# A Hybrid PAC Reinforcement Learning Algorithm

**Ashkan Zehfroosh**
Department of Mechanical Engineering
University of Delaware
Newark, DE 19716
ashkanz@udel.edu

**Herbert G. Tanner**
Department of Mechanical Engineering
University of Delaware
Newark, DE 19716
btanner@udel.edu

## Abstract

This paper offers a new hybrid probably approximately correct (PAC) reinforcement learning (RL) algorithm for Markov decision processes (MDPs) that intelligently maintains favorable features of its parents. The designed algorithm, referred to as the Dyna-Delayed Q-learning (DDQ) algorithm, combines model-free and model-based learning approaches while outperforming both in most cases. The paper includes a PAC analysis of the DDQ algorithm and a derivation of its sample complexity. Numerical results are provided to support the claim regarding the new algorithm's sample efficiency compared to its parents as well as the best known model-free and model-based algorithms in application.

## 1 Introduction

While several reinforcement learning (RL) algorithms can apply to a dynamical system modelled as a Markov decision process (MDP), few are probably approximately correct (PAC)—namely, they can guarantee how soon the algorithm will converge to the near-optimal policy. Existing PAC MDP algorithms can be broadly divided into two groups: model-based algorithms like [1–6], and model-free Delayed Q-learning algorithms [7–9]. Each group has its advantages and disadvantages. The goal here is to capture the advantages of both groups, while preserving PAC properties.

Model-free RL is a powerful approach for learning complex tasks. For many real-world learning problems, however, the approach is taxing in terms of the size of the necessary body of data—what is more formally referred to as its *sample complexity*. The reason is that model-free RL ignores rich information from state transitions and only relies on the observed rewards for learning the optimal policy [10]. A popular model-free PAC RL MDP algorithm is known as *Delayed Q-learning* [7]. The known upper-bound on the sample complexity of Delayed Q-learning suggests that it outperforms model-based alternatives only when the state-space size of the MDP is relatively large [11].

Model-based RL, on the other hand, utilizes all information from state transitions to learn a model, and then uses that model to compute an optimal policy. The sample complexity of model-based RL algorithms are typically lower than that of model-free ones [12]; the trade-off comes in the form of computational effort and possible bias [10].

A popular model-based PAC RL MDP algorithms is *R-max* [2]. The derived upper-bound for the sample complexity of the R-max algorithm [13] suggests that this model-based algorithm shines from the viewpoint of sample efficiency when the size of the state/action space is relatively small. This efficiency assessment can typically be generalized to most model-based algorithms.

Overall, R-max and Delayed Q-learning are incomparable in terms of their bound on the sample complexity. For instance, *for the same sample size*, R-max is bound to return a policy of higher accuracy compared to Delayed Q-learning, while the latter will converge much faster on problems with large state spaces.

Typically, model-free algorithms circumvent the model learning stage of the solution process, a move that by itself reduces complexity in problems of large size. In many applications, however, model learning is not the main

complexity bottleneck. Neurophysiologically-inspired hypotheses [14] have suggested that the brain approach toward complex learning tasks can be model-free (trial and error) or model-based (deliberate planning and computation) or even combination of both, depending on the amount and reliability of the available information. This intelligent combination is postulated to contribute in making the process efficient and fast. The design of the PAC MDP algorithm presented in this paper is motivated by such observations. Rather than following strictly one of the two prevailing directions, it orchestrates a marriage of a model-free (Delayed Q-learning) with a model-based (R-max) PAC algorithm, in order to give rise to a new PAC algorithm (Dyna-Delayed Q-learning (DDQ)) that combines the advantages of both.

The search for a connection between model-free and model-based RL algorithms begins with the Dyna-Q algorithm [15], in which synthetic generated experiences based on the learned model are used to expedite Q-learning. Some other examples that continued along this thread of research are partial model back propagation [16], training a goal condition Q function [17–20], and integrating model-based linear quadratic regulator based algorithm into model-free framework of path integral policy improvement [21]. The recently introduced Temporal Difference Models (TDM) provides a smooth(er) transition from model-free to model-based, during the learning process [10]. What is missed in the literature is a PAC combination of model-free and model-based frameworks.

Here the Dyna-Q idea is leveraged to combine two popular PAC algorithms, one model-free and one model-based, into a new one named DDQ, which is not only PAC like its parents, but also inherits the best of both worlds: it will intelligently behave more like a model-free algorithm on large problems, and operate more like a model-based algorithm on problems that require high accuracy, being content with the smallest among the sample sizes required by its parents. Specifically, the sample complexity of DDQ, in the worst case, matches the minimum bound between that of R-max and Delayed Q-learning, and often outperforms both. Note that DDQ algorithm is purely online and does not assume accessing to a generative model like in [22]. While the provable worst case upper bound on the sample complexity of DDQ algorithm is higher than the best known model-based [5] and model-free [8,9] algorithms, we can demonstrate (see Section 5) that the hybrid nature allows for superior performance of the DDQ algorithm in applications. The availability of a hybrid PAC algorithm like DDQ in hand, obviates the choice between a model-free and a model-based approach.

Our own motivation for developing of this new breed of RL algorithms comes from application problems in the area of early pediatric motor rehabilitation, where robots can be used as smart toys to socially interact with infants who have special needs, and engage with them socially in play-based activity that involves gross motion. There, MDP models can be constructed to capture the dynamics of the social interaction between infant and robot, and RL algorithms can guide the behavior of the robot as it interacts with the infant in order to achieve the maximum possible rehabilitation outcome—the latter possibly quantified by the overall length of infant displacement, or the frequency of infant motor transitions. Some early attempts at modeling such instances of human-robot interaction (HRI) did not result in models of particularly large state and action spaces, but were particularly complicated by the absence of sufficient data sets for learning [23, 24]. This is because every child is different, and the exposure of each particular infant to the smart robotic toys (during which HRI data can be collected) is usually limited to a few hours per month. There is a need, therefore, for reinforcement learning approaches that can maintain efficiency and accuracy even when the learning set is particularly small.

The paper starts with some technical preliminaries in Section 2. This section introduces the required properties of a PAC RL algorithm in the form of a well-known theorem. The DDQ algorithm is introduced in Section 3, with particular emphasis given on its update mechanism. Section 4 presents the mathematical analysis that leads the establishment of the algorithm's PAC properties, and the analytic derivation of its sample complexity. Finally, Section 5 offers numerical data to support the theoretical sample complexity claims, through an illustrative grid-world example. The data indicate that DDQ outperforms Delayed Q-learning and R-max in terms of the required samples to learn near-optimal policy. To promote readability, the proofs of most of the lemmas supporting the proof of our main result are included separately in the paper's Appendix.


## 2  Technical Preliminaries

A finite MDP $M$ is a tuple $\{S, A, R, T, \gamma\}$ with elements

| | |
|---|---|
| $S$ | a finite set of states |
| $A$ | a finite set of actions |
| $R : S \times A \to [0, 1]$ | the *reward* from executing $a$ at $s$ |
| $T : S \times A \times S \to [0, 1]$ | the *transition probabilities* |
| $\gamma \in [0, 1)$ | the *discount factor* |

A *policy* $\pi$ is a mapping $\pi : S \to A$ that selects an action $a$ to be executed at state $s$. A policy is *optimal* if it maximizes the expected sum of discounted rewards; if $t$ indexes the current time step and $a_t$, $s_t$ denote current action and state, respectively, then this expected sum is written $\mathbb{E}_M\left\{\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)\right\}$. The discount factor $\gamma$ here reflects the preference of immediate rewards over future ones. The *value* of state $s$ under policy $\pi$ in MDP $M$ is defined as

$$v_M^\pi(s) = \mathbb{E}_M\left\{R\big(s, \pi(s)\big) + \sum_{t=1}^{\infty} \gamma^t R\big(s_t, \pi(s_t)\big)\right\}$$

Note that an upper bound for the value at any state is $v_{\max} = \frac{1}{1-\gamma}$. Similarly defined is the value of *state-action pair* $(s,a)$ under policy $\pi$:

$$Q_M^\pi(s,a) = \mathbb{E}_M\left\{R(s,a) + \sum_{t=1}^{\infty} \gamma^t R\big(s_t, \pi(s_t)\big)\right\}$$

Every MDP $M$ has at least one optimal policy $\pi^*$ that results in an optimal value (or state-action value) assignment at all states; the latter is denoted $v_M^*(s)$ (or $Q_M^*(s,a)$, respectively).

The standard approach to finding the optimal values is through the search for a fix point of the Bellman equation

$$v_M^*(s) = \max_a \{R(s,a) + \gamma \sum_{s'} T(s, s', a) v_M^*(s')\}$$

which, after substituting $V_M^*(s') = \max_a Q_M^*(s', a)$, can equivalently be written in terms of state-action values

$$Q_M^*(s,a) = R(s,a) + \gamma \sum_{s'} T(s, s', a) v_M^*(s')$$

Reinforcement learning, (RL) is a procedure to obtain an optimal policy in an MDP, when the actual transition probabilities and/or reward function are not known. The procedure involves exploration of the MDP model. An RL algorithm usually maintains a table of state-action pair value estimates $Q(s,a)$ that are updated based on the exploration data. We denote $Q_t(s,a)$ the currently stored value for state-action pair $(s,a)$ at timestep $t$ during the execution of an RL algorithm. Consequently, $v_t(s) = \max_a Q_t(s,a)$. An RL algorithm is *greedy* if it at any timestep $t$, it always executes action $a_t = \arg\max_{a \in A} Q_t(s_t, a)$. The policy in force at time step $t$ is similarly denoted $\pi_t$. In what follows, we denote $|S|$ the cardinality of a set $S$.

