

---

# On the Sample Complexity of Learning Infinite-horizon Discounted Linear Kernel MDPs

---

Yuanzhou Chen<sup>\*1</sup> Jiafan He<sup>\*2</sup> Quanquan Gu<sup>2</sup>

## Abstract

We study reinforcement learning for infinite-horizon discounted linear kernel MDPs, where the transition probability function is linear in a predefined feature mapping. Existing UCLK (Zhou et al., 2021b) algorithm for this setting only has a regret guarantee, which cannot lead to a tight sample complexity bound. In this paper, we extend uniform-PAC sample complexity from the episodic setting to the infinite-horizon discounted setting, and propose a novel algorithm dubbed UPAC-UCLK that achieves an  $\tilde{O}(d^2/((1-\gamma)^4\epsilon^2) + 1/((1-\gamma)^6\epsilon^2))$  uniform-PAC sample complexity, where  $d$  is the dimension of the feature mapping,  $\gamma \in (0, 1)$  is the discount factor of the MDP and  $\epsilon$  is the accuracy parameter. To the best of our knowledge, this is the first  $\tilde{O}(1/\epsilon^2)$  sample complexity bound for learning infinite-horizon discounted MDPs with linear function approximation (without access to the generative model).

## 1. Introduction

In reinforcement learning (RL), the central goal is to design an efficient algorithm to learn the optimal (or near-optimal) policy through repeated interactions between the agent and an unknown environment. Markov decision processes (MDPs) are typical models that describe the environment formally. An MDP is described by a tuple of action space, state space, reward function and transition probability function. Based on the planning horizon and the transition probability kernel, MDPs can be categorized into three different types: (1) episodic MDPs (Azar et al., 2017; Dann et al., 2017; Jin et al., 2018; 2020; Ayoub et al.,

2020; Zhou et al., 2021a; Zhang et al., 2020a; Zanette et al., 2020; He et al., 2021b), (2) infinite-horizon average reward MDPs (Jaksch et al., 2010; Bartlett & Tewari, 2012; Ouyang et al., 2017; Agrawal & Jia, 2017; Fruit et al., 2018b;a; Talebi & Maillard, 2018; Zhang & Ji, 2019; Fruit et al., 2020; Ortner, 2020; Wei et al., 2021; Wu et al., 2022), and (3) infinite-horizon discounted MDPs (Kakade et al., 2003; Strehl et al., 2006; Dong et al., 2019; Zhang et al., 2020b; Szita & Szepesvári, 2010; Lattimore & Hutter, 2012; Liu & Su, 2020; He et al., 2021a; Zhou et al., 2021a;b). The theoretical studies of infinite-horizon discounted MDPs are still quite limited, compared with the other two types of MDPs. Without the restart process in episodic MDPs, or the bounded diameter condition in infinite-horizon average reward MDPs (Jaksch et al., 2010), the infinite-horizon discounted MDP poses a great challenge for both algorithm design and theoretical analysis.

Existing results for discounted MDPs mainly focus on the tabular setting, where the agent can access each state-action pair  $(s, a)$  and estimate the corresponding value functions individually. Without the help of a *generative model* (Kakade et al., 2003) that allows the agent to visit all state-action pairs  $(s, a)$ , Lattimore & Hutter (2012); Dong et al. (2019); Zhang et al. (2020b) and Liu & Su (2020); He et al. (2021a) showed that it is still possible to obtain a polynomial sample complexity or  $O(\sqrt{T})$ -regret for discounted tabular MDPs.

However, these tabular RL algorithms are computationally inefficient when the sizes of action and state spaces are large. A common approach to deal with high-dimensional (or even infinite) state space and action space is to use certain function classes, such as linear function class and neural networks, to approximate the transition probability kernel or the value function (Jin et al., 2020; Ayoub et al., 2020). For discounted MDPs, Zhou et al. (2021b) studied a linear kernel MDP, where the transition probability kernel can be represented by a linear function in some known feature mapping, and proposed the UCLK algorithm with an  $O(\sqrt{T})$ -regret guarantee. Unfortunately, unlike episodic MDPs initialized at the same state for different episodes, the regret guarantee cannot imply a finite sample complexity bound for discounted MDPs. Therefore, UCLK can fail to learn the  $\epsilon$ -optimal policy, even though its regret is sublinear. Thus, a natural question arises:

---

<sup>\*</sup>Equal contribution <sup>1</sup>School of Mathematical Sciences, Peking University, Beijing, China <sup>2</sup>Computer Science Department, University of California at Los Angeles, Los Angeles, California, USA. Correspondence to: Quanquan Gu <qgu@cs.ucla.edu>.

*Can we design provably efficient algorithms with polynomial sample complexity for discounted MDPs with linear function approximation?*

In this paper, we answer this question affirmatively by presenting a variant of the UCLK algorithm (Zhou et al., 2021b), namely UPAC-UCLK, and prove that this algorithm has a near-optimal sample complexity. In fact, our algorithm enjoys a stronger notion of sample complexity guarantee than the standard PAC sample complexity, which is called uniform-PAC sample complexity Dann et al. (2017). As discussed in Dann et al. (2017), the uniform-PAC guarantee is stronger than both PAC sample complexity and regret, and can further guarantee the convergence to the optimal policy up to an arbitrarily small error.

Our contributions are summarized as follows:

- We adapt the multi-level scheme from He et al. (2021b) to the UCLK algorithm (Zhou et al., 2021b), and propose a novel algorithm UPAC-UCLK for discounted linear kernel MDPs. Compared with the original multi-level scheme, we propose a novel *discounted data inheritance* technique, which adds a discounted portion of data at each level to the subsequent level. With this technique, newly added levels are non-empty at initialization due to data *inheritance*, and high levels contain low data information due to data *discounting*. This technique provides crucial guarantees for our algorithm, allowing us to uniformly bound the sample complexity.
- We show that our algorithm satisfies the uniform-PAC guarantee with sample complexity  $\tilde{O}(d^2/((1-\gamma)^4\epsilon^2) + 1/((1-\gamma)^6\epsilon^2))$ , where  $d$  is the dimension of the feature mapping,  $\gamma \in (0, 1)$  is the discount factor of the MDP and  $\epsilon$  is the accuracy parameter. This result immediately implies a high probability regret bound  $\tilde{O}(d\sqrt{T}/(1-\gamma)^2 + \sqrt{T}/(1-\gamma)^3)$ , where  $T$  is the number of interactions with the environment. The regret bound matches that of UCLK algorithm (Zhou et al., 2021b) up to a logarithmic factor, after ignoring the extra  $\tilde{O}(\sqrt{T}/(1-\gamma)^3)$  term in the regret. To the best of our knowledge, UPAC-UCLK is the first algorithm for infinite-horizon discounted RL with linear function approximation that enjoys a (uniform) sample complexity guarantee.

The remaining of this paper is organized as follows. Section 2 reviews the mostly related work. Section 3 provides the preliminaries of MDPs and reinforcement learning. Section 4 presents our main algorithms. Section 5 introduces the theoretical guarantees of our algorithm, and discusses its implications. Section 6 provides a proof outline for the main theorem, along with several key technical lemmas.

Section 7 concludes this work with discussions on the future work. The detailed proofs are deferred to the appendix.

**Notation** In this work, we use lower case letters to denote scalars, and use lower and upper case bold face letters to denote vectors in  $\mathbb{R}^d$  and  $d \times d$  matrices respectively. For a vector  $\mathbf{x} \in \mathbb{R}^d$ , we denote by  $\|\mathbf{x}\|_2$  the Euclidean norm. Furthermore, for a positive-definite matrix  $\Sigma \in \mathbb{R}^{d \times d}$ , we define the Mahalanobis norm of  $\mathbf{x}$  with respect to  $\Sigma$  to be  $\|\mathbf{x}\|_\Sigma = \sqrt{\mathbf{x}^\top \Sigma \mathbf{x}}$ . In addition, for a matrix  $\Sigma$  with real eigenvalues, we denote its largest and smallest eigenvalues by  $\lambda_{\max}(\Sigma)$  and  $\lambda_{\min}(\Sigma)$  respectively. For two sequences  $\{a_n\}$  and  $\{b_n\}$ , we write  $a_n = O(b_n)$  if there exists an integer subscript  $n_0$  and an absolute constant  $C$ , such that for all  $n \geq n_0$ ,  $a_n \leq Cb_n$ . We use  $\tilde{O}(\cdot)$  to further hide logarithmic factors except for the  $\log T$  term in all uniform-PAC guarantees.

## 2. Related Work

In this section, we will review the mostly related work to our work.

### 2.1. Infinite-horizon Discounted MDPs

For discounted tabular MDPs, there has been a series of work providing sample complexity guarantees without the help of a generative model (Kakade et al., 2003). For these results, the algorithms have access to every possible state-action pair. They can be mainly grouped into two categories: model-free algorithms and model-based algorithms. For model-free algorithms, the value function for each state-action pair is directly estimated. For instance, Strehl et al. (2006) first proposed the delayed Q-learning algorithm with polynomial sample complexity guarantee. Later, Dong et al. (2019) improved the sample complexity to  $\tilde{O}(SA/((1-\gamma)^7\epsilon^2))$  by proposing an infinite Q-learning with UCB algorithm, where  $S$  is the number of states and  $A$  is the number of actions. Recently, Zhang et al. (2020b) proposed the UCB-multistage-adv algorithm, which obtains a near-optimal sample complexity  $\tilde{O}(SA/((1-\gamma)^3\epsilon^2))$  that matches the theoretical lower bound. Model-based algorithms, on the other hand, estimate the underlining environment or the transition probability kernel and compute the value function from the estimated environment. For instance, Szita & Szepesvári (2010) introduced the MoRmax algorithm with  $\tilde{O}(SA/((1-\gamma)^6\epsilon^2))$  sample complexity. Later, Lattimore & Hutter (2012) proposed the UCRL- $\gamma$  algorithm and improved the sample complexity to  $\tilde{O}(SA/((1-\gamma)^3\epsilon^2))$ , under a strong assumption on the transition probability kernel. Our work also falls into the category of model-based algorithms, and considers the sample complexity as the performance measure.

Recently, Liu & Su (2020) extended the definition of re-

gret from episodic MDPs to the discounted MDPs, and named it  $\gamma$ -regret. Under this notion of regret, Liu & Su (2020) proposed the Double Q-Learning algorithm with an  $\tilde{O}(\sqrt{SAT}/(1-\gamma)^{2.5})$  regret guarantee for the first  $T$  steps. Later, He et al. (2021a) introduced the UCBVI- $\gamma$  algorithm and improved the regret bound to  $\tilde{O}(\sqrt{SAT}/(1-\gamma)^{1.5})$ . Additionally, He et al. (2021a) proved there exists a class of hard-to-learn discounted MDPs such that the regret for any algorithm is lower bounded by  $\tilde{\Omega}(\sqrt{SAT}/(1-\gamma)^{1.5})$ . This matches the regret upper bound of UCBVI- $\gamma$  algorithm up to a logarithmic factor.

In an attempt to design algorithms that can learn efficiently when the state and action spaces are high-dimensional or even infinite, Zhou et al. (2021b) focused on discounted MDPs with linear function approximation, and proposed the UCLK algorithm with a  $\tilde{O}(d\sqrt{T}/(1-\gamma)^2)$  regret guarantee, where  $d$  is the dimension of feature mapping. Later, Zhou et al. (2021a) introduced the Bernstein-type bonus and improved the regret bound to  $\tilde{O}(d\sqrt{T}/(1-\gamma)^{1.5})$  in their new UCLK<sup>+</sup> algorithm, which matches the regret lower bound proved in Zhou et al. (2021b) up to a logarithmic factor. However, both UCLK and UCLK<sup>+</sup> only have regret guarantees, and cannot provide any sample complexity guarantees. To the best of our knowledge, our UPAC-UCLK algorithm gives the first sample complexity guarantee for learning discounted MDPs with linear function approximation.

## 2.2. Uniform-PAC Guarantees in RL

Traditional analysis on reinforcement learning mainly focuses on the regret in the first  $T$  steps or the sample complexity with respect to a specific accuracy parameter  $\epsilon$ . Unfortunately, both of these performance measures fail to guarantee the convergence to optimal policy. To overcome this problem, Dann et al. (2017) first introduced a stronger performance guarantee named uniform-PAC for episodic MDPs. Different from traditional sample complexity guarantee, uniform-PAC requires bounding the sample complexity at all accuracy parameters  $\epsilon$  simultaneously, and thus guarantees that the number of steps with suboptimality larger than  $\epsilon$  is finite. Under this stronger performance measure, Dann et al. (2017) proposed the UBEV algorithm and obtained an  $\tilde{O}(SAH^4/\epsilon^2)$  uniform-PAC guarantee, which implies that the UBEV algorithm converges to the optimal policy. Here,  $H$  is the length of each episode. Later, He et al. (2021b) focused on episodic MDPs with linear function approximation and proposed the FLUTE algorithm. By introducing a minimax value function estimator and a multi-level scheme, the FLUTE algorithm obtains an  $\tilde{O}(d^3H^5/\epsilon^2)$  uniform-PAC guarantee for episodic linear MDPs (Jin et al., 2020). Our algorithm UPAC-UCLK also satisfies the stronger uniform-PAC guarantee, under the setting of discounted MDPs with linear function approximation.

## 3. Preliminaries

In this work, we consider the infinite-horizon discounted Markov decision processes (MDPs), denoted by  $M(\mathcal{S}, \mathcal{A}, \gamma, r, \mathbb{P})$ . In this tuple,  $\mathcal{S}, \mathcal{A}$  are the spaces of states and actions respectively;  $\gamma \in (0, 1)$  is the discount factor;  $r(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the reward function, which we assume to be deterministic and known to the agent;  $\mathbb{P}(\cdot | \cdot, \cdot) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^+$  is the transition probability function that satisfies the equations  $\int_{s' \in \mathcal{S}} \mathbb{P}(s' | s, a) d\mu(s') = 1$  for arbitrary  $s \in \mathcal{S}, a \in \mathcal{A}$ ; here  $d\mu$  is a probability measure on  $\mathcal{S}$ . We will abbreviate  $d\mu(s)$  as  $ds$  henceforth.

As an agent tries to learn an MDP via interaction with the environment, it repeats the following process: it takes an action  $a \in \mathcal{A}$  based on the current state  $s$ , then receives the next state  $s' \in \mathcal{S}$  with the reward  $r$ . We can thus label these states and actions in chronological order with time step subscripts  $t = 1, 2, \dots$ . For the agent, its choice of an action  $a_t$  is based only on the available observations before the action, namely a tuple  $(s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t)$ . Consequently, we can define the agent's (non-stationary) policy  $\pi$  on an MDP  $M(\mathcal{S}, \mathcal{A}, \gamma, r, \mathbb{P})$  to be  $\pi = \{\pi_t\}_{t=1}^\infty$ , where  $\pi_t : \{\mathcal{S} \times \mathcal{A}\}^{t-1} \times \mathcal{S} \rightarrow \mathcal{A}$  maps the tuple  $(s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t)$  to an action  $a_t$ . The agent's policy  $\pi$ , combined with the transition probability  $\mathbb{P}$ , give rise to the infinite-length random process consisting of states and actions  $\{s_t, a_t\}_{t=1}^\infty$ .

Now consider the expected total discounted reward of a given policy  $\pi$  after a certain point in the learning process, a.k.a., value functions. The action-value function and the value function for policy  $\pi$  at time step  $t$  are defined as

$$Q_t^\pi(s, a) := \mathbb{E} \left[ \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r(s_\tau, a_\tau) \middle| s_1, \dots, s_t = s, a_t = a \right],$$

$$V_t^\pi(s) := \mathbb{E} \left[ \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r(s_\tau, a_\tau) \middle| s_1, \dots, s_t = s \right],$$

where expectation is conditioned on the states and actions up to time step  $t$ . With these in place, we can define the optimal action-value function  $Q^*(s, a) := \sup_\pi Q_1^\pi(s, a)$  and optimal value function  $V^*(s) := \sup_\pi V_1^\pi(s)$ . We refer to the difference between the optimal value function and the value function of policy  $\pi$  at time step  $t$ ,  $\Delta_t = V^*(s_t) - V_t^\pi(s_t)$ , as the *suboptimality gap*. Now for any real function defined on the state space  $V : \mathcal{S} \rightarrow \mathbb{R}$ , we define  $\mathbb{P}V(s, a) := \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} V(s')$ . With this notation, we can now deduce and express the Bellman equations in discounted MDPs as follows:

$$Q_t^\pi(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{P}V_{t+1}^\pi(s_t, a_t) \quad (3.1)$$

$$Q^*(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{P}V^*(s_t, a_t). \quad (3.2)$$

We focus on linear kernel MDPs (Zhou et al., 2021b) in this

work, which is also known as linear mixture models (Modi et al., 2019) or linear mixture MDPs (Ayoub et al., 2020).

**Definition 3.1** (Linear Kernel MDP (Zhou et al., 2021b)). An MDP  $M(\mathcal{S}, \mathcal{A}, \gamma, r, \mathbb{P})$  is called a linear kernel MDP if there exists a *known* feature mapping  $\phi(\cdot|\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^d$  and an *unknown* vector  $\theta \in \mathbb{R}^d, \|\theta\|_2 \leq \sqrt{d}$ , such that:

1. For all  $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ , the transition probability function can be expressed as  $\mathbb{P}(s'|s, a) = \langle \phi(s'|s, a), \theta \rangle$ ;
2. For any bounded function  $V : \mathcal{S} \rightarrow [0, R]$ ,  $\|\phi_V(s, a)\|_2 \leq \sqrt{d}R$  holds for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , where  $\phi_V(s, a) := \int_{s' \in \mathcal{S}} \phi(s'|s, a)V(s')ds'$ .

We denote such an MDP by  $M(\theta; \phi)$  for simplicity.

From Definition 3.1, the following equation holds for any function  $V : \mathcal{S} \rightarrow [0, R]$ :

$$\langle \theta^*, \phi_V(s, a) \rangle = \mathbb{P}V(s, a). \quad (3.3)$$

Now we define the set of all feasible  $\theta$ 's satisfying the conditions in Definition 3.1 to be

$$\mathcal{B} := \left\{ \theta \in \mathbb{R}^d : \|\theta\|_2 \leq \sqrt{d} \text{ and } \langle \phi(\cdot|s, a), \theta \rangle \text{ is a probability measure on } \mathcal{S} \right\}, \quad (3.4)$$

which, as we will see in Lemma A.3 in the appendix, is a convex body when constrained to an affine subspace  $\mathcal{L}$  of  $\mathbb{R}^d$ .

For linear kernel MDPs, Zhou et al. (2021b) proposed the Upper-Confidence Linear Kernel (UCLK) algorithm that achieves near-optimal regret. However, as pointed out by Dann et al. (2017), algorithms with regret bounds do not generally converge to the optimal policy: as long as the regret tends to infinity as  $t \rightarrow \infty$ , high suboptimal gaps may occur infinitely often. To overcome this issue, Dann et al. (2017) proposed a stronger notion of performance measure called uniform-PAC that guarantees the convergence to the optimal policy with high probability. We now extend it to the discounted MDP setting.

**Definition 3.2** (Uniform-PAC). For an algorithm Alg on an MDP, denote its corresponding policy by  $\pi$ . For any  $\epsilon > 0$ , define  $N_\epsilon = \sum_{t=1}^{\infty} \mathbb{1}\{V^*(s_t) - V_t^\pi(s_t) > \epsilon\}$  to be the total number of time steps where the suboptimality gap is greater than  $\epsilon$ . Algorithm Alg is said to be uniform-PAC for some  $\delta \in (0, 1)$  with sample complexity  $\Gamma(1/\epsilon, \log(1/\delta))$  if

$$\mathbb{P}[\exists \epsilon > 0, N_\epsilon > \Gamma(1/\epsilon, \log(1/\delta))] \leq \delta, \forall \delta \in (0, 1),$$

where  $\Gamma(\cdot, \cdot)$  is a polynomial function dependent on the MDP itself.

**Remark 3.3.** Definition 3.2 is slightly different from that in Dann et al. (2017); He et al. (2021b) due to the different definitions of  $N_\epsilon$ . Specifically, Dann et al. (2017); He et al. (2021b) focused on episodic MDPs, where the suboptimality gap  $\Delta_k = V_1^*(s_1^k) - V_1^\pi(s_1^k)$ , and  $N_\epsilon$  is the number of episodes with suboptimality gap greater than  $\epsilon$ . Due to the difference between episodic MDPs and discounted MDPs, we define the suboptimality for each time step  $t$  as  $\Delta_t = V^*(s_t) - V_t^\pi(s_t)$  and define the sample complexity as  $N_\epsilon = \sum_{t=1}^{\infty} \mathbb{1}\{V^*(s_t) - V_t^\pi(s_t) > \epsilon\}$ . In addition, for linear kernel MDPs, the sample complexity bound should be expressed in the form  $\Gamma(1/\epsilon, \log(1/\delta); \gamma, d)$ , where  $\gamma, d$  are the parameters of  $M(\theta; \phi)$ , rather than  $\Gamma(1/\epsilon, \log(1/\delta); \gamma, |\mathcal{S}|, |\mathcal{A}|)$  for tabular MDPs. One would expect  $\Gamma$  to be polynomial in  $1/(1 - \gamma)$  and  $d$  as well. Though in Dann et al. (2017), uniform-PAC was defined in the finite-horizon setting with finite state and action spaces, the proof of its basic properties does not depend on the specific MDP model. Hence, we still have the conclusions that uniform-PAC sample complexity bound implies both the PAC sample complexity bound and the regret bound, and it guarantees the convergence to the optimal policy with high probability, as stated in Theorem 3 of Dann et al. (2017).

From now on, we assume the true  $\theta$  in the target linear kernel MDP  $M(\theta; \phi)$  to be  $\theta^*$ .

