicable actions at state s' . Then the estimation bias Z_{MN} for transition s, a, r, s' is defined as:

$$\begin{aligned} Z_{MN} &\stackrel{\text{def}}{=} \left(r + \gamma \max_{a'} Q_{s'a'}^{\min} \right) - \left(r + \gamma \max_{a'} Q_{s'a'}^* \right) \\ &= \gamma \left(\max_{a'} Q_{s'a'}^{\min} - \max_{a'} Q_{s'a'}^* \right) \end{aligned}$$

where

$$Q_{sa}^{\min} \stackrel{\text{def}}{=} \min_{i \in \{1, \dots, N\}} Q_{sa}^i = Q_{sa}^* + \min_{i \in \{1, \dots, N\}} e_{sa}^i$$

3.1 The Expected Estimation Bias And The Variance of Q

First of all, in Theorem 1, the paper wants to show that by controlling N in Maxmin Q-learning, the expectation of estimation bias $E[Z_{MN}]$ and the variance of Q_{sa}^{\min} will be affected.

Theorem 1. *Under the conditions stated above,*

(i) *the expected estimation bias is*

$$E[Z_{MN}] = \gamma\tau[1 - 2t_{MN}] \quad \text{where } t_{MN} = \frac{M(M-1) \cdots 1}{(M + \frac{1}{N})(M-1 + \frac{1}{N}) \cdots (1 + \frac{1}{N})}$$

$E[Z_{MN}]$ decreases as N increases: $E[Z_{M,N=1}] = \gamma\tau \frac{M-1}{M+1}$ and $E[Z_{M,N \rightarrow \infty}] = -\gamma\tau$.

(ii)

$$\text{Var}[Q_{sa}^{\min}] = \frac{4N\tau^2}{(N+1)^2(N+2)}$$

$\text{Var}[Q_{sa}^{\min}]$ decreases as N increases: $\text{Var}[Q_{sa}^{\min}] = \frac{\tau^2}{3}$ for $N = 1$ and $\text{Var}[Q_{sa}^{\min}] = 0$ for $N \rightarrow \infty$.

To prove it, the paper leverages the first lemma in [3].

Lemma 1. *Let X_1, \dots, X_N be N i.i.d. random variables from an absolutely continuous distribution with probability density function(PDF) $f(x)$ and cumulative distribution function (CDF) $F(x)$. Denote $\mu \stackrel{\text{def}}{=} E[X_i]$ and $\sigma^2 \stackrel{\text{def}}{=} \text{Var}[X_i] < +\infty$. Set $X_{1:N} \stackrel{\text{def}}{=} \min_{i \in \{1, \dots, N\}} X_i$ and $X_{N:N} \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, N\}} X_i$. Denote the PDF and CDF of $X_{1:N}$ as $f_{1:N}(x)$ and $F_{1:N}(x)$, respectively. Similarly, denote the PDF and CDF of $X_{N:N}$ as $f_{N:N}(x)$ and $F_{N:N}(x)$, respectively. We then have*

$$i. \mu - \frac{(N-1)\sigma}{\sqrt{2n-1}} \leq E[X_{1:N}] \leq \mu \text{ and } E[X_{1:N+1}] \leq E[X_{1:N}]$$

$$ii. F_{1:N}(x) = 1 - (1 - F(x))^N. f_{1:N}(x) = Nf(x)(1 - F(x))^{N-1}$$

$$iii. F_{N:N}(x) = (F(x))^N. f_{N:N}(x) = Nf(x)(F(x))^{N-1}$$

$$iv. \text{ If } X_1, \dots, X_N \sim U(-\tau, \tau), \text{ we have } \text{Var}(X_{1:N}) = \frac{4N\tau^2}{(N+1)^2(N+2)} \text{ and } \text{Var}(X_{1:N+1}) < \text{Var}(X_{1:N}) \leq \text{Var}(X_{1:1}) = \sigma^2 \text{ for any positive integer } N$$

Next, Theorem 1 is proved.