Reinforcement learning algorithms have been classified as *model-based* or *model-free*. Although the characterization is debatable, what is meant by calling an RL algorithm "model-based," is that $T$ and/or $R$ are estimated based on online observations (exploration data), and the resulting estimated model subsequently informs the computation of the the optimal policy. A model-free RL algorithm, on the other hand, would skip the construction of an estimated MDP model, and search directly for an optimal policy over the policy space. An RL algorithm is expected to converge to the optimal policy, practically reporting a near-optimal one at termination.

Probably approximately correct (PAC) analysis of RL algorithms deals with the question of how fast an RL algorithm converges to a near-optimal policy. An RL algorithm is PAC if there exists a probabilistic bound on the number of exploration steps that the algorithm can take before converging to a near-optimal policy.

**Definition 1.** *Consider that an RL algorithm $\mathcal{A}$ is executing on MDP $M$. Let $s_t$ be the visited state at time step $t$ and $\mathcal{A}_t$ denotes the (non-stationary) policy that the $\mathcal{A}$ executes at $t$. For a given $\epsilon > 0$ and $\delta > 0$, $\mathcal{A}$ is a PAC RL algorithm if there is an $N > 0$ such that with probability at least $1 - \delta$ and for all but $N$ time steps,*

$$v_M^{\mathcal{A}_t}(s_t) \geq v_M^*(s_t) - \epsilon \tag{1}$$

Equation (1) is known as the $\epsilon$-optimality condition and $N$ as the *sample complexity* of $\mathcal{A}$, which is a function of $\left(|S|, |A|, \frac{1}{\epsilon}, \frac{1}{\delta}, \frac{1}{1-\gamma}\right)$.

**Definition 2.** *Consider MDP $M = \{S, A, R, T, \gamma\}$ which at time $t$ has a set of state-action value estimates $Q_t(s,a)$, and let $K_t \subseteq S \times A$ be a set of state-action pairs labeled* known. *The known state-action MDP*

$$M_{K_t} = \left\{S \cup \{z_{s,a} | (s,a) \notin K_t\}, A, T_{K_t}, R_{K_t}, \gamma\right\}$$

*is an MDP derived from $M$ and $K_t$ by defining new states $z_{s,a}$ for each unknown state-action pair $(s,a) \notin K_t$, with self-loops for all actions, i.e., $T_{K_t}(z_{s,a}, \cdot, z_{s,a}) = 1$. For all $(s,a) \in K_t$, it is $R_{K_t}(s,a) = R(s,a)$ and $T_{K_t}(s, a, \cdot) =$*

$T(s, a, \cdot)$. When an unknown state-action pair $(s, a) \notin K_t$ is experienced, $R_{K_t}(s, a) = Q_t(s, a)(1 - \gamma)$ and the model jumps to $z_{s,a}$ with $T_{K_t}(s, a, z_{s,a}) = 1$; subsequently, $R_{K_t}(z_{s,a}, \cdot) = Q_t(s, a)(1 - \gamma)$.

Let $K_t$ be set of current known state-action pairs of an RL algorithm $\mathcal{A}$ at time $t$, and allow $K_t$ to be arbitrarily defined as long as it depends only on the history of exploration data up to $t$. Any $(s, a) \notin K_t$ experienced at time step $t$ marks an *escape event*.

**Theorem 1** ([11]). *Let $\mathcal{A}$ be a greedy RL algorithm for an arbitrary MDP $M$, and let $K_t$ be the set of current known state-action pairs, defined based only on the history of the exploration data up to timestep $t$. Assume that $K_t = K_{t+1}$ unless an update to some state-action value occurs or an escape event occurs at timestep $t$, and that $Q_t(s, a) \leq v_{\max}$ for all $(s, a)$ and $t$. Let $M_{K_t}$ be the known state-action MDP at timestep $t$ and $\pi_t(s) = \arg\max_a Q_t(s, a)$ denote the greedy policy that $\mathcal{A}$ executes. Suppose now that for any positive constant $\epsilon$ and $\delta$, the following conditions hold with probability at least $1 - \delta$ for all $s$, $a$ and $t$:*

**optimism:** $\quad v_t(s) \geq v_M^*(s) - \epsilon$

**accuracy:** $\quad v_t(s) - v_{M_{K_t}}^{\pi_t}(s) \leq \epsilon$

**complexity:** $\quad$ *sum of number of timesteps with Q-value updates plus number of timesteps with escape events is bounded by $\zeta(\epsilon, \delta) > 0$.*

*Then, executing algorithm $\mathcal{A}$ on any MDP $M$ will result in following a $4\epsilon$-optimal policy on all but*

$$\mathcal{O}\left(\frac{\zeta(\epsilon, \delta)}{\epsilon(1 - \gamma)^2} \ln\left(\frac{1}{\delta}\right) \ln\left(\frac{1}{\epsilon(1 - \gamma)}\right)\right) \simeq \mathcal{O}\left(\frac{\zeta(\epsilon, \delta)}{\epsilon(1 - \gamma)^2}\right) \tag{2}$$

*timesteps, with probability at least $1 - 2\delta$.*


## 3 DDQ Algorithm

This section presents Algorithm 1, the one we call DDQ and the main contribution of this paper. DDQ integrates elements of R-max and Delayed $Q$-learning, while preserving the implementation advantages of both. We refer to the assignment in line 31 of Algorithm 1 as a *type-1 update*, and to the one on line 52 as a *type-2 update*. Type-1 updates use the $m_1$ most recent experiences (occurances) of a state-action pair $(s, a)$ to update that pair's value, while a type-2 update is realized through a value iteration algorithm (lines $43 - 54$) and applies to state-action pairs experienced at least $m_2$ times. The outcome at timestep $t$ of the value iteration for a type-2 update is denoted $Q_t^{\text{vl}}(s, a)$. The value iteration is set to run for $\frac{\ln(1/(\epsilon_2(1-\gamma)))}{(1-\gamma)}$ iterations; parameter $\epsilon_2$ regulates the desired accuracy on the resulting estimate (Lemma 5). A type-1 update is successful only if the condition on line 30 of the algorithm is satisfied, and this condition ensures that the type-1 update necessarily decreases the value estimate by at least $\epsilon_1 = 3\epsilon_2$. Similarly, a type-2 update is successful only if the condition on line 51 of the algorithm holds. The DDQ algorithm maintains the following internal variables:

- $l(s, a)$: the number of samples gathered for the update type-1 of $Q(s, a)$ once $l(s, a) = m_1$.
- $U(s, a)$: the running sum of target values used for a type-1 update of $Q(s, a)$, once enough samples have been gathered.
- $b(s, a)$: the timestep at which the most recent or ongoing collection of $m_1$ $(s, a)$ experiences has started.
- $\text{learn}(s, a)$: a Boolean flag that indicates whether or not samples are being gathered for type-1 update of $Q(s, a)$. The flag is set to true initially, and is reset to true whenever some Q-value is updated. It flips to false when no updates to any Q-values occurs within a time window of $m_1$ experiences of $(s, a)$ in which attempted updates type-1 of $Q^i(s, a)$ fail.
- $n(s, a)$: variable that keeps track of the number of times $(s, a)$ is experienced.
- $n(s, a, s')$: variable that keeps track of the number of transitions to $s'$ on action $a$ at state $s$.
- $r(s, a)$: the accumulated rewards by doing $a$ in $s$.

The execution of the DDQ algorithm is tuned via the $m_1$ and $m_2$ parameters. One can practically reduce it to Delayed $Q$-learning by setting $m_2$ very large, and to R-max by setting $m_1$ large. The next section provides a formal

**Algorithm 1** The DDQ algorithm

---

1: **Inputs**: $S, A, \gamma, m_1, m_2, \epsilon_1, \epsilon_2$
2: **for all** $s, a, s'$ **do**
3:      $Q(s, a) \leftarrow v_{\max}$            $\triangleright$ initialize $Q$ values to its maximum
4:      $U(s, a) \leftarrow 0$            $\triangleright$ used for attempted updates of type-1
5:      $l(s, a) \leftarrow 0$            $\triangleright$ counters
6:      $b(s, a) \leftarrow 0$         $\triangleright$ beginning timestep of attempted update type-1
7:      $\mathsf{learn}(s, a) \leftarrow \mathrm{true}$          $\triangleright$ learn flags
8:      $n(s, a) \leftarrow 0$          $\triangleright$ number of times $(s, a)$ is tried
9:      $n(s, a, s') \leftarrow 0$        $\triangleright$ number of transitions to $s'$ by $a$ in $s$
10:     $r(s, a) \leftarrow 0$        $\triangleright$ accumulated reward by execution of $a$ in $s$
11: **end for**
12: $t^* \leftarrow 0$           $\triangleright$ time of the most recent successful timestep
13: **for** $t = 1, 2, 3, \ldots$ **do**
14:      let $s$ denotes the state at time $t$
15:      choose action $a = \arg\max_{a' \in A} Q(s, a')$
16:      observe immediate reward $r$ and next state $s'$
17:      $n(s, a) = n(s, a) + 1$
18:      $n(s, a, s') = n(s, a, s') + 1$
19:      $r(s, a) = r(s, a) + r$
20:      **if** $b(s, a) \leq t^*$ **then**
21:          $\mathsf{learn}(s, a) \leftarrow \mathrm{true}$
22:      **end if**
23:      **if** $\mathsf{learn}(s, a) = \mathrm{true}$ **then**
24:          **if** $l(s, a) = 0$ **then**
25:              $b(s, a) \leftarrow t$
26:          **end if**
27:          $l(s, a) \leftarrow l(s, a) + 1$
28:          $U(s, a) \leftarrow U(s, a) + r + \gamma \max_{a'} Q(s', a')$
29:          **if** $l(s, a) = m_1$ **then**
30:              **if** $Q(s, a) - U(s, a)/m_1 \geq 2\epsilon_1$ **then**
31:                  $Q(s, a) \leftarrow U(s, a)/m_1 + \epsilon_1 \ \text{ and } \ t^* \leftarrow t$
32:              **else if** $b(s, a) > t^*$ **then**
33:                  $\mathsf{learn}(s, a) \leftarrow \mathrm{false}$
34:              **end if**
35:              $U(s, a) \leftarrow 0 \ \text{ and } \ l(s, a) \leftarrow 0$
36:          **end if**
37:      **end if**
38:      **if** $n(s, a) = m_2$ **then**
39:          $t^* \leftarrow t$
40:          **for all** $(\bar{s}, \bar{a})$ **do**
41:              $Q_{\mathrm{vl}}(\bar{s}, \bar{a}) \leftarrow Q(\bar{s}, \bar{a})$
42:          **end for**
43:          **for** $i = 1, 2, 3, \ldots, \left( \frac{\ln\left(1/(\epsilon_2(1-\gamma))\right)}{(1-\gamma)} \right)$ **do**
44:              **for all** $(\bar{s}, \bar{a})$ **do**
45:                  **if** $n(\bar{s}, \bar{a}) \geq m_2$ **then**
46:                      $Q_{\mathrm{vl}}(\bar{s}, \bar{a}) \leftarrow \frac{r(\bar{s}, \bar{a})}{n(\bar{s}, \bar{a})} + \gamma \sum_{s''} \frac{n(\bar{s}, \bar{a}, s'')}{n(\bar{s}, \bar{a})} \max_{a'} Q_{\mathrm{vl}}(s'', a')$
47:                  **end if**
48:              **end for**
49:          **end for**
50:          **for all** $(\bar{s}, \bar{a})$ **do**
51:              **if** $Q_{\mathrm{vl}}(\bar{s}, \bar{a}) \leq Q(\bar{s}, \bar{a})$ **then**
52:                  $Q(\bar{s}, \bar{a}) \leftarrow Q_{\mathrm{vl}}(\bar{s}, \bar{a})$
53:              **end if**
54:          **end for**
55:      **end if**
56: **end for**

