



**HAL**  
open science

# Reinforcement Learning with a Near Optimal Rate of Convergence

Mohammad Gheshlaghi Azar, Rémi Munos, Mohammad Ghavamzadeh,  
Hilbert Kappen

► **To cite this version:**

Mohammad Gheshlaghi Azar, Rémi Munos, Mohammad Ghavamzadeh, Hilbert Kappen. Reinforcement Learning with a Near Optimal Rate of Convergence. [Technical Report] 2011. inria-00636615v1

**HAL Id: inria-00636615**

**<https://inria.hal.science/inria-00636615v1>**

Submitted on 27 Oct 2011 (v1), last revised 29 Nov 2011 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Reinforcement Learning with a Near Optimal Rate of Convergence

**Mohammad Gheshlaghi Azar**

M.AZAR@SCIENCE.RU.NL

*Department of Biophysics*

*Radboud University Nijmegen*

*6525 EZ Nijmegen, The Netherlands*

**Rémi Munos**

REMI.MUNOS@INRIA.FR

**Mohammad Ghavamzadeh**

MOHAMMAD.GHAVAMZADEH@INRIA.FR

*INRIA Lille, SequeL Project*

*40 avenue Halley*

*59650 Villeneuve d'Ascq, France*

**Hilbert J. Kappen**

B.KAPPEN@SCIENCE.RU.NL

*Department of Biophysics*

*Radboud University Nijmegen*

*6525 EZ Nijmegen, The Netherlands*

**Editor:** TBA

## Abstract

We consider the problem of model-free reinforcement learning in the Markovian decision processes (MDP) under the PAC (“probably approximately correct”) model. We introduce a new variant of Q-learning, called speedy Q-learning (SQL), to address the problem of the slow convergence in the standard Q-learning algorithm, and prove PAC bounds on the performance of SQL. The bounds show that for any MDP with  $n$  state-action pairs and the discount factor  $\gamma \in [0, 1)$  a total of  $O(n \log(n/\delta)/((1-\gamma)^4 \epsilon^2))$  step suffices for the SQL algorithm to converge to an  $\epsilon$ -optimal action-value function with probability  $1 - \delta$ . We also establish a lower-bound of  $\Omega(n \log(1/\delta)/((1-\gamma)^2 \epsilon^2))$  for all reinforcement learning algorithms, which matches the upper bound in terms of  $\epsilon$ ,  $\delta$  and  $n$  (up to a logarithmic factor). Further, our results have better dependencies on  $\epsilon$  and  $1 - \gamma$ , and thus, are tighter than the best available results for Q-learning. SQL also improves on existing results for the batch Q-value iteration, so far considered to be more efficient than the incremental methods like Q-learning.

## 1. Introduction

The Markovian decision process (MDP) problem is a classical problem in the fields of operations research and decision theory. When an explicit model of the MDP, transition probabilities and reward function, is known, one can rely on dynamic programming (DP) Bellman (1957) algorithms such as value iteration or policy iteration (see, e.g., Bertsekas, 2007a; Puterman, 1994), to compute the optimal policy. For example, value iteration computes the optimal value function by successive iterations of the Bellman operator. One can

show that, in the discounted infinite-horizon setting, the convergence of value iteration is exponentially fast since the Bellman operator is a contraction mapping (Bertsekas, 2007b). However, DP relies on an explicit knowledge of the MDP. In many real world problems the transition probabilities are not initially known, but one may observe transition samples using Monte-Carlo sampling, either as a single trajectory (rollout) obtained by following an exploration policy, or by simulating independent transition samples anywhere in the space by resorting to a *generative model*. The field of reinforcement learning (RL) is concerned with the problem of approximating the optimal policy, or the optimal value function, from the observed reward and transition samples (Szepesvári, 2010; Sutton and Barto, 1998).

One may characterize RL methods as model-based or model-free. In model-based RL we first learn a model of the MDP and then we use this model for computing an approximation of the value functions by dynamic programming techniques. Model-free methods, in contrast, compute directly an approximation of the value function by making use of a sample-based estimate of the Bellman operator, without resorting to learning a model. Q-learning (Watkins, 1989) is a well-known model-free reinforcement learning (RL) algorithm that, incrementally, finds an estimate of the optimal action-value function. The Q-learning algorithm can be seen as a combination of the value iteration algorithm and stochastic approximation. In finite state-action problems, it has been shown that Q-learning converges to the optimal action-value function (Jaakkola et al., 1994; Bertsekas and Tsitsiklis, 1996). However, it suffers from slow convergence, especially when the discount factor  $\gamma$  is close to one (Szepesvári, 1997; Even-Dar and Mansour, 2003). The main reason for the slow convergence of Q-learning is the combination of the sample-based stochastic approximation, which makes use of a decaying learning rate, and the fact that the Bellman operator propagates information throughout the whole space, especially when  $\gamma$  is close to 1.

In this paper, we focus on RL problems that are formulated as finite state-action discounted infinite horizon Markov decision processes (MDPs), and propose an algorithm, called *speedy Q-learning* (SQL), that addresses the problem of slow convergence of Q-learning. At each time step, SQL uses two successive estimates of the action-value function that makes its space complexity twice as the standard Q-learning. However, this allows SQL to use a more aggressive learning rate for one of the terms in its update rule and eventually achieves a faster convergence rate than the standard Q-learning, see Section 3.1 for a more detailed discussion.

Similar to Q-learning, SQL can be implemented in synchronous and asynchronous fashions. In the synchronous case, we update the action-value function of all state-action pairs simultaneously by resorting to a generative model. In the asynchronous case, in contrast, the set of samples is a sequence of state-action pairs generated by following a policy and at each step of learning the action-value function of only one state-action pair in the sequence is updated. In this paper, we report PAC (“probably approximately correct”) bounds on the performance of SQL for both synchronous and asynchronous SQL. The bounds show that only  $T = O(1/((1 - \gamma)^4 \epsilon^2))$  number of trials are required for SQL in order to guarantee an  $\epsilon$ -optimal action-value function with high probability. This is superior to the best result for the standard Q-learning by Even-Dar and Mansour (2003), both in terms of  $1/\epsilon$  and  $1/(1 - \gamma)$ . Also, our results for the asynchronous case, scales significantly better than the standard Q-learning with the *cover time* of MDP  $L$ , which measures the difficulty of exploring an MDP by following a policy. The rate for SQL is even better than that for the

*Phased Q-learning* algorithm, a model-free batch Q-value iteration algorithm proposed and analyzed by Kearns and Singh (1999). In addition, SQL’s rate is slightly better than the rate of the model-based batch Q-value iteration algorithm in Kearns and Singh (1999) and has a better computational and memory requirement, see Section 3.3.2 for more detailed comparisons.

The idea of using previous estimates of the action-values has already been used in order to improve the performance of Q-learning. A popular algorithm of this kind is  $Q(\lambda)$  (Peng and Williams, 1996; Watkins, 1989), which incorporates the concept of eligibility traces in Q-learning, and has been empirically shown to have a better performance than Q-learning, i.e.,  $Q(0)$ , for suitable values of  $\lambda$ . Another recent work in this direction is *Double Q-learning* (van Hasselt, 2010), which uses two estimators for the action-value function in order to alleviate the over-estimation of action-values in Q-learning. This over-estimation is caused by a positive bias introduced by using the maximum action value as an approximation for the maximum expected action value (van Hasselt, 2010).

The rest of the paper is organized as follows. After introducing the notations used in the paper in Section 2, we present our *Speedy Q-learning* algorithm in Section 3. We first describe the synchronous and asynchronous version of the algorithm in Section 3.1, then state our main theoretical result, i.e., high-probability bounds on the performance of SQL as well as a new lower bound for the sample complexity of reinforcement learning, in Section 3.2, and finally compare our bound with the previous results on Q-learning and Q-value iteration in Section 3.3. In Section 4, we evaluate the performance of SQL, numerically, on different problem domains. Section 5 contains the detailed proof of the performance bound of the SQL algorithm and the RL lower bound. Finally, we conclude the paper and discuss some future directions in Section 6.

## 2. Preliminaries

In this section, we introduce some concepts and definitions from the theory of Markov decision processes (MDPs) and stochastic processes that are used throughout the paper. We start by the definition of supremum norm ( $\ell_\infty$  norm). For a real-valued function  $g : \mathcal{Y} \mapsto \mathbb{R}$ , where  $\mathcal{Y}$  is a finite set, the supremum norm of  $g$  is defined as  $\|g\| \triangleq \max_{y \in \mathcal{Y}} |g(y)|$ .

We consider the standard reinforcement learning (RL) framework (Bertsekas and Tsitsiklis, 1996; Sutton and Barto, 1998) in which a learning agent interacts with a stochastic environment and this interaction is modeled as a discrete-time discounted MDP. A discounted MDP is a quintuple  $(\mathcal{X}, \mathcal{A}, P, \mathcal{R}, \gamma)$ , where  $\mathcal{X}$  and  $\mathcal{A}$  are the set of states and actions,  $P$  is the state transition distribution,  $\mathcal{R}$  is the reward function, and  $\gamma \in (0, 1)$  is a discount factor. We also denote the horizon of MDP by  $\beta$  defined by  $\beta = 1/(1 - \gamma)$ . We denote by  $P(\cdot|x, a)$  and  $r(x, a)$  the probability distribution over the next state and the immediate reward of taking action  $a$  at state  $x$ , respectively. To keep the representation succinct, we use  $\mathcal{Z}$  for the joint state-action space  $\mathcal{X} \times \mathcal{A}$ .

**Assumption A1 (MDP regularity)** *We assume  $\mathcal{Z}$  is a finite set with the cardinality  $n$ . We also assume that the immediate rewards  $r(x, a)$  are in the interval  $[0, 1]$ .*

A policy kernel  $\pi(\cdot|\cdot)$  determines the distribution of the control action given the past observations. The policy is called stationary if the distribution of the control action just

depends on the last state  $x$ . It is deterministic if this distribution concentrates over a single action. The *value* and the *action-value functions* of a policy  $\pi$ , denoted respectively by  $V^\pi : \mathcal{X} \mapsto \mathbb{R}$  and  $Q^\pi : \mathcal{Z} \mapsto \mathbb{R}$ , are defined as the expected sum of discounted rewards that are encountered when the policy  $\pi$  is executed. Given a MDP, the goal is to find a policy that attains the best possible values,  $V^*(x) \triangleq \sup_\pi V^\pi(x)$ ,  $\forall x \in \mathcal{X}$ . Function  $V^*$  is called the *optimal value function*. Similarly the *optimal action-value function* is defined as  $Q^*(x, a) = \sup_\pi Q^\pi(x, a)$ ,  $\forall (x, a) \in \mathcal{Z}$ . The optimal action-value function  $Q^*$  is the unique fixed-point of the *Bellman optimality operator*  $\mathcal{T}$  defined as:

$$(\mathcal{T}Q)(x, a) \triangleq r(x, a) + \gamma \sum_{y \in \mathcal{X}} P(y|x, a)(\mathcal{M}Q)(y), \quad \forall (x, a) \in \mathcal{Z}.$$

where the max operator  $\mathcal{M}$  over action-value functions is defined as  $(\mathcal{M}Q)(y) = \max_{a \in \mathcal{A}} Q(y, a)$ ,  $\forall y \in \mathcal{X}$ . It is important to note that  $\mathcal{T}$  is a  $\gamma$ -contraction mapping w.r.t. to the  $\ell_\infty$ -norm, i.e., for any pair of action-value functions  $Q$  and  $Q'$ , we have  $\|\mathcal{T}Q - \mathcal{T}Q'\| \leq \gamma \|Q - Q'\|$  (Bertsekas, 2007b, Chap. 1).

