
Is Q-Learning Provably Efficient? An Extended Analysis

Kushagra Rastogi *

University of California, Los Angeles
Los Angeles, CA 90095
krastogi@g.ucla.edu

Jonathan Lee *

University of California, Los Angeles
Los Angeles, CA 90095
jlee916@ucla.edu

Fabrice Harel-Canada *

University of California, Los Angeles
Los Angeles, CA 90095
fabricehc@cs.ucla.edu

Aditya Joglekar

University of California, Los Angeles
Los Angeles, CA 90095
adivj123@gmail.com

Abstract

This work extends the analysis of the theoretical results presented within the paper *Is Q-Learning Provably Efficient?* by Jin *et al.* [1]. We include a survey of related research to contextualize the need for strengthening the theoretical guarantees related to perhaps the most important threads of model-free reinforcement learning. We also expound upon the reasoning used in the proofs to highlight the critical steps leading to the main result showing that Q-learning with UCB exploration achieves a sample efficiency that matches the optimal regret that can be achieved by any model-based approach.

Introduction

State-of-the-art reinforcement learning (RL) has been dominated by model-free algorithms (like Q-learning) because they are online, more expressive and need less space. However, empirical work has shown that model-free algorithms have a higher sample complexity [2, 3], meaning that they require many more samples in order to perform well on a given task. Can we make model-free algorithms sample-efficient? This is one of the most fundamental questions in the reinforcement learning community that has yet to be answered definitely. As seen in the setting of multi-armed bandits, good sample efficiency is the result of aptly managing the exploration-exploitation trade-off. In our project, we aim to elaborate on the proofs establishing that Q-learning with Upper Confidence Bound (UCB) exploration, in an episodic MDP setting and without access to a “simulator”, matches the information-theoretic regret optimum, up to a single \sqrt{H} where H is the number of steps per episode. To do this, we will leverage our current understanding of Q-learning and survey existing literature related to sample efficiency and complexity of both model-free and model-based RL methods.

Related Work

This section reviews related work that compares Model-free (MF) and Model-based (MB) reinforcement learning (RL) in general before focusing on theoretical research into their respective sample efficiencies and complexities.

*Indicates equal contribution.

Model-free vs. Model-based RL

The study of reinforcement learning has given rise to two primary approaches for maximizing cumulative rewards while interacting with an unknown environment through time: model-based and model-free algorithms. MB algorithms are the “planners” that either learn or use a model of environmental dynamics to form a suitable control policy. On the other hand, MF algorithms make no attempt to model state transitions explicitly, instead updating their state and action value functions directly. Both fundamentally and in practice, the two approaches overlap substantially; indeed MF methods act as important building blocks for MB methods [4].

Despite the similarities, MF methods like classical Q-learning [5], DQNs [6] and their variants [7, 8], most policy gradient approaches [9, 10, 3], and many others dominate most of the progress in modern RL [1]. Table 1 highlights some of the pros and cons of both approaches and highlights why MF methods enjoy wide attention in the field. Of the cons, the most problematic is the tendency for MF approaches to be sample inefficient as they require many “experiences” to train. The current work we analyze by Jin *et al.* [1] establishes that this con does not apply to the entire class of MF algorithms by showing that not only is it possible to design MF algorithms that are sample efficient, but also that Q-learning with an upper confidence bound (UCB) exploration policy *is provably efficient*. However, before expanding on the illustrative process and proofs, we review other work related to sample efficiency and complexity in the next subsection.

	Model-free (MF)	Model-based (MB)
Pros	<ul style="list-style-type: none"> ◦ Computationally less complex than MB methods, requiring no model of the environment to be effective (which can be a bottleneck for MB methods) [4] ◦ Capable of functioning online (as opposed to working with batches) [1] ◦ Require less space (memory) [1] ◦ More expressive since specifying value functions / policies are more flexible than specifying a model for the environment [1] 	<ul style="list-style-type: none"> ◦ Tend to be more sample efficient [2, 11] ◦ More efficient handling of changing goals because it does not need “personal experience” with every state-action pair [11, 4]
Cons	<ul style="list-style-type: none"> ◦ Requires (repeated) “personal experience” with many state-action pairs in order to train, makes exploration more costly [4] ◦ Tend to be less sample efficient [4, 2, 11] 	<ul style="list-style-type: none"> ◦ Suffer from model bias, i.e., they inherently assume that the learned dynamics model sufficiently accurately resembles the real environment [2, 12, 13, 11] ◦ Computationally more complex than MF methods - can be difficult to learn a good model of state transitions / rewards [4]

Table 1: Pros & Cons of MF vs. MB RL Approaches

Sample Efficiency & Complexity

Within RL, *sample efficiency* $e(\cdot)$ measures the number of inputs an agent requires in order to achieve a given level of performance [14] on a particular task. For example, for any two agents A_1 and A_2 , $e(A_1) > e(A_2)$ if A_1 requires fewer inputs to achieve the *same* performance as A_2 on a given task. The related idea of *sample complexity* measures the minimum number of inputs required to guarantee a probably approximately correct (PAC) estimator [15]. Generally, the lower the sample complexity, the more efficient the class of estimators / agents.

In the MF setting, several recent works provide empirical evidence that MF algorithms generally require higher sample complexity [2, 3]. In these cases, the authors elected to measure the duration of interactions between the agent and the environment rather than the more literal count of inputs since there is a one-to-one correspondence between the two units of measure. As an illustrative example, the authors of PILCO [2] measure their MB approach against six MF approaches [16, 17, 18, 19, 20, 21, 22] and achieved up to 5x orders of magnitude reduction in time required to succeed at the classic *cart-pole* task.

In the MB setting, several publications [23, 24, 25, 26, 27] have been able to demonstrate asymptotically optimal sample efficiency by importing ideas from the bandit literature, such as the UCB variations that our selected paper also pairs with Q-learning to prove its efficiency. If the existence of a simulator is assumed, MF methods like Speedy Q-Learning [28] can be *almost* as efficient as the

best MB algorithms [29]. Unfortunately, the value of this work is undercut by the observation that simulators generally do not do a good job of representing real-world environments where exploration is significantly harder — i.e. using a uniformly random exploration policy is optimal for the simulator in question [29]. The only theoretical result for MF without using a simulator is that of “delayed Q-learning” by Strehl *et al.* [30], which achieves a total regret of $\mathcal{O}(T^{4/5})$ — ignoring factors in S , A , and H — compared to $\mathcal{O}(\sqrt{T})$ achieved by MB methods.