Proof. Let $f(x)$ and $F(x)$ be the cdf and pdf of e_{sa} , respectively. Similarly, Let $f_N(x)$ and $F_N(x)$ be the cdf and pdf of $\min_{i \in \{1, \dots, N\}} e_{sa}^i$. since e_{sa} is sampled from $\text{Uniform}(-\tau, \tau)$, it is easy to get $f(x) = \frac{1}{2\tau}$ and $F(x) = \frac{1}{2} + \frac{x}{2\tau}$. By Lemma 1, we have $f_N(x) = Nf(x)[1 - F(x)]^{N-1} =$

$\frac{N}{2\tau} \left(\frac{1}{2} - \frac{x}{2\tau}\right)^{N-1}$ and $F_N(x) = 1 - (1 - F(x))^N = 1 - \left(\frac{1}{2} - \frac{x}{2\tau}\right)^N$. The expectation of Z_{MN} is

$$\begin{aligned}
E[Z_{MN}] &= \gamma E \left[\left(\max_{a'} Q_{s'a'}^{\min} - \max_{a'} Q_{s'a'}^* \right) \right] \\
&= \gamma E \left[\max_{a'} \min_{i \in \{1, \dots, N\}} e_{sa'}^i \right] \\
&= \gamma \int_{-\tau}^{\tau} M x f_N(x) F_N(x)^{M-1} dx \\
&= \gamma \int_{-\tau}^{\tau} M N \frac{x}{2\tau} \left(\frac{1}{2} - \frac{x}{2\tau} \right)^{N-1} \left[1 - \left(\frac{1}{2} - \frac{x}{2\tau} \right)^N \right]^{M-1} dx \\
&= \gamma \int_{-\tau}^{\tau} x d \left[1 - \left(\frac{1}{2} - \frac{x}{2\tau} \right)^N \right]^M \\
&= \gamma \tau - \gamma \int_{-\tau}^{\tau} \left[1 - \left(\frac{1}{2} - \frac{x}{2\tau} \right)^N \right]^M dx \\
&= \gamma \tau \left[1 - 2 \int_0^1 (1 - y^N)^M dy \right] \quad \left(y \stackrel{\text{def}}{=} \frac{1}{2} - \frac{x}{2\tau} \right)
\end{aligned}$$

Let $t_{MN} = \int_0^1 (1 - y^N)^M dy$, so that $E[Z_{MN}] = \gamma \tau [1 - 2t_{MN}]$. Substitute y by t where $t = y^N$ then

$$\begin{aligned}
t_{MN} &= \frac{1}{N} \int_0^1 t^{\frac{1}{N}-1} (1-t)^M dt \\
&= \frac{1}{N} B \left(\frac{1}{N}, M+1 \right) \\
&= \frac{1}{N} \frac{\Gamma(M+1) \Gamma(\frac{1}{N})}{\Gamma(M + \frac{1}{N} + 1)} \\
&= \frac{\Gamma(M+1) \Gamma(1 + \frac{1}{N})}{\Gamma(M + \frac{1}{N} + 1)} \\
&= \frac{M(M-1) \cdots 1}{(M + \frac{1}{N})(M-1 + \frac{1}{N}) \cdots (1 + \frac{1}{N})}
\end{aligned}$$

Each term in the denominator decreases as N increases, because $1/N$ gets smaller. Therefore, $t_{M,N=1} = \frac{1}{M+1}$ and $t_{M,N \rightarrow \infty} = 1$. Using this, we conclude that $E[Z_{MN}]$ decreases as N increases and $E[Z_{M,N=1}] = \gamma \tau \frac{M-1}{M+1}$ and $E[Z_{M,N \rightarrow \infty}] = -\gamma \tau$.

By Lemma 1, the variance of Q_{sa}^{\min} is

$$\text{Var}[Q_{sa}^{\min}] = \frac{4N\tau^2}{(N+1)^2(N+2)}$$

$\text{Var}[Q_{sa}^{\min}]$ decreases as N increases. In particular, $\text{Var}[Q_{sa}^{\min}] = \frac{\tau^2}{3}$ for $N = 1$ and $\text{Var}[Q_{sa}^{\min}] = 0$ for $N \rightarrow \infty$ \square

Corollary 1. Assuming the n_{sa} samples are evenly allocated amongst the N estimators, then $\tau = \sqrt{3\sigma^2 N/n_{sa}}$ where σ^2 is the variance of samples for (s, a) and, for Q_{sa} the estimator that uses all n_{sa} samples for a single estimate,

$$\text{Var}[Q_{sa}^{\min}] = \frac{12N^2}{(N+1)^2(N+2)} \text{Var}[Q_{sa}]$$

Under this uniform random noise assumption, for $N \geq 8$, $\text{Var}[Q_{sa}^{\min}] < \text{Var}[Q_{sa}]$