Finally we define the cover time of MDP under the policy  $\pi$  as follows:

**Definition 1 (Cover time of MDP)** *Let  $t$  be an integer. Define  $\tau_\pi(x, t)$ , the cover time of the MDP under the (non-stationary) policy  $\pi$ , as the minimum number of steps required to visit all state-action pairs  $(x, a) \in \mathcal{Z}$  starting from state  $x \in \mathcal{X}$  at time-step  $0 \leq t$ . Also, the state-action space  $\mathcal{Z}$  is covered by the policy  $\pi$  if all the state-action pairs are performed at least once under the policy  $\pi$ .*

The following assumption which bounds the expected cover time of the MDP guarantees that, asymptotically, all the state-action pairs are visited infinitely many times under the policy  $\pi$ .

**Assumption A2 (Boundedness of the the expected cover time)** *Let  $0 < L < \infty$  and  $t$  be a integer. Then, under the policy  $\pi$  for all  $x \in \mathcal{X}$  and  $t > 0$  assume that:*

$$\mathbb{E}(\tau_\pi(x, t)) \leq L.$$

### 3. Speedy Q-Learning

In this section, we introduce a new RL method, called speedy Q-Learning (SQL), derive performance bounds for the asynchronous and synchronous variant of SQL, and compare these bounds with similar results on standard Q-learning. The derived performance bound shows that SQL has a rate of convergence of  $T = O(1/\epsilon^2)$ , which is optimal in the sense that it matches the lower bound of reinforcement learning.

#### 3.1 Algorithms

In this subsection, we introduce two variants of SQL algorithms, the synchronous SQL and asynchronous SQL. In the asynchronous version, at each time step, the action-value of only

one state-action pair, the current observed state-action, is updated, while the action-values of rest of the state-action pairs remain unchanged. For the convergence of this instance of the algorithm, it is required that all the states and actions are visited infinitely many times, which makes the analysis slightly more complicated. On the other hand, given a generative model, the algorithm may be also formulated in a synchronous fashion, in which we first generate a next state  $y \sim P(\cdot|x, a)$  for each state-action pair  $(x, a)$ , and then update the action-values of all the state-action pairs using these samples. The pseudo-code of the synchronous and asynchronous SQL are shown in Algorithm 1 and 2, respectively. One can show that asynchronous SQL is reduced to Algorithm 2 when the cover time  $\tau_\pi(x, t) = n$  for all  $x \in \mathcal{X}$  and  $t \geq 0$ , in which case the action-values of all state-action pairs are updated in a row. In other words, the Algorithm 1 can be seen as a special case of Algorithm 2. Therefore, in the sequel we only describe the more general asynchronous SQL algorithm.

As it can be seen from the update rule of Algorithm 2, at each time step, Algorithm 2 keeps track of the action-value functions of the two most recent iterations  $Q_k$  and  $Q_{k+1}$ , and its main update rule is of the following form for all  $(x, a) \in \mathcal{Z}$  at time step  $t$  and the iteration round  $k$ :

$$Q_{k+1}(X_t, A_t) = Q_k(X_t, A_t) + \alpha_k(\mathcal{T}_k Q_{k-1}(X_t, A_t) - Q_k(X_t, A_t)) + (1 - \alpha_k)(\mathcal{T}_k Q_k(X_t, A_t) - \mathcal{T}_k Q_{k-1}(X_t, A_t)), \quad (1)$$

where  $\mathcal{T}_k Q(X_t, A - t) = 1/|\mathcal{Y}_k| \sum_{y \in \mathcal{Y}_k} [r(X_t, A_t) + \gamma \mathcal{M}Q(y)]$  is the empirical Bellman optimality operator using the set of next state samples  $\mathcal{Y}_k = \{y|y \sim P(\cdot|X_t, A_t)\}$ , generated up to time step  $t$  at round  $k$  and  $|\mathcal{Y}_k|$  is the cardinality of  $\mathcal{Y}_k$ . At each time step  $t$ , Algorithm 2 works as follows: **(i)** It simulates the MDP for one-step, i.e., it draws the state-action pair  $(X_t, A_t) \in \mathcal{Z}$  from the distribution  $\mu$  and then make a transition to a new state  $y_k$ . **(ii)** it updates the two sample estimates  $\mathcal{T}_k Q_{k-1}(X_t, A_t)$  and  $\mathcal{T}_k Q_k(X_t, A_t)$  of the Bellman optimality operator (for state-action pair  $(X_t, A_t)$  using the next state  $y_k$ ) applied to the estimates  $Q_{k-1}$  and  $Q_k$  of the action-value function at the previous and current round  $k-1$  and  $k$ , **(iii)** it updates the action-value function of  $(X_t, A_t)$ , generates  $Q_{k+1}(X_t, A_t)$ , using the update rule of Eq. 1, and finally **(iv)** we check for the condition that if all  $(x, a) \in \mathcal{Z}$  have been visited at least one time at iteration  $k$ . If the condition is satisfied then we move to next round  $k+1$ , otherwise  $k$  remains unchanged. Moreover, we let  $\alpha_k$  decays linearly with the number of iterations  $k$ , i.e.,  $\alpha_k = 1/(k+1)$ , in Algorithm 2. Also, we notice that the update rule  $\mathcal{T}_k Q_k(X_t, A_t) := (1 - \eta_N)\mathcal{T}_k Q_k(X_t, A_t) + \eta_N(r(X_t, A_t) + \gamma \mathcal{M}Q_k(y_k))$  is being used to make an unbiased estimate of  $\mathcal{T}Q_k$  in an incremental fashion.

Now, let us consider the update rule of the standard Q-learning

$$Q_{k+1}(x, a) = Q_k(x, a) + \alpha_k(\mathcal{T}_k Q_k(x, a) - Q_k(x, a)),$$

which may be rewritten as

$$Q_{k+1}(x, a) = Q_k(x, a) + \alpha_k(\mathcal{T}_k Q_{k-1}(x, a) - Q_k(x, a)) + \alpha_k(\mathcal{T}_k Q_k(x, a) - \mathcal{T}_k Q_{k-1}(x, a)). \quad (2)$$

**Algorithm 1:** Synchronous speedy Q-learning

---

**Input:** Initial action-values  $Q_0$ , discount factor  $\gamma$  and number of steps  $T$

$Q_{-1} := Q_0;$  // Initialization  
 $t := k := 0;$

**repeat** // Main loop

$\alpha_k := \frac{1}{k+1};$

**foreach**  $(x, a) \in \mathcal{Z}$  **do** // Update the action-values for all  $(x, a) \in \mathcal{Z}$

Generate the next state sample  $y_k \sim P(\cdot|x, a);$

$\mathcal{J}_k Q_{k-1}(x, a) := r(x, a) + \gamma \mathcal{M} Q_{k-1}(y_k);$

$\mathcal{J}_k Q_k(x, a) := r(x, a) + \gamma \mathcal{M} Q_k(y_k);$

$Q_{k+1}(x, a) := Q_k(x, a) + \alpha_k (\mathcal{J}_k Q_{k-1}(x, a) - Q_k(x, a)) + (1 - \alpha_k) (\mathcal{J}_k Q_k(x, a) - \mathcal{J}_k Q_{k-1}(x, a));$

// SQL update rule

$t := t + 1;$

**end**

$k := k + 1;$

**until**  $t \geq T;$

**return**  $Q_k$

---

**Algorithm 2:** Asynchronous speedy Q-learning

---

**Input:** Initial action-values  $Q_0$ , The policy  $\pi(\cdot|·)$ , discount factor  $\gamma$ , number of step  $T$  and the initial state  $X_0$ .

$t := k := 0;$  // Initialization  
 $\alpha_0 = 1;$

**foreach**  $(x, a) \in \mathcal{Z}$  **do**

$Q_{-1}(x, a) := Q_0(x, a);$

$N_0(x, a) := 0;$

**end**

**repeat** // Main loop

Draw the action  $A_t \sim \pi(\cdot|X_t);$

Generate the next state sample  $y_k$  by simulating  $P(\cdot|X_t, A_t);$

$\eta_N := \frac{1}{N_k(X_t, A_t) + 1};$

**if**  $N_k(x, a) > 0$  **then**

$\mathcal{J}_k Q_{k-1}(X_t, A_t) := (1 - \eta_N) \mathcal{J}_k Q_{k-1}(X_t, A_t) + \eta_N (r(X_t, A_t) + \gamma \mathcal{M} Q_{k-1}(y_k));$

$\mathcal{J}_k Q_k(X_t, A_t) := (1 - \eta_N) \mathcal{J}_k Q_k(X_t, A_t) + \eta_N (r(X_t, A_t) + \gamma \mathcal{M} Q_k(y_k));$

**else**

$\mathcal{J}_k Q_{k-1}(X_t, A_t) := r(X_t, A_t) + \gamma \mathcal{M} Q_{k-1}(y_k);$

$\mathcal{J}_k Q_k(X_t, A_t) := r(X_t, A_t) + \gamma \mathcal{M} Q_k(y_k);$

**end**

$Q_{k+1}(X_t, A_t) :=$

$Q_k(X_t, A_t) + \alpha_k (\mathcal{J}_k Q_{k-1}(X_t, A_t) - Q_k(X_t, A_t)) + (1 - \alpha_k) (\mathcal{J}_k Q_k(X_t, A_t) - \mathcal{J}_k Q_{k-1}(X_t, A_t));$

// SQL update rule

$N_k(X_t, A_t) := N_k(X_t, A_t) + 1;$

$X_{t+1} = y_k;$

**if**  $\min_{(x,a) \in \mathcal{Z}} N_k(x, a) > 0$  **then** // Check if all  $(x, a) \in \mathcal{Z}$  have been visited at round  $k$

$k := k + 1;$

$\alpha_k := \frac{\alpha_0}{k+1};$

**foreach**  $(x, a) \in \mathcal{Z}$  **do**  $N_k(x, a) := 0;$

**end**

$t := t + 1;$

**until**  $t \geq T;$

**return**  $Q_k$

---



Comparing the Q-learning update rule of Eq. 2 with the one for SQL in Eq. 1, we first notice that the same terms:  $\mathcal{T}_k Q_{k-1} - Q_k$  and  $\mathcal{T}_k Q_k - \mathcal{T}_k Q_{k-1}$  appear on the right hand side of the update rules of both algorithms. However, while Q-learning uses the same conservative learning rate  $\alpha_k$  for both these terms, SQL uses  $\alpha_k$  for the first term and a bigger learning step  $1 - \alpha_k = k/(k+1)$  for the second one. Since the term  $\mathcal{T}_k Q_k - \mathcal{T}_k Q_{k-1}$  goes to zero as  $Q_k$  approaches its optimal value  $Q^*$ , it is not necessary that its learning rate approaches zero. As a result, using the learning rate  $\alpha_k$ , which goes to zero with  $k$ , is too conservative for this term. This might be a reason why SQL that uses a more aggressive learning rate  $1 - \alpha_k$  for this term has a faster convergence rate than Q-learning.