This general issue with MF methods suggests that it may be fruitful to combine key elements of MF and MB approaches to increase sample efficiency. While there is presently no theoretical basis for the benefits of this line of inquiry, several researchers have [31, 32] have demonstrated that there is at least some empirical evidence supporting the utility of blending both approaches. Nagabandi *et al.* [31] combine the expressiveness of deep neural networks with a model-based controller (MBC) to achieve $3 - 5\times$ efficiency improvement over MF baselines on the MuJoCo [33] locomotion benchmark. Similarly, Pong *et al.* [32] proposed the ideas of temporal difference models (TDMs), which are a family of goal-conditioned value functions trained with MF learning, but used for MB control. Their experimental results show substantial improvements in efficiency relative to *both* high performing MF methods like DDPG [34] and HER [35] as well as MB methods on a range of RL tasks.

Table 2 summarizes the regret of various algorithms discussed above and illustrates the comparative sample efficiency of the work done in our selected paper by Jin *et al.* [1].

	Algorithm	Regret	Time	Space
MB	UCRL2 [25]	$\geq \mathcal{O}(\sqrt{H^4 S^2 AT})$	$\Omega(TS^2A)$	$\mathcal{O}(S^2AH)$
	Agrawal & Jia [23]	$\geq \mathcal{O}(\sqrt{H^3 S^2 AT})$		
	UCBVI [24]	$\mathcal{O}(\sqrt{H^2 SAT})$	$\mathcal{O}(TS^2A)$	
	vUCQ [26]	$\mathcal{O}(\sqrt{H^2 SAT})$		
MF	Delayed Q-learning [30]	$\mathcal{O}_{S,A,H}(T^{4/5})$	$\mathcal{O}(T)$	$\mathcal{O}(SAH)$
	Q-learning (UCB-H) [1]	$\mathcal{O}(\sqrt{H^4 SAT})$		
	Q-learning (UCB-B) [1]	$\mathcal{O}(\sqrt{H^3 SAT})$		
	information theoretic lower bound [1]	$\Omega(\sqrt{H^2 SAT})$	—	—

Table 2: Regret comparisons for RL methods on Episodic MDP where $T = KH$ is the total number of steps, H is the steps per episode, S is the number of states, and A is the number of actions. NOTE: this table is presented for $T \geq \text{poly}(S, A, H)$, and thus omits the lower order terms.

Preliminary

The notation used in this paper is mostly adapted from [1]. We consider an episodic Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$, where \mathcal{S} is a finite set of states with $|\mathcal{S}| = S$, \mathcal{A} is a finite set of actions with $|\mathcal{A}| = A$, H is the number of steps in each episode, \mathbb{P} is the transition matrix where $\mathbb{P}_h(\cdot|x, a)$ is the distribution of states when action a is taken at state x at step $h \in [H]$ and $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a deterministic reward function at step h .

Each episode of the MDP begins with the agent at state x_1 . For each step $h \in [H]$, the agent observes state $x_h \in \mathcal{S}$, takes action $a_h \in \mathcal{A}$, receives reward $r_h(x_h, a_h)$ and subsequently transitions to the next state x_{h+1} that is drawn from $\mathbb{P}_h(\cdot|x_h, a_h)$. The episode ends when the agent reaches the terminal state x_{H+1} .

We define $V_h^\pi : \mathcal{S} \rightarrow \mathbb{R}$ as the agent’s state-value function at step h under policy π . We define $Q_h^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ as the agent’s Q-value function at step h under policy π .

$$V_h^\pi(x) = \mathbb{E} \left[\sum_{h'=h}^H r_{h'}(x_{h'}, \pi_{h'}(x_{h'})) | x_h = x \right]$$

$$Q_h^\pi(x, a) = r_h(x, a) + \mathbb{E} \left[\sum_{h'=h+1}^H r_{h'}(x_{h'}, \pi_{h'}(x_{h'})) | x_h = x, a_h = a \right]$$

For finite state and action spaces, we define the optimal state-value function as $V_h^*(x) = \max_\pi V_h^\pi(x)$ $\forall x \in \mathcal{S}$ and $h \in [H]$ with optimal policy π^* . Let the total number of episodes be K , initial

state be x_1^k for episode k and policy be π_k for the k th episode. Then, the total expected regret is $\text{Regret}(K) = \sum_{k=1}^K [V_1^*(x_1^k) - V_1^{\pi_k}(x_1^k)]$.

Main Results

We combine Q-learning with a UCB exploration strategy which has the following Q-value update: $Q_h(x, a) \leftarrow (1 - \alpha_t)Q_h(x, a) + \alpha_t[r_h(x, a) + V_{h+1}(x') + b_t]$ where t counts the number of times the algorithm has visited state-action pair (x, a) at step h , x' is the next state, b_t is the confidence bonus and $\alpha_t = \frac{H+1}{H+t}$ is the step-size (learning rate). This choice of α_t scales as $\mathcal{O}(H/t)$ which allows the regret to be sub-exponential in H , thus making Q-learning efficient.

Q-learning with Hoeffding bonus. Since $r_h \in [0, 1]$ and there are H steps in each episode, the Q-values are upper-bounded by H . By the Azuma-Hoeffding inequality, the Q-values confidence bound scales as $\mathcal{O}(1/\sqrt{t})$ if the state-action pair (x, a) is visited t times. Thus, a simple bonus would be $b_t = \mathcal{O}\left(\sqrt{\frac{H^3 \iota}{t}}\right)$ where $\iota = \log(SAT/p)$. We present Q-learning algorithm with UCB-Hoeffding bonus.

Algorithm 1 Q-learning with UCB-Hoeffding

```

1: Initialize  $Q_h(x, a) \leftarrow H, N_h(x, a) \leftarrow 0 \forall (x, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ 
2: for episode  $k = 1$  to  $K$  do
3:   get  $x_1$ 
4:   for step  $h = 1$  to  $H$  do
5:      $a_h \leftarrow \text{argmax}_{a'} Q_h(x_h, a')$ 
6:      $t = N_h(x, a) \leftarrow N_h(x, a) + 1$ 
7:      $b_t \leftarrow c\sqrt{H^3 \iota/t}$  where  $c > 0$  is a constant and  $\iota = \log(SAT/p)$ 
8:      $Q_h(x_h, a_h) \leftarrow (1 - \alpha_t)Q_h(x_h, a_h) + \alpha_t[r_h(x_h, a_h) + V_{h+1}(x_{h+1}) + b_t]$ 
9:      $V_h(x_h) \leftarrow \min\left(H, \max_{a' \in \mathcal{A}} Q_h(x_h, a')\right)$ 
10:   end for
11: end for

```

Theorem 1 (Hoeffding). *If $b_t = c\sqrt{H^3 \iota/t}$, then with probability $1 - p \forall p \in (0, 1)$, the total regret of Algorithm 1 is at most $\mathcal{O}(\sqrt{H^4 SAT \iota})$ where $c > 0$ is a constant and $\iota = \log(SAT/p)$.*

Algorithm 1 has a \sqrt{T} regret without having access to a simulator which makes it very efficient and comparable to model-based algorithms. As an online learning algorithm, Algorithm 1 only stores the Q-value table and has superior time and space complexities when $|\mathcal{S}|$ is large.