Proof. Because Q_{sa}^i is a sample mean, its variance is $\sigma^2 N/n_{sa}$ where σ^2 is the variance of samples for (s, a) and its mean is Q_{sa}^* (because it is an unbiased sample average). Consequently, e_{sa} has mean zero and variance $\sigma^2 N/n_{sa}$. Because e_{sa} is a uniform random variable which has variance $\frac{1}{3}\tau^2$, we know that $\tau = \sqrt{3\sigma^2 N/n_{sa}}$. Plugging this value into the variance formula in Theorem 1 we get that

$$\begin{aligned}\text{Var}[Q_{sa}^{\min}] &= \frac{4N\tau^2}{(N+1)^2(N+2)} \\ &= \frac{12N^2\sigma^2/n_{sa}}{(N+1)^2(N+2)} \\ &= \frac{12N^2}{(N+1)^2(N+2)} \text{Var}[Q_{sa}]\end{aligned}$$

because $\text{Var}[Q_{sa}] = \sigma^2/n_{sa}$ for the sample average Q_{sa} that uses all the samples for one estimator. Easy to verify that for $N \geq 8$, $\text{Var}[Q_{sa}^{\min}] < \text{Var}[Q_{sa}]$ \square

3.2 Convergence Analysis of Maxmin Q-Learning

Secondly, the convergence of Maxmin Q-Learning in the tabular setting will be shown. Here, the paper provides *Generalized Q-learning*: Q-learning where the bootstrap target uses a function G of N action values. The pros to use this framework to proof the convergence is that it can be used in many other variants of Q-learning. But the cons is that choosing G function may not be easy. G needs to be limited so that it maintains relative maximum values, as stated in Assumption 1.

Assumption 1. (Conditions on G) Let $G : \mathbb{R}^{nNK} \mapsto \mathbb{R}$ and $G(Q) = q$ where $Q = (Q_a^{ij}) \in \mathbb{R}^{nNK}$, $a \in \mathcal{A}$ and $|\mathcal{A}| = n$, $i \in \{1, \dots, N\}$, $j \in \{0, \dots, K-1\}$ and $q \in \mathbb{R}$.

- i. If $Q_a^{ij} = Q_a^{kl}$, $\forall i, k, \forall j, l$, and $\forall a$, then $q = \max_a Q_a^{ij}$
- ii. $\forall Q, Q' \in \mathbb{R}^{nNK}$, $|G(Q) - G(Q')| \leq \max_{a,i,j} |Q_a^{ij} - Q_a'^{ij}|$

This assumption guarantees that even though information can be outdated, any old information is eventually discarded.

Proof. We can verify that Assumption 1 holds for Maxmin Q-learning. Set $K = 1$ and set N to be a positive integer. Let $Q_s = (Q_s^1, \dots, Q_s^N)$ and define $G^{MQ}(Q_s) = \max_{a \in \mathcal{A}} \min_{i \in \{1, \dots, N\}} Q_s^i$. It is easy to check that part i. of Assumption 1 is satisfied. Part ii. is also satisfied because

$$|G(Q_s) - G(Q'_s)| \leq \left| \max_a \min_i Q_{sa}^i - \max_{a'} \min_{i'} Q_{sa'}^{i'} \right| \leq \max_{a,i} |Q_{sa}^i - Q_{sa}^{i'}|$$

\square

The target action-value of Generalized Q-learning Y^{GQ} is defined based on action-value estimates from both dimensions:

$$Y^{GQ} = r + \gamma Q_{s'}^{GQ}(t-1)$$

where t is the current time step and the action-value function $Q_s^{GQ}(t)$ is a function of $Q_s^1(t-K), \dots, Q_s^1(t-1), \dots, Q_s^N(t-K), \dots, Q_s^N(t-1)$:

$$Q_s^{GQ}(t) = G \begin{pmatrix} Q_s^1(t-K) & \dots & Q_s^1(t-1) \\ Q_s^2(t-K) & \dots & Q_s^2(t-1) \\ \vdots & \ddots & \vdots \\ Q_s^N(t-K) & \dots & Q_s^N(t-1) \end{pmatrix}$$