### 3.2 Main Theoretical Results

The main theoretical results of this paper are expressed as high-probability bounds over the performance of the SQL algorithms for both synchronous and asynchronous cases. We also report a lower bound on the number of transition samples required, for every reinforcement learning algorithm, to achieve an  $\epsilon$ -optimal estimate of  $Q^*$  with high probability.<sup>1</sup>

**Theorem 2** *Let A1 hold,  $T$  be a positive integer and  $Q_T$  be the estimate of  $Q^*$  by Algorithm 1 at time step  $T$ . Then, the uniform approximation error*

$$\|Q^* - Q_T\| \leq \beta^2 \left[ \frac{\gamma n}{T} + \sqrt{\frac{2n \log \frac{2n}{\delta}}{T}} \right],$$

with probability (w.p.) at least  $1 - \delta$ .

**Theorem 3** *Let A1 and A2 hold and  $T > 0$  be the number of time steps. Then, at step  $T$  of Algorithm 2 w.p. at least  $1 - \delta$*

$$\|Q^* - Q_T\| \leq \beta^2 \left[ \frac{\gamma e L \log \frac{2}{\delta}}{T} + \sqrt{\frac{2e L \log \left(\frac{2}{\delta}\right) \log \frac{4n}{\delta}}{T}} \right].$$

These results, combined with Borel-Cantelli lemma (Feller, 1968), guarantee that  $Q_T$  converges almost surely to  $Q^*$  with the rate  $\sqrt{1/T}$  for both Algorithm 1 and 2. Further, the following PAC bounds which quantify the number of steps  $T$  required to reach the error  $\epsilon > 0$  in estimating the optimal action-value function, w.p.  $1 - \delta$ , are immediate consequences of Theorem 2 and 3, respectively.

**Corollary 4 (Finite-time PAC bound of synchronous SQL)** *Under A1, after*

$$T = \frac{3.74n\beta^4 \log \frac{2n}{\delta}}{\epsilon^2}$$

steps (transitions) the uniform approximation error of Algorithm 1  $\|Q^* - Q_T\| \leq \epsilon$ , w.p. at least  $1 - \delta$ .

---

1. We report the proofs in Section 5.



**Corollary 5 (Finite-time PAC bound of asynchronous SQL)** *Under A1 and A2, after*

$$T = \frac{3.74eL\beta^4 \log\left(\frac{2}{\delta}\right) \log\frac{4n}{\delta}}{\epsilon^2}$$

*steps (transitions) the uniform approximation error of Algorithm 2  $\|Q^* - Q_T\| \leq \epsilon$ , w.p. at least  $1 - \delta$ .*

The following general result provides a lower bound on the number of transitions  $T$  for every reinforcement learning algorithm to achieve  $\epsilon$ -optimal performance w.p.  $1 - \delta$  under the assumption that the reinforcement learning algorithm is  $(\epsilon, \delta)$ -correct on the class of MDPs  $\mathbb{M}$ .

**Definition 6 ( $(\epsilon, \delta)$ -correct reinforcement learning)** *Define  $\mathbb{A}$  as the class of all reinforcement learning algorithms which rely on estimating the action-value function  $Q^*$ . Let  $Q_t^A$  be the estimate of  $Q^*$  using  $t \geq 0$  number of transition samples under the algorithm  $A \in \mathbb{A}$ . We then call  $A \in \mathbb{A}$   $(\epsilon, \delta)$ -correct on the class of MDP  $\mathbb{M}$ , if there exists some  $0 < T < +\infty$  such that, for all  $M \in \mathbb{M}$ ,  $\|Q^* - Q_t^A\| \leq \epsilon$  w.p. at least  $1 - \delta$  iff  $t \geq T$ .*

**Theorem 7 (Lower bound on the sample complexity of reinforcement learning)** *There exists some  $\epsilon_0, \delta_0, \gamma_0, c_1, c_2$  and a class of MDPs  $\mathbb{M}$  such that for all  $\epsilon \in (0, \epsilon_0)$ ,  $\delta \in (0, \delta_0)$  and  $\gamma \in (\gamma_0, 1)$  and for every  $(\epsilon, \delta)$ -correct reinforcement learning algorithm  $A \in \mathbb{A}$  on the class of MDPs  $\mathbb{M}$  the number of transitions*

$$T \geq \frac{n\beta^2}{c_1\epsilon^2} \log \frac{1}{c_2\delta}.$$

### 3.3 Relation to Existing Results

In this section, we first compare our results for SQL with the existing results on the convergence of standard Q-learning. This comparison indicates that SQL accelerates the convergence of Q-learning, especially for  $\gamma$  close to 1 and small  $\epsilon$ . We then compare SQL with batch Q-value iteration (QVI) in terms of sample and computational complexities, i.e., the number of samples and the computational cost required to achieve an  $\epsilon$ -optimal solution w.p.  $1 - \delta$ , as well as space complexity, i.e., the memory required at each step of the algorithm.

#### 3.3.1 A COMPARISON WITH THE CONVERGENCE RATE OF STANDARD Q-LEARNING

There are not many studies in the literature concerning the convergence rate of incremental model-free RL algorithms such as Q-learning. Szepesvári (1997) has provided the asymptotic convergence rate for Q-learning under the assumption that all the states have the same next state distribution. This result shows that the asymptotic convergence rate of Q-learning has an exponential dependency on  $1/\beta$ .

Even-Dar and Mansour (2003) thoroughly investigated the finite-time behavior of synchronous Q-learning for different time scales. Their main result indicates that by using the polynomial learning step  $\alpha_k = 1/(k+1)^\omega$ ,  $0.5 < \omega < 1$ , the asynchronous variant of

Q-learning achieves  $\epsilon$ -optimal performance w.p. at least  $1 - \delta$  after

$$T = O \left( n \left[ \left( \frac{\beta^4 \log \frac{n\beta}{\delta\epsilon}}{\epsilon^2} \right)^{\frac{1}{\omega}} + \left( \beta \log \frac{\beta}{\epsilon} \right)^{\frac{1}{1-\omega}} \right] \right), \quad (3)$$

steps. When  $\gamma \approx 1$ , one can argue that  $\beta$  becomes the dominant term in the bound of Eq. 3, and thus, the optimized bound w.r.t.  $\omega$  is obtained for  $\omega = 4/5$  and is of order  $\tilde{O}(\beta^5/\epsilon^{2.5})$ . On the other hand, SQL is guaranteed to achieve the same precision in only  $O(\beta^4/\epsilon^2)$  steps. The difference between these two bounds is significant for  $\gamma$ 's close to 1.

Even-Dar and Mansour (2003) also proved bounds for the asynchronous variant of Q-learning in the case that the cover time of MDP can be uniformly bounded from above by some finite constant. The extension of their results to the more realistic case that the expected value of the cover-time is bounded by some  $L > 0$  (assumption A2) leads to the following PAC bound:

**Proposition 8 (Even-Dar and Mansour, 2003)** *Under A1 and A2, for all  $\omega \in (0.5, 1)$ , after*

$$T = O \left( \left[ \frac{(L \log \frac{1}{\delta})^{1+3\omega} \beta^4 \log \frac{n\beta}{\delta\epsilon}}{\epsilon^2} \right]^{\frac{1}{\omega}} + \left[ L \beta \log \frac{1}{\delta} \log \frac{\beta}{\epsilon} \right]^{\frac{1}{1-\omega}} \right)$$

*steps (transitions) the uniform approximation error of asynchronous  $\|Q^* - Q_T\| \leq \epsilon$ , w.p. at least  $1 - \delta$ .*

The dependence on  $L$  of this algorithm is  $O(L^{3+\frac{1}{\omega}} + L^{\frac{1}{1-\omega}})$ , which leads, with the choice of  $\omega \approx 0.77$ , to the optimized dependency of  $O(L^{4.34})$ , whereas asynchronous SQL achieves the same accuracy after just  $O(L)$  steps. This result shows that for MDPs with large expected cover-time, i.e., slow mixing MDPs, asynchronous SQL may converge substantially faster to a near-optimal solution than its Q-learning counterpart.

### 3.3.2 SQL vs. Q-VALUE ITERATION

Finite sample bounds for both model-based and model-free (Phased Q-learning) QVI have been derived in Kearns and Singh (1999) and Even-Dar et al. (2002). These algorithms can be considered as the batch version of Q-learning. They show that to quantify  $\epsilon$ -optimal action-value functions with high probability, we need  $O(n\beta^5/\epsilon^2 \log(1/\epsilon)(\log(n\beta) + \log \log 1/\epsilon))$  and  $O(n\beta^4/\epsilon^2(\log(n\beta) + \log \log 1/\epsilon))$  samples in model-free and model-based QVI, respectively. A comparison between their results and the main result of this paper suggests that the sample complexity of SQL, which is of order  $O(n\beta^4/\epsilon^2 \log n)$ ,<sup>2</sup> is better than model-free QVI in terms of  $\beta$  and  $\log(1/\epsilon)$ . Although the sample complexities of SQL is only slightly tighter than the model-based QVI, SQL has a significantly better computational and space complexity than model-based QVI: SQL needs only  $2n$  memory space, while the space complexity of model-based QVI is given by either  $\tilde{O}(n/(\beta^4\epsilon^2))$  or  $n(|\mathcal{X}| + 1)$ , depending on whether the learned state transition matrix is sparse or not (see Kearns and Singh,

2. Note that at each round of SQL  $n$  new samples are generated. This combined with the result of Corollary 5 deduces the sample complexity of order  $O(n\beta^4/\epsilon^2 \log(n/\delta))$ .

1999). Also, SQL improves the computational complexity by a factor of  $\tilde{O}(\beta)$  compared to both model-free and model-based QVI.<sup>3</sup> Table 1 summarizes the comparisons between SQL and the other RL methods discussed in this section.