Theorem 2 (Bernstein). *For a specified b_t , with probability $1 - p \forall p \in (0, 1)$, the total regret of Q-learning with UCB-Bernstein exploration is at most $\mathcal{O}(\sqrt{H^3 SAT \iota} + \sqrt{H^9 S^3 A^3 \iota^4})$.*

Q-learning with UCB-Bernstein exploration improves the total regret by a factor of \sqrt{H} over Q-learning with UCB-Hoeffding exploration. Thus, the asymptotic regret of UCB-Bernstein is only a \sqrt{H} factor away from the optimal regret achieved by model-based algorithms. However, when T is small, total regret of UCB-Bernstein exploration is dominated by $\mathcal{O}(\sqrt{H^9 S^3 A^3 \iota^4})$.

Theorem 3 (Information-theoretic lower bound). *The total regret for any algorithm in an episodic MDP setting must be at least $\Omega(\sqrt{H^2 SAT})$.*

Note that the upper bounds mentioned in Theorem 1 and 2 differ from the optimal regret by a factor of H and \sqrt{H} respectively.

Proofs for Algorithm 1

Notation. We have (x_h^k, a_h^k) = the state-action pair observed and chosen at step h of episode k . $\mathbb{I}[A]$ is the indicator function for event A . We use Q_h^k, V_h^k, N_h^k to represent the Q_h, V_h, N_h functions at

the beginning of episode k . We get the following update rules for Algorithm 1:

$$V_h^k(x) \leftarrow \min \left(H, \max_{a' \in \mathcal{A}} Q_h^k(x, a') \right), \forall x \in \mathcal{S}$$

$$Q_h^{k+1}(x, a) = \begin{cases} (1 - \alpha_t)Q_h^k(x, a) + \alpha_t[r_h(x, a) + V_{h+1}^k(x_{h+1}^k) + b_t], & \text{if } (x, a) = (x_h^k, a_h^k) \\ Q_h^k(x, a), & \text{otherwise} \end{cases} \quad (1)$$

We have $[\mathbb{P}_h V_{h+1}](x, a) = \mathbb{E}_{x' \sim \mathbb{P}_h(\cdot | x, a)} V_{h+1}(x')$ and its empirical counterpart of episode k is $[\hat{\mathbb{P}}_h^k V_{h+1}](x, a) = V_{h+1}(x_{h+1}^k)$ which is only defined for $(x, a) = (x_h^k, a_h^k)$.

The learning rate is $\alpha_t = \frac{H+1}{H+t}$. Also, we present $\alpha_t^0 = \prod_{j=1}^t 1 - \alpha_j$ and $\alpha_t^i = \prod_{j=i+1}^t 1 - \alpha_j$. Since empty products are equal to 1 and empty summations equal to 0, we get $\alpha_t^0 = 1$ and $\sum_{i=1}^t \alpha_t^i = 0$ for $t = 0$. For $t \geq 1$, we get $\alpha_t^0 = \prod_{j=1}^t \frac{j-1}{H+j} = 0$ and $\sum_{i=1}^t \alpha_t^i = 1$. Combining these equations with (1), we get:

$$Q_h^k(x, a) = \alpha_t^0 H + \sum_{i=1}^t \alpha_t^i \left[r_h(x, a) + V_{h+1}^{k_i}(x_{h+1}^{k_i}) + b_i \right] \quad (2)$$

Lemma 1.1. *Properties of α_t^i :*

- (a) For every $t \geq 1$, $\frac{1}{\sqrt{t}} \leq \sum_{i=1}^t \frac{\alpha_t^i}{\sqrt{i}} \leq \frac{2}{\sqrt{t}}$.
- (b) For every $t \geq 1$, $\max_{i \in [t]} \alpha_t^i \leq \frac{2H}{t}$ and $\sum_{i=1}^t (\alpha_t^i)^2 \leq \frac{2H}{t}$.
- (c) For every $i \geq 1$, $\sum_{t=i}^{\infty} \alpha_t^i = 1 + \frac{1}{H}$.

Proof of Lemma 1.1. Our choice of the learning rate is crucial for Q-learning to be efficient. Property (c) is particularly important to bound the regret by a constant factor of $(1 + \frac{1}{H})^H$ for each step in each episode. We provide proofs for the properties.

(a) We use induction on t . For the base case $t = 1$, we get $\alpha_1^1 = 1$. Note that $\alpha_t^i = (1 - \alpha_t)\alpha_{t-1}^i$ for $i = 1, \dots, t-1$ and $t \geq 2$. This means $\sum_{i=1}^t \frac{\alpha_t^i}{\sqrt{i}} = \frac{\alpha_t}{\sqrt{t}} + (1 - \alpha_t) \sum_{i=1}^{t-1} \frac{\alpha_{t-1}^i}{\sqrt{i}}$. Recall that $H \geq 1$ for Q-learning to be meaningful. Using induction on both sides, we can show that $\frac{\alpha_t}{\sqrt{t}} + (1 - \alpha_t) \sum_{i=1}^{t-1} \frac{\alpha_{t-1}^i}{\sqrt{i}} \geq \frac{1}{\sqrt{t}}$ and $\frac{\alpha_t}{\sqrt{t}} + (1 - \alpha_t) \sum_{i=1}^{t-1} \frac{\alpha_{t-1}^i}{\sqrt{i}} \leq \frac{2}{\sqrt{t}}$.

(b) We have $\alpha_t^i = \frac{H+1}{H+i} \left(\frac{i}{H+i+1} \frac{i+1}{H+i+2} \dots \frac{t-1}{H+t} \right)$. By rearranging, we get $\alpha_t^i = \frac{H+1}{H+t} \prod_{i=1}^t \frac{i}{H+i} = \max_{i \in [t]} \alpha_t^i$. Each term in the product resembles $\frac{x}{x+y}$ with $y \geq 1$. Thus, $\frac{x}{x+y} \leq 1$ and hence $\alpha_t^i \leq \frac{H+1}{H+t}$. Since, $\frac{H+1}{H+t} \leq \frac{H+H}{H+t} \leq \frac{H+H}{t}$, then $\alpha_t^i \leq \frac{2H}{t}$. Thus, we have shown that $\max_{i \in [t]} \alpha_t^i \leq \frac{2H}{t}$. But $\sum_{i=1}^t (\alpha_t^i)(\alpha_t^i) \leq \sum_{i=1}^t \alpha_t^i (\max_{i \in [t]} \alpha_t^i)$ which implies $\sum_{i=1}^t (\alpha_t^i)^2 \leq \frac{2H}{t}$.