---

proof that DDQ is not only PAC but also *possesses the minimum sample complexity* between R-max and Delayed $Q$-learning in the worst case —often, it outperforms both.

# 4   PAC Analysis of DDQ Algorithm

In general, the sample complexity of R-max and Delayed $Q$-learning is incomparable [11]; the former is better in terms of the accuracy of the resulting policy while the latter is better in terms of scaling with the size of the state space. The sample complexity of R-max algorithm is $\frac{|S|^2|A|}{\epsilon^3(1-\gamma)^8}$ —note the power on $\epsilon$; the sample complexity of Delayed $Q$-learning algorithm is $\frac{|S||A|}{\epsilon^4(1-\gamma)^8}$ —note the linear scaling with $|S|$. It appears that DDQ can bring together the best of both worlds; its sample complexity is

$$\mathcal{O}\left(\min\left\{\mathcal{O}\left(\frac{|S|^2|A|}{\epsilon^3(1-\gamma)^8}\right), \mathcal{O}\left(\frac{|S||A|}{\epsilon^4(1-\gamma)^8}\right)\right\}\right)$$

Before formally stating the PAC properties of the DDQ algorithm and proving the bound on its sample complexity, some technical groundwork needs to be laid. To slightly simplify notation, let $\kappa \triangleq |S||A|(1 + \frac{1}{(1-\gamma)\epsilon_1})$. Moreover, subscript $t$ marks the value of a variable at the beginning of timestep $t$ (particularly line 23 of the algorithm).

**Definition 3.** *An event when* learn$(s, a)$ = true *and at the same time* $l(s, a) = m_1$ *or* $n(s, a) = m_2$*, is called an* attempted update.

**Definition 4.** *At any timestep* $t$ *in the execution of* DDQ *algorithm the set of* known state-action pairs *is defined as:*

$$K_t = \left\{(s, a) \mid n(s, a) \geq m_2 \quad \text{or} \quad Q_t(s, a) - \left(R(s, a) + \gamma\sum_{s'} T(s, a, s')v_t(s')\right) \leq 3\epsilon_1\right\}$$

In subsequent analysis, and to distinguish between the conditions that make a state-action pair $(s, a)$ known, the set $K_t$ will be partitioned into two subsets:

$$K_t^1 = \left\{(s, a) \mid Q_t(s, a) - \left(R(s, a) + \gamma\sum_{s'} T(s, a, s')v_t(s')\right) \leq 3\epsilon_1\right\}$$
$$K_t^2 = \left\{(s, a) \mid n(s, a) \geq m_2\right\}$$

**Definition 5.** *In the execution of* DDQ *algorithm a timestep* $t$ *is called a* successful timestep *if at that step any state-action value is updated or the number of times that a state-action pair is visited reaches* $m_2$*. Moreover, considering a particular state-action pair* $(s, a)$*, timestep* $t$ *is called a* successful timestep for $(s, a)$ *if at* $t$ *either update type-1 happens to* $Q(s, a)$ *or the number of times that* $(s, a)$ *is visited reaches* $m_2$*.*

Recall that a type-1 update necessarily decreases the Q-value by at least $\epsilon_1$. Defining rewards as positive quantities prevents the Q-values from becoming negative. At the same time, state-action pairs can initiate update type-2 only once they are experienced $m_2$ times. Together, these conditions facilitate the establishment of an upper-bound on the total number of successful timesteps during the execution of DDQ:

**Lemma 1.** *The number of successful timesteps for a particular state-action pair* $(s, a)$ *in a* DDQ *algorithm is at most* $1 + \frac{1}{(1-\gamma)\epsilon_1}$*. Moreover, the total number of successful timesteps is bounded by* $\kappa$*.*

*Proof.* See Appendix A.

**Lemma 2.** *The total number of attempted updates in* DDQ *algorithm is bounded by* $|S||A|(1 + \kappa)$*.*

*Proof.* See Appendix B.

**Lemma 3.** *Let* $M$ *be an* MDP *with a set of known state-action pairs* $K_t$*. If we assume that for all state-action pairs* $(s, a) \notin K_t$ *we have* $Q_t(s, a) \leq \frac{1}{1-\gamma}$*, then for all state-action pairs in the known state-action* MDP $M_{K_t}$ *it holds*

$$Q^*_{M_{K_t}}(s, a) \leq \frac{1}{1-\gamma}$$

6

*Proof.* See Appendix C.

Choosing $m_1$ big enough and applying Hoefding's inequality allows following conclusion (Lemma 4) for all type-1 updates, and paves the way for establishing the optimism condition of Theorem 1.

**Lemma 4.** *Suppose that at time $t$ during the execution of* DDQ *a state-action pair $(s,a)$ experiences a successful update of type-1 with its value changing from $Q(s,a)$ to $Q'(s,a)$, and that there exists $\exists \epsilon_2 \in (0, \frac{\epsilon_1}{2})$ such that $\forall s \in S$ and $\forall t' < t$, $v_{t'}(s) \geq v_M^*(s) - 2\epsilon_2$. If*

$$m_1 \geq \frac{\ln\left(\frac{8|S||A|(1+\kappa)}{\delta}\right)}{2(\epsilon_1 - 2\epsilon_2)^2(1-\gamma)^2} \simeq \mathcal{O}\left(\frac{\ln\left(\frac{|S|^2|A|^2}{\delta}\right)}{\epsilon_1^2(1-\gamma)^2}\right) \tag{3}$$

*for $\kappa = |S||A|(1 + 1/(1-\gamma)\epsilon_1)$, then $Q'(s,a) \geq Q_M^*(s,a)$ with probability at least $1 - \frac{\delta}{8}$.*

*Proof.* In Appendix D.

The following two lemmas are borrowed from [11] with very minor modifications, and inform on how to choose parameter $m_2$, and the number of iterations for the value iteration part of the DDQ algorithm in order to obtain a desired accuracy.

**Lemma 5.** *(cf. [11, Proposition 4]) Suppose the value-iteration algorithm runs on* MDP $M$ *for $\frac{\ln\left(1/\epsilon_2(1-\gamma)\right)}{1-\gamma}$ iterations, and each state-action value estimate $Q(s,a)$ is initialized to some value between $0$ and $v_{\max}$ for all states and actions. Let $Q'(s,a)$ be the state-action value estimate the algorithm yields. Then $\max_{s,a}\left\{|Q'(s,a) - Q_M^*(s,a)|\right\} \leq \epsilon_2$.*

**Lemma 6.** *Consider an* MDP $M$ *with reward function $R$ and transition probabilities $T$. Suppose another* MDP $\hat{M}$ *has the same state and action set as $M$, but maintains an maximum likelihood (*ML*) estimate of $R$ and $T$, with $n(s,a) \geq m_2$, in the form of $\hat{R}$ and $\hat{T}$ respectively. With $C$ a constant and for all state-action pairs, choosing*

$$m_2 \geq C\left(\frac{|S| + \ln\left(8|S||A|/\delta\right)}{\epsilon_2^2(1-\gamma)^4}\right) \simeq O\left(\frac{|S| + \ln\left(|S||A|/\delta\right)}{\epsilon_2^2(1-\gamma)^4}\right)$$

*guarantees*

$$|R(s,a) - \hat{R}(s,a)| \leq C\epsilon_2(1-\gamma)^2$$
$$\|T(s,a,\cdot) - \hat{T}(s,a,\cdot)\|_1 \leq C\epsilon_2(1-\gamma)^2$$

*with probability at least $1 - \frac{\delta}{8}$. Moreover, for any policy $\pi$ and for all state-action pairs,*

$$|Q_M^\pi(s,a) - Q_{\hat{M}}^\pi(s,a)| \leq \epsilon_2$$
$$|v_M^\pi(s) - v_{\hat{M}}^\pi(s)| \leq \epsilon_2$$

*with probability at least $1 - \frac{\delta}{8}$.*

*Proof.* Combine [11, Lemmas 12–15].