Table 1: Comparison between SQL, Q-learning, model-based and model-free Q-value iteration in terms of sample complexity, computational complexity, and space complexity.

| Method                   | SQL   | Q-learning (optimized)                                  | Model-based QVI                                     | Model-free QVI                                      |
|--------------------------|---|---|---|---|
| Sample complexity        | $\tilde{O}\left(\frac{n\beta^4}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{n\beta^5}{\epsilon^{2.5}}\right)$ | $\tilde{O}\left(\frac{n\beta^4}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{n\beta^5}{\epsilon^2}\right)$ |
| Computational complexity | $\tilde{O}\left(\frac{n\beta^4}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{n\beta^5}{\epsilon^{2.5}}\right)$ | $\tilde{O}\left(\frac{n\beta^5}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{n\beta^5}{\epsilon^2}\right)$ |
| Space complexity         | $\Theta(n)$   | $\Theta(n)$   | $\tilde{O}\left(\frac{n\beta^4}{\epsilon^2}\right)$ | $\Theta(n)$   |

## 4. Experiments

In this section, we analyze empirically the effectiveness of the proposed algorithms on different problem domains. We examine the convergence of synchronous SQL (Algorithm 1) as well as the asynchronous SQL (Algorithm 2) on several discrete state-action problems and compare it with Q-learning (Even-Dar and Mansour, 2003) (QL) and the model-based Q-value iteration (VI) of Kearns and Singh (1999). The source code of all tested algorithms are freely available at [http://www.mbfys.ru.nl/~mazar/Research\\_Top.html](http://www.mbfys.ru.nl/~mazar/Research_Top.html).

**Linear MDP:** this problem consists of states  $x_k \in \mathcal{X}, k = \{1, 2, \dots, 2500\}$  arranged in a one-dimensional chain (see Figure 1). There are two possible actions  $\mathcal{A} = \{-1, +1\}$  (left/right) and every state is accessible from any other state except for the two ends of the chain, which are absorbing states. A state  $x_k \in \mathcal{X}$  is called absorbing if  $P(x_k|x_k, a) = 1$  for all  $a \in \mathcal{A}$  and  $P(x_l|x_k, a) = 0, \forall l \neq k$ . Any transition to one of these two states has associated reward 1.

The transition probability for an interior state  $x_k$  to any other state  $x_l$  is inversely proportional to their distance in the direction of the selected action, and zero for all states corresponding to the opposite direction. Formally, consider the following

3. SQL has the sample and computational complexity of a same order since it performs only one Q-value update per sample, whereas, in the case of model-based QVI, the algorithm needs to iterate the action-value function of all state-action pairs at least  $\tilde{O}(\beta)$  times using Bellman operator, which leads to a computational complexity bound of order  $\tilde{O}(n\beta^5/\epsilon^2)$  given that only  $\tilde{O}(n\beta^4/\epsilon^2)$  entries of the estimated transition matrix are non-zero.

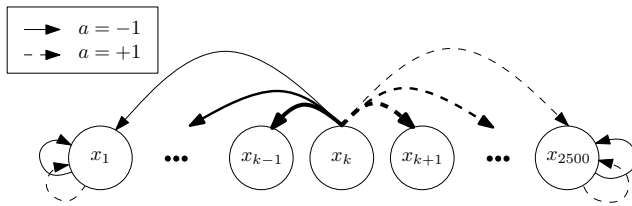


Figure 1: Linear MDP: Illustration of the linear MDP problem. Nodes indicate states. States  $x_1$  and  $x_{2500}$  are the two absorbing states and state  $x_k$  is an example of interior state. Arrows indicate possible transitions of **these three nodes only**. From  $x_k$  any other node is reachable with transition probability (arrow thickness) proportional to the inverse of the distance to  $x_k$  (see the text for details).

quantity  $n(x_l, a, x_k)$  assigned to all non-absorbing states  $x_k$  and to every  $(x_l, a) \in \mathcal{Z}$ :

$$n(x_l, a, x_k) = \begin{cases} \frac{1}{|l-k|} & \text{for } (l-k)a > 0 \\ 0 & \text{otherwise} \end{cases}.$$

We can write the transition probabilities as:

$$P(x_l|x_k, a) = \frac{n(x_l, a, x_k)}{\sum_{x_m \in \mathcal{X}} n(x_m, a, x_k)}.$$

Any transition that ends up in one of the interior states has associated reward  $-1$ .

The optimal policy corresponding to this problem is to reach the closest absorbing state as soon as possible.

**Combination lock:** the combination lock problem considered here is a stochastic variant of the reset state space models introduced in Koenig and Simmons (1993), where more than one reset state is possible (see Figure 2).

In our case we consider, as before, a set of states  $x_k \in \mathcal{X}, k \in \{1, 2, \dots, 2500\}$  arranged in a one-dimensional chain and two possible actions  $\mathcal{A} = \{-1, +1\}$ . In this problem, however, there is only one absorbing state (corresponding to the state *lock-opened*) with associated reward of 1. This state is reached if the all-ones sequence  $\{+1, +1, \dots, +1\}$  is entered correctly. Otherwise, if at some state  $x_k, k < 2500$ , action  $-1$  is taken, the lock automatically resets to some previous state  $x_l, l < k$  randomly (in the original combination lock problem, the reset state is always the initial state  $x_1$ ).

For every intermediate state, the rewards of actions  $-1$  and  $+1$  are set to 0 and  $-0.01$ , respectively. The transition probability upon taking the wrong action  $-1$  is, as before, inversely proportional to the distance of the states. That is

$$n(x_k, x_l) = \begin{cases} \frac{1}{k-l} & \text{for } l < k \\ 0 & \text{otherwise} \end{cases}, \quad P(x_l|x_k, 0) = \frac{n(k, l)}{\sum_{x_m \in \mathcal{X}} n(k, m)}.$$

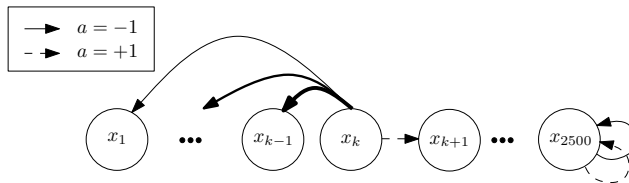


Figure 2: Combination lock: illustration of the combination lock MDP problem. Nodes indicate states. State  $x_{2500}$  is the goal (absorbing) state and state  $x_k$  is an example of interior state. Arrows indicate possible transitions of **these two nodes only**. From  $x_k$  any previous state is reachable with transition probability (arrow thickness) proportional to the inverse of the distance to  $x_k$ . Among the future states only  $x_{k+1}$  is reachable (arrow dashed).

Note that this problem is more difficult than the linear MDP since the goal state is only reachable from one state,  $x_{2499}$ .

**Grid world:** this MDP consists of a grid of  $50 \times 50$  states. A set of four actions {RIGHT, UP, DOWN, LEFT} is assigned to every state  $x \in \mathcal{X}$ . The location of each state  $x$  of the grid is determined by the coordinates  $c_x = (h_x, v_x)$ , where  $h_x$  and  $v_x$  are some integers between 1 and 50. There are 196 absorbing *firewall states* surrounding the grid and another one at the center of grid, for which a reward  $-1$  is assigned. The reward for the firewalls is

$$r(x, a) = -\frac{1}{\|c_x\|_2}, \quad \forall a \in \mathcal{A}.$$

Also, we assign reward 0 to all of the remaining (non-absorbing) states.

This means that both the top-left absorbing state and the central state have the least possible reward ( $-1$ ), and that the remaining absorbing states have reward which increases proportionally to the distance to the state in the bottom-right corner (but are always negative).

The transition probabilities are defined in the following way: taking action  $a$  from any non-absorbing state  $x$  results in a one-step transition in the direction of action  $a$  with probability 0.6, and a random move to a state  $y \neq x$  with probability inversely proportional to their Euclidean distance  $1/\|c_x - c_y\|_2$ .

The optimal policy then is to *survive* in the grid as long as possible by avoiding both the absorbing firewalls and the center of the grid. Note that because of the difference between the cost of firewalls, the optimal control prefers the states near the bottom-right corner of the grid, thus avoiding absorbing states with higher cost.

#### 4.0.3 EXPERIMENTAL SETUP AND RESULTS

We describe now our experimental setting. The convergence properties of SQL are compared with two other algorithms: a Q-learning (Even-Dar and Mansour, 2003) (QL) and

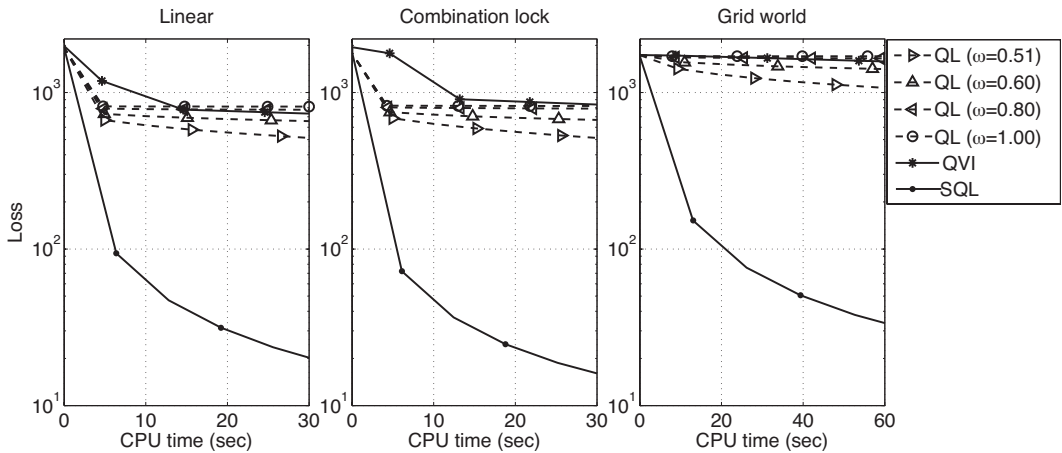


Figure 3: A comparison between SQL, Q-Learning and model-based QVI. Each plot compares the performance loss of the policy induced by the algorithms for a different MDP averaged over 50 different runs (see the text for details).

the model-based Q-value iteration (QVI) of Kearns and Singh (1999). VI is a batch reinforcement learning algorithm that first estimates the model using the whole data set and then performs value iteration on the learned model.

All algorithms are evaluated in terms of  $\ell_\infty$ -norm performance loss of the action-value function  $\|Q^* - Q_T\|$  at time-step  $T$ . We choose this performance measure in order to be consistent with the performance measure used in Section 3.2. The optimal action-value function  $Q^*$  is computed with high accuracy through value iteration.

We consider QL with polynomial learning step  $\alpha_k = 1/(k+1)^\omega$  where  $\omega \in \{0.51, 0.75\}$  and the linear learning step  $\alpha_k = 1/(k+1)$ . Note that  $\omega$  needs to be larger than 0.5, otherwise QL can asymptotically diverge (see Even-Dar and Mansour, 2003, for the proof).

To have a fair comparison of the three algorithms, since each algorithm requires different number of computations per iteration, we fix the total computational budget of the algorithms to the same value for each benchmark. The computation time is constrained to 30 seconds in the case of linear MDP and the combination lock problems. For the grid world, which has twice as many actions as the other benchmarks, the maximum run time is fixed to 60 seconds. We also fix the total number of samples, per state-action, to  $1 \times 10^5$  samples for all problems and algorithms. Significantly less number of samples leads to a dramatic decrease of the quality of the obtained solutions using all the approaches.