(c) We have

$$\begin{aligned} \sum_{t=1}^{\infty} \alpha_t^i &= \sum_{t=1}^{\infty} \alpha_i \prod_{j=i+1}^t (1 - \alpha_j) = \alpha_i \sum_{t=1}^{\infty} \prod_{j=i+1}^t (1 - \alpha_j) \\ &= \frac{H+1}{H+i} \left(1 + \frac{i}{H+i+1} + \frac{i}{H+i+1} \frac{i+1}{H+i+2} + \dots \right) \end{aligned}$$

To simplify the last equality, we conjecture the following identity and prove it by induction:

$$\frac{n}{k} = 1 + \frac{n-k}{n+1} + \frac{n-k}{n+1} \frac{n-k+1}{n+2} + \dots$$

where $n, k > 0$ and $n \geq k$.

Note that this is equivalent to induction on $\frac{n}{k} - \sum_{i=0}^t x_i = \frac{n-k}{k} \prod_{i=1}^t \frac{n-k+i}{n+i}$. For the base case $t = 1$, we get $\frac{n}{k} - 1 - \frac{n-k}{n+1} = \frac{n-k}{k} - \frac{n-k}{n+1}$ and $\frac{n-k}{k} \frac{n-k+1}{n+1} = \frac{n-k}{k} \left(1 - \frac{k}{n+1}\right) = \frac{n-k}{k} - \frac{n-k}{n+1}$. Assume the induction hypothesis holds for $t = m$ so $\frac{n}{k} - \sum_{i=0}^m x_i = \frac{n-k}{k} \prod_{i=1}^m \frac{n-k+i}{n+i}$. For $t = m + 1$,

$$\begin{aligned} \frac{n}{k} - \sum_{i=0}^m x_i - x_{m+1} &= \frac{n-k}{k} \prod_{i=1}^m \frac{n-k+i}{n+i} - x_{m+1} \\ &= \frac{n-k}{k} \prod_{i=1}^m \frac{n-k+i}{n+i} - \prod_{i=1}^{m+1} \frac{n-k+i-1}{n+i} \\ &= \left(\frac{n-k}{k} \prod_{i=1}^m \frac{n-k+i}{n+i} \right) \left(1 - \frac{k}{n+m+1} \right) \\ &= \left(\frac{n-k}{k} \prod_{i=1}^m \frac{n-k+i}{n+i} \right) \left(\frac{n-k+m+1}{n+m+1} \right) \\ &= \frac{n-k}{k} \prod_{i=1}^{m+1} \frac{n-k+i}{n+i} \end{aligned}$$

This finishes the induction. By taking $n = H + i$ and $k = H$, we get $\sum_{t=1}^{\infty} \alpha_t^i = \frac{H+1}{H+i} \frac{H+i}{H} = \frac{H+1}{H} = 1 + \frac{1}{H}$. This concludes the proof of Lemma 1.1.

Lemma 1.2. For any $(x, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and episode $k \in [K]$ let $t = N_h^k(x, a)$ and suppose (x, a) was previously taken at step h of episodes $k_1, k_2, \dots, k_t < k$. Then:

$$(Q_h^k - Q_h^*)(x, a) = \alpha_t^0 (H - Q_h^*(x, a)) + \sum_{i=1}^t \alpha_t^i \left[(V_{h+1}^{k_i} - V_{h+1}^*)(x_{h+1}^{k_i}) + [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h) V_{h+1}^*](x, a) + b_i \right]$$

Proof of Lemma 1.2. Recall that $\sum_{i=1}^t \alpha_t^i = 1$ and $[\hat{\mathbb{P}}_h^{k_i} V_{h+1}](x, a) = V_{h+1}(x_{h+1}^{k_i})$. The Bellman optimality equation is $Q_h^*(x, a) = (r_h + \mathbb{P}_h V_{h+1}^*)(x, a)$. Then, $\sum_{i=1}^t \alpha_t^i r_h(x, a) = r_h(x, a) \sum_{i=1}^t \alpha_t^i = r_h(x, a)$. Similarly, $[\mathbb{P}_h V_{h+1}^*](x, a) = [\mathbb{P}_h V_{h+1}^*](x, a) - [\hat{\mathbb{P}}_h^{k_i} V_{h+1}^*](x, a) + V_{h+1}^*(x_{h+1}^{k_i})$ and the same trick with $\sum_{i=1}^t \alpha_t^i$ applies here too. Furthermore, $Q_h^*(x, a) = (\alpha_t^0 Q_h^* + r_h + \mathbb{P}_h V_{h+1}^*)(x, a)$ where $\alpha_t^0 = \begin{cases} 0, & t \geq 1 \\ 1, & t = 0 \end{cases}$. This manipulation is valid since $t = 1$

represents the start of the episode so $Q_h^*(x, a)$ is technically just defined as itself at $t = 0$. By consolidating everything we get:

$$Q_h^*(x, a) = \alpha_t^0 Q_h^*(x, a) + \sum_{i=1}^t \alpha_t^i \left[r_h(x, a) + (\mathbb{P}_h - \hat{\mathbb{P}}_h^{k_i}) V_{h+1}^*(x, a) + V_{h+1}^*(x_{h+1}^{k_i}) \right] \quad (3)$$

We attain Lemma 1.2 by $Q_h^k(x, a) - Q_h^*(x, a)$ where $Q_h^k(x, a)$ comes from (2) and $Q_h^*(x, a)$ comes from (3). This concludes the proof of Lemma 1.2.

Lemma 1.3. There exists an absolute constant $c > 0$ such that, for any $p \in (0, 1)$, letting $b_t = c\sqrt{H^3 t}/t$, we have $\beta_t = 2 \sum_{i=1}^t \alpha_t^i b_i \leq 4c\sqrt{H^3 t}/t$ and, with probability at least $1 - p$, the following holds simultaneously $\forall (x, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$:

$$0 \leq (Q_h^k - Q_h^*)(x, a) \leq \alpha_t^0 H + \sum_{i=1}^t \alpha_t^i (V_{h+1}^{k_i} - V_{h+1}^*)(x_{h+1}^{k_i}) + \beta_t$$

Notation. The idea behind this lemma is to construct an upper confidence bound on the optimal state-action values, $Q_h^* \forall h \in \{1, 2, \dots, H\}$. Before going into the proof, we first define some notation.

For each state-action-step pair $(x, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, we denote k_i as the episode in which (x, a, h) occurs for the i^{th} time. Otherwise, $k_i, k_{i+1}, \dots, k_K = K + 1$ if (x, a, h) only occurs $i - 1$ times

over the K episodes. It is important to note that the K episodes are indexed based on the ordering in which they were observed, that is, $k = j$ indicates the j^{th} episode observed. Consequently, k_i is denoted as

$$k_i = \begin{cases} \min(k \in [K] \mid \{k > k_{i-1} \wedge (x_h^k, a_h^k)\} \cup \{K+1\}), & i \in [K] \\ 0, & i = 0 \end{cases}$$

The aforementioned notation will be utilized for the proofs of this lemma and Theorem 1.