## 4. Proposed Algorithms

In this section, we propose a new algorithm called UPAC-UCLK for learning infinite-horizon discounted linear kernel MDPs. Our algorithm UPAC-UCLK is inspired by the UCLK algorithm (Zhou et al., 2021b). To achieve the uniform-PAC guarantee, we incorporate the multi-level partition scheme proposed in He et al. (2021b) and separate all time steps and their corresponding data into different levels. The algorithm is described in Algorithm 1. It is worth noting that the direct application of multi-level scheme to the infinite-horizon discounted MDPs requires maintaining infinite number of levels. However, this is infeasible in practice. In order to overcome this problem, we will maintain  $L$  ‘‘effective’’ levels along with extra  $3L + C_0$  auxiliary levels in Algorithm 1, which can be seen as a finite-level approximation to the infinite number of levels.

### 4.1. Planning (Line 4 to 7 of Algorithm 1)

At the beginning of time step  $t$ , our algorithm calculates the estimates of  $\theta^*$  for all available levels  $l = 1, 2, \dots, 4L + C_0$ . By (3.3), we have  $\langle \theta^*, \phi_{V_t}(s_t, a_t) \rangle = \mathbb{P}V_t(s_t, a_t) = \mathbb{E}[V_t(s_{t+1})|s_t, a_t]$ , so a natural method to estimate  $\theta^*$  is using linear regression. In particular, we calculate the estimates by the closed-form solution of linear regression

**Algorithm 1** Uniform-PAC UCLK (UPAC-UCLK)

---

**Require:** Regularization parameter  $\lambda$ , exploration parameters  $\beta^l$  for  $l = 1, 2, \dots$ , number of value iteration rounds  $U_t$  for  $t = 1, 2, \dots$

- 1: **Initialize**  $C_0 \leftarrow 3 \lceil \log(\sqrt{d}/(1-\gamma))/\log 2 \rceil + 6$ ,  
 $L \leftarrow 0$ ,  $\Sigma^l \leftarrow \lambda \mathbf{I}$ ,  $\mathbf{b}^l \leftarrow \mathbf{0}$ ,  $l = 1, 2, \dots$
- 2: Receive initial state  $s_1$ .
- 3: **for** time step  $t = 1, 2, \dots$  **do**
- 4:   **for** all level  $l \in \{1, 2, \dots, 4L + C_0\}$  **do**
- 5:      $\hat{\theta}^l \leftarrow (\Sigma^l)^{-1} \mathbf{b}^l$ ,  $\mathcal{C}^l \leftarrow \{\theta : \|\theta - \hat{\theta}^l\|_{\Sigma^l} \leq \beta^l\}$ .
- 6:   **end for**
- 7:    $\{Q_t^l(\cdot, \cdot)\}_{l=1}^{4L+C_0}$ ,  $V_t(\cdot) \leftarrow \text{ML-EVI}(\{\mathcal{C}^l\}_{l=1}^{4L+C_0}, U_t)$ .
- 8:    $a_t \leftarrow \arg\max_a \min_{1 \leq l \leq L} Q_t^l(s_t, a)$ ,
- 9:   receive  $s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)$ .
- 10:    $l_t \leftarrow 1$ .
- 11:   **while**  $l_t \leq L$  **and**  
        $\|\phi_{V_t}(s_t, a_t)\|_{(\Sigma^{l_t})^{-1}} \leq 2^{-l_t} \sqrt{d}/(1-\gamma)$  **do**
- 12:      $l_t \leftarrow l_t + 1$
- 13:   **end while**
- 14:   **if**  $l_t = L + 1$  **then**
- 15:     **for** level  $l = 4L + C_0 + 1, \dots, 4L + C_0 + 4$  **do**
- 16:       $\Sigma^l \leftarrow \lambda \mathbf{I} + \frac{1}{2}(\Sigma^{l-1} - \lambda \mathbf{I})$ ,  $\mathbf{b}^l \leftarrow \frac{1}{2} \mathbf{b}^{l-1}$ .
- 17:     **end for**
- 18:      $L \leftarrow L + 1$ .
- 19:   **end if**
- 20:   **for** level  $l = l_t, \dots, 4L + C_0$  **do**
- 21:      $\Sigma^l \leftarrow \Sigma^l + 2^{l-t} \phi_{V_t}(s_t, a_t) \phi_{V_t}(s_t, a_t)^\top$ ,
- 22:      $\mathbf{b}^l \leftarrow \mathbf{b}^l + 2^{l-t} \phi_{V_t}(s_t, a_t) V_t(s_{t+1})$ .
- 23:   **end for**
- 24: **end for**

---

$\hat{\theta}^l = (\Sigma^l)^{-1} \mathbf{b}^l$  for each level, where  $\Sigma^l$  and  $\mathbf{b}^l$  are updated in Lines 16, 21 and 22. In Line 5, we then construct the confidence sets  $\mathcal{C}^l$  of  $\theta^*$  for all levels  $l$ , which are centered at  $\hat{\theta}^l$  with radius  $\beta^l$ . The Multi-Level Extended Value Iteration Algorithm (ML-EVI), as described in the next subsection, then calculates the estimates of value functions  $\{Q_t^l\}$  and  $V_t$ .

**4.2. Multi-Level Extended Value Iteration (Algorithm 2)**

The ML-EVI algorithm is described in Algorithm 2. It can be seen as an extension of the Extended Value Iteration algorithm (EVI, Algorithm 2 from Zhou et al. (2021b)). By alternating between performing the greedy policy (Line 3) and an approximate Bellman equation (Line 6), it finds an optimistic estimate for the optimal value functions  $Q^*(\cdot, \cdot)$  and  $V^*(\cdot)$ .

**Discounted Data Inheritance.** In He et al. (2021b), each data point (i.e., a pair of regression predictor and response) is only added to one single level. In contrast, in our setting,

**Algorithm 2** Multi-Level Extended Value Iteration (ML-EVI)

---

**Require:** Number of levels  $L$ , confidence sets  $\mathcal{C}_l$ ,  $l = 1, \dots, L$ , number of value iteration rounds  $U$ .

- 1:  $Q_l^{(0)}(\cdot, \cdot) \leftarrow 1/(1-\gamma)$  for  $l = 1, \dots, L$ .
- 2: **for**  $u = 1, \dots, U$  **do**
- 3:    $V^{(u-1)}(\cdot) \leftarrow \max_{a \in \mathcal{A}} \min_{1 \leq l \leq L} Q_l^{(u-1)}(\cdot, a)$ ,
- 4:   **for**  $l = 1, \dots, L$  **do**
- 5:     **if**  $\mathcal{B} \cap \mathcal{C}_l \neq \emptyset$  **then**
- 6:        $Q_l^{(u)}(\cdot, \cdot) \leftarrow r(\cdot, \cdot)$   
        $+ \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-1)}}(\cdot, \cdot) \rangle$ .
- 7:     **else**
- 8:        $Q_l^{(u)}(\cdot, \cdot) \leftarrow 1/(1-\gamma)$ .
- 9:     **end if**
- 10:   **end for**
- 11: **end for**
- 12:  $V^{(U)}(\cdot) \leftarrow \max_a \min_{1 \leq l \leq L} Q_l^{(U)}(\cdot, a)$ .
- 13: **return**  $\{Q_l^{(U)}(\cdot, \cdot)\}_{l=1}^L, V^{(U)}(\cdot)$ .

---

the data in one level are inherited by the subsequent levels after being discounted. Under this scheme, our estimate of  $\theta^*$  at level  $l$  at time step  $t$  regresses over the set of data

$$\{2^{(l_\tau - l)/2} (\phi_{V_\tau}(s_\tau, a_\tau), V_\tau(s_{\tau+1})) : \tau \leq t, l_\tau \leq l\},$$

which has a closed-form solution  $\hat{\theta}^l = (\Sigma^l)^{-1} \mathbf{b}^l$  with  $\Sigma^l$  and  $\mathbf{b}^l$  being defined as follows:

$$\Sigma^l = \lambda \mathbf{I} + \sum_{\tau=1}^t \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top,$$

$$\mathbf{b}^l = \sum_{\tau=1}^t \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}).$$

We now explain the reasons behind this algorithm design.

As we will see in (6.5) of Section 6, in order to bound the sample complexity, we need to control the term  $(V_t - V_{t+1})(s_{t+1})$ : the difference on  $s_{t+1}$  between two consecutive optimistic value functions output by Algorithm 2. When the total number of levels  $L$  is increased by 1 in Line 18 of Algorithm 1, we say a new level is activated. Without data inheritance, newly activated levels hold a very scarce amount of data and have large exploration radii. This implies that the corresponding confidence sets are large. If these levels were selected in the minimization operation of Line 3 or Line 12 in Algorithm 2 at a later iteration  $u$ , it may produce an inaccurate estimation of the value function. In contrast, with data inheritance, newly activated levels will have enough data to produce sufficiently small confidence sets and this problem can be alleviated.

In addition, we need to discount the data when inheriting data from lower levels. The main reason is that for those

levels much higher than  $L$ , they still have large exploration radii, so adding the undiscounted new data can result in big changes to their confidence sets. With data discounting, the addition of data to levels much higher than  $L$  is exponentially small and therefore negligible.

With the help of discounted data inheritance, the term  $(V_t - V_{t+1})(s_{t+1})$  can be roughly bounded by  $\tilde{O}(\sqrt{\phi_{V_t}(s_t, a_t)^\top (\Sigma^l)^{-1} \phi_{V_t}(s_t, a_t)}) = O(2^{-l_t})$  and will diminish as  $t \rightarrow \infty$ .

### 4.3. Execution (Line 8 to 9 in Algorithm 1)

The algorithm estimates the action-value function  $Q^*$  by taking a minimum on the upper bounds  $Q_t^l$ , and chooses the next action based on this estimate. The next state is then generated by the environment based on the underlying transition probability function  $\mathbb{P}$ .

### 4.4. Level Selection (Line 10 to 23 in Algorithm 1)

Since our new predictor at time step  $t$  is  $\phi_{V_t}(s_t, a_t)$ , its uncertainty at level  $l$  is characterized by  $\|\phi_{V_t}(s_t, a_t)\|_{(\Sigma^l)^{-1}}$ . Based on the uncertainty, our algorithm finds the minimum level  $l$  at which the uncertainty exceeds the threshold value at level  $l$ , i.e.,

$$l_t = \min_{l \leq L} \{ \|\phi_{V_t}(s_t, a_t)\|_{(\Sigma^l)^{-1}} > 2^{-l} \sqrt{d}/(1-\gamma) \}.$$

If this minimum value does not exist, or in other words the uncertainty is less than the thresholds at all effective levels  $l \leq L$ ,  $l_t$  is set to  $L+1$ ,  $L \leftarrow L+1$ , and a new effective level is activated along with 3 other auxiliary levels. Based on the discounted data inheritance scheme, the initialization of these new levels has the following form

$$\Sigma^l \leftarrow \lambda \mathbf{I} + \frac{1}{2} (\Sigma^{l-1} - \lambda \mathbf{I}), \mathbf{b}^l \leftarrow \frac{1}{2} \mathbf{b}^{l-1}.$$

Finally, Algorithm 1 updates the variables  $\Sigma, \mathbf{b}$  for all levels  $l \geq l_t$  by adding the new data using the discounted data inheritance scheme (Line 20 to 23).

## 5. Main Results

In this section, we present the theoretical guarantees for our UPAC-UCLK algorithm. Our main result is a theorem which certifies that our algorithm UPAC-UCLK has a uniform-PAC guarantee with an efficient sample complexity bound.

**Theorem 5.1.** For any linear mixture MDP, if we set  $\lambda \geq 1$

and the parameters  $\beta^l$  and  $U_t$  in Algorithm 1 as follows:

$$\begin{aligned} U_t &:= \lceil \log(t(t+1))/(1-\gamma) \rceil, \\ \beta^l &:= \frac{2}{1-\gamma} (3\sqrt{d \max\{l, l_0\}} + 2\sqrt{\log(1/\delta)}) + 2\sqrt{d\lambda}, \end{aligned} \quad (5.1)$$

where  $l_0 := \log(\sqrt{d}/(1-\gamma))/\log 2$ , then with probability at least  $1 - \delta$ , for all accuracy parameter  $\epsilon > 0$ , the number of time steps with suboptimality gap larger than  $\epsilon$  in Algorithm 1 is upper bounded by

$$\begin{aligned} &\Gamma(1/\epsilon, \log(1/\delta); \gamma, d) \\ &= \tilde{O}\left(\frac{1}{(1-\gamma)^6 \epsilon^2} + \frac{d^2 + d \log(1/\delta)}{(1-\gamma)^4 \epsilon^2}\right), \end{aligned}$$

where  $d$  is the dimension of feature mapping and  $\gamma$  is the discount factor.

**Remark 5.2.** Similar to the main results in Zhou et al. (2021b), our sample complexity bound does not depend on the size of state space  $\mathcal{S}$  or that of action space  $\mathcal{A}$ . This suggests that our algorithm is applicable to MDPs with large state and action spaces.

**Corollary 5.3.** Under the same conditions in Theorem 5.1, with probability at least  $1 - \delta$ , the regret of Algorithm 1 is bounded as follows

$$\text{Regret}(T) = \tilde{O}\left(\frac{d + \sqrt{d \log(1/\delta)}}{(1-\gamma)^2} \sqrt{T} + \frac{\sqrt{T}}{(1-\gamma)^3}\right).$$

**Remark 5.4.** Corollary 5.3 suggests that UPAC-UCLK enjoys an  $\tilde{O}(d\sqrt{T}/(1-\gamma)^2 + \sqrt{T}/(1-\gamma)^3)$  regret. In comparison, the UCLK algorithm in Zhou et al. (2021b) enjoys an  $\tilde{O}(d\sqrt{T}/(1-\gamma)^2)$  regret. Our algorithm suffers an extra  $\tilde{O}(\sqrt{T}/(1-\gamma)^3)$  term in the regret. This originates from a bound on the difference between ML-EVI outputs for two consecutive time steps, and we leave it as a future work to remove this term.

## 6. Proof Sketch of the Main Results

In this section, we will present the proof sketch of the main results. The detailed proof can be found in the appendix.

We first present several central lemmas, and then give a short sketch outlining the proof of Theorem 5.1. For the ease of presentation, we introduce new notations for the variables in Algorithm 1: denote the variables  $L, \Sigma^l, \mathbf{b}^l, \hat{\theta}^l$  and  $\mathcal{C}^l$  after Line 6 of time step  $t$  by  $L_t, \Sigma_t^l, \mathbf{b}_t^l, \hat{\theta}_t^l$  and  $\mathcal{C}_t^l$ .

For simplicity, we define  $k_t^l := \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l\} 2^{l_\tau - l}$  for  $t \geq 1$  and  $l \geq 0$ , which represents the total amount of data in level  $l$  at time step  $t$ . With this notation, we have the following lemma upper bounding this value:

**Lemma 6.1.** Let  $l_0 := \log(\sqrt{d}/(1-\gamma))/\log 2$ . Consider  $k_{\max}^l := \lim_{t \rightarrow \infty} k_t^l$ . For each level  $1 \leq l \leq l_0$ , we have  $k_{\max}^l \leq 20dl_0$ , while for  $l > l_0$ ,  $k_{\max}^l \leq 20dl4^{l-l_0}$ .

Lemma 6.1 is parallel to Lemma A.1 from He et al. (2021b). It suggests that the amount of data in level  $l$  is finite with an upper bound exponential in  $l$ . Note that when  $l$  reaches beyond  $l_0$ , the threshold value at level  $l$  namely  $2^{-l_0}\sqrt{d}/(1-\gamma)$  shrinks to below constant 1, which is considered the beginning of “useful” levels in our analysis.

**Lemma 6.2.** For the factor  $\beta^l$  defined in (5.1), with probability at least  $1 - \delta/2$ , for all steps  $t$  and levels  $l$ , we have  $\|\theta^* - \hat{\theta}_t^l\|_{\Sigma_t^l} \leq \beta^l/2$ .

Lemma 6.2 suggests that with the given factor  $\beta^l$ , the true vector  $\theta^*$  lies within halved confidence sets in each level and time step. Here the radii of our confidence sets are set to be twice as large as normally defined. This is to ensure that  $\theta^*$  stays relatively close to the center of the confidence sets, and will enable our later bounds concerning the difference between confidence sets. For simplicity, we denote the event in Lemma 6.2 as  $\mathcal{E}_1$ .

**Lemma 6.3.** Under  $\mathcal{E}_1$ , for all  $t, l$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we have

$$\begin{aligned} \frac{1}{1-\gamma} &\geq Q_t^l(s, a) \geq Q^*(s, a), \\ \frac{1}{1-\gamma} &\geq V_t(s) \geq V^*(s). \end{aligned} \quad (6.1)$$

Lemma 6.3 suggests that the results of Algorithm 2 are upper bounds of the optimal value functions  $Q^*$  and  $V^*$ .

**Lemma 6.4.** Under  $\mathcal{E}_1$ , for all  $t$  and  $l$ , there exists a certain  $\theta_t^l \in \mathcal{C}_t^l \cap \mathcal{B}$ , so that

$$Q_t^l(s_t, a_t) \leq r(s_t, a_t) + \gamma \langle \theta_t^l, \phi_{V_t}(s_t, a_t) \rangle + \gamma^{U_t}. \quad (6.2)$$

Lemma 6.4 characterizes the transition error of Algorithm 2. Inequality (6.2) is an approximate version of the Bellman equation (3.2), and suggests that in order to reach an error of  $O(1/\epsilon)$ , a number of iteration rounds of order  $O(\log(1/\epsilon))$  is needed.

We are now ready to begin our proof. Our analysis of sample complexity follows a similar road map as in Lemma 1 of Dong et al. (2019). We take the summation of large suboptimality over all time steps, which is trivially lower bounded, and then upper bound the sum with similar methods as in the proof of Theorem 1 of Zhou et al. (2021b). This yields an inequality centered around the sample complexity, which leads to the eventual conclusion.

*Proof sketch of Theorem 5.1.* In order to bound the sample complexity  $N_\epsilon = \sum_{t=1}^\infty \mathbb{1}\{V^*(s_t) - V_t^\pi(s_t) > \epsilon\}$ , consider the set

$$T_\epsilon = \{t : V^*(s_t) - V_t^\pi(s_t) > \epsilon\},$$

we instantly have

$$\epsilon N_\epsilon < \sum_{t \in T_\epsilon} (V^*(s_t) - V_t^\pi(s_t)). \quad (6.3)$$

We then bound the right hand side of (6.3). By (6.1) from Lemma 6.3, the suboptimality gap  $V^*(s_t) - V_t^\pi(s_t)$  is upper bounded by  $V_t(s_t) - V_t^\pi(s_t)$ . Using (6.2) from Lemma 6.4 and the policy Bellman equation (3.1), we can decompose the suboptimality and eventually obtain

$$\begin{aligned} (V_t - V_t^\pi)(s_t) &\leq \Gamma(t) + \gamma \Delta(t) + \gamma(V_t - V_{t+1})(s_{t+1}) \\ &\quad + \gamma(V_{t+1} - V_{t+1}^\pi)(s_{t+1}), \end{aligned} \quad (6.4)$$

where

$$\begin{aligned} \Gamma(t) &= 1/(t(t+1)) + 3 \cdot 2^{l_0-l_t} \gamma \beta^{l_t-1}, \\ \Delta(t) &= \mathbb{P}(V_t - V_{t+1}^\pi)(s_t, a_t) - (V_t - V_{t+1}^\pi)(s_{t+1}). \end{aligned}$$

By iterating (6.4) for  $T$  rounds, we can bound the suboptimality gap by

$$(V^* - V_t^\pi)(s_t) \leq \frac{\gamma^T}{1-\gamma} + A_t + B_t + C_t, \quad (6.5)$$

where

$$\begin{aligned} A_t &= \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t} \Gamma(\tau), \\ B_t &= \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t+1} \Delta(\tau), \\ C_t &= \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t+1} (V_\tau - V_{\tau+1})(s_{\tau+1}). \end{aligned}$$

However, when the level  $l_t$  is small, the bound above is too large; in fact, it could go much larger than the trivial bound  $1/(1-\gamma)$ . We thus define the condition  $\mathcal{Q}_t = \{l_\tau > l_0, \forall t \leq \tau < t+T\}$ , and separate the right hand side of (6.3) into two parts accordingly:

$$\begin{aligned} &\sum_{t \in T_\epsilon} (V^*(s_t) - V_t^\pi(s_t)) \\ &\leq \sum_{t \in T_\epsilon} \mathbb{1}\{\overline{\mathcal{Q}}_t\} \frac{1}{1-\gamma} \\ &\quad + \sum_{t \in T_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} \left( \frac{\gamma^T}{1-\gamma} + A_t + B_t + C_t \right). \end{aligned}$$

Since  $\mathcal{Q}_t$  is false only for a small amount of  $t$ , the first term on the right hand side here is easy to bound. We then seek to bound the sums  $\sum_{t \in T_\epsilon} \mathbb{1}\{\mathcal{Q}_t\}A_t$ ,  $\sum_{t \in T_\epsilon} \mathbb{1}\{\mathcal{Q}_t\}B_t$  and  $\sum_{t \in T_\epsilon} \mathbb{1}\{\mathcal{Q}_t\}C_t$  separately. For the sum of  $A_t$ , the individual terms are already explicit expressions. For  $B_t$ , we view  $\Delta(t)$  as a martingale sequence and use the Azuma-Hoeffding inequality for a bound that holds with probability at least  $1 - \delta/2$ . Bounding the sum of  $C_t$  is the hardest and most novel section of our result.

The term  $C_t$  comes down to the difference between ML-EVI outputs at two consecutive time steps. Unlike the regret analysis for the UCLK algorithm (Zhou et al., 2021b), where one can trivially bound the sum of the difference by  $1/(1 - \gamma)$ , we introduce a novel method that decomposes the difference and expresses it with terms representing changes in confidence ellipsoids. The effectiveness of this bound relies on the discounted data inheritance scheme we proposed earlier in Section 4.2.