Algorithms were implemented as MEX files (in C++) and ran on a Intel core i5 processor with 8 GB of memory. CPU time was acquired using the system function `times()` which provides process-specific CPU time. Randomization was implemented using `gsl_rng_uniform()` function of the GSL library, which is superior to the standard `rand()`.<sup>4</sup> Sampling time, which is the same for all algorithms, were not included in CPU time.

Figure 3 shows the performance-loss in terms of elapsed CPU time for the three problems and algorithms with the choice of  $\gamma = 0.999$ . The results are averages over 50 runs, where

4. <http://www.gnu.org/s/gsl>.

at the beginning of each run **(i)** the action-value functions are randomly initialized in the interval  $[-V_{\max}, V_{\max}]$ , and **(ii)** a new set of samples is generated from  $P(\cdot|x, a)$  for all  $(x, a) \in \mathcal{Z}$ . Results correspond to the average error computed after a small fixed amount of iterations.

First, we see that, in all cases, SQL improves the performance very fast, almost by an order of magnitude, in just a few seconds. SQL outperforms both QL and QVI in all the three benchmarks. The minimum and maximum errors are attained for the combination lock problem and the grid world, respectively. We also observe that the difference between the final outcome of SQL and Q-learning (second best method) is very significant, up to 30 times in the grid-world problem. QL shows the best performance for  $\omega = 0.51$  and the quality of the QL

These results show that, as proved before in Theorem 2 that SQL manages to average out the simulation noise caused by sampling and converges, rapidly, to a near optimal solution, which is very robust. In addition, we can conclude that SQL significantly improves the computational complexity of learning w.r.t. the standard QL and QVI in the three presented benchmarks for our choice of experimental setup.

## 5. Analysis

In this section, we give some intuition about the convergence of asynchronous variant of SQL and provide the full proof of the finite-time analysis reported in Theorem 2 and Theorem 3. Then, we prove a lower bound on the sample complexity of all  $(\epsilon, \delta)$ -correct reinforcement learning algorithms stated in Theorem 7. We start by introducing some notations.

Let  $y^k$  be the set of all samples drawn at round  $k$  of the SQL algorithms and  $\mathcal{F}_k$  be the filtration generated by the sequence  $\{y^0, y^1, y^2, \dots, y^k\}$ . We first notice that for all  $(x, a) \in \mathcal{Z}$ , the update rule of Eq. 1 may be rewritten in the following more compact form:

$$Q_{k+1}(x, a) = (1 - \alpha_k)Q_k(x, a) + \alpha_k \mathcal{D}_k[Q_k, Q_{k-1}](x, a),$$

where  $\mathcal{D}_k[Q_k, Q_{k-1}](x, a) \triangleq \frac{1}{\alpha_k}[(1 - \alpha_k)\mathcal{T}_k Q_k(x, a) - (1 - 2\alpha_k)\mathcal{T}_k Q_{k-1}(x, a)]$ .

We then define the operator  $\mathcal{D}[Q_k, Q_{k-1}]$  as the expected value of the empirical operator  $\mathcal{D}_k$  conditioned on  $\mathcal{F}_{k-1}$ :

$$\begin{aligned} \mathcal{D}[Q_k, Q_{k-1}](x, a) &\triangleq \mathbb{E}(\mathcal{D}_k[Q_k, Q_{k-1}](x, a) | \mathcal{F}_{k-1}) \\ &= \mathbb{E}\left(\frac{1 - \alpha_k}{\alpha_k} \mathcal{T}_k Q_k(x, a) - \frac{1 - 2\alpha_k}{\alpha_k} \mathcal{T}_k Q_{k-1}(x, a) \middle| \mathcal{F}_{k-1}\right) \\ &= \frac{1 - \alpha_k}{\alpha_k} \mathcal{T} Q_k(x, a) - \frac{1 - 2\alpha_k}{\alpha_k} \mathcal{T} Q_{k-1}(x, a), \end{aligned}$$

where the last line follows by the fact that, in both Algorithm 1 and 2,  $\mathcal{T}_k Q_k(x, a)$  and  $\mathcal{T}_k Q_{k-1}(x, a)$  are unbiased empirical estimates of the Bellman optimality operators  $\mathcal{T} Q_k(x, a)$  and  $\mathcal{T} Q_{k-1}(x, a)$ , respectively. Thus, the update rule of SQL can be re-expressed as

$$Q_{k+1}(x, a) = (1 - \alpha_k)Q_k(x, a) + \alpha_k (\mathcal{D}[Q_k, Q_{k-1}](x, a) - \epsilon_k(x, a)), \quad (4)$$

where the estimation error  $\epsilon_k$  is defined as the difference between the operator  $\mathcal{D}[Q_k, Q_{k-1}]$  and its sample estimate  $\mathcal{D}_k[Q_k, Q_{k-1}]$  for all  $(x, a) \in \mathcal{Z}$ :

$$\epsilon_k(x, a) \triangleq \mathcal{D}[Q_k, Q_{k-1}](x, a) - \mathcal{D}_k[Q_k, Q_{k-1}](x, a).$$



We have the property that  $\mathbb{E}[\epsilon_k(x, a) | \mathcal{F}_{k-1}] = 0$  which means that for all  $(x, a) \in \mathcal{Z}$  the sequence of estimation error  $\{\epsilon_1(x, a), \epsilon_2(x, a), \dots, \epsilon_k(x, a)\}$  is a martingale difference sequence w.r.t. the filtration  $\mathcal{F}_k$ . Finally, we define the martingale  $E_k(x, a)$  to be the sum of the estimation errors:

$$E_k(x, a) \triangleq \sum_{j=0}^k \epsilon_j(x, a), \quad \forall (x, a) \in \mathcal{Z}. \quad (5)$$

The following steps lead to the proof of Theorem 2 and Theorem 3 **(i)** Lemma 9 shows the stability of SQL (i.e., the sequence of  $Q_k$  stays bounded). **(ii)** Lemma 10 states the key property that the iterate  $Q_{k+1}$  is close to the Bellman operator  $\mathcal{T}$  applied to the previous iterate  $Q_k$  plus an estimation error term of order  $E_k/k$ . **(iii)** By induction, Lemma 11 provides a performance bound  $\|Q^* - Q_k\|$  in terms of a discounted sum of the cumulative estimation errors  $\{E_j\}_{j=0:k-1}$ . The above-mentioned results hold for both Algorithm 1 and Algorithm 2. Subsequently, **(iv)** we concentrate on proving Theorem 2 by making use of a maximal Azuma's inequality, stated in Lemma 13. **(v)** We then extend this result for the case of asynchronous SQL by making use of the result of Lemma 14.

For simplicity of the notations, we remove the dependence on  $(x, a)$  (e.g., writing  $Q$  for  $Q(x, a)$ ,  $E_k$  for  $E_k(x, a)$ ) when there is no possible confusion. Also, we notice that, for all  $k \geq 0$ , the following relations hold between  $\alpha_k$  and  $\alpha_{k+1}$  in Algorithm 1 and Algorithm 2:

$$\alpha_{k+1} = \frac{\alpha_k}{\alpha_k + 1} \quad \text{and} \quad \alpha_k = \frac{\alpha_{k+1}}{1 - \alpha_{k+1}}.$$

**Lemma 9 (Stability of SQL)** *Let A1 hold and assume that the initial action-value function  $Q_0 = Q_{-1}$  is uniformly bounded by  $V_{\max} = \beta$ , then we have*

$$\|Q_k\| \leq V_{\max}, \quad \|\epsilon_k\| \leq 2V_{\max}, \quad \text{and} \quad \|\mathcal{D}_k[Q_k, Q_{k-1}]\| \leq V_{\max} \quad \forall k \geq 0.$$

**Proof** We first prove that  $\|\mathcal{D}_k[Q_k, Q_{k-1}]\| \leq V_{\max}$  by induction. For  $k = 0$  we have:

$$\|\mathcal{D}_0[Q_0, Q_{-1}]\| = \|\mathcal{T}_0 Q_{-1}\| \leq \|r\| + \gamma \|\mathcal{M} Q_{-1}\| \leq R_{\max} + \gamma V_{\max} = V_{\max}.$$

Now for any  $k \geq 0$ , let us assume that the bound  $\|\mathcal{D}_k[Q_k, Q_{k-1}]\| \leq V_{\max}$  holds. Thus

$$\begin{aligned} \|\mathcal{D}_{k+1}[Q_{k+1}, Q_k]\| &= \left\| \frac{1-\alpha_{k+1}}{\alpha_{k+1}} \mathcal{T}_{k+1} Q_{k+1} - \frac{1-2\alpha_{k+1}}{\alpha_{k+1}} \mathcal{T}_{k+1} Q_k \right\| \\ &= \left\| \left[ \frac{1-\alpha_{k+1}}{\alpha_{k+1}} - \frac{1-2\alpha_{k+1}}{\alpha_{k+1}} \right] r \right\| + \gamma \left\| \frac{1-\alpha_{k+1}}{\alpha_{k+1}} \mathcal{M} Q_{k+1} - \frac{1-2\alpha_{k+1}}{\alpha_{k+1}} \mathcal{M} Q_k \right\| \\ &\leq \|r\| + \gamma \left\| \frac{1-\alpha_{k+1}}{\alpha_{k+1}} \mathcal{M}((1-\alpha_k)Q_k + \alpha_k \mathcal{D}_k[Q_k, Q_{k-1}]) - \frac{1-2\alpha_{k+1}}{\alpha_{k+1}} \mathcal{M} Q_k \right\| \\ &= \|r\| + \gamma \left\| \frac{1-\frac{\alpha_k}{\alpha_{k+1}}}{\frac{\alpha_k}{\alpha_{k+1}}} \mathcal{M}((1-\alpha_k)Q_k + \alpha_k \mathcal{D}_k[Q_k, Q_{k-1}]) - \frac{1-2\frac{\alpha_k}{\alpha_{k+1}}}{\frac{\alpha_k}{\alpha_{k+1}}} \mathcal{M} Q_k \right\| \\ &= \|r\| + \gamma \left\| \mathcal{M} \left( \frac{1-\alpha_k}{\alpha_k} Q_k + \mathcal{D}_k[Q_k, Q_{k-1}] \right) - \frac{1-\alpha_k}{\alpha_k} \mathcal{M} Q_k \right\| \\ &\leq \|r\| + \gamma \left\| \mathcal{M} \left( \frac{1-\alpha_k}{\alpha_k} Q_k + \mathcal{D}_k[Q_k, Q_{k-1}] - \frac{1-\alpha_k}{\alpha_k} Q_k \right) \right\| \\ &\leq \|r\| + \gamma \|\mathcal{D}_k[Q_k, Q_{k-1}]\| \leq R_{\max} + \gamma V_{\max} = V_{\max}, \end{aligned}$$

and by induction, we deduce that for all  $k \geq 0$ ,  $\|\mathcal{D}_k[Q_k, Q_{k-1}]\| \leq V_{\max}$ .