Proof of Lemma 1.3. For every fixed $(x, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, let $t = N_h^k(x, a)$, indicating the number of occurrences of (x, a, h) before the start of episode k . Moreover, let \mathcal{F}_i be a σ -field generated by all random variables up to episode k_i , step h . In the context of a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, $\{\mathcal{F}_i\}_{i=1}^K$ is defined as a filtration over $(\Omega, \mathcal{F}, \mathbb{P})$ consisting of an increasing family of sub- σ -fields [36] of the event space \mathcal{F} , where a σ -field \mathcal{F}_j can be interpreted as the accumulative information or collection of events generated from the observation of outcomes from the past episodes k_1, k_2, \dots, k_{j-1} and the current episode k_j . Note that we will only be concerned with episodes for which outcome (x, a, h) occurs by use of an indicator function $\mathbb{I}[k_i \leq K]$ in the latter half of the proof.

From the error $[(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a)$ of the empirical data in Lemma 1.2 along with the predefined notion of k_i and filtration, we now construct the sequence

$$\mathbb{E} \left[[(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \mid \mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_i \right] = 0 \quad \forall i \in \{1, 2, \dots, K\} \quad (4)$$

The result shown above stems from the fact that taking the expectation of $\hat{\mathbb{P}}_h^{k_i} V_{h+1}^*(x, a)$ conditioned on the past σ -fields $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_i$ provides knowledge of the probability transition matrix $\mathbb{P}_h(x' \mid x, a)$, which implies the following:

$$\begin{aligned} \mathbb{E} \left[[(\hat{\mathbb{P}}_h^{k_i} V_{h+1}^*](x, a) \mid \mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_i \right] &= \sum_{x' \in \mathcal{S}} \mathbb{P}_h(x' \mid x, a) \cdot V_{h+1}^*(x') \\ &= \mathbb{E} \left[[\mathbb{P}_h V_{h+1}^*](x, a) \mid \mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_i \right] \end{aligned}$$

Given that we assume the setting to be a tabular episodic finite-horizon MDP, $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$, where $|\mathcal{S}|, |\mathcal{A}|$, and H are finite with a finite amount of episodes K , then

$$\mathbb{E} \left[\left| [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| \right] < \infty \quad \forall i \in \{1, 2, \dots, K\} \quad (5)$$

Since (4) and (5) hold true, the sequence of empirical errors $\{[(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a)\}_{i=1}^K$ can be interpreted as a martingale difference sequence (MDS) with respect to the filtration $\{\mathcal{F}_i\}_{i=1}^K$ [37]. Therefore, we can use the Azuma-Hoeffding inequality to give a concentration result [38] for each index in the MDS, i.e., to construct confidence bounds for $Q_h^* \forall h \in \{1, 2, \dots, H\}$. Applying Azuma-Hoeffding and a union bound over all K episodes gives the following:

$$\left| \sum_{i=1}^{\tau} \alpha_{\tau}^i \cdot \mathbb{I}[k_i \leq K] \cdot [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| \leq \frac{cH}{2} \sqrt{\sum_{i=1}^{\tau} (\alpha_{\tau}^i)^2 \cdot \iota} \leq c \sqrt{\frac{H^3 \iota}{\tau}} \quad \forall \tau \in [K] \quad (6)$$

for some absolute constant c , with probability at least $1 - \frac{p}{5AH}$. Recall that $\mathbb{I}[k_i \leq K]$ is an indicator function that filters out episodes where (x, a) was not taken at step h . To prove the left inequality in (6), we consider a previously stated fact that $r_h \in [0, 1]$, implying $Q_h(x, a) \leq H$ and thus $V_h(x) \leq H$ for any x, a, h :

$$\begin{aligned} \left| [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| &\leq H \leq cH \leq \sqrt{2}cH \leq \sqrt{2}\alpha_{\tau}^i cH = c_i \\ &\forall i \in \{1, 2, \dots, K\}, c > 0 \end{aligned} \quad (7)$$

Note that c_i is the symmetric bound on the martingale difference $[(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a)$ and is used for the Azuma-Hoeffding inequality:

$$\mathbb{P} \left[\left| \sum_{i=1}^{\tau} \alpha_{\tau}^i \cdot \mathbb{I}(k_i \leq K) \cdot [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| \geq \epsilon \right] \leq 2 \exp \left(- \frac{2\epsilon^2}{\sum_{i=1}^K c_i^2} \right) \quad (8)$$

whose complementary event is

$$\mathbb{P} \left[\left| \sum_{i=1}^{\tau} \alpha_{\tau}^i \cdot \mathbb{I}(k_i \leq K) \cdot [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| \leq \epsilon \right] \geq 1 - 2 \exp \left(- \frac{2\epsilon^2}{\sum_{i=1}^K c_i^2} \right) \quad (9)$$

To find the proper choice of ϵ , we revisit the bound on the martingale difference:

$$\begin{aligned} \left| \sum_{i=1}^{\tau} \alpha_{\tau}^i \cdot \mathbb{I}(k_i \leq K) \cdot [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| &\leq \left| \sum_{i=1}^{\tau} \alpha_{\tau}^i \cdot \mathbb{I}(k_i \leq K) \cdot cH \right| \\ &= cH \sqrt{\sum_{i=1}^{\tau} (\alpha_{\tau}^i)^2 \cdot (\mathbb{I}(k_i \leq K))^2} \\ &\leq cH \sqrt{\sum_{i=1}^{\tau} (\alpha_{\tau}^i)^2} \leq cH \sqrt{\sum_{i=1}^{\tau} (\alpha_{\tau}^i)^2 \cdot \iota} = \epsilon \end{aligned} \quad (10)$$

With c_i in (7) and ϵ in (10), we can rewrite the right-hand side of the inequality in (9) as

$$1 - 2 \exp \left(- \frac{2\epsilon^2}{\sum_{i=1}^K c_i^2} \right) = 1 - \frac{2p}{SAH} \quad (11)$$

Therefore, results from (9) and (11) indicate an upper bound on the left-hand side of (6) with a probability of $1 - \frac{2p}{SAH}$. Rescaling p to $\frac{p}{2}$ finishes the proof of the left inequality of (6).

To remove the notation of learning rate as shown on the right-hand side of (6), we apply property (b) of Lemma 1.1, which gave an inclusive upper bound of $\frac{2H}{t}$ for $\sum_{i=1}^t (\alpha_t^i)^2$, $\forall t \geq 1$. Making the substitution on the middle term of (6), that is, $\frac{cH}{2} \sqrt{\sum_{i=1}^{\tau} (\alpha_{\tau}^i)^2} \cdot \iota$, concludes the proof of (6).