For simplicity, the vector $(Q_{sa}^{GQ}(t))_{a \in \mathcal{A}}$ is denoted as $Q_s^{GQ}(t)$, same for $Q_s^i(t)$. The corresponding update rule is

$$Q_{sa}^i(t) \leftarrow Q_{sa}^i(t-1) + \alpha_{sa}^i(t-1) (Y^{GQ} - Q_{sa}^i(t-1))$$

Assumption 2. (Conditions on the step-sizes) There exists some (deterministic) constant C such that for every $(s, a) \in \mathcal{S} \times \mathcal{A}$, $i \in \{1, \dots, N\}$, $0 \leq \alpha_{sa}^i(t) \leq 1$, and with probability 1,

$$\sum_{t=0}^{\infty} (\alpha_{sa}^i(t))^2 \leq C, \quad \sum_{t=0}^{\infty} \alpha_{sa}^i(t) = \infty$$

This assumption is needed on the step-sizes employed by the Generalized Q -learning algorithm which is standard for stochastic approximation algorithms. In particular, it requires that every state-action pair (s, a) is simulated an infinite number of times.

Theorem 2. Assume a finite MDP $(\mathcal{S}, \mathcal{A}, P, R)$ and that Assumption 1 and 2 hold. Then the action-value functions in Generalized Q -learning will converge to the optimal action-value function with probability 1, in either of the following cases: (i) $\gamma < 1$ or (ii) $\gamma = 1, \forall a \in \mathcal{A}, Q_{s_1 a}^i(t=0) = 0$ where s_1 is an absorbing state and all policies are proper.

Proof. To proof Theorem 2, some assumptions and theorems proposed in [4] are required.

Assumption 3. For any i and j , $\lim_{t \rightarrow \infty} \tau_j^i(t) = \infty$, with probability 1 where each $\tau_j^i(t)$ is an integer satisfying $0 \leq \tau_j^i(t) \leq t$.

This assumption of parameter is involved in the following assumptions.

Assumption 4. Let $\{\mathcal{F}(t)\}_{t=0}^{\infty}$ be an increasing sequence of subfields of \mathcal{F}

- i. $x(0)$ is $\mathcal{F}(0)$ -measurable.
- ii. For every i and t , $w_i(t)$ is $\mathcal{F}(t+1)$ -measurable.
- iii. For every i, j and t , $\alpha_i(t)$ and $\tau_j^i(t)$ are $\mathcal{F}(t)$ -measurable.
- iv. For every i and t , we have $E[w_i(t) \mid \mathcal{F}(t)] = 0$
- v. There exist (deterministic) constants A and B such that

$$E[w_i^2(t) \mid \mathcal{F}(t)] \leq A + B \max_j \max_{\tau \leq t} |x_j(\tau)|^2, \quad \forall i, t$$

The measure-theoretic terminology that "a random variable Z is $\mathcal{F}(t)$ -measurable" intuitively means that Z is completely determined by the history represented by $\mathcal{F}(t)$.

Assumption 5. There exists a vector $x^* \in \mathbb{R}^n$, a positive vector v , and a scalar $\beta \in [0, 1)$, such that

$$\|F(x) - x^*\|_v \leq \beta \|x - x^*\|_v, \quad \forall x \in \mathbb{R}^n$$

Theorem 3. Let Assumptions 3, 4, 5, 2 hold. Then $x(t)$ converges to x^* with probability 1 Detailed proofs can be found in the work of [4].

If Assumptions 3, 4, 5, 2 are all satisfied, then convergence is guaranteed according to Theorem 3. \square

4 Conclusion

It is inspiring that this paper not only proposed a novel method to control overestimation bias, it also provides an elegant framework to help prove convergence with theoretically support. However, the tuning of parameters requires lots of efforts. Therefore, the future work should investigate how best to select N .

References

- [1] Qingfeng Lan, Yangchen Pan, Alona Fyshe, and Martha White. Maxmin q-learning: Controlling the estimation bias of q-learning. 02 2020.
- [2] Hado Van Hasselt. Double q-learning. pages 2613–2621, 01 2010.
- [3] Sebastian Thrun and Anton Schwartz. Issues in using function approximation for reinforcement learning. 10 1993.
- [4] Dimitri Bertsekas and John Tsitsiklis. Parallel and distributed computation : numerical methods / dimitri p. bertsekas, john n. tsitsiklis. *SERBIULA (sistema Librum 2.0)*, 23, 01 1989.