Now the bound on  $\epsilon_k$  follows from  $\|\epsilon_k\| = \|\mathbb{E}(\mathcal{D}_k[Q_k, Q_{k-1}]|\mathcal{F}_{k-1}) - \mathcal{D}_k[Q_k, Q_{k-1}]\| \leq 2V_{\max}$ , and the bound  $\|Q_k\| \leq V_{\max}$  is deduced by noticing that  $Q_k = \frac{1}{k} \sum_{j=0}^{k-1} \mathcal{D}_j[Q_j, Q_{j-1}]$ .  
■

The next Lemma shows that  $Q_k$  is close to  $\mathcal{T}Q_{k-1}$ , up to a  $O(\frac{1}{k})$  term minus the cumulative estimation error  $\frac{1}{k}E_{k-1}$ .

**Lemma 10** *Under A1, for any  $k \geq 1$ :*

$$Q_k = \mathcal{T}Q_{k-1} + \frac{1}{k} (\mathcal{T}Q_0 - \mathcal{T}Q_{k-1} - E_{k-1}). \quad (6)$$

**Proof** We prove this result by induction. The result holds for  $k = 1$ , where (6) reduces to (4). We now show that if the property (6) holds for  $k \geq 1$  then it also holds for  $k + 1$ . Assume that (6) holds for  $k$ . Then, from (4) we have:

$$\begin{aligned} Q_{k+1} &= (1 - \alpha_k)Q_k + \alpha_k \left[ \frac{1 - \alpha_k}{\alpha_k} \mathcal{T}Q_k - \frac{1 - 2\alpha_k}{\alpha_k} \mathcal{T}Q_{k-1} - \epsilon_k \right] \\ &= (1 - \alpha_k) (\mathcal{T}Q_{k-1} + \alpha_{k-1} (\mathcal{T}Q_0 - \mathcal{T}Q_{k-1} - E_{k-1})) \\ &\quad + \alpha_k \left[ \frac{1 - \alpha_k}{\alpha_k} \mathcal{T}Q_k - \frac{1 - 2\alpha_k}{\alpha_k} \mathcal{T}Q_{k-1} - \epsilon_k \right] \\ &= (1 - \alpha_k) \left[ \mathcal{T}Q_{k-1} + \frac{\alpha_k}{1 - \alpha_k} (\mathcal{T}Q_0 - \mathcal{T}Q_{k-1} - E_{k-1}) \right] \\ &\quad + (1 - \alpha_k) \mathcal{T}Q_k - (1 - 2\alpha_k) \mathcal{T}Q_{k-1} - \alpha_k \epsilon_k \\ &= (1 - 2\alpha_k) \mathcal{T}Q_{k-1} + \alpha_k (\mathcal{T}Q_0 - E_{k-1}) + (1 - \alpha_k) \mathcal{T}Q_k - (1 - 2\alpha_k) \mathcal{T}Q_{k-1} - \alpha_k \epsilon_k \\ &= (1 - \alpha_k) \mathcal{T}Q_k + \alpha_k (\mathcal{T}Q_0 - E_{k-1} - \epsilon_k) = \mathcal{T}Q_k + \alpha_k (\mathcal{T}Q_0 - \mathcal{T}Q_k - E_k) \\ &= \mathcal{T}Q_k + \frac{1}{k+1} (\mathcal{T}Q_0 - \mathcal{T}Q_k - E_k). \end{aligned}$$

Thus (6) holds for  $k + 1$ , and is thus true for all  $k \geq 1$ . ■

Now we bound the difference between  $Q^*$  and  $Q_k$  in terms of the discounted sum of the cumulative estimation errors  $\{E_0, E_1, \dots, E_{k-1}\}$ .

**Lemma 11 (Error Propagation of SQL)** *Let A1 hold and assume that the initial action-value function  $Q_0 = Q_{-1}$  is uniformly bounded by  $V_{\max} = \beta$ , then for all  $k \geq 1$ , we have:*

$$\|Q^* - Q_k\| \leq \frac{1}{k} \left[ 2\gamma\beta^2 + \sum_{j=1}^k \gamma^{k-j} \|E_{j-1}\| \right], \quad (7)$$

$$\|Q^* - Q_k\| \leq \frac{\beta}{k} \left[ 2\gamma\beta + \max_{j=1:k} \|E_{j-1}\| \right]. \quad (8)$$

**Proof**

We first notice that  $\sum_{j=1}^k \gamma^{k-j} \|E_{j-1}\| \leq \beta \max_{j=1:k} \|E_{j-1}\|$  for any sequence of cumulative errors  $\{E_0, E_1, E_2, \dots, E_{k-1}\}$ . Therefore, we only need to prove (7) and (8) follows this result.

Again we prove this lemma by induction. The result holds for  $k = 1$  as:

$$\|Q^* - Q_1\| = \|\mathcal{T}Q^* - \mathcal{T}Q_0 - \epsilon_0\| \leq \|Q^* - Q_0\| + \|\epsilon_0\| \leq 2\gamma V_{\max} + \|\epsilon_0\| \leq 2\gamma\beta^2 + \|E_0\|.$$

We now show that if the bound holds for  $k$ , then it also holds for  $k + 1$ . Thus, assume that (7) holds for  $k$ . By using Lemma 10:

$$\begin{aligned} \|Q^* - Q_{k+1}\| &= \left\| Q^* - \mathcal{T}Q_k - \frac{1}{k+1} (\mathcal{T}Q_0 - \mathcal{T}Q_k - E_k) \right\| \\ &\leq \|\alpha_k (\mathcal{T}Q^* - \mathcal{T}Q_0) + (1 - \alpha_k) (\mathcal{T}Q^* - \mathcal{T}Q_k)\| + \alpha_k \|E_k\| \\ &\leq \alpha_k \|\mathcal{T}Q^* - \mathcal{T}Q_0\| + (1 - \alpha_k) \|\mathcal{T}Q^* - \mathcal{T}Q_k\| + \alpha_k \|E_k\| \\ &\leq 2\gamma\alpha_k V_{\max} + \gamma(1 - \alpha_k) \|Q^* - Q_k\| + \alpha_k \|E_k\|. \end{aligned}$$

By plugging (7) into (9) we then deduce:

$$\begin{aligned} \|Q^* - Q_{k+1}\| &\leq 2\alpha_k \gamma V_{\max} + \gamma(1 - \alpha_k) \alpha_{k-1} \left[ 2\gamma\beta^2 + \sum_{j=1}^k \gamma^{k-j} \|E_{j-1}\| \right] + \alpha_k \|E_k\| \\ &\leq 2\alpha_k \gamma \beta + \gamma(1 - \alpha_k) \frac{\alpha_k}{1 - \alpha_k} \left[ 2\gamma\beta^2 + \sum_{j=1}^k \gamma^{k-j} \|E_{j-1}\| \right] + \alpha_k \|E_k\| \\ &\leq 2\alpha_k \gamma \beta + \gamma \alpha_k 2\gamma\beta^2 + \gamma \alpha_k \sum_{j=1}^k \gamma^{k-j} \|E_{j-1}\| + \alpha_k \|E_k\| \\ &= \frac{1}{k+1} \left[ 2\gamma\beta^2 + \sum_{j=1}^{k+1} \gamma^{k+1-j} \|E_{j-1}\| \right]. \end{aligned}$$

Thus, the result holds for  $k + 1$  thus for all  $k \geq 1$  by induction.  $\blacksquare$

In the next lemma, we prove a high probability bound over the estimation error term of Lemma 11

**Lemma 12** *Let A1 hold and assume that the initial action-value function  $Q_0 = Q_{-1}$  is uniformly bounded by  $V_{\max} = \beta$  then for all  $0 < \delta < 1$  and  $k \geq 1$ , we have:*

$$\max_{j=1:k} \|E_{j-1}\| \leq \beta \sqrt{8k \log \frac{2n}{\delta}}, \quad w.p. 1 - \delta. \quad (9)$$

**Proof** We start by providing a high probability bound for  $\max_{1 \leq j \leq k} |E_{j-1}(x, a)|$  for a given  $(x, a)$ . First notice that:

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq j \leq k} |E_{j-1}(x, a)| > \epsilon\right) &= \mathbb{P}\left(\max\left[\max_{1 \leq j \leq k} (E_{j-1}(x, a)), \max_{1 \leq j \leq k} (-E_{j-1}(x, a))\right] > \epsilon\right) \\ &= \mathbb{P}\left(\left\{\max_{1 \leq j \leq k} (E_{j-1}(x, a)) > \epsilon\right\} \cup \left\{\max_{1 \leq j \leq k} (-E_{j-1}(x, a)) > \epsilon\right\}\right) \\ &\leq \mathbb{P}\left(\max_{1 \leq j \leq k} (E_{j-1}(x, a)) > \epsilon\right) + \mathbb{P}\left(\max_{1 \leq j \leq k} (-E_{j-1}(x, a)) > \epsilon\right), \end{aligned} \tag{10}$$

and each term is now bounded by using a maximal Azuma-Hoeffding's inequality, reminded now (see e.g., Cesa-Bianchi and Lugosi (2006)).

**Lemma 13 (Maximal Hoeffding-Azuma's Inequality)** *Let  $\mathcal{V} = \{V_1, V_2, \dots, V_k\}$  be a martingale difference sequence w.r.t. a sequence of random variables  $\{X_1, X_2, \dots, X_k\}$  (i.e.,  $\mathbb{E}(V_{j+1}|X_1, \dots, X_j) = 0$  for all  $0 < j \leq k$ ) such that  $\mathcal{V}$  is uniformly bounded by  $L > 0$ . If we define  $S_k = \sum_{i=1}^k V_i$ , then for any  $\epsilon > 0$ , we have:*

$$\mathbb{P}\left(\max_{1 \leq j \leq k} S_j > \epsilon\right) \leq \exp\left(\frac{-\epsilon^2}{2kL^2}\right).$$

As mentioned earlier, the sequence of random variables  $\{\epsilon_0(x, a), \epsilon_1(x, a), \dots, \epsilon_k(x, a)\}$  is a martingale difference sequence w.r.t. the filtration  $\mathcal{F}_k$  (generated by the random samples  $\{y_0, y_1, \dots, y_k\}(x, a)$  for all  $(x, a)$ ), i.e.,  $\mathbb{E}[\epsilon_k(x, a)|\mathcal{F}_{k-1}] = 0$ . It follows from Lemma 13 that for any  $\epsilon > 0$  we have:

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq j \leq k} (E_{j-1}(x, a)) > \epsilon\right) &\leq \exp\left(\frac{-\epsilon^2}{8kV_{\max}^2}\right) \\ \mathbb{P}\left(\max_{1 \leq j \leq k} (-E_{j-1}(x, a)) > \epsilon\right) &\leq \exp\left(\frac{-\epsilon^2}{8kV_{\max}^2}\right). \end{aligned} \tag{11}$$