Because the inequality in (6) holds for all fixed $\tau \in [K]$ uniformly, it also holds for $\tau = t = N_h^k(x, a) \leq [K]$. As a result, we can rewrite (6) in a way that removes the indicator function:

$$\left| \sum_{i=1}^t \alpha_t^i \cdot [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x, a) \right| \leq c \sqrt{\frac{H^3 \iota}{t}} \quad \text{where } t = N_h^k(x, a) \quad (12)$$

If we choose the Hoeffding-style bonus b_t to be $c\sqrt{\frac{H^3 \iota}{t}}$ from the equation above, then from property (a) in Lemma 1.1,

$$\sum_{i=1}^t \alpha_t^i b_i = \sum_{i=1}^t \alpha_t^i \cdot c \sqrt{\frac{H^3 \iota}{t}} \in \left[c \sqrt{\frac{H^3 \iota}{t}}, 2c \sqrt{\frac{H^3 \iota}{t}} \right] \quad (13)$$

For notational convenience, we introduce $\frac{\beta}{2} = \sum_{i=1}^t \alpha_t^i b_i$. The final step is putting everything together to yield an upper confidence bound for Q_h^* :

$$\begin{aligned} (Q_h^k - Q_h^*)(x, a) &\leq \alpha_t^0 \cdot (H - Q_h^*(x, a)) + \sum_{i=1}^t \alpha_t^i \cdot \left((V_{h+1}^{k_i} - V_{h+1}^*)(x_{h+1}^{k_i}) + [(\hat{\mathbb{P}}_h^{k_i} - \mathbb{P}_h)V_{h+1}^*](x_{h+1}^{k_i}, a_{h+1}^{k_i}) + b_i \right) \\ &\leq \alpha_t^0 \cdot (H - Q_h^*(x, a)) + \sum_{i=1}^t \alpha_t^i \cdot \left((V_{h+1}^{k_i} - V_{h+1}^*)(x_{h+1}^{k_i}) + b_i \right) + c \sqrt{\frac{H^3 \iota}{t}} \\ &\leq \alpha_t^0 \cdot H + \sum_{i=1}^t \alpha_t^i \cdot \left((V_{h+1}^{k_i} - V_{h+1}^*)(x_{h+1}^{k_i}) \right) + \beta_t \end{aligned}$$

where the first inequality stems immediately from Lemma 1.2. The right inequality in (6) is then applied as an inclusive upper bound for the next step. Lastly, the definition of β and the fact that $\sum_{i=1}^t \alpha_i^i \leq 1$ are utilized to construct the final inequality, thus completing the proof of Lemma 1.3.

Proof of Theorem 1.

The proof of Theorem 1 uses Lemma 1.3 and the Azuma-Hoeffding inequality to produce a recursive formulation for the upper bound of the regret. Figure 1 illustrates the high-level flow of the proof to follow.

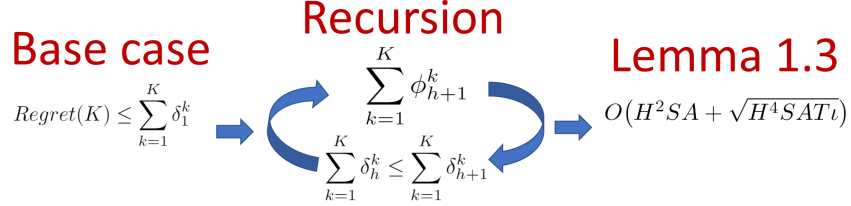


Figure 1: Flowchart for the Proof of Theorem 1.

We define $\delta_h^k := (V_h^k - V_h^{\pi^k})(x_h^k)$ and $\phi_h^k := (V_h^k - V_h^*)(x_h^k)$. Using Lemma 1.3, the regret can be upper bounded as $Regret(K) \leq \sum_{k=1}^K \delta_1^k$.

The main idea is to upper bound $\sum_{k=1}^K \delta_h^k$ by the next step $\sum_{k=1}^K \delta_{h+1}^k$ which gives a recursive relation for the total regret. For any fixed $(k, h) \in [K] \times [H]$, let $t = N_h^k(x_h^k, a_h^k)$ and suppose (x_h^k, a_h^k) was previously taken at step h of episodes $k_1, k_2, \dots, k_t < k$. Then,

$$\begin{aligned} \delta_h^k &\leq (Q_h^k - Q_h^{\pi^k})(x_h^k, a_h^k) = (Q_h^k - Q_h^*)(x_h^k, a_h^k) + (Q_h^* - Q_h^{\pi^k})(x_h^k, a_h^k) \\ &\leq \alpha_t^0 H + \sum_{i=1}^t \alpha_i^i \phi_{h+1}^{k_i} + \beta_t + [\mathbb{P}_h(V_{h+1}^* - V_{h+1}^{\pi^k})](x_h^k, a_h^k) \\ &= \alpha_t^0 H + \sum_{i=1}^t \alpha_i^i \phi_{h+1}^{k_i} + \beta_t - \phi_{h+1}^k + \delta_{h+1}^k + \epsilon_{h+1}^k \end{aligned} \quad (14)$$

where $\beta_t = 2 \sum \alpha_i^i b_i \leq \mathcal{O}(1) \sqrt{H^3 t}$ and $\epsilon_{h+1}^k = [(\mathbb{P}_h - \hat{\mathbb{P}}_h^k)(V_{h+1}^* - V_{h+1}^{\pi^k})](x_h^k, a_h^k)$ is a martingale difference sequence. Line 1 uses the definition of Q-value function and $V_h^k(x_h^k) \leq \max_{a' \in A} Q_h^k(x_h^k, a') = Q_h^k(x_h^k, a_h^k)$. Line 2 follows from Lemma 1.3, the Bellman equation $Q_h^{\pi^k}(x, a) = (r_h + \mathbb{P}_h V_{h+1}^{\pi^k})(x, a)$ and Bellman optimality equation $Q_h^*(x, a) = (r_h + \mathbb{P}_h V_{h+1}^*)(x, a)$. Finally, Line 3 holds by definition of $\delta_{h+1}^k - \phi_{h+1}^k = (V_{h+1}^* - V_{h+1}^{\pi^k})(x_{h+1}^k)$.

Now, we use (14) to compute $\sum_{k=1}^K \delta_h^k$. Hence, we get:

$$\sum_{k=1}^K \delta_h^k \leq \sum_{k=1}^K \alpha_t^0 H + \sum_{k=1}^K \sum_{i=1}^{n_h^k} \alpha_{n_h^k}^i \phi_{h+1}^{k_i(x_h^k, a_h^k)} + \sum_{k=1}^K \delta_{h+1}^k + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \quad (15)$$

Let $n_h^k = t = N_h^k(x_h^k, a_h^k)$. The first term of (15) is $\sum_{k=1}^K \alpha_{n_h^k}^0 H = \sum_{k=1}^K H \cdot \mathbb{I}[n_h^k = 0] \leq SAH$.

The equality follows from $\alpha_t^0 = \begin{cases} 0, & t \geq 1 \\ 1, & t = 0 \end{cases}$. The inequality stems from the fact that, in the worst case, $n_h^k = 0$ for all state-action pairs (x, a) which results in an upper bound of SAH .