**Lemma 7.** *Let $t_1 < t_2$ be two timesteps during the execution of the* DDQ *algorithm. If*

$$Q_{t_1}(s,a) \geq Q_{M_{\kappa_{t_1}^2}}^*(s,a) - 2\epsilon_2 \quad \forall(s,a) \in S \times A$$

*then with probability at least $1 - \frac{\delta}{8}$*

$$Q_{M_{\kappa_{t_1}^2}}^*(s,a) \geq Q_{M_{\kappa_{t_2}^2}}^*(s,a) \quad \forall(s,a) \in S \times A$$

*Proof.* See Appendix E.

Lemmas 5 and 6 together have as a consequence the following Lemma, which contributes to establishing the accuracy condition of Theorem 1 for the DDQ algorithm.

**Lemma 8.** *During the execution of* DDQ, *for all $t$ and $(s, a) \in S \times A$, we have:*

$$Q^*_{M_{K^2_t}}(s, a) - 2\epsilon_2 \leq Q_t(s, a) \leq Q^*_{M_{K^2_t}}(s, a) + 2\epsilon_2 \tag{4}$$

*with probability at least $1 - \frac{3\delta}{8}$.*

*Proof.* See Appendix F.

Lemma 1 has already offered a bound on the number of updates in DDQ; however, for the complexity condition of Theorem 1 to be satisfied, one needs to show that during the execution of Algorithm 1 the number of escape events is also bounded. The following Lemma is the first step: it states that by picking $m_1$ as in (3), and under specific conditions, an escape event necessarily results in a successful type-1 update. With the number of updates bounded, Lemma 9 can be utilized to derive a bound on the number of escape events.

**Lemma 9.** *With the choice of $m_1$ as in (3), and assuming the DDQ algorithm at timestep $t$ with $(s, a) \notin K_t$, $l(s, a) = 0$ and $\mathrm{learn}(s, a) = \mathrm{true}$, we know that an attempted type-1 update of $Q(s, a)$ will necessarily occur within $m_1$ occurrences of $(s, a)$ after $t$, say at timestep $t_{m_1}$. If $(s, a)$ has been visited fewer than $m_2$ till $t_{m_1}$, then the attempted type-1 update at $t_{m_1}$ will be successful with probability at least $1 - \frac{\delta}{8}$.*

*Proof.* See Appendix G.

**Lemma 10.** *Let $t$ be the timestep when $(s, a)$ has been visited for $m_1$ times after the conditions of Lemma 9 were satisfied. If the update at timestep $t$ is unsuccessful and at timestep $t + 1$ it is $\mathrm{learn}(s, a) = \mathrm{false}$, then $(s, a) \in K_{t+1}$.*

*Proof.* See Appendix H.

A bound on the number the escape events of DDQ algorithm can be derived in a straightforward way. Note that a state-action pair that is visited $m_2$ times becomes a permanent member of set $K_t$. Therefore, the number of escape events is bounded by $|S||A|m_2$. On the other hand, Lemma 9 and the learn flag mechanism (i.e. Lemma 10) suggest another upper bound on escape events. The following Lemma simply states an upper bound for escape events in DDQ as the minimum among the two bounds.

**Lemma 11.** *During the execution of* DDQ, *with the assumption that Lemma 9 holds, the total number of timesteps with $(s_t, a_t) \notin K_t$ (i.e. escape events) is at most $\min\{2m_1\kappa, |S||A|m_2\}$.*

*Proof.* See Appendix I.

Next comes the main result of this paper. The statement that follows establishes the PAC properties of the DDQ algorithm and provides a bound on its sample complexity.

**Theorem 2.** *Consider an* MDP $M = \{S, A, T, R, \gamma\}$, *and let $\epsilon \in (0, \frac{1}{1-\gamma})$, and $\delta \in (0, 1)$. There exist $m_1 = \mathcal{O}\left(\ln(|S|^2|A|^2/\delta)/\epsilon_1^2(1-\gamma)^2\right)$ and $m_2 = \mathcal{O}\left(|S| + \ln(|S||A|/\delta)/\epsilon_2^2(1-\gamma)^4\right)$ with $\frac{1}{\epsilon_1} = \frac{3}{(1-\gamma)\epsilon} = \mathcal{O}\left(1/\epsilon(1-\gamma)\right)$ and $\epsilon_2 = \frac{\epsilon_1}{3}$, such that if DDQ algorithm is executed, $M$ follows a $4\epsilon$-optimal policy with probability at least $1 - 2\delta$ on all but*

$$\mathcal{O}\left(\min\left\{\mathcal{O}(|S|^2|A|/\epsilon^3(1-\gamma)^8), \mathcal{O}(|S||A|/\epsilon^4(1-\gamma)^8)\right\}\right)$$

*timesteps (logarithmic factors ignored).*

*Proof.* We intend to apply Theorem 1. To satisfy the *optimism* condition, we start by proving that $Q_t(s, a) \geq Q^*_M(s, a) - 2\epsilon_2$ by strong induction for all state-action pairs: (i) At $t = 1$, the value of all state-action pairs are set to the maximum possible value in MDP $M$. This implies that $Q_1(s, a) \geq Q^*_M(s, a) \geq Q^*_M(s, a) - 2\epsilon_2$, therefore $v_t(s) \geq v^*_M(s) - 2\epsilon_2$. (ii) Assume that $Q_t(s, a) \geq Q^*_M(s, a) - 2\epsilon_2$ holds for all timesteps before or equal to $t = n - 1$. (iii) At timestep $t = n$, all $(s, a) \notin K^2_n$ can only be updated by a type-1 update before or at $t = n$. For these state-action pairs, Lemma 4 implies that it will be $Q_n(s, a) \geq Q^*_M(s, a)$ with probability $1 - \frac{\delta}{8}$.

For all $(s, a) \in K^2_n$, on the other hand, by Lemma 8 and with probability $1 - \frac{3\delta}{8}$:

$$Q_n(s, a) \geq Q^*_{M_{K^2_n}}(s, a) - 2\epsilon_2 \geq Q^*_M(s, a) - 2\epsilon_2$$

8

Note that $Q^*_{M_{K_n^2}}(s,a) \geq Q^*_M(s,a)$ since $M_{K_n^2}$ is similar to $M$ exept for $(s,a) \notin K_n^2$ which their values are set to be more than or equal to $Q^*_M(s,a)$. Therefore, $Q_t(s,a) \geq Q^*_M(s,a) - 2\epsilon_2$ holds for all timesteps $t$ and all state-action pairs, which directly implies $v_t(s) \geq v^*_M(s) - 2\epsilon_2 \geq v^*_M(s) - \epsilon$.

To establish the *accuracy* condition, first write

$$Q_t(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_t(s',a') + \beta(s,a) \tag{5}$$

If $(s,a) \in K_t$, there can be two cases: either $(s,a) \in K_t^1$ or $(s,a) \in K_t^2$. If $(s,a) \in K_t^1$, then by Definition 4 $\beta(s,a) \leq 3\epsilon_1$. If $(s,a) \in K_t^2$, then Lemma 8 (right-hand side inequality) implies that with probability at least $1 - \frac{3\delta}{8}$

$$2\epsilon_2 \geq Q_t(s,a) - Q^*_{M_{K_t^2}}(s,a) \tag{6}$$

Meanwhile,

$$Q^*_{M_{K_t^2}}(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q^*_{M_{K_t^2}}(s',a') \tag{7}$$

and substituting from (7) and (5) into (6) yields

$$\gamma \sum_{s'} T(s,a,s') \left( \max_{a'} Q_t(s',a') - \max_{a'} Q^*_{M_{K_t^2}}(s',a') \right) + \beta(s,a) \leq 2\epsilon_2 \tag{8}$$

Let $a_1 := \arg\max_{a'} Q_{M_{K_t^2}}(s',a')$ and bound the difference

$$\max_{a'} Q_t(s',a') - \max_{a'} Q^*_{M_{K_t^2}}(s',a') = \max_{a'} Q_t(s',a') - Q^*_{M_{K_t^2}}(s',a_1)$$
$$\geq Q_t(s',a_1) - Q^*_{M_{K_t^2}}(s',a_1)$$

Apply Lemma 8 (left-hand side inequality) to the latter expression to get

$$\max_{a'} Q_t(s',a') - \max_{a'} Q^*_{M_{K_t^2}}(s',a') \geq -2\epsilon_2$$

which implies for (8) that

$$2\epsilon_2 \geq \beta(s,a) - 2\gamma\epsilon_2 \implies \beta(s,a) \leq 2(1+\gamma)\epsilon_2 \leq 3\epsilon_2$$

Thus in any case when $(s,a) \in K_t$, $\beta(s,a) \leq 3\epsilon_1$ with probability at least $1 - \frac{3\delta}{8}$. In light of this, considering a policy dictating actions $a = \pi_t(s)$ and mirroring (5)–(7) we write for the values of states in which $(s, \pi_t(s)) \in K_t$

$$v^{\pi_t}_{M_{K_t}}(s) = R\big(s, \pi_t(s)\big) + \gamma \sum_{s'} T\big(s, \pi_t(s), s'\big) v^{\pi_t}_{M_{K_t}}(s')$$
$$v_t(s) = R\big(s, \pi_t(s)\big) + \gamma \sum_{s'} T\big(s, \pi_t(s), s'\big) v_t(s') + \beta(s,a)$$

while for those in which $(s, \pi_t(s)) \notin K_t$, we already know that

$$v^{\pi_t}_{M_{K_t}}(s) = Q_t\big(s, \pi_t(s)\big)$$
$$v_t(s) = Q_t\big(s, \pi_t(s)\big)$$

So now if one denotes

$$\alpha := \max_s \big( v_t(s) - v^{\pi_t}_{M_{K_t}}(s) \big) = v_t(s^*) - v^{\pi_t}_{M_{K_t}}(s^*)$$

then either $\alpha = 0$ (when $(s, \pi_t(s)) \notin K_t$) or it affords an upper bound

$$\gamma \sum_{s'} T\big(s^*, \pi_t(s^*), s'\big) \big( v_t(s') - v^{\pi_t}_{M_{K_t}}(s') \big) + \beta\big(s^*, \pi_t(s^*)\big)$$
$$\leq \gamma \sum_{s'} T\big(s^*, \pi_t(s^*), s'\big) \big( v_t(s') - v^{\pi_t}_{M_{K_t}}(s') \big) + 3\epsilon_1 \leq \gamma\alpha + 3\epsilon_1$$

from which it follows that $\alpha \leq \gamma\alpha + 3\epsilon_1 \implies \alpha \leq \frac{3\epsilon}{1-\gamma} = \epsilon$.