By combining (11) with (10) we deduce that  $\mathbb{P}(\max_{1 \leq j \leq k} |E_{j-1}(x, a)| > \epsilon) \leq 2 \exp\left(\frac{-\epsilon^2}{8kV_{\max}^2}\right)$ , and by a union bound over the state-action space, we deduce that

$$\mathbb{P}\left(\max_{1 \leq j \leq k} \|E_{j-1}\| > \epsilon\right) \leq 2n \exp\left(\frac{-\epsilon^2}{8kV_{\max}^2}\right). \tag{12}$$

This bound can be rewritten, for any  $\delta > 0$ , as:

$$\mathbb{P}\left(\max_{1 \leq k \leq T} \|E_{j-1}\| \leq V_{\max} \sqrt{8k \log \frac{2n}{\delta}}\right) \geq 1 - \delta,$$

which completes the proof. ■

The result of Theorem 2 then follows by plugging (9) into (8) and taking into account that after  $T$  step of Algorithm 1 we have  $T \leq n(k+1)$ , where  $k$  is the number of iterations at time step  $T$ . For the proof of Theorem 3 we rely on the following result which bound the number of steps required to visit all states-actions  $k$  times with high probability :

**Lemma 14** *Under A2, from any initial state  $x_0$  and for any integer  $k > 0$ , after a run of  $T = keL \log \frac{1}{\delta}$  steps the state-action space  $\mathcal{Z}$  is covered at least  $k$  times under the policy  $\pi$  w.p. at least  $1 - \delta$ .*

**Proof**

We begin by defining the random event  $\mathcal{Q}_k$  as the number of steps required to cover the whole MDP for  $k$  times starting from any state  $x_0 \in \mathcal{X}$  at any time  $t > 0$ . We then bound  $\mathcal{Q}_k$  for any  $x_0$  in high probability using Markov inequality (Feller, 1968):

$$\mathbb{P}(\mathcal{Q}_k > ekL) \leq \frac{\mathbb{E}(\mathcal{Q}_k)}{ekL} \leq \frac{k \sup_{t>0} \max_{x \in \mathcal{X}} \mathbb{E}(\tau_\pi(x, t))}{ekL} \leq \frac{kL}{ekL} = \frac{1}{e}.$$

In words after a run of length  $ekL$  the probability that the entire state-action space is not covered for at least  $k$  times is less than  $\frac{1}{e}$ . The fact that the bound holds for any initial state and time implies that after  $m > 0$  intervals of length  $ekL$  the chance of not covering  $\mathcal{Z}$  for  $k$  times is less than  $\frac{1}{e^m}$ :

$$\mathbb{P}(\mathcal{Q}_k > mekL) \leq \frac{1}{e^m}.$$

By the choice of  $m = \log \frac{1}{\delta}$  we deduce:

$$\mathbb{P}\left(\mathcal{Q}_k > keL \log \frac{1}{\delta}\right) \leq \delta.$$

The bound can be then rewritten as follows:

$$\mathbb{P}\left(\mathcal{Q}_k \leq keL \log \frac{1}{\delta}\right) \geq 1 - \delta,$$

which concludes the proof. ■

Plugging the results of Lemma 14 and Lemma 12 into (8) concludes the proof of theorem 3.

### 5.1 Proof of Theorem 7

In this subsection, we provide the full proof of Theorem 7. In our analysis, we rely on the likelihood-ratio method, which has been perviously used to prove a lower-bound for multi-armed bandits, (see Mannor and Tsitsiklis, 2004), and extend this approach for reinforcement learning in Markovian decision problems. We begin by defining a class of MDPs for which the proposed lower-bound holds. The class of MDPs  $\mathbb{M}$  (see Figure 4) is defined as a set of all MDPs with  $n = 3n_1 \geq 3$  state-action pairs, in which the set of state-action pairs  $\mathcal{Z}$  is made up of three smaller sets  $\mathcal{Z}_0$ ,  $\mathcal{Z}_1$  and  $\mathcal{Z}_2$  each of size  $n_1$ . We first make the assumption that, for all  $M \in \mathbb{M}$ , both  $\mathcal{Z}_1$  and  $\mathcal{Z}_2$  consist of only absorbing states, i.e.,  $\mathbb{P}(z|z) = 1$  for all  $z \in \mathcal{Z}_1 \cup \mathcal{Z}_2$ . We then assume that for all  $M \in \mathbb{M}$ , every entry  $z_0^l$ ,  $1 \leq l \leq n_1$ , of the set  $\mathcal{Z}_0$  is connected to only two other state-action pairs  $z_1^l \in \mathcal{Z}_1$  and

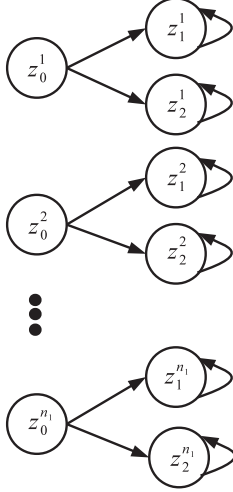


Figure 4: Illustration of the class of MDP considered for the lower-bounds. Nodes indicate state-(action) pairs. Arrows indicate possible transitions of these nodes. (see the text for details).

$z_2^l \in \mathcal{Z}_2$ , such that the probability of moving to  $z_1^l$  by taking the state-action pair  $z_0^l$  to the corresponding state in  $\mathcal{Z}_1$  is given by  $p_M(z_1^l)$ , for all  $1 \leq l \leq n_1$ . Also, we assume that the instant reward  $r(z)$  is set to 1 for all  $z \in \mathcal{Z}_1$  and 0 for the rest of the state-actions. One can then solve the Bellman equation, For all  $z \in \mathcal{Z}$ , in close-form and compute the optimal action-value function  $Q^*$  as follows:

$$Q^*(z) = \begin{cases} \gamma\beta p_M(z) & z \in \mathcal{Z}_0 \\ \beta & z \in \mathcal{Z}_1 \\ 0 & z \in \mathcal{Z}_2 \end{cases}$$

Here, we only consider a subset of  $n_1 + 1$  MDPs in the class  $\mathbb{M}$  denoted by  $\mathbb{M}^* = \{M_0, M_1, \dots, M_{n_1}\}$  and show that for an RL algorithm to be  $(\epsilon, \delta)$ -correct on  $\mathbb{M}^*$ , the number of samples should be at least of the order  $O(\frac{n\beta^2}{c_1} \log(n/c_2\delta)/\epsilon^2)$ . For all  $M_m \in \mathbb{M}^*$  and  $z \in \mathcal{Z}_1$ , the state transition probability  $p_{M_m}(z)$  is given by :

$$p_{M_m}(z) = \begin{cases} \frac{1 + 8\epsilon/\beta}{2} & z = z_m \\ \frac{1}{2} & \text{otherwise} \end{cases}$$

In words, for every MDP  $M_m \in \mathbb{M}^*$ ,  $m \geq 1$ , the transition probability  $p_{M_m}(z) = 1/2$  for all  $z \in \mathcal{Z}_0$ , except  $z_0^m$  in which case it is more likely to move to  $z_1^m$  in the high-value

set  $\mathcal{Z}_1$  rather than  $z_2^m$  in the low-value set  $\mathcal{Z}_2$ . Also, we notice that for  $M_0$  the transition probability  $p_{M_0}(z) = 1/2$  for all  $z \in \mathcal{Z}_1$ .

In the remainder of this subsection, we define  $\epsilon_0 = 1/8$ ,  $\delta_0 = \frac{e^{-2}}{4}$  and  $\gamma_0 = 1/2$ , and fix  $\epsilon \in (0, \epsilon_0)$ ,  $\delta \in (0, \delta_0)$  and  $\gamma \in (\gamma_0, 1)$ . We denote by  $\mathbb{E}_m$  and  $\mathbb{P}_m$  the expectation and the probability under the model  $M_m \in \mathbb{M}$ . we define  $t^*$  by

$$t^* \triangleq \frac{\beta^2}{c\epsilon^2} \log \frac{1}{\theta},$$

where  $\theta = 4\delta$  and  $c_1$  and  $c_2$  are absolute constants whose values will be define later.

Define the event  $\Omega_1(z) \triangleq \{\|Q_t^A(z) - \frac{\gamma}{2(1-\gamma)}\| \leq \epsilon\}$  for all  $z \in \mathcal{Z}$ , where  $Q_t^A$  is the estimation of the optimal value function after making  $t$  transition by the RL algorithm  $A$ . Also, define  $t(z)$  as the number of times that each state-action pair  $z \in \mathcal{Z}$  is sampled from before the algorithm achieves an  $\epsilon$ -optimal estimate of the value function. Finally, we define  $K_t(z) = r_1(z) + r_2(z) + \dots + r_t(z)$  as the sum of rewards of trying  $z \in \mathcal{Z}$  for  $t$  times. Based on these definitions we introduce the event  $\Omega_2(z)$ , for all  $z \in \mathcal{Z}$ :

$$\Omega_2(z) \triangleq \left\{ \max_{1 \leq t \leq t^*} \left( \frac{1}{2}t(z) - K_t(z) \right) \leq \sqrt{\frac{t^* \log \frac{1}{\theta}}{2}} \right\}. \quad (13)$$

In the sequel, we restrict ourselves to the case that the only unknown factor is the state-transition probability. Therefore, the RL task is to learn only those value functions which depend on  $p_M$ , i.e.,  $Q^*(z)$  for all  $z \in \mathcal{Z}_0$ . We follow the following steps in the proof: **(i)** We assume that  $t(z) \leq t^*$  for every  $z_l^0 \in \mathcal{Z}_0$ ,  $l = 1 : n_1$  and that there exists an  $(\epsilon, \delta)$ -correct reinforcement learning  $A$  for all  $M \in \mathbb{M}^*$  which achieves the  $\epsilon$ -optimal estimate of  $Q^*$  after making  $t \leq t^*$  transition per  $z \in \mathcal{Z}_0$ . Then, **(ii)** we show that under these assumptions the probability of  $\mathbb{P}_l(\Omega_1(z_l^0)) \geq \delta$ , for all  $z_l^0 \in \mathcal{Z}_0$ , which violates the  $(\epsilon, \delta)$ -correctness of the algorithm  $A$ . It then **(iii)** follows, by a contradiction argument, that  $t(z) > \frac{\beta^2}{c\epsilon^2} \log \frac{1}{\theta}$  for all  $z \in \mathcal{Z}_0$ , which proves the theorem.

We state the following lemmas required for our analysis.<sup>5</sup>

**Lemma 15 (Mannor and Tsitsiklis, 2004)**

$$\mathbb{P}_0(\Omega_2(z)) > \frac{3}{4}, \quad \forall z \in \mathcal{Z}. \quad (14)$$

**Lemma 16** *if  $x \leq 1$  and  $y \geq 0$  then*

$$(1 - x)^y \geq \exp(-dxy), \quad (15)$$

where  $d = 1.78$ .