Next, we bound the second term of (15): $\sum_{k=1}^K \sum_{i=1}^{n_h^k} \alpha_{n_h^k}^i \phi_{h+1}^{k_i(x_h^k, a_h^k)}$ where $k_i(x_h^k, a_h^k)$ is the episode in which (x_h^k, a_h^k) was taken at step h for the i th time. We first reorder the $\alpha_{n_h^k}^i$ and $\phi_{h+1}^{k_i(x_h^k, a_h^k)}$ terms. Note that $n_h^k = n_h^{k'} + j$ where $j = 1, 2, \dots$ is the j th time $\phi_{h+1}^{k'}$ appears in the summand

due to the fact that $\forall k' \in [K]$, the term $\phi_{h+1}^{k'}$ appears in the summand with $k > k'$ if and only if $(x_h^k, a_h^k) = (x_h^{k'}, a_h^{k'})$. This results in the following simplification:

$$\sum_{k=1}^K \sum_{i=1}^{n_h^k} \alpha_{n_h^k}^i \phi_{h+1}^{k_i(x_h^k, a_h^k)} \leq \sum_{k=1}^K \phi_{h+1}^{k'} \sum_{t=n_h^{k'}+1}^{n_h^k} \alpha_{n_h^k}^t \leq \left(1 + \frac{1}{H}\right) \sum_{k=1}^K \phi_{h+1}^k$$

where the first inequality uses the reasoning above and the final inequality uses property (c) of Lemma 1.1. Plugging the above inequalities into (15) results in:

$$\begin{aligned} \sum_{k=1}^K \delta_h^k &\leq SAH + \left(1 + \frac{1}{H}\right) \sum_{k=1}^K \phi_{h+1}^k - \sum_{k=1}^K \phi_{h+1}^k + \sum_{k=1}^K \delta_{h+1}^k + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \\ &= SAH + \frac{1}{H} \sum_{k=1}^K \phi_{h+1}^k + \sum_{k=1}^K \delta_{h+1}^k + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \\ &\leq SAH + \left(1 + \frac{1}{H}\right) \sum_{k=1}^K \delta_{h+1}^k + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \end{aligned} \quad (16)$$

where the last inequality is true because $\phi_{h+1}^k \leq \delta_{h+1}^k$ since $V^* \geq V_k^\pi$.

Inequality (16) recursively upper bounds $\sum_{k=1}^K \delta_h^k$ by $\sum_{k=1}^K \delta_{h+1}^k$. Applying recursion for steps $h \in \{1, 2, \dots, H\}$ and using $\delta_{H+1}^K = 0$ (the algorithm reaches the terminal state so $V_{H+1}^K = V_{H+1}^{\pi_K} = 0$) gives:

$$\begin{aligned} \sum_{k=1}^K \delta_1^k &\leq SAH + \left(1 + \frac{1}{H}\right) \sum_{k=1}^K \delta_2^k + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \\ &\leq SAH + \left(1 + \frac{1}{H}\right) \left[SAH + \left(1 + \frac{1}{H}\right) \sum_{k=1}^K \delta_3^k + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \right] + \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k) \\ &= SAH + \left(1 + \frac{1}{H}\right) SAH + \left(1 + \frac{1}{H}\right)^2 SAH + \dots + \left(1 + \frac{1}{H}\right)^{H-1} SAH \\ &\quad + \mathcal{O}\left(\sum_{h=1}^H \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k)\right) \\ &= \mathcal{O}\left(H^2 SA + \sum_{h=1}^H \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k)\right) \end{aligned} \quad (17)$$

Overall, we achieve $\sum_{k=1}^K \delta_1^k \leq \mathcal{O}\left(H^2 SA + \sum_{h=1}^H \sum_{k=1}^K (\beta_{n_h^k} + \epsilon_{h+1}^k)\right)$ from (17).

By definition of β , we have $\sum_{k=1}^K \beta_{n_h^k} \leq \mathcal{O}(1) \cdot \sum_{k=1}^K \sqrt{\frac{H^3 \iota}{n_h^k}}$. Applying the pigeon-hole principle to the inequality would mean the following: Suppose we play $1/\sqrt{n}$ at a state-action pair (x, a) . If we visit (x, a) again, then we only need to play $1/\sqrt{n+1}$ since we cannot include $1/\sqrt{n}$ twice in the summation for the same (x, a) . Thus for every (x, a) , we have $\sum_{n=1}^{N_h^K(x, a)} \sqrt{\frac{1}{n}}$. Hence we get:

$$\sum_{k=1}^K \beta_{n_h^k} \leq \mathcal{O}(1) \cdot \sum_{k=1}^K \sqrt{\frac{H^3 \iota}{n_h^k}} = \mathcal{O}(1) \cdot \sum_{x, a} \sum_{n=1}^{N_h^K(x, a)} \sqrt{\frac{H^3 \iota}{n}}$$

Note that $\sum_{x, a} N_h^K(x, a) = K$ because we are summing all occurrences of state-action pairs that occur at step h over all episodes. Since there are K episodes, there are K occurrences of state-action pairs occurring at step h .

Now, we have $\sqrt{H^3\iota} \sum_{x,a} \sum_{n=1}^{N_h^K} \frac{1}{\sqrt{n}} \leq \sqrt{H^3\iota} \sum_{x,a} \sqrt{N_h^K} = \sqrt{H^3\iota} \mathbf{1}^T v$ where $v = [N_h^K(x_1, a_1), N_h^K(x_2, a_2), \dots, N_h^K(x_{SA}, a_{SA})]^T$. Using the Cauchy-Schwarz inequality, we get $\sqrt{H^3\iota} \mathbf{1}^T v \leq \sqrt{H^3\iota} \sqrt{SA \sum_{x,a} N_h^K} = \sqrt{H^3SAK\iota} = \sqrt{H^2SAT\iota}$ by realizing that $T = KH$. Consolidating everything in one place, we get the following:

$$\sum_{k=1}^K \beta_{n_h^k} \leq \mathcal{O}(1) \cdot \sum_{k=1}^K \sqrt{\frac{H^3\iota}{n_h^k}} = \mathcal{O}(1) \cdot \sum_{x,a} \sum_{n=1}^{N_h^K(x,a)} \sqrt{\frac{H^3\iota}{n}} \leq \mathcal{O}(H^3SAK\iota) = \mathcal{O}(\sqrt{H^2SAT\iota}) \quad (18)$$

By the Azuma-Hoeffding inequality, with probability $1 - p$, we get:

$$\left| \sum_{h=1}^H \sum_{k=1}^K \epsilon_{h+1}^k \right| = \left| \sum_{h=1}^H \sum_{k=1}^K [(\mathbb{P}_h - \hat{\mathbb{P}}_h^k)(V_{h+1}^* - V_{h+1}^k)](x_h^k, a_h^k) \right| \leq cH\sqrt{T\iota} \quad (19)$$

Substituting (18) and (19) in (17) gives the following with probability $1 - p$:

$$\begin{aligned} \sum_{k=1}^K \delta_1^k &\leq \mathcal{O}(H^2SA + H\sqrt{H^2SAT\iota} + cH\sqrt{T\iota}) \\ &= \mathcal{O}(H^2SA + \sqrt{H^4SAT\iota} + c\sqrt{H^2T\iota}) \\ &= \mathcal{O}(H^2SA + \sqrt{H^4SAT\iota}) \end{aligned}$$

where the final equality is valid since $c\sqrt{H^2T\iota}$ is the smallest of the three terms. This concludes the proof of Theorem 1.