9

Finally, to analyze *complexity* invoke Lemmas 1 and 11 to see that the learning complexity $\zeta(\epsilon, \delta)$ is bounded by $\kappa + \min\left(2m_1\kappa, |S||A|m_2\right)$ with probability $1 - \frac{\delta}{8}$.

In conclusion, the conditions of Theorem 2 are satisfied with probability $1 - \delta$ and therefore the DDQ algorithm is PAC. Substituting $\zeta(\epsilon, \delta)$ into (2) completes the proof.

## 5 Numerical Results

This section opens with a comparison of the DDQ algorithm to its parent technologies. It proceeds with additional comparisons to the state-of-the-art in both model-based [5] as well as model-free [9] RL algorithms. For this comparison, the algorithms with the currently best sample complexity are implemented on a type of MDP which has been proposed and used in literature as a model which is objectively difficult to learn [11].

The first round of comparisons start with R-max, Delayed Q-learning and DDQ being implemented on a small-scale grid-world example (Fig. 1). This example test case has nine states, with the initial state being the one labeled 1, and the terminal (goal) state labeled 9. Each state is assigned a reward of 0 except for the terminal state which has 1. For this example, $\gamma := 0.8$. In all states but the terminal one, the system has four primitive actions available: down (d), left (l), up (u), and right (r). The grid-world of Fig. 1 includes cells with two types of boundaries: the boundaries marked with a single-line afford transition probabilities of 0.9 through them; the boundaries marked with a double line afford transitions through them at probability 0.1. The optimal policy for this grid-world example is shown in Fig. 2.



Figure 1: The grid-world example.



Figure 2: The actual optimal policy in the grid-world example.

Initializing the three PAC algorithms with parameters $m_1 = 65$, $m_2 = 175$ and $\varepsilon = 0.06$, yields the performance metrics shown in Table 1, which are measured in terms of the number of samples needed to reach at $4\varepsilon$ optimality, averaged over 10 algorithm runs. Parameters $m_1$ and $m_2$ are intentionally chosen to enable a fair comparison, in the sense that the sample complexity of the model-free Delayed Q-learning, and the model-based R-max algorithms is almost identical. In this case, and with these same tuning parameters, DDQ yields a modest but notable sample complexity improvement.

The lowest known bound on the sample complexity of a model-based RL algorithm on a infinite-horizon MDP is $\frac{|S||A|}{\epsilon^2(1-\gamma)^6}$ (by the Mormax algorithm [5]). For the model-free case (again on a infinite-horizon MDP), the lowest bound on the sample complexity is $\frac{|S||A|}{\epsilon^2(1-\gamma)^7}$, achieved by UCB Q-learning [9] (the extended version of [8] which is for finite-horizon MDP).

Table 1: Average # of samples for reaching $4\varepsilon$ optimality

| Algorithms | # of samples |
|:---:|:---:|
| Delayed Q-learning | 6622 |
| R-max | 6727 |
| DDQ | 5960 |

To perform a fair and meaningful comparison of these algorithms to DDQ, consider a family of "difficult-to-learn" MDP as Fig. 3. The MDP has $N + 2$ states as $S = \{1, 2, \ldots, N, +, -\}$, and $A$ different actions. Transitions from each state $i \in \{1, \ldots, N\}$ are the same, so only transitions from state 1 are shown. One of the actions (marked by solid line) deterministically transports the agent to state + with reward $0.5 + \epsilon'$ (with $\epsilon' > 0$). Let $a$ be any of the other $A - 1$ actions (represented by dashed lines). From any state $i \in \{1, \ldots, N\}$, taking action $a$ will trigger a transition to state + with reward 1 and probability $p_{ia}$, or to state − with reward 0 and probability $1 - p_{ia}$, where $p_{ia} \in \{0.5, 0.5 + 2\epsilon'\}$ are numbers very close to $0.5 + \epsilon'$. For each state $i \in \{1, \ldots, N\}$, there is at most one $a$ such that $p_{ia} = 0.5 + 2\epsilon'$. Transitions from states + and − are identical; they simply reset the agent to one of the states $\{1, \ldots, N\}$ uniformly at random.
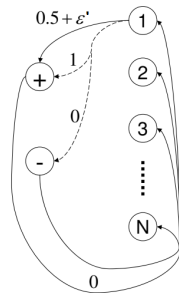


Figure 3: A family of difficult-to-learn MDPs. [11]

For an MDP such as the one shown in Fig. 3, the optimal action in any state $i \in \{1, \ldots, N\}$ is independent of the other states; specifically, it is the action marked by the solid arrow if $p_{ia} = 0.5$ for all dashed actions $a$, or the action marked by the dashed arrow for which $p_{ia} = 0.5 + 2\epsilon'$, otherwise. Intuitively, this MDP is hard to learn for exactly the same reason that a biased coin is hard to be recognized as such if its bias (say, the probability of landing on head) is close to $0.5$ [11].

We thus try to learn such an MDP $M$ with $N = 2$, $A = 2$, and $\epsilon' = 0.04$. The accuracy that the learned policy should satisfy is set to $\epsilon = 0.0025$, and the probability of failure is set to $\delta = 0.01$. Results are averaged over 50 runs of each algorithm running on MDP $M$.

We empirically fine-tune the parameters of Mormax and UCB Q-learning algorithms to maximize their performance on learning the near optimal ($4\epsilon$-optimal) policy of $M$ in terms of the required samples. As expected, the required samples decrease (almost linearly) in $m$ (Fig. 4) until the necessary condition for the convergence of the algorithm is violated (at around $m = 600$). For that reason, we cap $m$ at 600 which requires 7770 samples on average for Mormax to learn the optimal policy. Yet another important performance metric to record for a model-based RL algorithm is the number of times it needs resolve the learned model through value-iteration, since the associated computational effort is highly dependent on this number. For Mormax, the average number of times it performs model resolution is 12.06.

The performance of the UCB Q-learning algorithm appears to be very sensitive to its $c_2$ parameter. The value of $4\sqrt{2}$ that has been suggested for $c_2$ [9] proved very conservative, with the algorithm sometimes requiring millions of data for converging to the optimal policy on $M$. The reason is that values of $c_2$ that high cause the effective updates to start when the learning rate has already become very small, thus slowing down the convergence speed. We therefore tune the UCB Q-learning algorithm to achieve maximum performance on $M$ by setting its parameter $c_2 = 1/50$ (see Fig. 5); with this setting, the algorithm requires 8097 samples to learn the optimal policy on average. Setting $c_2 < 1/50$ may cause the algorithm to lie outside the upper confidence interval, and as a result, the algorithm either requires an actual higher number of samples or it fails to convege altogether to the optimal policy after $10^6$ samples.
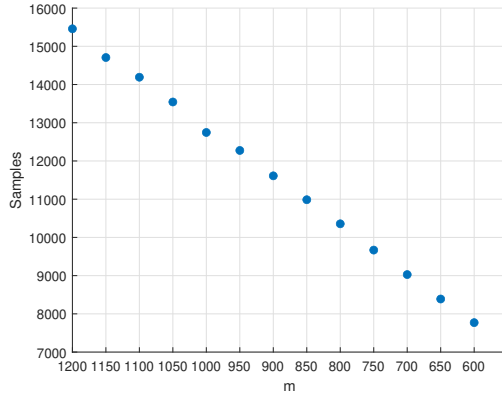
Figure 4: The number of samples required by the Mormax algorithm.
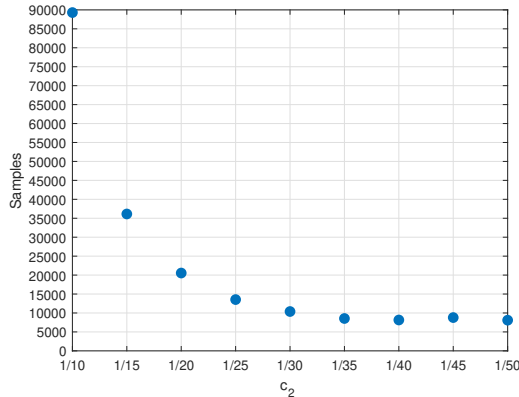


Figure 5: The required samples by UCB Q-learning algorithm

We compare the best performance we could achieve with Mormax and UCB Q-learning with that of DDQ which we tune with $m_1 = 150$ and $m_2 = 750$. The average required samples required by DDQ for learning the $4\epsilon$-optimal policy on $M$ is 5662, while the number of times that the R-max component of the algorithm resolves the model through value-iteration part is 3.76 on average.

Thus, although the provable worst-case bound on the sample complexity of DDQ algorithm appears higher than that of Mormax and UCB Q-learning, DDQ can outperform both algorithms in terms of the required data samples, especially in difficult learning tasks. What is more, the hybrid nature of DDQ algorithm enables significant savings in terms of computational effort —the latter captured by the number of times when the algorithm resorts to model resolution— compared to model-based algorithms like Mormax. Table 2 summarizes the results of this comparison.