Now let define  $\mathcal{S}(z) \triangleq \Omega_1(z) \cap \Omega_2(z)$ , for all  $z \in \mathcal{Z}$ . Then, under an  $(\epsilon, \delta)$ -correct RL algorithm, we have  $\mathbb{P}_0(\Omega_1(z)) \geq 1 - \delta \geq 1 - e^{-2}/4 > 3/4$ . This combined with Lemma 15 implies that  $\mathbb{P}_0(\mathcal{S}(z)) > 1/2$ , for all  $z \in \mathcal{Z}$ . Based on these results, we show that for all  $0 \leq l \leq n_1$  and MDP  $M_l \in \mathbb{M}^*$ :

5. For the proofs see Mannor and Tsitsiklis (2004).



**Lemma 17** *If  $t(z_0^l) < t^*$  and  $c \geq 0$ , then  $\mathbb{P}_l(\Omega(z_0^l)) > \delta$ .*

Let  $W_l$  be the history of all the outcomes of trying  $z_0^l$  for  $t(z_0^l)$  times. We define the likelihood function of  $L_m^l(w)$  as:

$$L_m^l(w) \triangleq \mathbb{P}_m(W_l = w),$$

for every possible history  $w$  and all  $1 \leq l \leq n_1$  and  $0 \leq m \leq n_1$ . This function can be used to define a random variable  $L_l(W_m)$ . We also let  $k$  be a shorthand notation for  $K_t(z_l^0)$ , the total sum of rewards obtained by trying  $z_l$  for  $t(z_l^0)$  times and  $t_l$  for  $t(z_l^0)$  itself. The likelihood ratio of the event  $w$  for every  $1 \leq l \leq n_1$  and  $0 \leq m \leq n_1$  is given by:

$$\begin{aligned} \frac{L_l^m(W)}{L_l^0(W)} &= \frac{(\frac{1}{2} + 4\epsilon/\beta)^k (\frac{1}{2} - 4\epsilon/\beta)^{t_l}}{\frac{1}{2}^{t_l+k}} \\ &= (1 + 8\epsilon/\beta)^k (1 - 8\epsilon/\beta)^k (1 - 8\epsilon/\beta)^{t_l-2k} \\ &= (1 - 64\epsilon^2/\beta^2)^k (1 - 8\epsilon/\beta)^{t_l-2k} \end{aligned} \quad (16)$$

We will now determine the lower bounds of the terms in the RHS of (16) when the event  $\mathcal{S} = \Omega_1 \cap \Omega_2$  occurs:

If  $\mathcal{S}$  occurs, then  $\Omega_1$  occurs, and we have  $k \leq t_l \leq t^*$ . Therefore, we deduce:

$$\begin{aligned} (1 - 64\epsilon^2/\beta^2)^k &\geq (1 - 64\epsilon^2/\beta^2)^{t^*} = (1 - 64\epsilon^2/\beta^2)^{\frac{\beta^2}{c\epsilon^2} \log(\frac{1}{\theta})} \\ &\geq e^{-64\frac{d}{c} \log(\frac{1}{\theta})} \\ &= \theta^{64\frac{d}{c}} \end{aligned} \quad (17)$$

Likewise, if the event  $\mathcal{S}$  occurs, then  $\Omega_2$  occurs which leads to:

$$t_l - 2k \leq 2\sqrt{\frac{t^* \log(1/\theta)}{2}} = \frac{2\beta}{\epsilon\sqrt{2c}} \log \frac{1}{\theta},$$

which implies:

$$\begin{aligned} (1 - 8\epsilon/\beta)^{t-2k} &\geq (1 - 8\epsilon/\beta)^{\frac{2\beta}{\epsilon\sqrt{2c}} \log \frac{1}{\theta}} \\ &\geq e^{\frac{-16d}{\sqrt{2c}} \log(\frac{1}{\theta})} \\ &= \theta^{\frac{16d}{\sqrt{2c}}} \end{aligned} \quad (18)$$

By plugging (17) and (18) into (16), we deduce:

$$\frac{L_l^m(W)}{L_l^0(W)} \geq \theta^{16d(\frac{1}{\sqrt{2c}} + \frac{4}{c})}. \quad (19)$$

Then, by the choice of  $c = 700$ , we obtain:

$$\frac{L_l^m(W)}{L_l^0(W)} \mathbf{1}_{\mathcal{S}} \leq 2\delta \mathbf{1}_{\mathcal{S}}, \quad (20)$$

where  $\mathbf{1}_{\mathcal{S}}$  is the indicator function of the event  $\mathcal{S}$ . Then,

$$\mathbb{P}_l(\Omega_1) \geq \mathbb{P}_l(\mathcal{S}) = \mathbb{E}_l[\mathbf{1}_\mathcal{S}] = \mathbb{E}_0 \left( \frac{L_l^m(W)}{L_l^0(W)} \mathbf{1}_\mathcal{S} \right) \geq \mathbb{E}_0 [2\delta \mathbf{1}_\mathcal{S}] = 2\delta \mathbb{P}_0(\mathcal{S}) > \delta, \quad (21)$$

where we make use of the fact that  $\mathbb{P}_0(\mathcal{S}) > \frac{1}{2}$ . This result shows that if, for all  $z_l^0 \in \mathcal{Z}_0$ ,  $t_l \leq (\frac{\beta^2}{700\epsilon^2}) \log(\frac{1}{4\delta})$  then  $\mathbb{P}_l(\Omega_1(z_l^0)) > \delta$ , which violates the  $(\epsilon, \delta)$ -correctness of the algorithm  $A$ . This is due the fact  $\Omega_1(z_l^0)$  and  $|Q^*(z_l) - Q_t^A(z_l)| \leq \epsilon$  are two separate events under the MDP  $M_l$  and therefore for the algorithm  $A$  to be  $(\epsilon, \delta)$ -correct  $\mathbb{P}_l(\Omega_1(z_l^0))$  needs to be smaller than  $\delta$ . This result implies that every  $(\epsilon, \delta)$ -correct algorithm on the class of MDPs  $\mathbb{M}$  needs to make at least  $t^* = \beta^2/(700\epsilon^2) \log \frac{1}{4\delta}$  transitions from every state-action pair  $z \in \mathcal{Z}_0$ . The proof then follows by summing up the number of samples for all  $z \in \mathcal{Z}_0$ .

## 6. Conclusions and Future Work

In this paper, we introduced a new reinforcement learning algorithm, called speedy Q-learning (SQL). We analyzed the finite-time behavior of this algorithm as well as its asymptotic convergence to the optimal action-value function. Our result is in the form of high probability bound on the performance loss of SQL, which suggests that the algorithm converges to the optimal action-value function in a faster rate than the standard Q-learning. Overall, SQL is a simple, efficient and theoretically well-founded reinforcement learning algorithm, which improves on existing RL algorithms such as Q-learning and the sample-based value iteration.

In this work, we are only interested in the estimation of the optimal action-value function and not the problem of exploration. Therefore, we did not compare our result to the PAC-MDP methods (Szita and Szepesvári, 2010; Strehl et al., 2009) and the upper-confidence bound based algorithms (Jaksch et al., 2010; Bartlett and Tewari, 2009), in which the choice of the exploration policy impacts the behavior of the learning algorithms. However, we believe that it would be possible to gain w.r.t. the state of the art in PAC-MDPs, by combining the asynchronous version of SQL with a smart exploration strategy. This is mainly due to the fact that the bound for SQL has been proved to be tighter than the RL algorithms that have been used for estimating the value function in PAC-MDP methods, especially in the model-free case. Also, SQL has a better computational requirement in comparison to the standard RL methods. We consider this as a subject for future research.

Another possible direction for future work is to scale up SQL to large (possibly continuous) state and action spaces where function approximation is needed. We believe that it would be possible to extend our current SQL analysis to the continuous case along the same path as in the fitted value iteration analysis by Munos and Szepesvári (2008) and Antos et al. (2007). This would require extending the error propagation result of Lemma 11 to a  $\ell_2$ -norm analysis and combining it with the standard regression bounds.

## References

- A. Antos, R. Munos, and Cs. Szepesvári. Fitted Q-iteration in continuous action-space MDPs. In *Proceedings of the 21st Annual Conference on Neural Information Processing Systems*, 2007.
- P. L. Bartlett and A. Tewari. REGAL: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence*, 2009.
- R. Bellman. *Dynamic Programming*. Princeton Univ. Press, 1957.
- D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume I. Athena Scientific, Belmont, Massachusetts, third edition, 2007a.
- D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume II. Athena Scientific, Belmont, Massachusetts, third edition, 2007b.
- D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, Massachusetts, 1996.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006.
- E. Even-Dar and Y. Mansour. Learning rates for Q-learning. *Journal of Machine Learning Research*, 5:1–25, 2003.
- E. Even-Dar, S. Mannor, and Y. Mansour. PAC bounds for multi-armed bandit and Markov decision processes. In *15th Annual Conference on Computational Learning Theory*, pages 255–270, 2002.
- W. Feller. *An Introduction to Probability Theory and Its Applications*, volume 1. Wiley, January 1968.
- T. Jaakkola, M. I. Jordan, and S. Singh. On the convergence of stochastic iterative dynamic programming. *Neural Computation*, 6(6):1185–1201, 1994.
- T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- M. Kearns and S. Singh. Finite-sample convergence rates for Q-learning and indirect algorithms. In *Advances in Neural Information Processing Systems 12*, pages 996–1002. MIT Press, 1999.
- S. Koenig and R. G. Simmons. Complexity analysis of real-time reinforcement learning. In *Proceedings of the Eleventh National Conference on Artificial Intelligence*. AAAI Press, 1993.
- S. Mannor and J. N. Tsitsiklis. The sample complexity of exploration in the multi-armed bandit problem. *Journal of Machine Learning Research*, 5:623–648, 2004.

- R. Munos and Cs. Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9:815–857, 2008.
- J. Peng and R. J. Williams. Incremental multi-step Q-learning. *Machine Learning*, 22(1-3): 283–290, 1996.
- M. L. Puterman. *Markov Decision Processes, Discrete Stochastic Dynamic Programming*. A Wiley-Interscience Publication, 1994.
- A. L. Strehl, L. Li, and M. L. Littman. Reinforcement learning in finite MDPs: PAC analysis. *Journal of Machine Learning Research*, 10:2413–2444, 2009.
- R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, Massachusetts, 1998.
- Cs. Szepesvári. The asymptotic convergence-rate of Q-learning. In *Advances in Neural Information Processing Systems 10, Denver, Colorado, USA, 1997*, 1997.
- Cs. Szepesvári. *Algorithms for Reinforcement Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2010.
- I. Szita and Cs. Szepesvári. Model-based reinforcement learning with nearly tight exploration complexity bounds. In *Proceedings of the 27th International Conference on Machine Learning*, pages 1031–1038. Omnipress, 2010.
- H. van Hasselt. Double Q-learning. In *Advances in Neural Information Processing Systems 23*, pages 2613–2621, 2010.
- C.J.C.H. Watkins. *Learning from Delayed Rewards*. PhD thesis, Kings College, Cambridge, England, 1989.