We backtrack the extended value iteration for this bound. Denote the variables  $Q_t^{(u)}$  and  $V^{(u)}$  in the run of Algorithm 2 at time step  $t$  as  $Q_{l,t}^{(u)}$  and  $V_t^{(u)}$ . Next define  $D_u := \max_{s \in \mathcal{S}} (V_t^{(U_t-u)} - V_{t+1}^{(U_{t+1}-u)})(s)$ . By expanding the value functions according to the iteration rules (specified in Lines 3 and 6 of Algorithm 2), we can prove the inequality

$$D_{u-1} \leq \gamma D_u + \gamma \max_{\phi_V, l} (\max_{\theta_1} \langle \theta_1, \phi_V \rangle - \max_{\theta_2} \langle \theta_2, \phi_V \rangle), \quad (6.6)$$

where  $\theta_1$  takes the maximum in  $\mathcal{C}_t^l \cap \mathcal{B}$ , and  $\theta_2$  from  $\mathcal{C}_{t+1}^l \cap \mathcal{B}$ . We can iterate (6.6) with respect to  $u$ , and the problem boils down to bounding the difference between confidence sets.

Suppose  $\tilde{\theta}_1 := \operatorname{argmax}_{\theta_1} \langle \theta_1, \phi_V \rangle$ . For the maximum concerning  $\theta_2$  in (6.6), we need to find a suitable  $\tilde{\theta}_2$  that can be used to lower bound  $\max_{\theta_2} \langle \theta_2, \phi_V \rangle$  while being comparable to  $\tilde{\theta}_1$ . Starting from  $\theta^*$ , we go in the direction of  $\tilde{\theta}_1$  until we reach the boundaries of  $\mathcal{C}_{t+1}^l \cap \mathcal{B}$  and name our ending point  $\tilde{\theta}_2$ . In other words, consider  $\mu_{\max} = \sup\{\mu : \theta^* + \mu(\tilde{\theta}_1 - \theta^*) \in \mathcal{C}_{t+1}^l \cap \mathcal{B}\}$ , and let  $\tilde{\theta}_2 = \theta^* + \mu_{\max}(\tilde{\theta}_1 - \theta^*)$ . We have

$$\begin{aligned} & \max_{\theta_1} \langle \theta_1, \phi_V \rangle - \max_{\theta_2} \langle \theta_2, \phi_V \rangle \\ & \leq (1 - \mu_{\max}) \langle \tilde{\theta}_1 - \theta^*, \phi_V \rangle. \end{aligned} \quad (6.7)$$

Next, using the fact that  $\tilde{\theta}_2$  belongs to the border of  $\mathcal{C}_{t+1}^l \cap \mathcal{B}$ , we can obtain an upper bound for  $1 - \mu_{\max}$ . Intuitively, when  $\tilde{\theta}_2 \in \partial \mathcal{C}_{t+1}^l$ , we have  $\|\tilde{\theta}_2 - \hat{\theta}_t^l\|_{\Sigma_t^l} \approx \|\tilde{\theta}_2 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} = \beta^l$ , and since  $\|\tilde{\theta}_1 - \hat{\theta}_t^l\|_{\Sigma_t^l} = \beta^l$ , we further have  $\tilde{\theta}_1 \approx \tilde{\theta}_2$  and  $\mu_{\max} \approx 1$ . Combining this with (6.7) and (6.6), we

obtain an upper bound for  $D_0 = \max_{s \in \mathcal{S}} (V_t - V_{t+1})(s)$ , and consequently a bound for  $C_t$ .

Putting the bounds for  $A_t$ ,  $B_t$  and  $C_t$  together, we ultimately acquire an inequality of the form

$$\epsilon N_\epsilon \leq C_1 N_\epsilon + C_2 \sqrt{N_\epsilon} + C_3,$$

where  $C_1$  can be arbitrarily small with proper parameter choices. Solving this inequality with  $N_\epsilon$  as the main variable, we can finally obtain an upper bound on  $N_\epsilon$  that holds with probability at least  $1 - \delta/2 - \delta/2 = 1 - \delta$ .  $\square$

## 7. Conclusions

We proposed a new algorithm UPAC-UCLK for learning linear kernel MDPs and proved that it is uniform-PAC with sample complexity  $\tilde{O}(d^2/((1 - \gamma)^4 \epsilon^2) + 1/((1 - \gamma)^6 \epsilon^2))$ . This is the first algorithm with uniform PAC guarantee in the discounted MDPs setting. Our result shows that the optimal policy in linear kernel MDPs can be learned efficiently. Our sample complexity bound has a term  $\tilde{O}(1/((1 - \gamma)^6 \epsilon^2))$  that stems from the ML-EVI in our analysis. Achieving minimax optimal sample complexity bound is left as an open question for future work.

## Acknowledgements

We thank the anonymous reviewers for their helpful comments. JH and QG are partially supported by the National Science Foundation CAREER Award 1906169. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.

## References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Agrawal, S. and Jia, R. Posterior sampling for reinforcement learning: worst-case regret bounds. In *Advances in Neural Information Processing Systems*, pp. 1184–1194, 2017.
- Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. F. Model-based reinforcement learning with value-targeted regression. *arXiv preprint arXiv:2006.01107*, 2020.
- Azar, M. G., Osband, I., and Munos, R. Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 263–272. JMLR. org, 2017.



- Bartlett, P. L. and Tewari, A. Regal: A regularization based algorithm for reinforcement learning in weakly communicating mdps. *arXiv preprint arXiv:1205.2661*, 2012.
- Cesa-Bianchi, N. and Lugosi, G. *Prediction, learning, and games*. Cambridge university press, 2006.
- Dann, C., Lattimore, T., and Brunskill, E. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5713–5723, 2017.
- Dong, K., Wang, Y., Chen, X., and Wang, L. Q-learning with ucb exploration is sample efficient for infinite-horizon mdp. In *International Conference on Learning Representations*, 2019.
- Fruit, R., Pirotta, M., and Lazaric, A. Near optimal exploration-exploitation in non-communicating markov decision processes. *arXiv preprint arXiv:1807.02373*, 2018a.
- Fruit, R., Pirotta, M., Lazaric, A., and Ortner, R. Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *International Conference on Machine Learning*, pp. 1578–1586. PMLR, 2018b.
- Fruit, R., Pirotta, M., and Lazaric, A. Improved analysis of ucl2 with empirical bernstein inequality. *arXiv preprint arXiv:2007.05456*, 2020.
- He, J., Zhou, D., and Gu, Q. Nearly minimax optimal reinforcement learning for discounted mdps. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021a.
- He, J., Zhou, D., and Gu, Q. Uniform-PAC bounds for reinforcement learning with linear function approximation. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021b.
- Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. Is q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pp. 4868–4878, 2018.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pp. 2137–2143, 2020.
- Kakade, S. M. et al. *On the sample complexity of reinforcement learning*. PhD thesis, University of London London, England, 2003.
- Lattimore, T. and Hutter, M. Pac bounds for discounted mdps. In *International Conference on Algorithmic Learning Theory*, pp. 320–334. Springer, 2012.
- Liu, S. and Su, H. Regret bounds for discounted mdps, 2020.
- Modi, A., Jiang, N., Tewari, A., and Singh, S. Sample complexity of reinforcement learning using linearly combined model ensembles. *arXiv preprint arXiv:1910.10597*, 2019.
- Ortner, R. Regret bounds for reinforcement learning via markov chain concentration. *Journal of Artificial Intelligence Research*, 67:115–128, 2020.
- Ouyang, Y., Gagrani, M., Nayyar, A., and Jain, R. Learning unknown markov decision processes: A thompson sampling approach. *arXiv preprint arXiv:1709.04570*, 2017.
- Strehl, A. L., Li, L., Wiewiora, E., Langford, J., and Littman, M. L. Pac model-free reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pp. 881–888. ACM, 2006.
- Szita, I. and Szepesvári, C. Model-based reinforcement learning with nearly tight exploration complexity bounds. In *ICML*, pp. 1031–1038, 2010.
- Talebi, M. S. and Maillard, O.-A. Variance-aware regret bounds for undiscounted reinforcement learning in mdps. In *Algorithmic Learning Theory*, pp. 770–805. PMLR, 2018.
- Wei, C.-Y., Jahromi, M. J., Luo, H., and Jain, R. Learning infinite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pp. 3007–3015. PMLR, 2021.
- Wu, Y., Zhou, D., and Gu, Q. Nearly minimax optimal regret for learning infinite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, 2022.
- Zanette, A., Lazaric, A., Kochenderfer, M., and Brunskill, E. Learning near optimal policies with low inherent bellman error. *arXiv preprint arXiv:2003.00153*, 2020.
- Zhang, Z. and Ji, X. Regret minimization for reinforcement learning by evaluating the optimal bias function. *arXiv preprint arXiv:1906.05110*, 2019.
- Zhang, Z., Zhou, Y., and Ji, X. Almost optimal model-free reinforcement learning via reference-advantage decomposition. *arXiv preprint arXiv:2004.10019*, 2020a.

Zhang, Z., Zhou, Y., and Ji, X. Model-free reinforcement learning: from clipped pseudo-regret to sample complexity. *arXiv preprint arXiv:2006.03864*, 2020b.

Zhou, D., Gu, Q., and Szepesvari, C. Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pp. 4532–4576. PMLR, 2021a.

Zhou, D., He, J., and Gu, Q. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*. PMLR, 2021b.

## A. Proof of the Main Results

We first give a few topological definitions. Under universal set  $\mathcal{U} \subset \mathbb{R}^d$ , for  $x \in \mathbb{R}^d$ , define  $B_{\mathcal{U}}(x, r) := \{y \in \mathcal{U} : \|y - x\|_2 < r\}$ . For  $\mathcal{X} \subset \mathcal{U}$ , a point  $x \in \mathcal{X}$  is its inner point (when constrained to  $\mathcal{U}$ ) if  $\exists r > 0, B_{\mathcal{U}}(x, r) \in \mathcal{X}$ ; a point  $x \in \mathcal{U}$  is its boundary point (when constrained to  $\mathcal{U}$ ) if  $\forall r > 0, B_{\mathcal{U}}(x, r) \cap \mathcal{X} \neq \emptyset, B_{\mathcal{U}}(x, r) \cap \mathcal{X}^C \neq \emptyset$ . Denote the boundary of  $\mathcal{X}$  as  $\partial_{\mathcal{U}}\mathcal{X}$ , which is defined as the set of all boundary points of  $\mathcal{X}$  (when constrained to  $\mathcal{U}$ ). If  $\mathcal{U} = \mathbb{R}^d$ , we will not mention  $\mathcal{U}$  when using these terminologies, and we define the corresponding notations in this scenario as  $B(x, r) := B_{\mathbb{R}^d}(x, r), \partial := \partial_{\mathbb{R}^d}$ .

An affine subspace in  $\mathbb{R}^d$  is defined as a subset of  $\mathbb{R}^d$  of the form  $\{x \in \mathbb{R}^d : L_i(x) = \alpha_i, i = 1, \dots, d_0\}$ , where  $d_0 \leq d$  is a positive integer,  $L_1, \dots, L_{d_0} : \mathbb{R}^d \rightarrow \mathbb{R}$  are mutually independent linear functions of  $x \in \mathbb{R}^d$ , and  $\alpha_1, \dots, \alpha_{d_0} \in \mathbb{R}$ .

Before diving into the main proof, we provide a few more lemmas in aid of our proof.

The first lemma sets up alternative definitions for the variables  $\Sigma_t^l$  and  $\mathbf{b}_t^l$ , which clarify their utilities and drop the restraint  $l \leq 4L_t + C_0$  between  $l$  and  $t$ .

**Lemma A.1.** For each time step  $t$  and level  $l$ , define the alternative and extended expressions for the variables  $\Sigma_t^l := \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbb{1}\{l_{\tau} \leq l\} 2^{l_{\tau}-l} \phi_{V_{\tau}}(s_{\tau}.a_{\tau}) \phi_{V_{\tau}}(s_{\tau}.a_{\tau})^{\top}$  and  $\mathbf{b}_t^l := \sum_{\tau=1}^{t-1} \mathbb{1}\{l_{\tau} \leq l\} 2^{l_{\tau}-l} \phi_{V_{\tau}}(s_{\tau}.a_{\tau}) V_{\tau}(s_{\tau+1})$ . Then for  $l \leq 4L_t + C_0$ , the definitions given here yield the same values as the variables given in Algorithm 1.

From now on, for the variables  $\Sigma_t^l$  and  $\mathbf{b}_t^l$ , we use the definitions given in Lemma A.1. We also define  $\widehat{\theta}_t^l := (\Sigma_t^l)^{-1} \mathbf{b}_t^l$ ,  $\mathcal{C}_t^l := \{\theta \in \mathbb{R}^d : \|\theta - \widehat{\theta}_t^l\|_{\Sigma_t^l} \leq \beta^l\}$  for arbitrary  $t, l$ .

The next lemma does the same thing for the ML-EVI algorithm. More specifically, it shows that one can extend the total number of levels in Algorithm 2 to infinity and obtain essentially the same results.

**Lemma A.2.** For each time step  $t$ , consider an alternative run of the ML-EVI algorithm with  $L = \infty, \mathcal{C}_l = \mathcal{C}_t^l$  for  $l = 1, 2, \dots$  and  $U = U_t$ . We compare this newly constructed run to the canonical run at time step  $t$ .

To do so, we write the variables in this run with an additional suffix  $*$ , resulting in the notations  $Q_{l,*}^{(u)}, V_*^{(u)}$  for iteration round  $u = 0, 1, \dots, U_t$ . Then we have

$$\begin{aligned} Q_{l,*}^{(u)} &= Q_l^{(u)}, \forall 1 \leq l \leq 4L_t + C_0, \forall 0 \leq u \leq U_t, \\ Q_{l,*}^{(u)} &\geq Q_{4L_t+C_0}^{(u)}, \forall l > 4L_t + C_0, \forall 0 \leq u \leq U_t, \\ V_*^{(u)} &= V^{(u)}, \forall 0 \leq u \leq U_t. \end{aligned}$$

We will use this alternative run of ML-EVI instead of the canonical run in future analysis. Thus, we are completely rid of the constraint  $l \leq 4L_t + C_0$ .

**Lemma A.3.** The set of  $\theta \in \mathbb{R}^d$  that satisfies the conditions for linear kernel MDPs in Definition 3.1, denoted by  $\mathcal{B}$ , lies within an affine subspace  $\mathcal{L} \subset \mathbb{R}^d$ ; furthermore, it is a convex body (closed and bounded with non-empty interior) when constrained to  $\mathcal{L}$ .

This lemma describes the general shape of  $\mathcal{B}$  and prepares us for future geometrical analysis on the parameter space.

### A.1. Proof of Theorem 5.1

We are now ready to prove our main theorem in detail.

*Proof of Theorem 5.1.* We will work under the event  $\mathcal{E}_1$  in this proof, which holds with probability at least  $1 - \delta/2$  according to the definition of  $\mathcal{E}_1$  and Lemma 6.2. For convenience, our proof is segmented into four parts.

#### Part I: Sample Complexity Analysis

In the first part, we give an iterable upper bound on the suboptimality gap, and deduce from it an inequality for the sample complexity  $N_{\epsilon}$ . Without losing any generality, we assume  $\epsilon \leq 1/(1 - \gamma)$ .

We start by analyzing the suboptimality gap. From Lemma 6.3, we have  $V^*(s_t) - V_t^\pi(s_t) \leq V_t(s_t) - V_t^\pi(s_t)$ . In the following, we seek to deduce an upper bound for  $(V_t - V_t^\pi)(s_t)$  that contains the next term  $(V_{t+1} - V_{t+1}^\pi)(s_{t+1})$ , so we can iterate this bound over  $t$ .

As a first step, notice that according to the definitions of value functions for policy  $\pi$ , we have  $V_t^\pi(s_t) = Q_t^\pi(s_t, a_t)$ ; according to Line 12 of Algorithm 2 and Line 8 of Algorithm 1, we have  $V_t(s_t) = \max_{a \in \mathcal{A}} \min_{l \geq 1} Q_t^l(s_t, a) = \min_{l \geq 1} Q_t^l(s_t, a_t)$ . Combining these equations with Lemma 6.4 and the Bellman equation (3.2), we can deduce

$$\begin{aligned} V_t(s_t) - V_t^\pi(s_t) &= \min_{l \geq 1} Q_t^l(s_t, a_t) - Q_t^\pi(s_t, a_t) \\ &\leq Q_t^{l_t-1}(s_t, a_t) - Q_t^\pi(s_t, a_t) \\ &\leq r(s_t, a_t) + \gamma \langle \theta_t^{l_t-1}, \phi_{V_t}(s_t, a_t) \rangle + \gamma^{U_t} - r(s_t, a_t) - \gamma \langle \theta^*, \phi_{V_{t+1}^\pi}(s_t, a_t) \rangle \\ &= \gamma^{U_t} + \gamma \left[ \langle \theta_t^{l_t-1}, \phi_{V_t}(s_t, a_t) \rangle - \langle \theta^*, \phi_{V_{t+1}^\pi}(s_t, a_t) \rangle \right]. \end{aligned} \quad (\text{A.1})$$

Now consider the bracketed part of (A.1), we decompose the difference into two parts, namely

$$\langle \theta_t^{l_t-1}, \phi_{V_t}(s_t, a_t) \rangle - \langle \theta^*, \phi_{V_{t+1}^\pi}(s_t, a_t) \rangle = \underbrace{\langle \theta_t^{l_t-1} - \theta^*, \phi_{V_t}(s_t, a_t) \rangle}_{I_1} + \underbrace{\langle \theta^*, \phi_{V_t}(s_t, a_t) - \phi_{V_{t+1}^\pi}(s_t, a_t) \rangle}_{I_2}. \quad (\text{A.2})$$

For  $I_1$ , under the event  $\mathcal{E}_1$ , we have from Lemma 6.2 that for all  $l$ ,  $\|\theta^* - \hat{\theta}_t^l\|_{\Sigma_t^l} \leq \beta^l/2$ ; also, from the condition  $\theta_t^l \in \mathcal{C}_t^l \cap \mathcal{B}$  in Lemma 6.4, we have  $\|\theta_t^l - \hat{\theta}_t^l\|_{\Sigma_t^l} \leq \beta^l$ . On the other hand, from the selection rule of  $l_t$  specified in Lines 10 to 13 of Algorithm 1, inequality  $\|\phi_{V_t}(s_t, a_t)\|_{(\Sigma_t^{l_t-1})^{-1}} \leq 2^{-(l_t-1)}\sqrt{d}/(1-\gamma)$  holds. Together with the Cauchy-Schwartz Inequality  $\|x\|_{\Sigma} \cdot \|y\|_{\Sigma^{-1}} \geq \langle x, y \rangle$ , we have that

$$\begin{aligned} I_1 &= \langle \theta_t^{l_t-1} - \theta^*, \phi_{V_t}(s_t, a_t) \rangle \\ &\leq \|\theta_t^{l_t-1} - \theta^*\|_{\Sigma_t^{l_t-1}} \|\phi_{V_t}(s_t, a_t)\|_{(\Sigma_t^{l_t-1})^{-1}} \\ &\leq (\|\theta_t^{l_t-1} - \hat{\theta}_t^{l_t-1}\|_{\Sigma_t^{l_t-1}} + \|\hat{\theta}_t^{l_t-1} - \theta^*\|_{\Sigma_t^{l_t-1}}) \cdot \|\phi_{V_t}(s_t, a_t)\|_{(\Sigma_t^{l_t-1})^{-1}} \\ &\leq (\beta^{l_t-1} + \frac{1}{2}\beta^{l_t-1}) \cdot 2^{-(l_t-1)}\sqrt{d}/(1-\gamma) \\ &= 3 \cdot 2^{l_0-l_t} \beta^{l_t-1}, \end{aligned} \quad (\text{A.3})$$

where in the second inequality we used the triangle inequality, and in the final equation we used the definition of  $l_0$  from Lemma 6.1.

For  $I_2$ , recall from (3.3) that  $\langle \theta^*, \phi_V(s_t, a_t) \rangle = \mathbb{P}V(s_t, a_t)$  for arbitrary  $V$ , therefore

$$\begin{aligned} I_2 &= \langle \theta^*, \phi_{V_t}(s_t, a_t) - \phi_{V_{t+1}^\pi}(s_t, a_t) \rangle \\ &= \mathbb{P}(V_t - V_{t+1}^\pi)(s_t, a_t). \end{aligned} \quad (\text{A.4})$$

Additionally, from the definition of  $U_t$  in (5.1), we have

$$\begin{aligned} \gamma^{U_t} &\leq \gamma^{\log[t(t+1)]/(1-\gamma)} \\ &= \left[ (1 - (1-\gamma))^{1/(1-\gamma)} \right]^{\log[t(t+1)]} \\ &\leq [e^{-1}]^{\log[t(t+1)]} \\ &= \frac{1}{t(t+1)}. \end{aligned} \quad (\text{A.5})$$

Now substitute (A.3) and (A.4) into (A.2) and then into (A.1) along with the above inequality, to obtain

$$\begin{aligned}
 V_t(s_t) - V_t^\pi(s_t) &\leq \gamma^{U_t} + \gamma[3 \cdot 2^{l_0 - l_t} \beta^{l_t - 1} + \mathbb{P}(V_t - V_{t+1}^\pi)(s_t, a_t)] \\
 &\leq \frac{1}{t(t+1)} + 3\gamma \cdot 2^{l_0 - l_t} \beta^{l_t - 1} + \gamma[\mathbb{P}(V_t - V_{t+1}^\pi)(s_t, a_t) - (V_t - V_{t+1}^\pi)(s_{t+1})] + \gamma(V_t - V_{t+1}^\pi)(s_{t+1}) \\
 &= \Gamma(t) + \gamma\Delta(t) + \gamma(V_t - V_{t+1}^\pi)(s_{t+1}) + \gamma(V_{t+1} - V_{t+1}^\pi)(s_{t+1}), \tag{A.6}
 \end{aligned}$$

where in the final equation we define  $\Gamma(t) := 1/(t(t+1)) + 3\gamma \cdot 2^{l_0 - l_t} \beta^{l_t - 1}$  and  $\Delta(t) := \mathbb{P}(V_t - V_{t+1}^\pi)(s_t, a_t) - (V_t - V_{t+1}^\pi)(s_{t+1})$ .