Table 2: The best possible performance on learning MDP $M$

| Algorithms | # of samples | # of model resolution |
|---|---|---|
| Mormax | 7770 | 12.06 |
| UCB Q-learning | 8097 | 0 (model-free) |
| DDQ | 5662 | 3.76 |

# 6 Conclusion

It is possible to build an RL algorithm that captures favorable features of both model-based and model-free learning and most importantly preserves the PAC property. One such algorithm is the DDQ; this algorithm leverages the idea of Dyna-Q to combine two existing PAC algorithms, namely the model-based R-max and the model-free Delayed Q-learning, in a way that achieves the best (complexity results) of both. Theoretical analysis establishes that DDQ enjoys a sample complexity that is at worst as high as the smallest of its constituent technologies; yet, in practice, as the numerical example included suggests, DDQ can outperform them both. In addition, numerical example on the comparison of DDQ to the state of the art in model-based and model free RL suggests clear advantages in practical implementations.

# References

[1] Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232, 2002.

[2] Ronen I Brafman and Moshe Tennenholtz. R-max a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.

[3] Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.

[4] Alexander L Strehl, Lihong Li, and Michael L Littman. Incremental model-based learners with formal learning-time guarantees. *arXiv preprint arXiv:1206.6870*, 2012.

[5] István Szita and Csaba Szepesvári. Model-based reinforcement learning with nearly tight exploration complexity bounds. In *International Conference on Machine Learning*, 2010.

[6] Tor Lattimore and Marcus Hutter. Near-optimal pac bounds for discounted mdps. *Theoretical Computer Science*, 558:125–143, 2014.

[7] 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*, pages 881–888. ACM, 2006.

[8] Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.

[9] Kefan Dong, Yuanhao Wang, Xiaoyu Chen, and Liwei Wang. Q-learning with UCB exploration is sample efficient for infinite-horizon MDP. *arXiv preprint arXiv:1901.09311*, 2019.

[10] Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model-free deep RL for model-based control. *arXiv preprint arXiv:1802.09081*, 2018.

[11] Alexander L Strehl, Lihong Li, and Michael L Littman. Reinforcement learning in finite MDPs: PAC analysis. *Journal of Machine Learning Research*, 10(Nov):2413–2444, 2009.

[12] Anusha Nagabandi, Gregory Kahn, Ronald S Fearing, and Sergey Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *2018 IEEE International Conference on Robotics and Automation*, pages 7559–7566, 2018.

[13] Sham Machandranath Kakade et al. *On the sample complexity of reinforcement learning*. PhD thesis, University of London London, England, 2003.

[14] Sang Wan Lee, Shinsuke Shimojo, and John P O'Doherty. Neural computations underlying arbitration between model-based and model-free learning. *Neuron*, 81(3):687–699, 2014.

[15] Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.

[16] Nicolas Heess, Gregory Wayne, David Silver, Timothy Lillicrap, Tom Erez, and Yuval Tassa. Learning continuous control policies by stochastic value gradients. In *Advances in Neural Information Processing Systems*, pages 2944–2952, 2015.

[17] Ronald Parr, Lihong Li, Gavin Taylor, Christopher Painter-Wakefield, and Michael L Littman. An analysis of linear models, linear value-function approximation, and feature selection for reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pages 752–759. ACM, 2008.

[18] Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, Adam White, and Doina Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pages 761–768. International Foundation for Autonomous Agents and Multiagent Systems, 2011.

[19] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In *International Conference on Machine Learning*, pages 1312–1320, 2015.

[20] 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 *Advances in Neural Information Processing Systems*, pages 5048–5058, 2017.

[21] Yevgen Chebotar, Karol Hausman, Marvin Zhang, Gaurav Sukhatme, Stefan Schaal, and Sergey Levine. Combining model-based and model-free updates for trajectory-centric reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 703–711. JMLR. org, 2017.

[22] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.

[23] A. Zehfroosh, E. Kokkoni, H. G. Tanner, and J. Heinz. Learning models of human-robot interaction from small data. In *2017 25th IEEE Mediterranean Conference on Control and Automation*, pages 223–228, July 2017.

[24] A. Zehfroosh, H. G. Tanner, and J. Heinz. Learning option mdps from small data. In *2018 IEEE American Control Conference*, pages 252–257, 2018.

# Appendices

## A  Proof of Lemma 1

Consider a fixed state-action pair $(s, a)$. Its value $Q(s, a)$ is initially set to $v_{\max} = \frac{1}{1-\gamma}$. When an update of type-1 (Algorithm 1 line 30) is successful $Q(s, a)$ is reduced by at least $\epsilon_1$. Since the reward function $R(s, a)$ is non-negative, we must have $Q(s, a) \geq 0$ in all timesteps, which means that there can be at most $\frac{1}{\epsilon_1(1-\gamma)}$ updates of type-1 for $(s, a)$. On the other hand, a type-2 update (Algorithm 1 line 51) can occur only once when $n(s, a) = m_2$. Therefore, the total number of successful timesteps for $(s, a)$ is at most $1 + \frac{1}{\epsilon_1(1-\gamma)}$ times. With $|S||A|$ total state-action pairs, the total number of successful timesteps is bounded by $\kappa = |S||A| + \frac{|S||A|}{(1-\gamma)\epsilon_1}$.

## B  Proof of Lemma 2

Suppose an attempted update occurs at timestep $t$ to some $(s, a)$. By definition, for a subsequent attempted update to $(s, a)$ to occur at timestep $t' > t$, at least one successful timestep must occur between $t$ and $t'$. Lemma 1 ensures that there can be no more than $\kappa$ successful timesteps. In other words, the most frequent occurrence of attempted updates is interlaced between successful updates, which implies that at most $1 + \kappa$ attempted updates are possible for $(s, a)$. Scaling this argument to all state-action pairs we arrive at the $|S||A|(1 + \kappa)$ upper bound.

## C  Proof of Lemma 3

Let $Q^*_{M_{K_t}}(s^*, a^*)$ denote $\max_{(s,a)} Q^*_{M_{K_t}}(s, a)$. If $(s^*, a^*) \notin K_t$, we are done since $Q^*_{M_{K_t}}(s^*, a^*) = Q_t(s^*, a^*) \leq \frac{1}{1-\gamma}$. Otherwise, for $(s^*, a^*) \in K_t$ write

$$Q^*_{M_{K_t}}(s^*, a^*) = R(s^*, a^*) + \gamma \sum_s T(s^*, a^*, s) \max_a Q^*_{M_{K_t}}(s, a)$$
$$\leq R(s^*, a^*) + \gamma Q^*_{M_{K_t}}(s^*, a^*) \sum_s T(s^*, a^*, s)$$
$$= R(s^*, a^*) + \gamma Q^*_{M_{K_t}}(s^*, a^*)$$
$$\leq 1 + \gamma Q^*_{M_{K_t}}(s^*, a^*) \implies Q^*_{M_{K_t}}(s^*, a^*) \leq \frac{1}{1-\gamma}$$

14

## D  Proof of Lemma 4

Let an update of type-1 occur for $(s,a)$ at timestep $t$. Suppose that the latest $m_1$ experiences of $(s,a)$ happened at timesteps $t_1 < t_2 < \cdots < t_{m_1} = t$, when the system was rewarded $r[1], r[2], \ldots, r[m_1]$ and jumped to states $s[1], s[2], \ldots, s[m_1]$, respectively. Define the random variable $Y = r[i] + \gamma v_M^*(s[i])$ for $1 \le i \le m_1$ and note that $0 \le Y \le \frac{1}{1-\gamma}$. Then a direct application of the Hoeffding inequality for bounded random variables and with the choice of $m_1$ as in (3) implies that

$$\frac{1}{m_1} \sum_{i=1}^{m_1} \big( r[i] + \gamma v_M^*(s[i]) \big) > \mathbb{E}\{Y\} - (\epsilon_1 - 2\epsilon_2) = Q_M^*(s,a) - \epsilon_1 + 2\epsilon_2$$

with probability $1 - \delta/8(|S||A|(1+\kappa))$.

Now we have:

$$
\begin{aligned}
Q'(s,a) &= \frac{1}{m_1}\Big( \sum_{i=1}^{m_1} r[i] + \gamma v_{t_i}(s[i]) \Big) + \epsilon_1 \\
&\ge \frac{1}{m_1}\Big( \sum_{i=1}^{m_1} r[i] + \gamma v_M^*(s[i]) \Big) - 2\gamma\epsilon_2 + \epsilon_1 \\
&\ge Q_M^*(s,a) - \epsilon_1 + 2\epsilon_2 - 2\gamma\epsilon_2 + \epsilon_1 \ge Q_M^*(s,a)
\end{aligned}
$$

Finally, we want this fact to be true for all possible attempted updates of type-1. According to Lemma 2, an upper bound for all possible attempted updates is $|S||A|(1+\kappa)$. Therefore, the above fact is true with probability at least $\big(1 - \delta/8(|S||A|(1+\kappa))\big)^{|S||A|(1+\kappa)}$. An induction argument can now be employed to show that $1 - \frac{\delta}{8}$ bounds the above expression from below.

## E  Proof of Lemma 7

First note that $K_{t_1}^2 \subseteq K_{t_2}^2$. For all $(s,a) \notin K_{t_2}^2$

$$Q_{M_{K_{t_1}^2}}^*(s,a) = Q_{t_1}(s,a) \ge Q_{t_2}(s,a) = Q_{M_{K_{t_2}^2}}^*(s,a) \tag{9}$$

while for all $(s,a) \in K_{t_1}^2$

$$Q_{M_{K_{t_1}^2}}^*(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{M_{K_{t_1}^2}}^*(s',a')$$

$$Q_{M_{K_{t_2}^2}}^*(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{M_{K_{t_2}^2}}^*(s',a')$$

implying

$$Q_{M_{K_{t_1}^2}}^*(s,a) - Q_{M_{K_{t_2}^2}}^*(s,a) = \gamma \sum_{s'} T(s,a,s') \times \Big( \max_{a'} Q_{M_{K_{t_1}^2}}^*(s',a') - \max_{a'} Q_{M_{K_{t_2}^2}}^*(s',a') \Big) \tag{10}$$

Every $(s,a) \in K_{t_2}^2 \setminus K_{t_1}^2$ falls in one of the following categories:

- $(s,a)$ is a state-action pair that has not been updated ever before or at timestep $t_1$. The Lemma 3 implies

$$Q_{M_{K_{t_1}^2}}^*(s,a) = Q_{t_1}(s,a) = v_{\max} = \frac{1}{1-\gamma} \ge Q_{M_{K_{t_2}^2}}^*(s,a)$$

  which completes the proof.