Conclusion

In this paper, we showed that a subset of model-free reinforcement learning algorithms can be made sample efficient. Specifically, we proved that, in an episodic setting, Q-learning with UCB-Hoeffding exploration strategy achieves a regret of $\mathcal{O}(\sqrt{H^4SAT\iota})$. This is the first time a regret analysis features a \sqrt{T} factor for model-free algorithms that do not require access to a "simulator". Thus, the key takeaways from the paper are:

- Use UCB exploration over ε -greedy in the model-free setting for better treatment of uncertainties in different states and actions.
- Use dynamic learning rates $\alpha_t = \mathcal{O}(H/t)$ such as $\frac{H+1}{H+t}$ instead of the commonly used $1/t$ for updates at time step t . This applies more weight to more recent updates and is critical for sample-efficiency guarantees.

We can build upon our current work by examining and unfolding the proof of Q-learning with the more sophisticated UCB-Berstein exploration strategy. Lastly, we can attempt to apply the theoretical framework used in this paper to analyze the pairing of Q-learning with another kind of exploration strategy, such as optimistic initial values.

References

- [1] Chi Jin, Zeyuan Allen-Zhu, Sébastien Bubeck, and Michael I. Jordan. Is q-learning provably efficient? In *NeurIPS*, 2018.
- [2] Marc Deisenroth and Carl Rasmussen. Pilco: A model-based and data-efficient approach to policy search., 01 2011.
- [3] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization, 2015.
- [4] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.
- [5] Christopher John Cornish Hellaby Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Cambridge, UK, May 1989.
- [6] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013.
- [7] Hado V. Hasselt. Double q-learning. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 2613–2621. Curran Associates, Inc., 2010.
- [8] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay, 2015. cite arxiv:1511.05952Comment: Published at ICLR 2016.
- [9] Richard S Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. A. Solla, T. K. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 1057–1063. MIT Press, 2000.
- [10] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. *CoRR*, abs/1602.01783, 2016.
- [11] Christopher G. Atkeson and Juan Carlos Santamaria. A comparison of direct and model-based reinforcement learning. In *IN INTERNATIONAL CONFERENCE ON ROBOTICS AND AUTOMATION*, pages 3557–3564. IEEE Press, 1997.
- [12] Jeff G. Schneider. Exploiting model uncertainty estimates for safe dynamic control learning. In *Proceedings of the 9th International Conference on Neural Information Processing Systems, NIPS’96*, page 1047–1053, Cambridge, MA, USA, 1996. MIT Press.
- [13] Stefan Schaal. Learning from demonstration. In *Proceedings of the 9th International Conference on Neural Information Processing Systems, NIPS’96*, page 1040–1046, Cambridge, MA, USA, 1996. MIT Press.
- [14] Brian Everitt. *The Cambridge dictionary of statistics*. Cambridge University Press, Cambridge, UK; New York, 2002.
- [15] Leslie Valiant. *Probably Approximately Correct: Nature’s Algorithms for Learning and Prospering in a Complex World*. Basic Books, Inc., USA, 2013.
- [16] Hajime Kimura and Shigenobu Kobayashi. Efficient non-linear control by combining q-learning with local linear controllers. In *Proceedings of the Sixteenth International Conference on Machine Learning, ICML ’99*, page 210–219, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [17] Kenji Doya. Reinforcement learning in continuous time and space. *Neural Computation*, 12(1):219–245, 2000.
- [18] Rémi Coulom. *Reinforcement Learning Using Neural Networks, with Applications to Motor Control*. PhD thesis, Institut National Polytechnique de Grenoble, 2002.
- [19] Pawel Wawrzynski and Andrzej Pacut. Model-free off-policy reinforcement learning in continuous environment, 08 2004.

- [20] Martin Riedmiller. Neural fitted q iteration – first experiences with a data efficient neural reinforcement learning method. In *Proceedings of the 16th European Conference on Machine Learning*, ECML'05, page 317–328, Berlin, Heidelberg, 2005. Springer-Verlag.
- [21] Tapani Raiko and Matti Törnio. Variational bayesian learning of nonlinear hidden state-space models for model predictive control. *Neurocomputing*, 72:3704–3712, 10 2009.
- [22] Hado Philip van Hasselt. *Insights in Reinforcement Learning: formal analysis and empirical evaluation of temporal-difference learning algorithms*. PhD thesis, Universiteit Utrecht, January 2011.
- [23] Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. *ArXiv*, abs/1705.07041, 2017.
- [24] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *ICML*, 2017.
- [25] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *J. Mach. Learn. Res.*, 11:1563–1600, August 2010.
- [26] Sham Kakade, Mengdi Wang, and Lin Yang. Variance reduction methods for sublinear reinforcement learning, 02 2018.
- [27] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions, 2016.
- [28] Mohammad Gheshlaghi Azar, Remi Munos, Mohammad Ghavamzadeh, and Hilbert J. Kappen. Speedy q-learning. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, NIPS'11, page 2411–2419, Red Hook, NY, USA, 2011. Curran Associates Inc.
- [29] Mohammad Gheshlaghi Azar, Rémi Munos, and Bert Kappen. On the sample complexity of reinforcement learning with a generative model . In *ICML*. icml.cc / Omnipress, 2012.
- [30] Alexander L. Strehl, Lihong Li, Eric Wiewiora, John Langford, and Michael L. Littman. Pac model-free reinforcement learning. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, page 881–888, New York, NY, USA, 2006. Association for Computing Machinery.
- [31] Anusha Nagabandi, Gregory Kahn, Ronald Fearing, and Sergey Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning, 05 2018.
- [32] Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model-free deep rl for model-based control. *ArXiv*, abs/1802.09081, 2018.
- [33] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *IROS*, pages 5026–5033. IEEE, 2012.
- [34] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Manfred Otto Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *CoRR*, abs/1509.02971, 2015.
- [35] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5048–5058. Curran Associates, Inc., 2017.
- [36] Takis Konstantopoulos. *Conditional Expectation and Probability*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011.
- [37] Ying-Xia Chen, Shui-Li Zhang, and Fu-Qiang Ma. On the complete convergence for martingale difference sequence. *Communications in Statistics - Theory and Methods*, 46(15):7603–7611, 2017.
- [38] David Williams. *Probability with Martingales*. Cambridge University Press, 1991.