Iterate Inequality (A.6) for  $T$  rounds (through  $t, t+1, \dots, t+T$ ), and we get the following suboptimality bound:

$$\begin{aligned}
 V^*(s_t) - V_t^\pi(s_t) &\leq V_t(s_t) - V_t^\pi(s_t) \\
 &\leq \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t} [\Gamma(\tau) + \gamma\Delta(\tau) + \gamma(V_\tau - V_{\tau+1}^\pi)(s_{\tau+1})] + \gamma^T (V_{t+T} - V_{t+T}^\pi)(s_{t+T}) \\
 &\leq \underbrace{\sum_{\tau=t}^{t+T-1} \gamma^{\tau-t} \Gamma(\tau)}_{A_t} + \underbrace{\sum_{\tau=t}^{t+T-1} \gamma^{\tau-t+1} \Delta(\tau)}_{B_t} + \underbrace{\sum_{\tau=t}^{t+T-1} \gamma^{\tau-t+1} (V_\tau - V_{\tau+1}^\pi)(s_{\tau+1})}_{C_t} + \frac{\gamma^T}{1-\gamma}, \tag{A.7}
 \end{aligned}$$

where the last inequality is due to  $0 \leq V_{t+T}^\pi(s_{t+T}), V_{t+T}(s_{t+T}) \leq 1/(1-\gamma)$ .

Now suppose the set of all time steps where the suboptimality gap is greater than  $\epsilon$  is  $\mathcal{T}_\epsilon := \{t : V^*(s_t) - V_t^\pi(s_t) > \epsilon\}$ , with a total number of  $N_\epsilon = \#\mathcal{T}_\epsilon$  elements. Define the condition  $\mathcal{Q}_t = \{l_\tau > l_0, \forall t \leq \tau < t+T\}$ . Since we also have trivially that  $V^*(s_t) - V_t^\pi(s_t) \leq 1/(1-\gamma)$ , we can obtain by summing (A.7) over all time steps  $t$  in set  $\mathcal{T}_\epsilon$  that

$$\begin{aligned}
 \epsilon N_\epsilon &\leq \sum_{t \in \mathcal{T}_\epsilon} (V^*(s_t) - V_t^\pi(s_t)) \\
 &\leq \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} \left( \frac{\gamma^T}{1-\gamma} + A_t + B_t + C_t \right) + \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\overline{\mathcal{Q}_t}\} \frac{1}{1-\gamma} \\
 &\leq \frac{\gamma^T}{1-\gamma} N_\epsilon + \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} (A_t + B_t + C_t) + \frac{1}{1-\gamma} \sum_{t=1}^{\infty} \mathbb{1}\{\overline{\mathcal{Q}_t}\}.
 \end{aligned}$$

Consider now the indicator function for  $\overline{\mathcal{Q}_t}$ . Only when  $l_\tau \leq l_0$  holds for some  $\tau \in \{t, \dots, t+T-1\}$  is  $\mathcal{Q}_t$  false, therefore  $\mathbb{1}\{\overline{\mathcal{Q}_t}\} \leq \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau \leq l_0\} = \sum_{l=1}^{\lfloor l_0 \rfloor} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau = l\}$ . Combine this with the definition of  $k_{\max}^l = \lim_{t \rightarrow \infty} k_t^l = \sum_{\tau=1}^{\infty} \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l}$  in Lemma 6.1, the last addend above can be bounded by

$$\begin{aligned}
 \frac{1}{1-\gamma} \sum_{t=1}^{\infty} \mathbb{1}\{\overline{\mathcal{Q}_t}\} &\leq \frac{1}{1-\gamma} \sum_{t=1}^{\infty} \sum_{l=1}^{\lfloor l_0 \rfloor} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau = l\} \\
 &\leq \frac{T}{1-\gamma} \sum_{l=1}^{\lfloor l_0 \rfloor} \sum_{\tau=1}^{\infty} \mathbb{1}\{l_\tau = l\} \\
 &= \frac{T}{1-\gamma} \sum_{l=1}^{\lfloor l_0 \rfloor} \left( k_{\max}^l - \frac{1}{2} k_{\max}^{l-1} \right) \\
 &= \frac{T}{1-\gamma} \left( k_{\max}^{\lfloor l_0 \rfloor} + \frac{1}{2} \sum_{l=1}^{\lfloor l_0 \rfloor - 1} k_{\max}^l \right) \\
 &\leq \frac{T}{1-\gamma} \cdot \frac{l_0 + 1}{2} \cdot 20dl_0 \\
 &= \frac{10dT}{1-\gamma} l_0(l_0 + 1),
 \end{aligned}$$

where the second inequality is based on the observation that each  $\tau$  is visited in the last summation at most  $T$  times, and the third inequality is according to the conclusions of Lemma 6.1. Hence we have the inequality

$$\epsilon N_\epsilon \leq \frac{\gamma^T}{1-\gamma} N_\epsilon + \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} A_t + \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} B_t + \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} C_t + \frac{10dT}{1-\gamma} l_0(l_0+1). \quad (\text{A.8})$$

## Part II: Upper Bounds for the Summation of $A_t$ and $B_t$

For this next part, we bound the second and third terms on the right hand side of (A.8). First, for the sum of the term  $A_t$ , we expand the expression to the following:

$$\begin{aligned} \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} A_t &= \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{\mathcal{Q}_t\} \gamma^{\tau-t} \Gamma(\tau) \\ &\leq \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t} \Gamma(\tau) \\ &= \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t} \left[ \frac{1}{\tau(\tau+1)} + 3\gamma \cdot 2^{l_0-l_\tau} \beta^{l_\tau-1} \right] \\ &\leq \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \frac{\gamma^{\tau-t}}{\tau(\tau+1)} + 3 \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t+1} 2^{l_0-l_\tau} \beta^{l_\tau-1}. \end{aligned} \quad (\text{A.9})$$

Regarding the first summation in (A.9), we shrink the  $\tau$ 's in the denominator to  $t$ :

$$\begin{aligned} \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \frac{\gamma^{\tau-t}}{\tau(\tau+1)} &\leq \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \frac{\gamma^{\tau-t}}{t(t+1)} \\ &= \sum_{t \in \mathcal{T}_\epsilon} \frac{1}{t(t+1)} \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t} \\ &\leq \sum_{t \in \mathcal{T}_\epsilon} \frac{1}{t(t+1)} \cdot \frac{1}{1-\gamma} \\ &\leq \frac{1}{1-\gamma}. \end{aligned} \quad (\text{A.10})$$

We hence obtain by substituting this into (A.9) that

$$\sum_{t \in \mathcal{T}_\epsilon} A_t \leq \frac{1}{1-\gamma} + 3 \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t+1} 2^{l_0-l_\tau} \beta^{l_\tau-1}. \quad (\text{A.11})$$

We leave the leftover part to be tackled in Part IV.

Second, we discard the indicator term in the sum of  $B_t$  and expand it to

$$\begin{aligned} \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} B_t &\leq \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t+1} \Delta(\tau) \\ &= \sum_{t \in \mathcal{T}_\epsilon} \sum_{n=0}^{T-1} \gamma^{n+1} \Delta(t+n) \\ &= \sum_{n=0}^{T-1} \left[ \sum_{t \in \mathcal{T}_\epsilon} \gamma^{n+1} \left[ \mathbb{P}(V_{t+n} - V_{t+n+1}^\pi)(s_{t+n}, a_{t+n}) - (V_{t+n} - V_{t+n+1}^\pi)(s_{t+n+1}) \right] \right], \end{aligned}$$

where in the first equation we substitute  $n = \tau - t$ , and in the second we exchange the ordering of summation.

Write  $\mathcal{T}_\epsilon = \{t_1, t_2, \dots, t_{N_\epsilon}\}$ , then for  $k = 1, \dots, N_\epsilon$ , the random variable  $t_k$  is the  $k$ -th time step where the suboptimality gap is larger than  $\epsilon$  and hence a stopping time. Now for each individual  $n \geq 0$ , let  $\mathcal{F}_{k,n}$  be the  $\sigma$ -field generated by all random variables before the generation of  $s_{t_k+n+1}$  in Line 9 of Algorithm 1 at time step  $t_k + n$ . Then

$$\left\{ \mathbb{P}(V_{t_k+n} - V_{t_k+n+1}^\pi)(s_{t_k+n}, a_{t_k+n}) - (V_{t_k+n} - V_{t_k+n+1}^\pi)(s_{t_k+n+1}) \right\}_{k=1}^{N_\epsilon}$$

is a martingale sequence with respect to the filtration  $\{\mathcal{F}_{k,n}\}_{k=0}^{N_\epsilon}$ . Now notice the term  $(V_{t+n} - V_{t+n+1}^\pi)(s)$  is bounded within  $[-1/(1-\gamma), 1/(1-\gamma)]$ , and therefore  $\mathbb{P}(V_{t+n} - V_{t+n+1}^\pi)(s_{t+n}, a_{t+n}) - (V_{t+n} - V_{t+n+1}^\pi)(s_{t+n+1})$  is bounded within  $[-2/(1-\gamma), 2/(1-\gamma)]$ . Consequently, from Lemma D.1, with probability at least  $1 - \delta/(2(n+1)(n+2))$ , the following holds:

$$\sum_{t \in \mathcal{T}_\epsilon} [\mathbb{P}(V_{t+n} - V_{t+n+1}^\pi)(s_{t+n}, a_{t+n}) - (V_{t+n} - V_{t+n+1}^\pi)(s_{t+n+1})] \leq \frac{2}{1-\gamma} \sqrt{2N_\epsilon \log(2(n+1)(n+2)/\delta)}.$$

Taking a union bound, we have that with probability at least  $1 - \sum_{n=0}^{\infty} \delta/(2(n+1)(n+2)) = 1 - \delta/2$ , the subsequent upper bound holds:

$$\begin{aligned} \sum_{t \in \mathcal{T}_\epsilon} B_t &= \sum_{n=0}^{T-1} \left[ \gamma^{n+1} \sum_{t \in \mathcal{T}_\epsilon} [\mathbb{P}(V_{t+n} - V_{t+n+1}^\pi)(s_{t+n}, a_{t+n}) - (V_{t+n} - V_{t+n+1}^\pi)(s_{t+n+1})] \right] \\ &\leq \sum_{n=0}^{T-1} \gamma^{n+1} \cdot \frac{2}{1-\gamma} \sqrt{2N_\epsilon \log(2(n+1)(n+2)/\delta)} \\ &\leq \sum_{n=0}^{T-1} \frac{2\gamma^{n+1}}{1-\gamma} \sqrt{2N_\epsilon \log(2T(T+1)/\delta)} \\ &\leq \frac{2\gamma}{(1-\gamma)^2} \sqrt{2N_\epsilon \log(2T(T+1)/\delta)}, \end{aligned} \tag{A.12}$$

where in the second inequality, we amplified the  $n+1$  under the square roots to  $T$ . We denote the event where this inequality holds as  $\mathcal{E}_2$ , then with probability at least  $1 - \delta/2 - \delta/2 = 1 - \delta$ , the joint event  $\mathcal{E} = \mathcal{E}_1 \cup \mathcal{E}_2$  holds.

### Part III: Upper bound for the Summation of $C_t$

For the third part, we provide an upper bound for  $\sum_{t \in \mathcal{T}_\epsilon} C_t$ . Since this process is rather lengthy, we partition it into four steps.

**Step 1**, we seek to obtain a bound for  $(V_t - V_{t+1})(s_{t+1})$ . According to Line 7 of Algorithm 1, the two value functions  $V_t, V_{t+1}$  are results of Algorithm 2 from two consecutive time steps  $t$  and  $t+1$ .

Denote the variables  $Q_l^{(u)}, V^{(u)}$  in Algorithm 2 at time step  $t$  by  $Q_{l,t}^{(u)}, V_t^{(u)}$  for iteration round  $u = 1, \dots, U_t$  and  $t = 1, 2, \dots$ . Next for a fixed  $t$ , define  $D_u := \max_{s \in \mathcal{S}} [V_t^{(U_t-u)} - V_{t+1}^{(U_{t+1}-u)}](s)$ . By the update rules of  $V$  and  $Q$  in

respectively Line 3 and Line 6 of Algorithm 2, we now give the following derivations:

$$\begin{aligned}
 D_{u-1} &= \max_{s \in \mathcal{S}} \left[ V_t^{(U_t-u+1)} - V_{t+1}^{(U_{t+1}-u+1)} \right] (s) \\
 &= \max_{s \in \mathcal{S}} \left[ \max_{a \in \mathcal{A}} \min_{l \geq 1} Q_{l,t}^{(U_t-u+1)}(s, a) - \max_{a \in \mathcal{A}} \min_{l \geq 1} Q_{l,t+1}^{(U_{t+1}-u+1)}(s, a) \right] \\
 &\leq \max_{s \in \mathcal{S}, a \in \mathcal{A}} \left[ \min_{l \geq 1} Q_{l,t}^{(U_t-u+1)} - \min_{l \geq 1} Q_{l,t+1}^{(U_{t+1}-u+1)} \right] (s, a) \\
 &\leq \max_{s \in \mathcal{S}, a \in \mathcal{A}, l \geq 1} \left[ Q_{l,t}^{(U_t-u+1)} - Q_{l,t+1}^{(U_{t+1}-u+1)} \right] (s, a) \\
 &= \max_{s, a, l} \left[ \left( r(s, a) + \gamma \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_{V_t^{(U_t-u)}}(s, a) \rangle \right) - \left( r(s, a) + \gamma \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_{V_{t+1}^{(U_{t+1}-u)}}(s, a) \rangle \right) \right] \\
 &= \gamma \max_{s, a, l} \left[ \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_{V_t^{(U_t-u)}}(s, a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_{V_{t+1}^{(U_{t+1}-u)}}(s, a) \rangle \right], \tag{A.13}
 \end{aligned}$$

where we used the trivial facts  $\max_{x \in \mathcal{X}} F(x) - \max_{x \in \mathcal{X}} G(x) \leq \max_{x \in \mathcal{X}} (F - G)(x)$  and  $\min_{x \in \mathcal{X}} F(x) - \min_{x \in \mathcal{X}} G(x) \leq \max_{x \in \mathcal{X}} (F - G)(x)$  in the inequalities above.

Next we denote for ease of expression that  $V_1 = V_t^{(U_t-u)}$ ,  $V_2 = V_{t+1}^{(U_{t+1}-u)}$ , and then separate the difference inside the outer maximum in (A.13) into two parts:

$$\begin{aligned}
 &\max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_{V_1}(s, a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_{V_2}(s, a) \rangle \\
 &\leq \underbrace{\max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_{V_1}(s, a) \rangle - \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_{V_2}(s, a) \rangle}_{J_1} + \underbrace{\max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_{V_2}(s, a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_{V_2}(s, a) \rangle}_{J_2}.
 \end{aligned}$$

Recall the definition of  $\phi_V$  from Definition 3.1, the first part

$$\begin{aligned}
 J_1 &\leq \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, (\phi_{V_1} - \phi_{V_2})(s, a) \rangle \\
 &= \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \int_{s' \in \mathcal{S}} \phi(s'|s, a)(V_1 - V_2)(s') ds' \rangle \\
 &= \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \int_{s' \in \mathcal{S}} \langle \theta_1, \phi(s'|s, a) \rangle (V_1 - V_2)(s') ds'. \tag{A.14}
 \end{aligned}$$

By the definition of  $\mathcal{B}$  in (3.4), for  $\theta \in \mathcal{B}$ ,  $\langle \theta, \phi(s'|s, a) \rangle$  is a probability distribution on  $\mathcal{S}$ , which we now denote by  $\tilde{\mathbb{P}} = \langle \theta, \phi \rangle$ . We then continue from (A.14):

$$\begin{aligned}
 J_1 &\leq \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}, \tilde{\mathbb{P}} = \langle \theta_1, \phi \rangle} \int_{s' \in \mathcal{S}} \tilde{\mathbb{P}}(s'|s, a)(V_1 - V_2)(s') ds' \\
 &\leq \max_{\theta_1 \in \mathcal{B}, \tilde{\mathbb{P}} = \langle \theta_1, \phi \rangle} \mathbb{E}_{s' \sim \tilde{\mathbb{P}}(\cdot|s, a)} [(V_1 - V_2)(s') | s, a] \\
 &\leq \max_{\theta_1 \in \mathcal{B}, \tilde{\mathbb{P}} = \langle \theta_1, \phi \rangle} \max_{s' \in \mathcal{S}} (V_1 - V_2)(s') \\
 &= \max_{s' \in \mathcal{S}} \left[ V_t^{(U_t-u)} - V_{t+1}^{(U_{t+1}-u)} \right] (s'), \tag{A.15}
 \end{aligned}$$

where in the second inequality we relaxed  $\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}$  to  $\theta_1 \in \mathcal{B}$ , and the third inequality is due to the fact that the expected value of a random variable is no greater than its maximum possible value. After that, the equation follows because the variable  $\theta_1$  is no longer tied to the value of  $(V_1 - V_2)(s')$ .

For  $J_2$ , we relax  $V_2 = V_{t+1}^{(U_{t+1}-u)}$  to an arbitrary function  $V : \mathcal{S} \rightarrow [0, 1/(1-\gamma)]$ :

$$J_2 \leq \max_{V: \mathcal{S} \rightarrow [0, 1/(1-\gamma)]} \left[ \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_V(s, a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_V(s, a) \rangle \right]. \tag{A.16}$$



Now substitute (A.15) and (A.16) into (A.13):

$$\begin{aligned}
 D_{u-1} &\leq \gamma \max_{s,a,l} (J_1 + J_2) \\
 &\leq \gamma \max_{s,a,l} \left[ \max_{s' \in \mathcal{S}} \left( V_t^{(U_t-u)} - V_{t+1}^{(U_{t+1}-u)} \right) (s') \right. \\
 &\quad \left. + \max_{V: \mathcal{S} \rightarrow [0,1/(1-\gamma)]} \left( \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_V(s,a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_V(s,a) \rangle \right) \right] \\
 &= \gamma \max_{s' \in \mathcal{S}} \left[ V_t^{(U_t-u)} - V_{t+1}^{(U_{t+1}-u)} \right] (s') + \gamma \max_{s,a,l,V} \left[ \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_V(s,a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_V(s,a) \rangle \right] \\
 &= \gamma D_u + \gamma \max_{s,a,l,V} \left[ \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_V(s,a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_V(s,a) \rangle \right],
 \end{aligned}$$

where the first equation holds because the first term in the big bracket is independent of the variables  $s, a, l$ , and the final equation holds based on the definition of  $D_u$ . Further note that according to Algorithm 1, the confidence sets at levels lower than  $l_t$  are not updated from  $t$  to  $t+1$ , and hence these confidence sets are identical for the two runs of ML-EVI. Therefore, we may exclude these smaller levels in the outside maximization above, and deduce

$$D_{u-1} \leq \gamma D_u + \gamma \max_{s,a,l \geq l_t, V} \left[ \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \phi_V(s,a) \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \phi_V(s,a) \rangle \right]. \quad (\text{A.17})$$

**Step 2**, for a given  $\psi = \phi_V(s,a)$ , we consider the expression in (A.17) inside the outer maximum, namely  $\max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \psi \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \psi \rangle$ .

Suppose  $\tilde{\theta}_1 := \operatorname{argmax}_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \psi \rangle$ . If  $\tilde{\theta}_1 \in \mathcal{C}_{t+1}^l$ , we immediately have  $\max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \psi \rangle = \langle \tilde{\theta}_1, \psi \rangle \leq \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \psi \rangle$  and the expression of interest is upper bounded by 0. Otherwise when  $\tilde{\theta}_1 \notin \mathcal{C}_{t+1}^l$ , consider  $\mu_{\max} = \sup\{\mu : \theta^* + \mu(\tilde{\theta}_1 - \theta^*) \in \mathcal{C}_{t+1}^l \cap \mathcal{B}\}$ .

According to Lemma A.3,  $\mathcal{B}$  is a convex body with non-empty interior when constrained to the affine subspace  $\mathcal{L} \subset \mathbb{R}^d$ . Combine this with  $\theta^*, \tilde{\theta}_1 \in \mathcal{B}$ , we see the segment  $\mathcal{J} \subset \mathbb{R}^d$  with  $\theta^*$  and  $\tilde{\theta}_1$  as endpoints is contained by  $\mathcal{B}$ . Since  $\mathcal{C}_{t+1}^l$  is a closed ellipsoid and hence a convex body in  $\mathbb{R}^d$ , its intersection with  $\mathcal{B}$  must also be a convex body when constrained to  $\mathcal{L}$ . Furthermore, by the conclusions of Lemma 6.2,  $\theta^*$  is an inner point in  $\mathcal{C}_{t+1}^l$  under  $\mathcal{E}_1$ , so the intersection of  $\mathcal{C}_{t+1}^l$  and  $\mathcal{J}$  must be a nontrivial segment  $\mathcal{J}'$ . Since the set  $\{\theta = \theta^* + \mu(\tilde{\theta}_1 - \theta^*) \in \mathcal{C}_{t+1}^l \cap \mathcal{B}\}$ , as the intersection of a line and a convex set, contains  $\mathcal{J}'$ , the set  $\{\mu : \theta^* + \mu(\tilde{\theta}_1 - \theta^*) \in \mathcal{C}_{t+1}^l \cap \mathcal{B}\}$  must be a non-trivial closed interval  $\mathcal{I}$  on  $\mathbb{R}$ .