- $(s,a)$ is a state-action pair that has experienced an type-1 update before or at $t_1$. Assume that the most recent type-1 update of $(s,a)$ occurred at some timestep $t \le t_1$. Suppose that the $m_1$ visits to $(s,a)$ that triggered this update occurred at instances $t^1 < t^2 < \ldots < t^{m_1} = t \le t_1$, and the observed rewards and next states were $r[1], r[2], \ldots, r[m_1]$ and $s[1], s[2], \ldots, s[m_1]$, respectively. For the random variable $Z = r[i] + \gamma v_t(s[i])$,

$$\mathbb{E}\{Z\} = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_t(s',a')$$

Then

$$Q^*_{M_{K^2_{t_1}}}(s,a) = Q_{t_1}(s,a) = Q_t(s,a) = \frac{\sum_{i=1}^{m_1} r[i] + \gamma v_{t^i}(s[i])}{m_1} + \epsilon_1 \geq \frac{\sum_{i=1}^{m_1} r[i] + \gamma v_t(s[i])}{m_1} + \epsilon_1$$

and applying Hoeffding inequality to the right hand side

$$Q^*_{M_{K^2_{t_1}}}(s,a) > \mathbb{E}\{Z\} - \epsilon_1 + 2\epsilon_2 + \epsilon_1 = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_t(s',a') + 2\epsilon_2$$

$$\geq R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{t_1}(s',a') + 2\epsilon_2$$

$$\overset{(9)}{\geq} R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q^*_{M_{K^2_{t_1}}}(s',a')$$

with probability $1 - \frac{\delta}{8|S||A|(1+\kappa)}$. Then — following the final steps of Lemma 4 — with probability at least $1 - \frac{\delta}{8}$ after all possible attempted updates,

$$Q^*_{M_{K^2_{t_1}}}(s,a) - Q^*_{M_{K^2_{t_2}}}(s,a) \geq \gamma \sum_{s'} T(s,a,s') \Big( \max_{a'} Q^*_{M_{K^2_{t_1}}}(s',a') - \max_{a'} Q^*_{M_{K^2_{t_2}}}(s',a') \Big) \qquad (11)$$

In any case, therefore, i.e., either when $(s,a) \notin K^2_{t_2}$ or when $(s,a) \in K^2_{t_2} \smallsetminus K^2_{t_1}$, one can define

$$\alpha := \min_{(s,a)} \Big( Q^*_{M_{K^2_{t_1}}}(s,a) - Q^*_{M_{K^2_{t_2}}}(s,a) \Big) := Q^*_{M_{K^2_{t_1}}}(s^*,a^*) - Q^*_{M_{K^2_{t_2}}}(s^*,a^*)$$

and if $\alpha \geq 0$ recognize that the proof is completed. Assume for the sake of argument that $\alpha < 0$; then still either (10) is true if $(s,a) \notin K^2_{t_2}$, or (11) if $(s,a) \in K^2_{t_2} \smallsetminus K^2_{t_1}$. Let $a_{s'} := \arg\max_{a'} Q^*_{M_{K^2_{t_2}}}(s',a')$, then in either case,

$$\alpha = Q^*_{M_{K^2_{t_1}}}(s^*,a^*) - Q^*_{M_{K^2_{t_2}}}(s^*,a^*)$$

$$\geq \gamma \sum_{s'} T(s^*,a^*,s') \Big( \max_{a'} Q^*_{M_{K^2_{t_1}}}(s',a') - \max_{a'} Q^*_{M_{K^2_{t_2}}}(s',a') \Big)$$

$$= \gamma \sum_{s'} T(s^*,a^*,s') \Big( \max_{a'} Q^*_{M_{K^2_{t_1}}}(s',a') - Q^*_{M_{K^2_{t_2}}}(s',a_{s'}) \Big)$$

$$\geq \gamma \sum_{s'} T(s^*,a^*,s') \Big( Q^*_{M_{K^2_{t_1}}}(s',a_{s'}) - Q^*_{M_{K^2_{t_2}}}(s',a_{s'}) \Big)$$

$$\geq \gamma \alpha \implies \alpha \geq 0$$

which is a contradiction. Therefore $\alpha$ cannot be negative and therefore $Q^*_{M_{K^2_{t_1}}}(s,a) - Q^*_{M_{K^2_{t_2}}}(s,a) \geq 0$.

# F   Proof of Lemma 8

For all $(s,a) \notin K^2_t$

$$Q_t(s,a) = Q^*_{M_{K^2_t}}(s,a) \leq Q^*_{M_{K^2_t}}(s,a) + 2\epsilon_2 \qquad (12a)$$

Now for $(s,a) \in K^2_t$, and referring to line 50 of Algorithm 1 one sees that for timestep $t$ it is $Q_t(s,a) \leq Q_{\text{vl}}(s,a)$. Meanwhile, for timestep $t$ Lemma 5 ensures

$$Q_{\text{vl}}(s,a) \leq Q^*_{\hat{M}_{K^2_t}}(s,a) + \epsilon_2 \qquad (12b)$$

while Lemma 6 implies

$$Q^*_{\hat{M}_{K^2_t}}(s,a) + \epsilon_2 \leq Q^*_{M_{K^2_t}}(s,a) + 2\epsilon_2 \qquad (12c)$$

with probability $1 - \frac{\delta}{8}$. Combining (12) one obtains the right hand side of (4). Establishing the left hand side of (4) is done by strong induction. At $t = 1$, we have $K^2_1 = \varnothing$ and thus

$$Q_1(s,a) = Q^*_{M_{K^2_1}}(s,a) \geq Q^*_{M_{K^2_1}}(s,a) - 2\epsilon_2$$

16

Assume that $Q_t(s,a) = Q^*_{M_{K_t^2}}(s,a) \geq Q^*_{M_{K_t^2}}(s,a) - 2\epsilon_2$ for $t \leq n-1$. If timestep $t = n$ is not a successful timestep (Definition 5), nothing happens so equality holds; thus let us assume that $t = n$ is successful. Then, and for all $(s,a) \notin K_n^2$ we have automatically

$$Q_n(s,a) = Q^*_{M_{K_n^2}}(s,a) \geq Q^*_{M_{K_n^2}}(s,a) - 2\epsilon_2$$

Just as before, for $(s,a) \in K_n^2$ for which a type-2 update succeeded at timestep $t$

$$Q_n(s,a) = Q_{vl}(s,a) \geq Q^*_{\hat{M}_{K_n^2}}(s,a) - \epsilon_2 \tag{13a}$$

as a result of Lemma 5, and

$$Q^*_{\hat{M}_{K_n^2}}(s,a) - \epsilon_2 \geq Q^*_{M_{K_n^2}}(s,a) - 2\epsilon_2 \tag{13b}$$

with probability $1 - \frac{\delta}{8}$, due to Lemma 6. For those $(s,a) \in K_n^2$ for which a type-2 update did *not* succeed at timestep $t$, it is $Q_n(s,a) = Q_{n-1}(s,a)$ and there are three distinct possibilities:

- Value $Q_{n-1}(s,a)$ has never been updated before. Then,

$$Q_n(s,a) = \frac{1}{1-\gamma} \overset{\text{Lemma 3}}{\geq} Q^*_{M_{K_n^2}}(s,a) \geq Q^*_{M_{K_n^2}}(s,a) - 2\epsilon_2$$

- The most recent update for $(s,a)$ was of type-2 and occured at some $t \leq n-1$. Then,

$$Q_n(s,a) \overset{\text{Lemmas 5\&6}}{\geq} Q^*_{M_{K_t^2}}(s,a) - 2\epsilon_2$$

with probability $1 - \frac{\delta}{8}$, and

$$Q^*_{M_{K_t^2}}(s,a) - 2\epsilon_2 \overset{\text{Lemma 7}}{\geq} Q^*_{M_{K_n^2}}(s,a) - 2\epsilon_2$$

also with with probability $1 - \frac{\delta}{8}$, so

$$Q_n(s,a) \geq Q^*_{M_{K_n^2}}(s,a) - 2\epsilon_2$$

with probability at least $1 - \frac{2\delta}{8} \leq (1 - \frac{\delta}{8})^2$.

- The most recent update for $(s,a)$ was of type-1 and occured at some $t' \leq n-1$. Then suppose that the $m_1$ collection of visits of $(s,a)$ for this update occurred at timesteps $t^1 < t^2 < \cdots < t^{m_1} = t' \leq n-1$, with the corresponding observed reward and next states being $r[1], r[2], \ldots, r[m_1]$ and $s[1], s[2], \ldots, s[m_1]$, respectively. The expectation of the random variable $F = r[i] + \gamma v_{t^{m_1}}(s[i])$ is

$$\mathbb{E}\{F\} = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{t^{m_1}}(s',a')$$

which, with the use of Hoeffding inequality, bounds the sum in

$$Q_n(s,a) = \frac{1}{m_1} \left( \sum_{i=1}^{m_1} r[i] + \gamma v_{t^i}(s[i]) \right) + \epsilon_1$$

$$\geq \frac{1}{m_1} \left( \sum_{i=1}^{m_1} r[i] + \gamma v_{t^{m_1}}(s[i]) \right) + \epsilon_1 > \mathbb{E}\{F\} - \epsilon_1 + 2\epsilon_2 + \epsilon_1$$

$$= R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{t'}(s',a') + 2\epsilon_2$$

and yields

$$Q_n(s,a) \geq R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_n(s',a')$$

with probability $1 - \delta/8(|S||A|(1+\kappa))$. Following the steps in the proof of Lemma 4 when thinking of all possible attempted updates, one states the above with probability $1 - \frac{\delta}{8}$. Subtracting now $Q^*_{M_{K_n^2}}(s,a)$ from both sides yields

$$\gamma \sum_{s'} T(s,a,s') \left( \max_{a'} Q_n(s',a') - \max_{a'} Q^*_{M_{K_n^2}}(s',a') \right) \leq Q_n(s,a) - Q^*_{M_{K_n^2}}(s,a) \tag{14}$$

17

and if one denotes
$$\alpha := \min_{(s,a)} \left( Q_n(s,a) - Q^*_{M_{K^2_n}}(s,a) \right) = Q_n(s^*,a^*) - Q^*_{M_{K^2_n}}(s^*,a^*)$$

then we want to show $\alpha \geq -2\epsilon_2$. Let $a_{s'} := \arg\max_{a'} Q^*_{M_{K^2_n}}(s',a')$, then (14) implies

$$\begin{aligned}
\alpha &= Q_n(s^*,a^*) - Q^*_{M_{K^2_n}}(s^*,a^*) \\
&\geq \gamma \sum_{s'} T(s^*,a^*,s') \left( \max_{a'} Q_n(s',a') - \max_{a'} Q^*_{M_{K^2_n}}(s',a') \right) \\
&= \gamma \sum_{s'} T(s^*,a^*,s') \left( \max_{a'} Q_n(s',a') - Q^*_{M_{K^2_n}}(s',a_{s'}) \right) \\
&\geq \gamma \sum_{s'} T(s^*,a^*,s') \left( Q_n(s',a_{s'}) - Q^*_{M_{K^2_n}}(s',a_{s'}) \right) \\
&\geq \gamma\alpha \implies \alpha \geq 0 \geq -2\epsilon_2
\end{aligned}$$

Summing up, the right side of (4) holds with probability $1 - \frac{\delta}{8}$, while the left side is true with probability at least $(1 - \frac{\delta}{12})^2$. Together, both inequalities are true with probability at least $(1 - \frac{\delta}{12})^3 \geq 1 - \frac{3\delta}{8}$.

# G   Proof of Lemma 9

Assume that at timestep $t$, $(s,a) \notin K_t$, $l(s,a) = 0$ and $\mathsf{learn}(s,a) = \mathrm{true}$, and suppose that $m_1$ experiences of $(s,a)$ after $t$ happen at timesteps $t \leq t_1 < t_2 < \cdots < t_{m_1}$. Let $r[1], r[2], \ldots, r[m_1]$ and $s[1], s[2], \ldots, s[m_1]$ be the rewards and next states observed for the $m_1$ experiences of $(s,a)$. Then define the random variable $X = r[i] + \gamma v_{t_1}(s[i])$ letting $i$ range in $\{1, \ldots, m_1\}$, and note that $0 \leq X \leq 1$.

A direct application of the Hoeffding inequality with the choice of $m_1$ as in (3) yields

$$\frac{1}{m_1} \left( \sum_{i=1}^{m_1} r[i] + \gamma v_{t_1}(s[i]) \right) - \mathbb{E}\{X\} < \epsilon_1 - 2\epsilon_2 < \epsilon_1$$

with probability $1 - \frac{\delta}{8|S||A|(1+\kappa)}$. Since the DDQ algorithm only allows for updates that decrease the value estimate for any stat-action pairs, we can write:

$$\begin{aligned}
Q_t(s,a) - \frac{1}{m_1} \left( \sum_{i=1}^{m_1} r[i] + \gamma v_{t_i}(s[i]) \right) &\geq Q_t(s,a) - \frac{1}{m_1} \left( \sum_{i=1}^{m_1} r[i] + \gamma v_{t_1}(s[i]) \right) \\
&> Q_t(s,a) - \mathbb{E}\{X\} - \epsilon_1
\end{aligned}$$

and because $(s,a) \notin K_t$ meaning $Q_t(s,a) - \mathbb{E}\{X\} > 3\epsilon_1$,

$$Q_t(s,a) - \mathbb{E}\{X\} - \epsilon_1 > 2\epsilon_1$$

guaranteeing success for the type-1 update at timestep $t_{m_1}$. Since for the case that $l(s,a) = 0$ and $\mathsf{learn}(s,a) = \mathrm{true}$, an attempted update will necessarily happen; there can be at most $|S||A|(1 + \kappa)$ instances of such an event. Working in a fashion similar to the proof of Lemma 4, one concludes that the lemma's statement holds with probability at least $1 - \frac{\delta}{8}$.

# H   Proof of Lemma 10

We will assume that $(s,a)$ has not already been visited $m_2$ times before timestep $t$, because then it is obvious that $(s,a) \in K_{t+1}$. Thus we work under the assumption that $(s,a)$ has been visited fewer than $m_2$ times up until $t$, at which time an unsuccessful update of $(s,a)$ occurs, while right after at $t+1$ we see $\mathsf{learn}(s,a) = \mathrm{false}$. Now set up a contradiction argument: under those conditions, *assume that* $(s,a) \notin K_{t+1}$. Since the update at $t$ was unsuccessful, $K_{t+1} = K_t$, which would also imply that $(s,a) \notin K_t$. Now label the times of the most recent $m_1$ experiences of $(s,a)$ as $b(s,a) \triangleq t_1 < t_2 < \cdots < t_{m_1} = t$. The contrapositive of the statement proved in Lemma 9, suggests that since the update at $t$ is unsuccessful, it must be $(s,a) \in K_{t_1}$. Since $(s,a) \notin K_t$, some timestep between $t_1$ and $t$ must have been successful. Let us denote that successful timestep $t^* > b(s,a)$. But then the condition $t_1 = b(s,a) < t^*$ would not allow

the learn flag to be set to $\mathsf{false}$ in between these two timesteps, and we know from the statement of the lemma that this is true. Therefore, we have a contradiction; the assumption made is invalid, and therefore $(s, a) \in K_t = K_{t+1}$.

## I   Proof of Lemma 11

Fix a state-action pair $(s, a)$, We begin by showing that if $(s, a) \notin K_t$ at timestep $t$, then within at most $2m_1$ more experiences of $(s, a)$ after $t$, a successful timestep for $(s, a)$ must occur. Toward that end, we analyse the worst case where $m_2$-th visit of $(s, a)$ will not occur within $2m_1$ more experiences of $(s, a)$ after timestep $t$. For $(s, a) \notin K_t$, distinguish two possible cases at the beginning of timestep $t$: either $\mathsf{learn}(s, a) = \mathsf{false}$ or $\mathsf{learn}(s, a) = \mathsf{true}$. Consider first the case where $\mathsf{learn}(s, a) = \mathsf{false}$. Assume that the most recent attempted update of $(s, a)$ occurred at some timestep $t'$ which was unsuccessful and set the flag $\mathsf{learn}(s, a)$ to false. Then, according to Lemma 10, it will be $(s, a) \in K_{t'+1}$. However, now it is $(s, a) \notin K_t$, which implies that a successful timestep must have occurred at some $t^*$ with $t' + 1 < t^* < t$. Thus the flag $\mathsf{learn}(s, a)$ will set to $\mathsf{true}$ during timestep $t$. Then, at $t$ we have all conditions of Lemma 9 (i.e. $\mathsf{learn}(s, a) = \mathsf{true}$, $(s, a) \notin K_t$ and $l(s, a) = 0$) and thus the type-1 update upon $m_1$-th visit of $(s, a)$ after $t$ will be successful.

Take now the case where $\mathsf{learn}(s, a) = \mathsf{true}$. We know that an attempted type-1 update for $(s, a)$ will occur in at most $m_1$ experiences of $(s, a)$, and those are assumed occurring at timesteps $t_1 < \cdots < t_{m_1}$, then $t_1 \leq t \leq t_{m_1}$. Consider the two possibilities: $(s, a) \notin K_{t_1}$ or $(s, a) \in K_{t_1}$. In the former case, Lemma 9 indicates that the attempted update type-1 at $t_{m_1}$ will be successful. In the latter case, given that $(s, a) \notin K_t$, a successful timestep $t^*$ must have taken place between $t_1$ and $t$ (since $K_{t_1} \neq K_t$). Thus, however the attempted update at $t_{m_1}$ is unsuccessful, $\mathsf{learn}(s, a)$ will remain $\mathsf{true}$ and at timestep $t_{m_1} + 1$ we will have $\mathsf{learn}(s, a) = \mathsf{true}$, $l(s, a) = 0$, and $(s, a) \notin K_{t_{m_1}+1}$; this would trigger Lemma 9, and the attempted update type-1 upon $m_1$-th visit of $(s, a)$ after timestep $t_{m_1} + 1$ (which is within at most $2m_1$ more experiences of $(s, a)$ after $t$), will be successful.

Thus far, we showed that after $(s, a) \notin K_t$, within at most $2m_1$ more experiences of $(s, a)$, at least one successful timestep for $(s, a)$ must occur. According to lemma 1, the total number of successful timesteps for $(s, a)$ are bounded by $1 + \frac{1}{(1-\gamma)\epsilon_1}$. This means that the total number of timesteps with $(s, a) \notin K_t$ is bounded by $2m_1(1 + \frac{1}{(1-\gamma)\epsilon_1})$. On the other hand, once a state-action pair $(s, a)$ is experienced for $m_2$-th time at any timestep $t$, it will become a member of $K_t$ and will never leave $K_t$ anymore. So, $m_2$ is another upper-bound for the number of timesteps with $(s, a) \notin K_t$.

Generalizing the above fact for all state-action pairs, we conclude that the total number of escape events (timesteps $t$ with $(s_t, a_t) \notin K_t$) is bounded by $\min(2m_1\kappa, |S||A|m_2)$.