Since  $\theta^* \in \mathcal{C}_{t+1}^l \cap \mathcal{B}$  under  $\mathcal{E}_1$  and  $\tilde{\theta}_1 \notin \mathcal{C}_{t+1}^l$  by our earlier assumption, we have  $0 \in \mathcal{I}$  and  $1 \notin \mathcal{I}$ , and therefore  $0 \leq \mu_{\max} < 1$ . Now define  $\tilde{\theta}_2 := \theta^* + \mu_{\max}(\tilde{\theta}_1 - \theta^*) \in \mathcal{C}_{t+1}^l \cap \mathcal{B}$ , and obtain:

$$\begin{aligned}
 \max_{\theta_1 \in \mathcal{C}_t^l \cap \mathcal{B}} \langle \theta_1, \psi \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \psi \rangle &= \langle \tilde{\theta}_1, \psi \rangle - \max_{\theta_2 \in \mathcal{C}_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \psi \rangle \\
 &\leq \langle \tilde{\theta}_1, \psi \rangle - \langle \theta^* + \mu_{\max}(\tilde{\theta}_1 - \theta^*), \psi \rangle \\
 &= (1 - \mu_{\max}) \langle \tilde{\theta}_1 - \theta^*, \psi \rangle.
 \end{aligned} \quad (\text{A.18})$$

It is worth noting that the second multiplier  $\langle \tilde{\theta}_1 - \theta^*, \psi \rangle$  is non-negative based on the definition of  $\tilde{\theta}_1$ .

Next, we start with giving a bound for the second multiplier. Since  $\tilde{\theta}_1 \in \mathcal{B}$ , we have  $\tilde{\mathbb{P}}(\cdot | s, a) = \langle \tilde{\theta}_1, \phi(\cdot | s, a) \rangle$  is a probability function for arbitrary  $s \in \mathcal{S}, a \in \mathcal{A}$ , and therefore

$$\begin{aligned}
 \langle \tilde{\theta}_1 - \theta^*, \psi \rangle &= \langle \tilde{\theta}_1, \psi \rangle - \langle \theta^*, \psi \rangle \\
 &= \tilde{\mathbb{P}}V(s, a) - \mathbb{P}V(s, a) \\
 &\leq \frac{1}{1-\gamma},
 \end{aligned} \quad (\text{A.19})$$

where we used  $0 \leq V(s) \leq 1/(1-\gamma)$  for arbitrary  $s \in \mathcal{S}$ .

Now we focus on  $\mu_{\max}$ . By its definition, we have  $\tilde{\theta}_2 = \theta^* + \mu_{\max}(\tilde{\theta}_1 - \theta^*) \in \partial_{\mathcal{L}}(\mathcal{C}_{t+1}^l \cap \mathcal{B})$ . If  $\tilde{\theta}_2 \in \partial_{\mathcal{L}}\mathcal{B}$ , since for  $0 \leq \mu < \mu_{\max}$ ,  $\theta^* + \mu(\tilde{\theta}_1 - \theta^*) \in \mathcal{B}$ , we can deduce that for any  $\mu > \mu_{\max}$  that  $\theta^* + \mu(\tilde{\theta}_1 - \theta^*) \notin \mathcal{B}$ . Because  $\mu_{\max} < 1$ , this leads to  $\tilde{\theta}_1 = \theta^* + 1 \cdot (\tilde{\theta}_1 - \theta^*) \notin \mathcal{B}$ , a contradiction. Otherwise  $\tilde{\theta}_2 \in \partial_{\mathcal{L}}\mathcal{C}_{t+1}^l \subset \partial\mathcal{C}_{t+1}^l$ , which implies

$$\begin{aligned}
 \beta^l &= \|\tilde{\theta}_2 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} \\
 &= \|\mu_{\max}\tilde{\theta}_1 + (1 - \mu_{\max})\theta^* - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} \\
 &= \|\mu_{\max}(\tilde{\theta}_1 - \hat{\theta}_{t+1}^l) + (1 - \mu_{\max})(\theta^* - \hat{\theta}_{t+1}^l)\|_{\Sigma_{t+1}^l} \\
 &\leq \mu_{\max}\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} + (1 - \mu_{\max})\|\theta^* - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} \\
 &\leq \mu_{\max}\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} + (1 - \mu_{\max}) \cdot \frac{1}{2}\beta^l,
 \end{aligned} \tag{A.20}$$

where we used the triangle inequality in the first inequality, and the conclusion of Lemma 6.2 in the second inequality.

Moving the second addend on the right hand side of (A.20) to the left hand side, we get

$$\begin{aligned}
 \frac{1}{2}(1 + \mu_{\max})\beta^l &\leq \mu_{\max}\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} \\
 &\leq \frac{1}{2}(1 + \mu_{\max})\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l}, \quad (\text{since } \mu_{\max} \leq 1)
 \end{aligned}$$

which leads to the conclusion  $\beta^l \leq \|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l}$ . We thus detract from  $\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l}$  both sides of (A.20) and obtain

$$\begin{aligned}
 \|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} - \beta^l &\geq (1 - \mu_{\max}) \times (\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} - 0.5\beta^l) \\
 &\geq (1 - \mu_{\max}) \cdot 0.5\beta^l,
 \end{aligned}$$

which yields

$$\begin{aligned}
 1 - \mu_{\max} &\leq \frac{\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} - \beta^l}{0.5\beta^l} \\
 &= 2\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} / \beta^l - 2.
 \end{aligned} \tag{A.21}$$

**Step 3**, we now set out to bound  $\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l}$ , aiming to prove it is close to  $\beta^l$ . In preparation, we first list the following relations between the variables based on the expressions given in Lemma A.1:  $\Sigma_{t+1}^l = \Sigma_t^l + 2^{t-l}\phi_0\phi_0^\top$ ,  $\mathbf{b}_{t+1}^l = \mathbf{b}_t^l + 2^{t-l}\phi_0V_0$ , where  $\phi_0 := \phi_{V_t}(s_t, a_t)$ ,  $V_0 := V_t(s_{t+1})$ . We also have  $\hat{\theta}_\tau^l = (\Sigma_\tau^l)^{-1}\mathbf{b}_\tau^l$ ,  $\tau = t, t+1$ . Now we break up the value using the triangle inequality:

$$\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} \leq \|\tilde{\theta}_1 - \hat{\theta}_t^l\|_{\Sigma_{t+1}^l} + \|\hat{\theta}_t^l - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l}. \tag{A.22}$$

We expand the first norm in order to use the relation  $\|\tilde{\theta}_1 - \hat{\theta}_t^l\|_{\Sigma_t^l} \leq \beta^l$  for an upper bound, specifically:

$$\begin{aligned}
 \|\tilde{\theta}_1 - \hat{\theta}_t^l\|_{\Sigma_{t+1}^l}^2 &= (\tilde{\theta}_1 - \hat{\theta}_t^l)^\top \Sigma_{t+1}^l (\tilde{\theta}_1 - \hat{\theta}_t^l) \\
 &= (\tilde{\theta}_1 - \hat{\theta}_t^l)^\top (\Sigma_t^l + 2^{t-l}\phi_0\phi_0^\top) (\tilde{\theta}_1 - \hat{\theta}_t^l) \\
 &= (\tilde{\theta}_1 - \hat{\theta}_t^l)^\top \Sigma_t^l (\tilde{\theta}_1 - \hat{\theta}_t^l) + 2^{t-l}[\phi_0^\top (\tilde{\theta}_1 - \hat{\theta}_t^l)]^2 \\
 &\leq \|\tilde{\theta}_1 - \hat{\theta}_t^l\|_{\Sigma_t^l}^2 + 2^{t-l}[\|\tilde{\theta}_1 - \hat{\theta}_t^l\|_{\Sigma_t^l} \cdot \|\phi_0\|_{(\Sigma_t^l)^{-1}}]^2 \quad (\text{Cauchy-Schwartz}) \\
 &\leq (\beta^l)^2 + 2^{t-l}(\beta^l)^2 \|\phi_0\|_{(\Sigma_t^l)^{-1}}^2.
 \end{aligned} \tag{A.23}$$

For the second norm, the vector is the difference between two confidence set centers. We tackle this difference by extracting  $(\Sigma_{t+1}^l)^{-1}$  from both centers:

$$\begin{aligned}
 \widehat{\theta}_t^l - \widehat{\theta}_{t+1}^l &= (\Sigma_t^l)^{-1} \mathbf{b}_t^l - (\Sigma_{t+1}^l)^{-1} \mathbf{b}_{t+1}^l \\
 &= (\Sigma_{t+1}^l)^{-1} [\Sigma_{t+1}^l (\Sigma_t^l)^{-1} \mathbf{b}_t^l - \mathbf{b}_{t+1}^l] \\
 &= (\Sigma_{t+1}^l)^{-1} [(\Sigma_t^l + 2^{l-t} \phi_0 \phi_0^\top) (\Sigma_t^l)^{-1} \mathbf{b}_t^l - (\mathbf{b}_t^l + 2^{l-t} \phi_0 V_0)] \\
 &= (\Sigma_{t+1}^l)^{-1} [\mathbf{b}_t^l + 2^{l-t} (\phi_0^\top (\Sigma_t^l)^{-1} \mathbf{b}_t^l) \phi_0 - \mathbf{b}_t^l - 2^{l-t} V_0 \phi_0] \\
 &= 2^{l-t} (\phi_0^\top (\Sigma_t^l)^{-1} \mathbf{b}_t^l - V_0) (\Sigma_{t+1}^l)^{-1} \phi_0.
 \end{aligned} \tag{A.24}$$

Substitute the results (A.23) and (A.24) back into (A.22):

$$\begin{aligned}
 \|\widehat{\theta}_1 - \widehat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} &\leq \sqrt{(\beta^l)^2 + 2^{l-t} (\beta^l)^2 \|\phi_0\|_{(\Sigma_t^l)^{-1}}^2} + \|2^{l-t} (\phi_0^\top (\Sigma_t^l)^{-1} \mathbf{b}_t^l - V_0) (\Sigma_{t+1}^l)^{-1} \phi_0\|_{\Sigma_{t+1}^l} \\
 &= \beta^l \sqrt{1 + 2^{l-t} \|\phi_0\|_{(\Sigma_t^l)^{-1}}^2} + 2^{l-t} |\phi_0^\top (\Sigma_t^l)^{-1} \mathbf{b}_t^l - V_0| \|(\Sigma_{t+1}^l)^{-1} \phi_0\|_{\Sigma_{t+1}^l} \\
 &\leq \beta^l (1 + 2^{l-t-1} \|\phi_0\|_{(\Sigma_t^l)^{-1}}) + 2^{l-t} |\phi_0^\top \widehat{\theta}_t^l - V_0| \|\phi_0\|_{(\Sigma_{t+1}^l)^{-1}},
 \end{aligned} \tag{A.25}$$

where we used the trivial inequality  $\sqrt{1+x} \leq 1+x/2$  for any  $x \geq 0$  in the last relation.

To tackle the right hand side of (A.25), we need to provide an upper bound for  $\|\phi_0\|_{(\Sigma_t^l)^{-1}}$ , where  $\tau = t, t+1$ . First we have trivially that  $\Sigma_{t+1}^l \succeq \Sigma_t^l$ . Furthermore, by the alternative definition of  $\Sigma_t^l$  given in Lemma A.1, we can see based on  $l \geq l_t$  that when  $l_t > 1$ ,

$$\begin{aligned}
 \Sigma_t^l &= \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l\} 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\
 &\succeq \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l_t - 1\} 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\
 &= \lambda \mathbf{I} + 2^{l_t - l - 1} \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l_t - 1\} 2^{l_\tau - l_t + 1} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\
 &\succeq 2^{l_t - l - 1} \Sigma_t^{l_t - 1}.
 \end{aligned}$$

The above relations lead to  $\|\phi_0\|_{(\Sigma_{t+1}^l)^{-1}} \leq \|\phi_0\|_{(\Sigma_t^l)^{-1}} \leq \|\phi_0\|_{(2^{l-t-t-1} \Sigma_t^{l_t-1})^{-1}} = 2^{(l-l_t+1)/2} \|\phi_0\|_{(\Sigma_t^{l_t-1})^{-1}} \leq 2^{(l-l_t+1)/2} \cdot 2^{-l_t+1} \sqrt{d}/(1-\gamma) = 2^{(l-3l_t+3)/2} \sqrt{d}/(1-\gamma)$ , where the  $l_t$  selection rule in Lines 10 to 13 of Algorithm 1 yields the last inequality. Recall our definition of  $l_0$  from Lemma 6.1, we have shown that when  $l_t > 1$ ,

$$\|\phi_0\|_{(\Sigma_{t+1}^l)^{-1}} \leq 2^{(l+2l_0-3l_t+3)/2}. \tag{A.26}$$

Note that in the special case where  $l_t = 1$ , we have trivially  $\|\phi_0\|_{(\Sigma_{t+1}^l)^{-1}} \leq \|\phi_0\|_{\lambda^{-1} \mathbf{I}} = \lambda^{-1/2} \|\phi_0\|_2 \leq \sqrt{d}/(1-\gamma) = 2^{l_0}$ , and (A.26) still holds.

Apart from this, we also need to bound  $|\phi_0^\top \widehat{\theta}_t^l - V_0|$  in (A.25). Since  $\mathbb{P}V_t(s_t, a_t) = \langle \theta^*, \phi_{V_t}(s_t, a_t) \rangle$  from (3.3), we have

$$\begin{aligned}
 |\phi_0^\top \widehat{\theta}_t^l - V_0| &= |\langle \widehat{\theta}_t^l, \phi_0 \rangle - V_t(s_{t+1})| \\
 &= |\langle \widehat{\theta}_t^l, \phi_0 \rangle - \langle \theta^*, \phi_0 \rangle + \mathbb{P}V_t(s_t, a_t) - V_t(s_{t+1})| \\
 &\leq |\langle \widehat{\theta}_t^l - \theta^*, \phi_0 \rangle| + |\mathbb{P}V_t(s_t, a_t) - V_t(s_{t+1})| && \text{(triangle inequality)} \\
 &\leq \|\theta^* - \widehat{\theta}_t^l\|_{\Sigma_t^l} \cdot \|\phi_0\|_{(\Sigma_t^l)^{-1}} + |\mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_t, a_t)} V_t(s') - V_t(s_{t+1})| && \text{(Cauchy-Schwartz)} \\
 &\leq \frac{1}{2} \beta^l \|\phi_0\|_{(\Sigma_t^l)^{-1}} + \frac{1}{1-\gamma},
 \end{aligned} \tag{A.27}$$

where in the last inequality we used the conclusion of Lemma 6.2 and  $0 \leq V_t(s) \leq 1/(1-\gamma)$ .

Now we turn back to (A.25). Implementing the inequalities (A.26) and (A.27) in the right hand side, we obtain

$$\begin{aligned}
 \|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} &\leq \beta^l (1 + 2^{l-l-1} \|\phi_0\|_{(\Sigma_t^l)^{-1}}^2) + 2^{l-l} \left( \frac{1}{2} \beta^l \|\phi_0\|_{(\Sigma_t^l)^{-1}} + \frac{1}{1-\gamma} \right) \|\phi_0\|_{(\Sigma_t^l)^{-1}} \\
 &= \beta^l + 2^{l-l} \left( \frac{1}{2} \beta^l \|\phi_0\|_{(\Sigma_t^l)^{-1}}^2 + \frac{1}{2} \beta^l \|\phi_0\|_{(\Sigma_t^l)^{-1}} + \frac{1}{1-\gamma} \|\phi_0\|_{(\Sigma_t^l)^{-1}} \right) \\
 &\leq \beta^l + 2^{l-l} \left( \beta^l 2^{2l_0-3l_t+3} + \frac{1}{1-\gamma} 2^{(l+2l_0-3l_t+3)/2} \right) \\
 &= \beta^l \left( 1 + 2^{2l_0-2l_t+3} + (1/(1-\gamma)\beta^l) \cdot 2^{l_0-(l+l_t-3)/2} \right) \\
 &\leq \beta^l \left( 1 + 2^{2l_0-2l_t+3} + (1/6\sqrt{d}) \cdot 2^{l_0-l_t+3/2}\sqrt{d} \right) \\
 &\leq \beta^l \left( 1 + 2^{2l_0-2l_t+3} + 2^{l_0-l_t}/2\sqrt{d} \right), \tag{A.28}
 \end{aligned}$$

where the third inequality is due to  $l \geq l_t$  and the definition of  $\beta^l$  from (5.1), which leads to  $(1-\gamma)\beta^l \geq 6\sqrt{d} \geq 6\sqrt{d}$ .

**Step 4**, we are putting everything together. Substituting (A.28) into (A.21), and then into (A.18) together with (A.19), we have

$$\begin{aligned}
 \max_{\theta_1 \in C_t^l \cap \mathcal{B}} \langle \theta_1, \psi \rangle - \max_{\theta_2 \in C_{t+1}^l \cap \mathcal{B}} \langle \theta_2, \psi \rangle &\leq (1 - \mu_{\max}) \langle \tilde{\theta}_1 - \theta^*, \psi \rangle \\
 &\leq (2\|\tilde{\theta}_1 - \hat{\theta}_{t+1}^l\|_{\Sigma_{t+1}^l} / \beta^l - 2) \cdot \frac{1}{1-\gamma} \\
 &\leq (2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}) \cdot \frac{1}{1-\gamma}.
 \end{aligned}$$

Now substitute the above inequality into (A.17) to get  $D_{u-1} \leq \gamma D_u + (\gamma/(1-\gamma)) \cdot [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}]$ . Recalling from (5.1) that  $U_{t+1} \geq U_t$ , and from (A.5) that  $\gamma^{U_t} \leq 1/(t(t+1))$ , along with the trivial bound for  $D_{U_t} \leq 1/(1-\gamma)$ , we iterate this inequality for  $D_u$  over  $u = 1, \dots, U_t$  to obtain:

$$\begin{aligned}
 D_0 &\leq \gamma^{U_t} D_{U_t} + \sum_{u=1}^{U_t} \frac{\gamma^u}{1-\gamma} \cdot [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}] \\
 &\leq \frac{1}{t(t+1)} \cdot \frac{1}{1-\gamma} + \frac{1}{1-\gamma} \cdot [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}] \sum_{u=1}^{U_t} \gamma^u \\
 &\leq \frac{1}{(1-\gamma)t(t+1)} + \frac{\gamma}{(1-\gamma)^2} [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}].
 \end{aligned}$$

Combining this with  $V_t(s_{t+1}) - V_{t+1}(s_{t+1}) = V_t^{(U_t)}(s_t) - V_{t+1}^{(U_{t+1})}(s_t) \leq D_0$ , we can now bound the sum of  $C_t$  by

$$\begin{aligned}
 \sum_{t \in \mathcal{T}_\epsilon} \mathbb{1}\{\mathcal{Q}_t\} C_t &= \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{\mathcal{Q}_t\} \gamma^{\tau-t+1} (V_\tau - V_{\tau+1})(s_{\tau+1}) \\
 &\leq \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{\mathcal{Q}_t\} \gamma^{\tau-t+1} \left[ \frac{1}{(1-\gamma)\tau(\tau+1)} + \frac{\gamma}{(1-\gamma)^2} [2^{2l_0-2l_\tau+4} + 2^{l_0-l_\tau+1.5}/3\sqrt{d}] \right] \\
 &\leq \frac{\gamma}{1-\gamma} \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \frac{\gamma^{\tau-t}}{\tau(\tau+1)} + \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \frac{\gamma^{\tau-t+2}}{(1-\gamma)^2} [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}] \\
 &\leq \frac{\gamma}{(1-\gamma)^2} + \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \frac{\gamma^{\tau-t+2}}{(1-\gamma)^2} [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}], \tag{A.29}
 \end{aligned}$$

where in the third inequality we relaxed the indicator function  $\mathbb{1}\{\mathcal{Q}_t\}$  to 1 and  $\mathbb{1}\{l_\tau > l_0\}$  for the two summations respectively. The final inequality is due to (A.10).

#### Part IV: Final Arrangements

We now combine the three bounds for the sums of  $A_t, B_t, C_t$ , namely (A.11), (A.12) and (A.29), and substitute them into (A.8) to obtain

$$\begin{aligned}
 \epsilon N_\epsilon &\leq \frac{\gamma^T}{1-\gamma} N_\epsilon + \left[ \frac{1}{1-\gamma} + 3 \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t+1} 2^{l_0-l_\tau} \beta^{l_\tau-1} \right] + \frac{2\gamma}{(1-\gamma)^2} \sqrt{2N_\epsilon \log(2T(T+1)/\delta)} \\
 &\quad + \left[ \frac{\gamma}{(1-\gamma)^2} + \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \frac{\gamma^{\tau-t+2}}{(1-\gamma)^2} [2^{2l_0-2l_t+4} + 2^{l_0-l_t+1.5}/3\sqrt{d}] \right] + \frac{10dT}{1-\gamma} l_0(l_0+1) \\
 &= \frac{\gamma^T}{1-\gamma} N_\epsilon + \frac{2\gamma}{(1-\gamma)^2} \sqrt{2N_\epsilon \log(2T(T+1)/\delta)} + \frac{1}{(1-\gamma)^2} + \frac{10dT}{1-\gamma} l_0(l_0+1) \\
 &\quad + \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t} \left[ 3\gamma 2^{l_0-l_\tau} \beta^{l_\tau-1} + \frac{16\gamma^3}{(1-\gamma)^2} \cdot 2^{2l_0-2l_t} + \frac{2\sqrt{2}\gamma^2}{3(1-\gamma)^2\sqrt{d}} \cdot 2^{l_0-l_t} \right]. \tag{A.30}
 \end{aligned}$$

We will now focus on the second half of (A.30). More generally, we present an upper bound for the term  $\sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t} \eta^{l_0-l_\tau} f(l_\tau - 1)$  in the following, where  $f: \mathbb{R}_+ \rightarrow \mathbb{R}$  is of the form  $f(x) = c_1 \sqrt{x} + c_2$ ,  $c_1, c_2 \geq 0$ , and  $2 \leq \eta \leq 4$ . We do this by using indicator functions to break up the summation, namely:

$$\begin{aligned}
 &\sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t} \eta^{l_0-l_\tau} f(l_\tau - 1) \\
 &= \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \sum_{l=[l_0]+1}^{\infty} \gamma^{\tau-t} \eta^{l_0-l} f(l-1) \mathbb{1}\{l = l_\tau\} \\
 &= \sum_{n=0}^{T-1} \sum_{\tau=1+n}^{\infty} \sum_{l=[l_0]+1}^{\infty} \gamma^n \eta^{l_0-l} f(l-1) \mathbb{1}\{l = l_\tau\} \mathbb{1}\{\tau - n \in \mathcal{T}_\epsilon\} \\
 &= \sum_{n=0}^{T-1} \gamma^n \sum_{l=[l_0]+1}^{\infty} \eta^{l_0-l} f(l-1) \sum_{\tau=1+n}^{\infty} \mathbb{1}\{l = l_\tau\} \mathbb{1}\{\tau - n \in \mathcal{T}_\epsilon\} \\
 &\leq \sum_{n=0}^{T-1} \gamma^n \sum_{l=[l_0]+1}^{\infty} \eta^{l_0-l} f(l-1) \min \left\{ \sum_{\tau=1+n}^{\infty} \mathbb{1}\{l = l_\tau\}, \sum_{\tau=1+n}^{\infty} \mathbb{1}\{\tau - n \in \mathcal{T}_\epsilon\} \right\} \\
 &\leq \sum_{n=0}^{T-1} \gamma^n \sum_{l=[l_0]+1}^{\infty} \eta^{l_0-l} f(l-1) \min\{k_{\max}^l, N_\epsilon\} \\
 &\leq \frac{1}{1-\gamma} \sum_{l=[l_0]+1}^{\infty} \eta^{l_0-l} f(l-1) \min\{20dl^{l-l_0}, N_\epsilon\}, \tag{A.31}
 \end{aligned}$$

where: in the first equation we expanded  $\mathbb{1}\{l_\tau > l_0\}$  into  $\sum_{l>l_0} \mathbb{1}\{l = l_\tau\}$ ; in the second equation we used the variable  $n = \tau - t$  to replace  $t$ , which adds a multiplier  $\mathbb{1}\{t := \tau - n \in \mathcal{T}_\epsilon\}$  to the expression under summation; in the second to last inequality we used the definition of  $k_{\max}^l$  from Lemma 6.1 and the definition  $N_\epsilon = \#\mathcal{T}_\epsilon$ ; the last inequality is due to Lemma 6.1.

Ideally, we can determine an exact threshold integer value  $L'_0 = \max\{l > l_0 : 20dl^{l-l_0} \leq N_\epsilon\}$ , and take the correct minimization in (A.31) for all  $l \in \mathbb{Z}$ , yet this will intensely complicate our deductions. Instead, we will set a undetermined threshold value  $L_0$ , under which the minimization takes the upper bound  $20dl^{l-l_0}$ , and above which the minimization takes

the upper bound  $N_\epsilon$ :

$$\begin{aligned}
 & \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbf{1}\{l_\tau > l_0\} \gamma^{\tau-t} \eta^{l_0-l_\tau} f(l_\tau - 1) \\
 & \leq \frac{1}{1-\gamma} \sum_{l=[l_0]+1}^{L_0} \eta^{l_0-l} f(l-1) \cdot 20dL_0 4^{l-l_0} + \frac{1}{1-\gamma} \sum_{l=L_0+1}^{\infty} \eta^{l_0-l} f(l-1) N_\epsilon \\
 & \leq \frac{1}{1-\gamma} \cdot f(L_0 - 1) \cdot 20dL_0 \sum_{l=[l_0]+1}^{L_0} \eta^{l_0-l} 4^{l-l_0} + \frac{N_\epsilon}{1-\gamma} \sum_{l=L_0+1}^{\infty} \eta^{l_0-l} f(l-1) \\
 & = \frac{20d}{1-\gamma} L_0 f(L_0 - 1) \sum_{l=[l_0]+1}^{L_0} (4/\eta)^{l-l_0} + \frac{N_\epsilon}{1-\gamma} \sum_{l=L_0+1}^{\infty} \eta^{l_0-l} (c_1 \sqrt{l-1} + c_2), \tag{A.32}
 \end{aligned}$$

where the second inequality holds because  $f(l-1) \cdot 20dl$  is increasing in  $l$ , and in the equation we plugged in  $f(x) = c_1 \sqrt{x} + c_2$ .

To tackle the second half of (A.32), gather from the first conclusion of Lemma D.3 that  $\sum_{l=L}^{\infty} \eta^{l_0-l} \sqrt{l} \leq \sqrt{L+1} \eta^{l_0-L+1} / (\eta - 1)$ , which leads to the upper bound

$$\begin{aligned}
 \sum_{l=L_0+1}^{\infty} \eta^{l_0-l} (c_1 \sqrt{l-1} + c_2) &= c_2 \sum_{l=L_0+1}^{\infty} \eta^{l_0-l} + c_1 \sum_{l=L_0}^{\infty} \eta^{l_0-l-1} \sqrt{l} \\
 &\leq (c_2 + c_1 \sqrt{L_0+1}) \eta^{l_0-L_0} / (\eta - 1) \\
 &= f(L_0+1) \eta^{l_0-L_0} / (\eta - 1). \tag{A.33}
 \end{aligned}$$

Substituting (A.33) into (A.32), we have

$$\begin{aligned}
 \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \gamma^{\tau-t} \eta^{l_0-l_\tau} f(l_\tau - 1) &\leq \frac{20d}{1-\gamma} L_0 f(L_0 - 1) \sum_{l=[l_0]+1}^{L_0} (4/\eta)^{l-l_0} + \frac{N_\epsilon}{1-\gamma} f(L_0+1) \eta^{l_0-L_0} / (\eta - 1) \\
 &\leq \frac{f(L_0+1)}{1-\gamma} \left[ 20dL_0 \sum_{l=[l_0]+1}^{L_0} (4/\eta)^{l-l_0} + N_\epsilon \eta^{l_0-L_0} / (\eta - 1) \right], \tag{A.34}
 \end{aligned}$$

where for the second inequality, we used the fact  $f(L_0 - 1) \leq f(L_0 + 1)$  due to  $c_1 \geq 0$ .

Now we go back to consider the second half of (A.30). With  $\beta^l = 2(3\sqrt{dl} + 2\sqrt{\log(1/\delta)}) / (1-\gamma) + 2\sqrt{d\lambda}$ , we can take  $\eta = 2, 4$  in (A.34) to get

$$\begin{aligned}
 & \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbf{1}\{l_\tau > l_0\} \gamma^{\tau-t} \left[ 3\gamma 2^{l_0-l_\tau} \beta^{l_\tau-1} + \frac{16\gamma^3}{(1-\gamma)^2} \cdot 2^{2l_0-2l_\tau} + \frac{2\sqrt{2}\gamma^2}{3(1-\gamma)^2\sqrt{d}} \cdot 2^{l_0-l_\tau} \right] \\
 &= \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbf{1}\{l_\tau > l_0\} \gamma^{\tau-t} 2^{l_0-l_\tau} \left[ 3\gamma \beta^{l_\tau-1} + \frac{2\sqrt{2}\gamma^2}{3(1-\gamma)^2\sqrt{d}} \right] + \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbf{1}\{l_\tau > l_0\} \gamma^{\tau-t} 4^{l_0-l_\tau} \cdot \frac{16\gamma^3}{(1-\gamma)^2} \\
 &\leq \frac{1}{1-\gamma} \left( 3\gamma \beta^{L_0+1} + \frac{2\sqrt{2}\gamma^2}{3(1-\gamma)^2\sqrt{d}} \right) \left( 20dL_0 \sum_{l=[l_0]+1}^{L_0} 2^{l-l_0} + N_\epsilon 2^{l_0-L_0} \right) \\
 &\quad + \frac{16\gamma^3}{(1-\gamma)^3} \left( 20dL_0 \sum_{l=[l_0]+1}^{L_0} 1 + N_\epsilon 4^{l_0-L_0} / 3 \right) \\
 &\leq \left( \frac{3\beta^{L_0+1}}{1-\gamma} 2^{l_0-L_0} + \frac{2\sqrt{2}}{3(1-\gamma)^3\sqrt{d}} 2^{l_0-L_0} + \frac{16}{3(1-\gamma)^3} 4^{l_0-L_0} \right) N_\epsilon \\
 &\quad + \left( \frac{3\beta^{L_0+1}}{1-\gamma} 2^{L_0-l_0+1} + \frac{2\sqrt{2}}{3(1-\gamma)^3\sqrt{d}} 2^{L_0-l_0+1} + \frac{16}{(1-\gamma)^3} L_0 \right) \cdot 20dL_0, \tag{A.35}
 \end{aligned}$$

where in the last inequality we used  $\gamma \leq 1$  to get rid of all the excess  $\gamma$ 's.

Now take

$$T := \left\lceil \log \frac{3}{(1-\gamma)\epsilon} / (1-\gamma) \right\rceil,$$

and

$$L_0 := \left\lceil l_0 + \log \left( M + K \sqrt{2 \log(M+K) / \log 2} + 2(l_0 + 1) \right) / \log 2 \right\rceil,$$

where

$$M = \frac{3}{\epsilon} \cdot \left[ \frac{1}{(1-\gamma)^3 \sqrt{d}} + \frac{12 \sqrt{\log(1/\delta)}}{(1-\gamma)^2} + \frac{6 \sqrt{d\lambda}}{1-\gamma} \right],$$

$$K = \frac{54 \sqrt{d}}{(1-\gamma)^2 \epsilon}.$$

It is not hard to check that

$$L_0 = O\left( \log \frac{d \log(1/\delta)}{(1-\gamma)} \right) + O\left( \log \frac{1}{\epsilon} \right),$$

where we separated the term  $\log(1/\epsilon)$  since the maximum possible value of  $\epsilon \leq 1 - \gamma$  is dependent on the parameter  $\gamma$ . In the third conclusion of Lemma D.3, by taking  $a = 2^{l_0+1}M$  and  $b = 2^{l_0+1}K$ , we have  $L_0 + 1 = \lceil \log(a + b\sqrt{2 \log(a+b)/\log 2}) / \log 2 \rceil$ , and hence

$$\begin{aligned} 2^{L_0-l_0} &\geq 2^{-l_0-1}(a + b\sqrt{L_0+1}) \\ &= M + K\sqrt{L_0+1} \\ &= \frac{3}{\epsilon} \cdot \left[ \frac{1}{(1-\gamma)^3 \sqrt{d}} + \frac{12 \sqrt{\log(1/\delta)}}{(1-\gamma)^2} + \frac{6 \sqrt{d\lambda}}{1-\gamma} + \frac{18 \sqrt{d}}{(1-\gamma)^2} \sqrt{L_0+1} \right] \\ &= \frac{3}{\epsilon} \cdot \left[ \frac{1}{(1-\gamma)^3 \sqrt{d}} + \frac{3\beta^{L_0+1}}{1-\gamma} \right], \end{aligned}$$

where the final equation is due to the definition of  $\beta^l$  in (5.1). As a side effect of the second to last expression above, we have the relation  $2^{L_0-l_0} \geq 3(6 + 18\sqrt{2})\sqrt{d}$  due to  $\epsilon \leq 1/(1-\gamma)$ ,  $\lambda \geq 1$  and  $L_0 \geq 1$ . On the other hand, we have an upper bound of similar composition:

$$\begin{aligned} 2^{L_0-l_0} &\leq 2 \lceil M + K \sqrt{2 \log(M+K) / \log 2} + 2l_0 + 2 \rceil \\ &\leq 2 \lceil M + K \sqrt{2L_0 + 2} \rceil \\ &= \frac{6}{\epsilon} \cdot \left[ \frac{1}{(1-\gamma)^3 \sqrt{d}} + \frac{12 \sqrt{\log(1/\delta)}}{(1-\gamma)^2} + \frac{6 \sqrt{d\lambda}}{1-\gamma} + \frac{18 \sqrt{2d}}{(1-\gamma)^2} \sqrt{L_0 + 1} \right], \end{aligned}$$

where in the second inequality we used the relation  $L_0 \geq \log(M+K)/\log 2$ .

We now substitute the two inequalities above into (A.35):

$$\begin{aligned}
 & \sum_{t \in \mathcal{T}_\epsilon} \sum_{\tau=t}^{t+T-1} \mathbb{1}\{l_\tau > l_0\} \gamma^{\tau-t} \left[ 3\gamma 2^{l_0-l_\tau} \beta^{l_\tau-1} + \frac{16\gamma^3}{(1-\gamma)^2} \cdot 2^{2l_0-2l_t} + \frac{2\sqrt{2}\gamma^2}{3(1-\gamma)^2\sqrt{d}} \cdot 2^{l_0-l_t} \right] \\
 & \leq \left( \frac{3\beta^{L_0+1}}{1-\gamma} + \frac{2\sqrt{2}}{3(1-\gamma)^3\sqrt{d}} + \frac{16}{3(1-\gamma)^3} \cdot \frac{1}{3(6+18\sqrt{2})\sqrt{d}} \right) 2^{l_0-L_0} N_\epsilon \\
 & \quad + \left( \frac{3\beta^{L_0+1}}{1-\gamma} + \frac{2}{3(1-\gamma)^3\sqrt{d}} \right) \cdot 20dL_0 2^{L_0-l_0} + \frac{320dL_0^2}{(1-\gamma)^3} \\
 & \leq \frac{\epsilon}{3} N_\epsilon + 20dL_0 \cdot \frac{6}{\epsilon} \left[ \frac{1}{(1-\gamma)^3\sqrt{d}} + \frac{12\sqrt{\log(1/\delta)}}{(1-\gamma)^2} + \frac{6\sqrt{d\lambda}}{1-\gamma} + \frac{18\sqrt{2d}}{(1-\gamma)^2} \sqrt{L_0+1} \right] \\
 & \quad \times \left( \frac{2\sqrt{2}}{3(1-\gamma)^3\sqrt{d}} + \frac{12\sqrt{\log(1/\delta)}}{(1-\gamma)^2} + \frac{6\sqrt{d\lambda}}{1-\gamma} + \frac{18\sqrt{d}}{(1-\gamma)^2} \sqrt{L_0+1} \right) + \frac{320dL_0^2}{(1-\gamma)^3} \\
 & \leq \frac{\epsilon}{3} N_\epsilon + 120dL_0/\epsilon \cdot \sqrt{1^2+12^2+6^2+(18\sqrt{2})^2} \cdot \sqrt{(2\sqrt{2}/3)^2+12^2+6^2+18^2} \\
 & \quad \times \left[ \frac{1}{(1-\gamma)^6d} + \frac{\log(1/\delta)}{(1-\gamma)^4} + \frac{dL_0}{(1-\gamma)^4} + \frac{d\lambda}{(1-\gamma)^2} \right] + \frac{320dL_0^2}{(1-\gamma)^3} \\
 & \leq \frac{\epsilon}{3} N_\epsilon + CdL_0/\epsilon \cdot \left[ \frac{1}{(1-\gamma)^6d} + \frac{\log(1/\delta)}{(1-\gamma)^4} + \frac{dL_0}{(1-\gamma)^4} + \frac{d\lambda}{(1-\gamma)^2} \right], \tag{A.36}
 \end{aligned}$$

where: the second inequality is partly due to  $2\sqrt{2}/3 + 16/9(6+18\sqrt{2}) \leq 1$ ; we used the standard Cauchy-Schwartz inequality in the third inequality; for the last expression,  $C$  is an absolute constant.

Besides, by the definition of  $T$ ,

$$\begin{aligned}
 \gamma^T & \leq (\gamma^{1/(1-\gamma)})^{\log[3/(1-\gamma)\epsilon]} \\
 & \leq \exp[-\log(3/(1-\gamma)\epsilon)] \\
 & \leq (1-\gamma)\epsilon/3.
 \end{aligned}$$

We then bring this and (A.36) back to (A.30):

$$\begin{aligned}
 \epsilon N_\epsilon & \leq \frac{\epsilon}{3} N_\epsilon + \frac{2\gamma}{(1-\gamma)^2} \sqrt{2N_\epsilon \log(2T(T+1)/\delta)} + \frac{1}{(1-\gamma)^2} + \frac{10dT}{1-\gamma} l_0(l_0+1) \\
 & \quad + \frac{\epsilon}{3} N_\epsilon + \frac{CdL_0}{\epsilon} \left[ \frac{1}{(1-\gamma)^6d} + \frac{\log(1/\delta)}{(1-\gamma)^4} + \frac{dL_0}{(1-\gamma)^4} + \frac{d\lambda}{(1-\gamma)^2} \right].
 \end{aligned}$$

We go on to move the two terms  $(\epsilon/3)N_\epsilon$  to the left hand side, then multiply both sides by the factor  $3/\epsilon$ , thus transforming the above inequality into (after erasing yet another  $\gamma$  from the first term):

$$\begin{aligned}
 N_\epsilon & \leq \frac{6}{(1-\gamma)^2\epsilon} \sqrt{2\log(2T(T+1)/\delta)} \sqrt{N_\epsilon} + \frac{3}{(1-\gamma)^2\epsilon} + \frac{30dT}{(1-\gamma)\epsilon} l_0(l_0+1) \\
 & \quad + \frac{3C}{\epsilon^2} \cdot \left[ \frac{L_0}{(1-\gamma)^6} + \frac{d\log(1/\delta)L_0 + d^2L_0^2}{(1-\gamma)^4} + \frac{d^2\lambda L_0}{(1-\gamma)^2} \right].
 \end{aligned}$$

Using the second conclusion in Lemma D.3, which states that  $x - a\sqrt{x} \leq b \Rightarrow x \leq a^2 + 2b$  for positive real numbers  $a, b$ ,



we see the above inequality implies

$$\begin{aligned}
 N_\epsilon &\leq \left[ \frac{6}{(1-\gamma)^2\epsilon} \sqrt{2\log(2T(T+1)/\delta)} \right]^2 + \frac{6}{(1-\gamma)^2\epsilon} + \frac{60dT}{(1-\gamma)\epsilon} l_0(l_0+1) \\
 &\quad + \frac{6C}{\epsilon^2} \cdot \left[ \frac{L_0}{(1-\gamma)^6} + \frac{d\log(1/\delta)L_0 + d^2L_0^2}{(1-\gamma)^4} + \frac{d^2\lambda L_0}{(1-\gamma)^2} \right] \\
 &= O\left( \frac{L_0}{(1-\gamma)^6\epsilon^2} + \frac{d\log(1/\delta)L_0 + d^2L_0^2 + \log T}{(1-\gamma)^4\epsilon^2} + \frac{d^2\lambda L_0}{(1-\gamma)^2\epsilon^2} + \frac{dl_0^2 \log(1/(1-\gamma)\epsilon)}{(1-\gamma)^2\epsilon} \right) \\
 &= \tilde{O}\left( \frac{1}{(1-\gamma)^6\epsilon^2} + \frac{d^2 + d\log(1/\delta)}{(1-\gamma)^4\epsilon^2} \right),
 \end{aligned}$$

which concludes the proof.  $\square$

## A.2. Proof of Corollary 5.3

*Proof.* Theorem 3 of [Dann et al. \(2017\)](#) states that if an algorithm is uniform-PAC with sample complexity  $\Gamma(\epsilon, \delta) = \tilde{O}(C_1/\epsilon + C_2/\epsilon^2)$  for some  $\delta > 0$ , then with probability at least  $1 - \delta$ , its regret is bounded by  $\tilde{O}(\sqrt{C_2T} + \max\{C_1, C_2\})$ . Taking  $C_1 = 0$ ,  $C_2 = 1/(1-\gamma)^3 + (d + \sqrt{d\log(1/\delta)})/(1-\gamma)^2$ , and discarding the non-dominating terms without  $\sqrt{T}$ , this suggests that our algorithm has a high probability regret bound of

$$\text{Regret}(T) = \tilde{O}\left( \frac{d + \sqrt{d\log(1/\delta)}}{(1-\gamma)^2} \sqrt{T} + \frac{\sqrt{T}}{(1-\gamma)^3} \right). \quad (\text{A.37})$$

$\square$

## B. Proof of Lemmas in Section 6

In this section, we focus on the lemmas presented in Section 6 and give detailed proofs for each of them.

### B.1. Proof of Lemma 6.1

*Proof of Lemma 6.1.* By Definition 3.1, we have  $\|\phi_{V_t}(s_t, a_t)\|_2 \leq \sqrt{d} \cdot 1/(1-\gamma)$ . Take  $\{X_t\} = \{2^{(l_\tau-1)/2} \phi_{V_\tau}(s_\tau, a_\tau) : \tau \geq 1, l_\tau \leq l\}$ ,  $V = \lambda \mathbf{I}$  and  $L = \sqrt{d}/(1-\gamma)$  which, considering Lemma A.1, means that

$$\begin{aligned}
 \bar{V}_{t'} &= \lambda \mathbf{I} + \sum_{\tau=1}^t \mathbf{1}\{l_\tau \leq l\} 2^{l_\tau-l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\
 &= \Sigma_{t+1}^l,
 \end{aligned}$$

where  $t' = \sum_{\tau=1}^t \mathbf{1}\{l_\tau \leq l\}$ . We can thus deduce from Lemma 11 of [Abbasi-Yadkori et al. \(2011\)](#) the following:

$$\sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l\} \min \{2^{l_\tau-l} \|\phi_{V_\tau}(s_\tau, a_\tau)\|_{(\Sigma_\tau^l)^{-1}}^2, 1\} \leq 2 \log \frac{\det(\Sigma_t^l)}{\det(\lambda \mathbf{I})}. \quad (\text{B.1})$$

On one hand, from Lines 10 to 13 of algorithm 1, we have when  $l_t = l$  and  $l_t \neq L_t + 1$ ,  $\|\phi_{V_t}(s_t, a_t)\|_{(\Sigma_t^l)^{-1}} \geq 2^{-l} \sqrt{d}/(1-\gamma)$ . Combine this with the definition  $k_t^l = \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l\} 2^{l_\tau-l}$  and the observation that  $l_t = l = L_t + 1$  happens exactly once for each level  $l$ , we have that when  $t$  is large enough, the left hand side of inequality (B.1) can be lower bounded by

$$\begin{aligned}
 \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau \leq l\} \min \{2^{l_\tau-l} \|\phi_{V_\tau}(s_\tau, a_\tau)\|_{(\Sigma_\tau^l)^{-1}}^2, 1\} &\geq \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau = l, L_\tau + 1 \neq l\} \min \{2^{l_\tau-l} \cdot 2^{-2l} d/(1-\gamma)^2, 1\} \\
 &= \left[ \sum_{\tau=1}^{t-1} \mathbf{1}\{l_\tau = l, L_\tau + 1 \neq l\} \right] \min \{4^{-l} d/(1-\gamma)^2, 1\} \\
 &= (k_t^l - \frac{1}{2} k_t^{l-1} - 1) \min \{4^{l_0-l}, 1\}.
 \end{aligned}$$

On the other hand, for the upper bound of the right hand side of (B.1), notice that

$$\begin{aligned}
 \det(\Sigma_t^l) &\leq \left[ \frac{\text{tr}(\Sigma_t^l)}{d} \right]^d \\
 &= \left[ \frac{d\lambda + \sum_{\tau=1}^{t-1} \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l} \text{tr} [\phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top]}{d} \right]^d \\
 &\leq \left[ \frac{d\lambda + k_t^l \cdot d/(1-\gamma)^2}{d} \right]^d \\
 &= (\lambda + k_t^l/(1-\gamma)^2)^d,
 \end{aligned} \tag{B.2}$$

where the second inequality holds because  $\text{tr} [\phi_{V_t}(s_t, a_t) \phi_{V_t}(s_t, a_t)^\top] = \|\phi_{V_t}(s_t, a_t)\|_2^2 \leq d/(1-\gamma)^2$ . Substituting these bounds into (B.1), we have

$$\begin{aligned}
 (k_t^l - \frac{1}{2}k_t^{l-1} - 1) \min \{4^{l_0-l}, 1\} &\leq 2 \log \frac{(\lambda + k_t^l/(1-\gamma)^2)^d}{\lambda^d} \\
 &\leq 2d \log (1 + k_t^l/(1-\gamma)^2),
 \end{aligned} \tag{B.3}$$

where in the second inequality we used  $\lambda \geq 1$ .

The above inequality immediately implies that  $k_t^l$  is finite for arbitrary  $t, l$ . Take  $t \rightarrow \infty$ , we can replace the suffixes  $t$  in (B.3) by  $\infty$ . Next we consider two separate cases of  $l$ .

First consider the case  $l \leq l_0$ , where (B.3) becomes

$$k_{\max}^l - \frac{1}{2}k_{\max}^{l-1} - 1 \leq 2d \log (1 + k_{\max}^l/(1-\gamma)^2).$$

We will prove through induction on  $l$  that

$$k_{\max}^l \leq 20dl_0. \tag{B.4}$$

Note that (B.4) holds trivially for  $l = 0$ . Now suppose (B.4) already holds for  $l - 1$ . Then

$$k_{\max}^l - 2d \log (1 + k_{\max}^l/(1-\gamma)^2) \leq \frac{1}{2} \cdot 20dl_0 + 1.$$

By viewing  $k_{\max}^l$  as a variable and differentiating on it, we see the left hand side is increasing for  $k_{\max}^l > 2d$ . It is evident from (B.4) that the desired upper bound is indeed greater than  $2d$ , hence we only need to prove that the above inequality does not hold when  $k_{\max}^l$  is above this upper bound, in other words

$$20dl_0 - 2d \log (1 + 20dl_0/(1-\gamma)^2) \geq 10dl_0 + 1.$$

Substitute  $4^{l_0} = d/(1-\gamma)^2$  and move the first term on the right to the left, this is equivalent to

$$10dl_0 - 2d \log (1 + 20l_0 \cdot 4^{l_0}) \geq 1.$$

The left hand side

$$\begin{aligned}
 d[10l_0 - 2 \log (1 + 20l_0 \cdot 4^{l_0})] &\geq 10l_0 - 2 \log(1/4 + 20) - 2 \log (l_0 \cdot 4^{l_0}) \\
 &= 10l_0 - 2 \log(81/4) - 2(2 \log 2 \cdot l_0 + \log l_0) \\
 &= (10 - 4 \log 2)l_0 - 2 \log l_0 - 2 \log(81/4) \\
 &\geq 10 - 4 \log 2 - 2 \log(81/4) \\
 &> 1,
 \end{aligned}$$

where in the second inequality we shrank  $l_0$  to 1 based on monotony. This concludes the induction, and hence we have proven (B.4).

Now for  $l > l_0$ , the bound (B.3) becomes

$$(k_{\max}^l - \frac{1}{2}k_{\max}^{l-1} - 1) \cdot 4^{l_0-l} \leq 2d \log(1 + k_{\max}^l/(1-\gamma)^2). \quad (\text{B.5})$$

Write  $k_{\max}^l = \alpha_l d l 4^{l-l_0}$ . Notice that from (B.4),  $k_{\max}^{\lfloor l_0 \rfloor} \leq 20d(\lfloor l_0 \rfloor + 1) \cdot 4^{\lfloor l_0 \rfloor + 1 - l_0}$ . We further define  $\alpha_{\lfloor l_0 \rfloor} = 80$ , then plug these into (B.5), to obtain

$$\alpha_l d l - \frac{1}{8} \alpha_{l-1} d l - 4^{l_0-l} \leq 2d \log(1 + \alpha_l 4^l l).$$

This inequality holds for arbitrary  $l > l_0$ : for  $l = \lfloor l_0 \rfloor + 1$ , this is the exact result of substituting the bound for  $k_{\max}^{\lfloor l_0 \rfloor}$  into (B.5); for  $l > l_0 + 1$ , we further expanded  $l - 1$  on the left hand side to  $l$  to obtain this inequality. We then go on to divide both sides by  $dl$ :

$$\alpha_l - \frac{1}{8} \alpha_{l-1} - \frac{1}{dl} \leq \frac{2}{l} \log(1 + \alpha_l 4^l l). \quad (\text{B.6})$$

From here we only need to verify through induction that  $\alpha_l \leq 20$  for  $l > l_0$ , which yields  $k_t^l \leq 20dl4^{l-l_0}$  as required. For the induction basis, by shrinking  $d$  and  $l$  to 1, we have  $\alpha_{\lfloor l_0 \rfloor + 1} \leq 11 + 2 \log(1 + 4\alpha_{\lfloor l_0 \rfloor + 1})$ . Because the function  $x - 2 \log(1 + 4x)$  is increasing in  $x$  when  $x \geq 2$ , and  $20 > 11 + 2 \log(1 + 4 \times 20)$ , we can deduce from here  $\alpha_{\lfloor l_0 \rfloor + 1} \leq 20$ .

Suppose we already have  $\alpha_{l-1} \leq 20$  in (B.6). We can then obtain  $\alpha_l \leq 5/2 + 1/l + (2/l) \log(1 + \alpha_l 4^l l)$ . Combine with the trivial lower bound  $\alpha_l 4^l l \geq k_{\max}^l \geq 1$ , we have

$$\begin{aligned} \alpha_l &\leq \frac{5}{2} + \frac{1}{l} + \frac{2}{l} \log(2\alpha_l 4^l l) \\ &\leq \frac{7}{2} + 2 \log 2 + 2 \log \alpha_l + 2 \log 4 + \frac{2}{l} \log l \\ &\leq \left[ \frac{7}{2} + \log 64 + \frac{2}{3} \log 3 \right] + 2 \log \alpha_l \\ &< 9 + 2 \log \alpha_l. \end{aligned}$$

Therefore, in light of the fact that  $x - 2 \log x$  is increasing for  $x \geq 2$ , and that  $20 > 9 + 2 \log 20$ , we have  $\alpha_l \leq 20$ , which concludes the proof.  $\square$

## B.2. Proof of Lemma 6.2

*Proof of Lemma 6.2.* By Lemma A.1, we have the following expression for  $\widehat{\theta}_t^l$ :

$$\widehat{\theta}_t^l = \left( \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \right)^{-1} \left( \sum_{\tau=1}^{t-1} \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}) \right).$$

Next, we have from (3.3) that  $\langle \theta^*, \phi_{V_t}(s_t, a_t) \rangle = \mathbb{P}V_t(s_t, a_t) = \mathbb{E}_{s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)} [V_t(s_{t+1})]$ , which means  $\{V_t(s_{t+1}) - \langle \theta^*, \phi_{V_t}(s_t, a_t) \rangle\}$  forms a martingale difference sequence. Furthermore, since  $V_t(s)$  is bounded by  $1/(1-\gamma)$ , this sequence is  $1/(1-\gamma)$ -subgaussian. Combine this with the fact  $\|\theta^*\| \leq \sqrt{d}$  from Definition 3.1, we may deduce from Theorem 2 of Abbasi-Yadkori et al. (2011) that with probability at least  $\delta/(2l(l+1))$ , the following holds for all  $t \geq 0$ :

$$\begin{aligned} \|\widehat{\theta}_t^l - \theta^*\|_{\Sigma_t} &\leq \frac{1}{1-\gamma} \sqrt{2 \log \left[ \frac{\det(\Sigma_t^l)^{1/2} \det(\lambda \mathbf{I})^{-1/2}}{\delta/(2l(l+1))} \right]} + \sqrt{\lambda} \cdot \sqrt{d} \\ &\leq \frac{1}{1-\gamma} \sqrt{2 \log \left[ \frac{[\lambda + k_t^l/(1-\gamma)^2]^{d/2}}{\lambda^{d/2}} \cdot \frac{2l(l+1)}{\delta} \right]} + \sqrt{\lambda d} \\ &= \frac{1}{1-\gamma} \sqrt{d \log(1 + k_t^l/(1-\gamma)^2 \lambda) + 2 \log(2l(l+1)/\delta)} + \sqrt{\lambda d} \\ &\leq \frac{1}{1-\gamma} \sqrt{d \log(1 + 20 \cdot 4^{l'}) + 2 \log(2l'(l'+1)/\delta)} + \sqrt{\lambda d}, \end{aligned} \quad (\text{B.7})$$

where in the second inequality we borrowed (B.2) from the proof of Lemma 6.1, and in the final inequality we used the conclusion of Lemma 6.1 and substituted  $l' = \max\{l, l_0\}$ .

We now seek to prove the right hand side of (B.7) is no greater than half of  $\beta^l$ . First, for the expression inside the big square root, we have

$$\begin{aligned}
 & d \log(1 + 20 \cdot 4^l l) + 2 \log(2l(l+1)/\delta) \\
 & \leq d[\log(1 + 20 \cdot 4^l l) + 2 \log(2l(l+1))] + 2 \log(1/\delta) \\
 & \leq d[\log((1 + 1/80) \cdot 20 \cdot 4^l l) + \log 4 + 2 \log l + 2 \log(l+1)] + 2 \log(1/\delta) \\
 & = d[\log 81 + (\log 4) \cdot l + 3 \log l + 2 \log((l+1)/2) + \log 4] + 2 \log(1/\delta) \\
 & \leq d[\log(81 \cdot 4) + (\log 4) \cdot l + 3(l-1) + 2((l+1)/2 - 1)] + 2 \log(1/\delta) \\
 & = d[(4 + \log 4)l + (\log 324 - 4)] + 2 \log(1/\delta) \\
 & < 9dl + 2 \log(1/\delta),
 \end{aligned}$$

where we used the relation  $\log x \leq x - 1$  for real number  $x > 0$  twice in the third inequality.

Second, we take the square root in the above to obtain as a continuation of (B.7) that

$$\begin{aligned}
 \|\hat{\theta}_t^l - \theta^*\|_{\Sigma_t^l} & \leq \frac{1}{1-\gamma} \sqrt{9dl' + 2 \log(1/\delta)} + \sqrt{d\lambda} \\
 & \leq \frac{1}{1-\gamma} (3\sqrt{dl'} + 2\sqrt{\log(1/\delta)}) + \sqrt{d\lambda} \\
 & = \frac{1}{2} \beta^l.
 \end{aligned}$$

Now take a uniform bound, with probability at least  $1 - \sum_{l=1}^{\infty} \delta/(2l(l+1)) = 1 - \delta/2$ ,  $\|\theta^* - \hat{\theta}_t^l\|_{\Sigma_t^l} \leq \beta^l/2$  holds for all  $l \geq 1$ , which finishes the proof.  $\square$

### B.3. Proof of Lemma 6.3

*Proof of Lemma 6.3.* Focusing on one specific run of the ML-EVI algorithm, we use induction on  $u$  to prove that  $1/(1-\gamma) \geq Q_l^{(u)}(s, a) \geq Q^*(s, a)$ . When  $u = 0$ ,  $1/(1-\gamma) = Q_l^{(0)}(s, a) \geq Q^*(s, a)$ ,  $\forall s, a$ .

Suppose the conclusion holds for  $u - 1$ . First recall that under the event  $\mathcal{E}_1$  in Lemma 6.2,  $\theta^* \in \mathcal{B}$  and  $\mathcal{B} \cap \mathcal{C}_t^l \neq \emptyset$  for all  $t, l$ , so the special update rule in Line 8 is never called. Based on the update rule for  $V$  in Line 3 of Algorithm 2,  $V^{(u-1)}(s) = \max_a \min_l Q_l^{(u-1)}(s, a)$ , which right hand side is clearly upper bounded by  $1/(1-\gamma)$  and lower bounded by  $\max_a Q^*(s, a) = V^*(s)$ . Now combine this with the update rule of  $Q$  in Line 6, we see that

$$\begin{aligned}
 Q_l^{(u)}(s, a) & = r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_t^l} \langle \theta, \phi_{V^{(u-1)}}(s, a) \rangle \\
 & = r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_t^l} \int_{\mathcal{S}} V^{(u-1)}(s') \langle \theta, \phi(s'|s, a) \rangle ds' \\
 & \geq r(s, a) + \gamma \int_{\mathcal{S}} V^{(u-1)}(s') \langle \theta^*, \phi(s'|s, a) \rangle ds' \\
 & \geq r(s, a) + \gamma \int_{\mathcal{S}} V^*(s') \langle \theta^*, \phi(s'|s, a) \rangle ds' \\
 & = Q^*(s, a),
 \end{aligned}$$

where the second equation is based on the expression of  $\phi_V$  given in Definition 3.1, and the last equation is the Bellman equation for optimal value functions in (3.2), combined with  $\mathbb{P}(s'|s, a) = \langle \theta^*, \phi(s'|s, a) \rangle$  from Definition 3.1.

On the other hand, for the lower bound, the definition of  $\mathcal{B}$  in (3.4) tells us  $\langle \theta, \phi(\cdot|s, a) \rangle$  is a probability measure on  $\mathcal{S}$  for

arbitrary  $\theta \in \mathcal{B}$ , so there exists such a probability measure  $\tilde{\mathbb{P}}$  such that

$$\begin{aligned}
 Q_l^{(u)}(s, a) &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \int_{\mathcal{S}} V^{(u-1)}(s') \langle \theta, \phi(s'|s, a) \rangle ds' \\
 &= r(s, a) + \gamma \int_{\mathcal{S}} \tilde{\mathbb{P}}(s'|s, a) V^{(u-1)}(s') ds' \\
 &= r(s, a) + \gamma \mathbb{E}_{s' \sim \tilde{\mathbb{P}}(s'|s, a)} V^{(u-1)}(s') \\
 &\leq 1 + \gamma \cdot \frac{1}{1 - \gamma} \\
 &= \frac{1}{1 - \gamma},
 \end{aligned} \tag{B.8}$$

where in the inequality we used  $r(s, a) \leq 1$  and the proven conclusion  $V^{(u-1)} \leq 1/(1 - \gamma)$  from the induction hypothesis.

Since  $Q_t^l, V_t$  are the returns of Algorithm 2 at time step  $t$ , we have that  $Q_t^l(s, a)$  is the value function from the final iteration, and so  $1/(1 - \gamma) \geq Q_t^l(s, a) \geq Q^*(s, a)$ ; furthermore, from Line 12 of Algorithm 2, we have that  $1/(1 - \gamma) \geq V_t(s) = \max_a \min_{1 \leq l \leq L} Q_t^l(s, a) \geq \max_a Q^*(s, a) = V^*(s)$ . These are the desired results, and our proof is completed.  $\square$

#### B.4. Proof of Lemma 6.4

*Proof of Lemma 6.4.* Again focusing on one single run of ML-EVI, we first prove the following inequality:

$$\max_{s, a, l} |Q_l^{(u)}(s, a) - Q_l^{(u-1)}(s, a)| \leq \gamma^{u-1}, \forall u \in \{1, 2, \dots, U\}, \tag{B.9}$$

where in the maximization, the variables are taken from the following sets:  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ , and  $l \in \{1, 2, \dots, L\}$ .

We use induction on  $u$  for the proof. For  $u = 1$ , since  $Q_l^{(0)}(s, a) = 1/(1 - \gamma)$ , we have  $V^{(0)}(s) = \max_{a \in \mathcal{A}} \min_{1 \leq l \leq L} 1/(1 - \gamma) = 1/(1 - \gamma)$ , and hence from (B.8) we see

$$\begin{aligned}
 Q_l^{(1)}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim \tilde{\mathbb{P}}(s'|s, a)} V^{(0)}(s') \\
 &= r(s, a) + \frac{\gamma}{1 - \gamma} \\
 &\in \left[ \frac{\gamma}{1 - \gamma}, \frac{1}{1 - \gamma} \right],
 \end{aligned}$$

which leads to  $|Q_l^{(0)}(s, a) - Q_l^{(1)}(s, a)| \leq |1/(1 - \gamma) - \gamma/(1 - \gamma)| = 1$  for all  $s, a, l$ .

Now assume inequality (B.9) holds for  $u - 1$ . From the update rule of  $Q$  in Line 6, we have

$$\begin{aligned}
 Q_l^u(s, a) &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-1)}}(s, a) \rangle, \\
 Q_l^{u-1}(s, a) &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-2)}}(s, a) \rangle.
 \end{aligned}$$

We therefore obtain by subtracting the two equations that

$$\begin{aligned}
 |Q_l^u(s, a) - Q_l^{u-1}(s, a)| &= \gamma \left| \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-1)}}(s, a) \rangle - \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-2)}}(s, a) \rangle \right| \\
 &\leq \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} |\langle \theta, \phi_{V^{(u-1)}}(s, a) - \phi_{V^{(u-2)}}(s, a) \rangle| \\
 &= \gamma |\tilde{\mathbb{P}}(V^{(u-1)} - V^{(u-2)})(s, a)| \\
 &\leq \gamma \max_s |(V^{(u-1)} - V^{(u-2)})(s)|,
 \end{aligned} \tag{B.10}$$

where as before, we denote  $\tilde{\mathbb{P}} = \langle \tilde{\theta}, \phi_V \rangle$  as the probability measure corresponding to  $\tilde{\theta}$ , which attains the maximum in the second line above.

Next, the update rule for  $V$  in Line 3 of Algorithm 2 suggests that

$$\begin{aligned} V^{(u-1)}(s) &= \max_a \min_{1 \leq l \leq L} Q_l^{(u-1)}(s, a), \\ V^{(u-2)}(s) &= \max_a \min_{1 \leq l \leq L} Q_l^{(u-2)}(s, a). \end{aligned}$$

Plugging these into (B.10), we get

$$\begin{aligned} |Q_l^{(u)}(s, a) - Q_l^{(u-1)}(s, a)| &\leq \gamma \max_s \left| \max_a \min_{1 \leq l \leq L} Q_l^{(u-1)}(s, a) - \max_a \min_{1 \leq l \leq L} Q_l^{(u-2)}(s, a) \right| \\ &\leq \gamma \max_{s, a, l} |Q_l^{(u-1)}(s, a) - Q_l^{(u)}(s, a)| \\ &\leq \gamma \cdot \gamma^{u-2} \\ &= \gamma^{u-1}, \end{aligned}$$

where we used the induction hypothesis in the third inequality.

Now consider the result of Algorithm 2 at time step  $t$ , namely  $Q_t^l = Q_l^{(U)}$  and  $V_t = V^{(U)}$ . There again exists a probability measure  $\tilde{\mathbb{P}} = \langle \tilde{\boldsymbol{\theta}}, \phi \rangle$  to support the following:

$$\begin{aligned} Q_t^l(s_t, a_t) &= Q_l^{(U)}(s_t, a_t) \\ &= r(s_t, a_t) + \gamma \max_{\boldsymbol{\theta} \in \mathcal{B} \cap \mathcal{C}_l} \langle \boldsymbol{\theta}, \phi_{V^{(U-1)}}(s_t, a_t) \rangle \\ &= r(s_t, a_t) + \gamma \tilde{\mathbb{P}} V^{(U-1)}(s_t, a_t) \\ &= r(s_t, a_t) + \gamma [\tilde{\mathbb{P}} V_t(s_t, a_t) + \tilde{\mathbb{P}}(V^{(U-1)} - V^{(U)})(s_t, a_t)] \\ &\leq r(s_t, a_t) + \gamma \tilde{\mathbb{P}} V_t(s_t, a_t) + \gamma \max_{s \in \mathcal{S}} |V^{(U-1)} - V^{(U)}|(s) \\ &\leq r(s_t, a_t) + \gamma \max_{\boldsymbol{\theta} \in \mathcal{B} \cap \mathcal{C}_l} \langle \boldsymbol{\theta}, \phi_{V_t}(s_t, a_t) \rangle + \gamma \max_{s \in \mathcal{S}} \left| \max_a \min_l Q^{(U-1)}(s, a) - \max_a \min_l Q^{(U)}(s, a) \right| \\ &\leq r(s_t, a_t) + \gamma \max_{\boldsymbol{\theta} \in \mathcal{B} \cap \mathcal{C}_l} \langle \boldsymbol{\theta}, \phi_{V_t}(s_t, a_t) \rangle + \gamma \max_{s, a, l} |Q^{(U-1)}(s, a) - Q^{(U)}(s, a)| \\ &\leq r(s_t, a_t) + \gamma \max_{\boldsymbol{\theta} \in \mathcal{B} \cap \mathcal{C}_l} \langle \boldsymbol{\theta}, \phi_{V_t}(s_t, a_t) \rangle + \gamma^U, \end{aligned}$$

where in the last inequality we implemented (B.9). This concludes our proof of this lemma.  $\square$

## C. Proof of Lemmas in Appendix A

### C.1. Proof of Lemma A.1

*Proof of Lemma A.1.* We use induction on time step  $t$ . When  $t = 1$ , according to the initial values  $\boldsymbol{\Sigma}_t^l = \lambda \mathbf{I}$  and  $\mathbf{b}_t^l = \mathbf{0}$ , the alternative expressions hold.

Suppose the conclusion holds for  $t$ . We compare the variables at time steps  $t$  and  $t + 1$  of Algorithm 1, respectively after the execution of Line 7. According to the update rules in Lines 14 to 19, if  $l_t = L_t + 1$  and thus  $L_{t+1} = L_t + 1$ , new dormant levels for  $l = 4L_t + C_0 + 1, \dots, 4L_t + C_0 + 4$  are created. By the rule specified in Line 16, we use the induction hypothesis on  $(t, l)$  and  $(t, 4L_t + C_0)$ , to see that for  $l > 4L_t + C_0$ , the additional variables  $\boldsymbol{\Sigma}$  are

$$\begin{aligned} \boldsymbol{\Sigma}^l - \lambda \mathbf{I} &= \frac{1}{2} (\boldsymbol{\Sigma}^{l-1} - \lambda \mathbf{I}) \\ &= (1/2)^{l-(4L_t+C_0)} (\boldsymbol{\Sigma}_t^{4L_t+C_0} - \lambda \mathbf{I}) \\ &= 2^{(4L_t+C_0)-l} \sum_{\tau=1}^{t-1} \mathbb{1}(l_\tau \leq 4L_t + C_0) 2^{l_\tau - (4L_t+C_0)} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\ &= \sum_{\tau=1}^{t-1} \mathbb{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\ &= \boldsymbol{\Sigma}_t^l - \lambda \mathbf{I}, \end{aligned}$$

where in the second to last equation we used the fact  $l_\tau \leq L_t + 1 < 4L_t + C_0$  for arbitrary  $\tau < t$ . Similarly, for the new variables  $\mathbf{b}$ ,

$$\begin{aligned}
 \mathbf{b}^l &= \frac{1}{2} \mathbf{b}^{l-1} \\
 &= (1/2)^{l-(4L_t+C_0)} \mathbf{b}_t^{4L_t+C_0} \\
 &= 2^{(4L_t+C_0)-l} \sum_{\tau=1}^{t-1} \mathbf{1}(l_\tau \leq 4L_t + C_0) 2^{l_\tau - (4L_t+C_0)} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}) \\
 &= \sum_{\tau=1}^{t-1} \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}) \\
 &= \mathbf{b}_t^l.
 \end{aligned}$$

In other words, these variables for the newly created dormant levels are the same as the respective extended variables at time step  $t$ , and so this update rule can be ignored since we are using the alternative expressions for  $\Sigma_t^l$  and  $\mathbf{b}_t^l$  for all  $t, l$ .

Now we go on to the update rules in Lines 20 to 23. For levels  $l < l_t$ , the two variables of interest do not change, hence

$$\begin{aligned}
 \Sigma_{t+1}^l &= \Sigma_t^l \\
 &= \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top \\
 &= \lambda \mathbf{I} + \sum_{\tau=1}^t \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top,
 \end{aligned}$$

and

$$\begin{aligned}
 \mathbf{b}_{t+1}^l &= \mathbf{b}_t^l \\
 &= \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}) \\
 &= \lambda \mathbf{I} + \sum_{\tau=1}^t \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}).
 \end{aligned}$$

For levels  $l \geq l_t$ , a single term is added to the two variables respectively. Namely,

$$\begin{aligned}
 \Sigma_{t+1}^l &= \Sigma_t^l + 2^{l_t - l} \phi_{V_t}(s_t, a_t) \phi_{V_t}(s_t, a_t)^\top \\
 &= \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top + 2^{l_t - l} \phi_{V_t}(s_t, a_t) \phi_{V_t}(s_t, a_t)^\top \\
 &= \lambda \mathbf{I} + \sum_{\tau=1}^t \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) \phi_{V_\tau}(s_\tau, a_\tau)^\top,
 \end{aligned}$$

and that

$$\begin{aligned}
 \mathbf{b}_{t+1}^l &= \mathbf{b}_t^l + 2^{l_t - l} \phi_{V_t}(s_t, a_t) V_t(s_{t+1}) \\
 &= \sum_{\tau=1}^{t-1} \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}) + 2^{l_t - l} \phi_{V_t}(s_t, a_t) V_t(s_{t+1}) \\
 &= \sum_{\tau=1}^t \mathbf{1}(l_\tau \leq l) 2^{l_\tau - l} \phi_{V_\tau}(s_\tau, a_\tau) V_\tau(s_{\tau+1}).
 \end{aligned}$$

These expressions suggest the conclusion holds for arbitrary  $l \leq 4L_t + C_0$  at time step  $t$ .  $\square$

### C.2. Proof of Lemma A.2

*Proof of Lemma A.2.* We use induction on  $u$  for this proof. For  $u = 0$ , by the initialization rules  $Q_{l_1,*}^{(0)} = Q_{l_2}^{(0)} = 1/(1 - \gamma)$ ,  $\forall l_1 \in \{1, 2, \dots\}, l_2 \in \{1, \dots, 4L_t + C_0\}$ , and hence  $V_*^{(0)} = V^{(0)} = 1/(1 - \gamma)$ .

Suppose the conclusions already hold for  $u - 1$ . Consider the conclusions of Lemma D.2, which states the confidence sets of levels above  $4L_t + C_0$  trivially contain the entirety of  $\mathcal{B}$ , or  $\mathcal{B} \cap \mathcal{C}_l = \mathcal{B} \supset \mathcal{B} \cap \mathcal{C}_{4L_t+C_0}$ . Then for any level  $l > 4L_t + C_0$  and state-action pair  $(s, a)$ , by the update rule for  $Q$  in Line 6 of Algorithm 2:

$$\begin{aligned} Q_{l,*}^{(u)}(s, a) &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V_*^{(u-1)}}(s, a) \rangle \\ &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-1)}}(s, a) \rangle \\ &\geq r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_{4L_t+C_0}} \langle \theta, \phi_{V^{(u-1)}}(s, a) \rangle \\ &= Q_{4L_t+C_0}^{(u)}(s, a). \end{aligned}$$

Moreover, for any  $l \leq 4L_t + C_0$ , we have

$$\begin{aligned} Q_{l,*}^{(u)}(s, a) &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V_*^{(u-1)}}(s, a) \rangle \\ &= r(s, a) + \gamma \max_{\theta \in \mathcal{B} \cap \mathcal{C}_l} \langle \theta, \phi_{V^{(u-1)}}(s, a) \rangle \\ &= Q_l^{(u)}(s, a) \end{aligned}$$

We now give the relation for the value functions  $V$ . For any state  $s$ , by the update rule in Line 3 of Algorithm 2,

$$\begin{aligned} V_*^{(u)}(s) &= \max_{a \in \mathcal{A}} \min_{l \geq 1} Q_{l,*}^{(u)}(s, a) \\ &= \max_{a \in \mathcal{A}} \min \left\{ \min_{1 \leq l \leq 4L_t+C_0} Q_{l,*}^{(u)}(s, a), \min_{l > 4L_t+C_0} Q_{l,*}^{(u)}(s, a) \right\} \\ &= \max_{a \in \mathcal{A}} \min \left\{ \min_{1 \leq l \leq 4L_t+C_0} Q_l^{(u)}(s, a), \min_{l > 4L_t+C_0} Q_{l,*}^{(u)}(s, a) \right\} \\ &= \max_{a \in \mathcal{A}} \min_{1 \leq l \leq 4L_t+C_0} Q_l^{(u)}(s, a) \\ &= V^{(u)}(s), \end{aligned}$$

where we used the relations above in the third and fourth equations. This concludes the induction and our proof.  $\square$

### C.3. Proof of Lemma A.3

*Proof of Lemma A.3.* Consider the definition in (3.4):

$$\mathcal{B} := \{ \theta \in \mathbb{R}^d : \|\theta\|_2 \leq \sqrt{d} \text{ and } \langle \phi(\cdot|s, a), \theta \rangle \text{ is a probability measure on } \mathcal{S} \}.$$

We separate  $\mathcal{B}$  into  $\mathcal{B}_1 \cap \mathcal{B}_2$ , where  $\mathcal{B}_1 := \{ \theta : \|\theta\|_2 \leq \sqrt{d} \}$  and  $\mathcal{B}_2 := \{ \theta : \langle \phi(\cdot|s, a), \theta \rangle \text{ is a probability measure on } \mathcal{S} \}$ . Obviously  $\mathcal{B}_1$  is a convex body in  $\mathbb{R}^d$ . For  $\mathcal{B}_2$ , consider the various constraints. For an arbitrary state-action pair  $(s, a)$ ,  $\langle \phi(\cdot|s, a), \theta \rangle$  being a probability measure is equivalent to the following conditions:

$$\begin{aligned} \langle \phi(s'|s, a), \theta \rangle &\geq 0, \forall s' \in \mathcal{S}, \\ \int_{s' \in \mathcal{S}} \langle \phi(s'|s, a), \theta \rangle ds' &= 1. \end{aligned}$$

The first constraint defines a closed half-space  $\mathcal{H}_{s'|s,a} := \{ \theta : \langle \phi(s'|s, a), \theta \rangle \geq 0 \}$  in  $\mathbb{R}^d$  for all states  $s' \in \mathcal{S}$ , while the second defines a hyperplane  $\mathcal{P}_{s,a} := \{ \theta : \langle \int_{s' \in \mathcal{S}} \phi(s'|s, a) ds', \theta \rangle = 1 \}$  in  $\mathbb{R}^d$ . With these notations, we can now write

$$\mathcal{B}_2 = \left( \bigcap_{s', s, a} \mathcal{H}_{s'|s,a} \right) \cap \left( \bigcap_{s, a} \mathcal{P}_{s,a} \right).$$



Since the sets  $\mathcal{H}_{s'|s,a}$  and  $\mathcal{P}_{s,a}$  are all closed convex sets,  $\mathcal{B}_2$  being the intersection of these sets must also be closed and convex, and finally so must the intersection of all the sets  $\mathcal{B}$ . Since  $\theta^* \in \mathcal{B}$ , this intersection is non-empty.

Now consider the affine subspace  $\mathcal{L} := \{\theta^* + \lambda(\theta_1 - \theta_2) : \theta_1, \theta_2 \in \mathcal{B}, \lambda \in \mathbb{R}\}$  of  $\mathbb{R}^d$  with dimension  $k$ , which is basically the subspace passing through  $\theta^*$  and spanned by all the vectors insider  $\mathcal{B}$ . To see that  $\mathcal{L}$  is indeed an affine subspace of  $\mathbb{R}^d$ , we take  $\kappa_1, \kappa_2 \in \mathcal{L}$ , where  $\kappa_1 = \theta^* + \lambda_1(\theta_1 - \theta_2)$  and  $\kappa_2 = \theta^* + \lambda_2(\theta_3 - \theta_4)$ ,  $\theta_i \in \mathcal{B}$ ,  $i = 1, 2, 3, 4$ . Then for any  $\mu \in \mathbb{R}$ ,

$$\begin{aligned} \mu\kappa_1 + (1 - \mu)\kappa_2 &= \theta^* + (\lambda_1\mu\theta_1 - \lambda_1\mu\theta_2 + \lambda_2(1 - \mu)\theta_3 - \lambda_2(1 - \mu)\theta_4) \\ &= \theta^* + (|\lambda_1\mu| + |\lambda_2(1 - \mu)|) \cdot [(\mu'\theta'_1 + (1 - \mu')\theta'_3) - (\mu'\theta'_2 + (1 - \mu')\theta'_4)] \\ &\in \mathcal{L}, \end{aligned}$$

where  $\mu' = |\lambda_1\mu| / (|\lambda_1\mu| + |\lambda_2(1 - \mu)|)$ , and we used the fact  $\mathcal{B}$  is convex in the final relation.

By taking  $\lambda = 1$  and  $\theta_2 = \theta^*$  in the definition of  $\mathcal{B}$ , it is clear that  $\mathcal{B} \subset \mathcal{L}$ . On the other hand, this subspace can be generated by the vectors  $\{\theta - \theta^* : \theta \in \mathcal{B}\}$  starting from  $\theta^*$ . From this vector set, we select a basis for the subspace, namely  $\{\theta_i^* - \theta^* : i = 1, \dots, k\}$ , and denote  $\theta_0^* = \theta^*$ .

Now consider the simplex  $\{\sum_{i=0}^k \lambda_i \theta_i^* : \lambda_i \geq 0, \forall i, \sum_{i=0}^k \lambda_i = 1\}$ , which, since it has  $k + 1$  linearly independent vertices, is  $k$ -dimensional with inner points when constrained to  $\mathcal{L}$ . Furthermore, because  $\mathcal{B}$  is convex and contains all its vertices, this simplex must be a subset of  $\mathcal{B}$ , so  $\mathcal{B}$  itself has inner points when constrained to  $\mathcal{L}$ . Combine this with the fact that  $\mathcal{B}_1$  is bounded, and so must  $\mathcal{B}$  be, we have the conclusion that  $\mathcal{B}$  is a convex body when constrained to  $\mathcal{L}$ .  $\square$

## D. Auxiliary Lemmas

We will give a few auxiliary lemmas in this section. The first is the well-known Azuma-Hoeffding inequality.

**Lemma D.1.** (Azuma-Hoeffding Inequality (Cesa-Bianchi & Lugosi, 2006)). Let  $\{x_i\}_{i=1}^n$  be a real-valued martingale difference sequence with respect to the filtration  $\{\mathcal{F}_i\}_{i=1}^n$ , which suggests  $\mathbb{E}[x_i | \mathcal{F}_i] = 0$  and  $x_i$  is  $\mathcal{F}_{i+1}$ -measurable. Further assume  $|x_i| \leq K$  for some positive constant  $K$ . Then with probability at least  $1 - \delta$ , we have

$$\left| \sum_{i=1}^n x_i \right| \leq K \sqrt{2n \log \frac{1}{\delta}}.$$

This next lemma proves that the fictional confidence sets with levels greater than  $4L_t + C_0$  are trivial.

**Lemma D.2.** For time step  $t$  and level  $l > 4L_t + C_0$ , we have  $\mathcal{C}_t^l \supset \mathcal{B}$ .

*Proof of Lemma D.2.* We first give an upper bound for the term  $k_t^l$  when  $l > 4L_t + C_0$ . The definition of  $k_t^l$  tells us  $k_t^l = \sum_{\tau=1}^{t-1} \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - l} = 2^{L_t - l} \sum_{\tau=1}^{t-1} \mathbb{1}\{l_\tau \leq l\} 2^{l_\tau - L_t} = 2^{L_t - l} k_t^{L_t}$ . Recall from Algorithm 1 that  $C_0 = 4\lceil l_0 \rceil + 6$ , we go on to apply the conclusion of Lemma 6.1:

$$\begin{aligned} k_t^l &\leq 2^{L_t - l} \cdot 20d \max\{L_t 2^{2L_t - 2l_0}, l_0\} \\ &= 2^{L_t - l} \cdot 20(1 - \gamma)^2 \max\{L_t 2^{2L_t}, l_0 2^{2l_0}\} \\ &\leq 20(1 - \gamma)^2 \max\{L_t 2^{-L_t - C_0}, l_0 2^{-3L_t + 2l_0 - C_0}\} \\ &\leq 20/64 \cdot 2^{-l_0} \max\{L_t 2^{-L_t}, l_0 2^{-l_0}\} \\ &\leq 5/16 \cdot (1 - \gamma) / \sqrt{d} \\ &< \log 1.5 \cdot (1 - \gamma) / \sqrt{d}, \end{aligned}$$

where we used the definition for  $l_0$  in Lemma 6.1 twice above, and used the trivial relation  $2^x \geq x$  for arbitrary real number  $x$  in the third inequality.

Now we target the maximum eigenvalue of  $\Sigma_t^l$  when  $l > 4L_t + C_0$ . Borrowing from the determinant upper bound (B.2)

from the proof of Lemma 6.1, we have

$$\begin{aligned}
 \lambda_{\max}(\boldsymbol{\Sigma}_t^l) &\leq \frac{\det(\boldsymbol{\Sigma}_t^l)}{\lambda^{d-1}} \\
 &\leq \frac{(\lambda + \lambda k_t^l/d)^d}{\lambda^{d-1}} \\
 &= \lambda(1 + k_t^l/d)^d \\
 &\leq \lambda \exp k_t^l \\
 &\leq 1.5\lambda,
 \end{aligned}$$

where in the first inequality we used the fact all eigenvalues of  $\boldsymbol{\Sigma}_t^l$  are no less than  $\lambda$ , and in the third inequality we used the elementary inequality  $1 + x \leq \exp(x)$ .

Next we upper bound the 2-norm of  $\mathbf{b}_t^l$  when  $l > 4L_t + C_0$ . By the alternative expression in Lemma A.1:

$$\begin{aligned}
 \|\mathbf{b}_t^l\|_2 &\leq \sum_{\tau=1}^{t-1} \mathbb{1}(l \geq l_\tau) 2^{l_\tau - l} \|\phi_{V_\tau}(s_\tau, a_\tau)\|_2 \cdot |V_\tau(s_{\tau+1})| \\
 &\leq k_t^l \cdot \frac{\sqrt{d}}{1-\gamma} \cdot \frac{1}{1-\gamma} \\
 &\leq \frac{\log 1.5}{1-\gamma},
 \end{aligned}$$

where in the second inequality we used  $\|\phi_{V_\tau}(s_\tau, a_\tau)\|_2 \leq \sqrt{d}/(1-\gamma)$  from Definition 3.1 and  $V_\tau(s_{\tau+1}) \leq 1/(1-\gamma)$  from Lemma 6.3.

Now for any  $\boldsymbol{\theta} \in \mathcal{B}$ , we substitute these two bounds in the following deductions:

$$\begin{aligned}
 \|\boldsymbol{\theta} - (\boldsymbol{\Sigma}_t^l)^{-1} \mathbf{b}_t^l\|_{\boldsymbol{\Sigma}_t^l} &\leq \|\boldsymbol{\theta}\|_{\boldsymbol{\Sigma}_t^l} + \|(\boldsymbol{\Sigma}_t^l)^{-1} \mathbf{b}_t^l\|_{\boldsymbol{\Sigma}_t^l} \\
 &= \|\boldsymbol{\theta}\|_{\boldsymbol{\Sigma}_t^l} + \|\mathbf{b}_t^l\|_{(\boldsymbol{\Sigma}_t^l)^{-1}} \\
 &\leq \|\boldsymbol{\theta}\|_2 \cdot \sqrt{\lambda_{\max}(\boldsymbol{\Sigma}_t^l)} + \|\mathbf{b}_t^l\|_2 \cdot \frac{1}{\sqrt{\lambda_{\min}(\boldsymbol{\Sigma}_t^l)}} \\
 &\leq \sqrt{d} \cdot \sqrt{1.5\lambda} + \frac{\log 1.5}{1-\gamma} \cdot \frac{1}{\sqrt{\lambda}} \\
 &= \sqrt{1.5d\lambda} + \frac{\log 1.5}{1-\gamma} \\
 &< \beta^l,
 \end{aligned}$$

where in the second inequality we used the relation  $\|x\|_{\boldsymbol{\Sigma}}^2 = x^\top \boldsymbol{\Sigma} x \leq x^\top [\lambda_{\max}(\boldsymbol{\Sigma}) \mathbf{I}] x = \|x\|_2^2 \cdot \lambda_{\max}(\boldsymbol{\Sigma})$ , and in the third inequality we used  $\|\boldsymbol{\theta}\|_2 \leq \sqrt{d}$  for any  $\boldsymbol{\theta} \in \mathcal{B}$ . Thus when  $l > 4L_t + C_0$ ,  $\boldsymbol{\theta} \in \mathcal{C}_t^l$  for any  $\boldsymbol{\theta} \in \mathcal{B}$ , and hence  $\mathcal{B} \subset \mathcal{C}_t^l$ .  $\square$

Finally, we list a few useful elementary inequalities and relations that were implemented in our proof of the main theory, and gather them into one single lemma below.

**Lemma D.3.** The following conclusions hold:

1. For  $\eta \geq 2$  and integer  $L > 0$ , we have  $\sum_{l=L}^{\infty} \eta^{-l} \sqrt{l} \leq \sqrt{L+1} \eta^{-L+1} / (\eta - 1)$ .
2. For any positive real numbers  $a, b$  and  $x$ , the inequality  $x - a\sqrt{x} \leq b$  implies the inequality  $x \leq a^2 + 2b$ .
3. For  $a, b \geq 1$ , any real number  $x \geq \log(a + b\sqrt{2\log(a+b)/\log 2}) / \log 2$  satisfies the equation  $2^x \geq a + b\sqrt{x}$ .

*Proof of Lemma D.3.* For the first inequality:

$$\begin{aligned}
 \sum_{l=L}^{\infty} \eta^{-l} \sqrt{l} &= \sum_{l=L}^{\infty} \eta^{-l} \left( \sqrt{L} + \sum_{l'=L+1}^l (\sqrt{l'} - \sqrt{l'-1}) \right) \\
 &= \sqrt{L} \sum_{l=L}^{\infty} \eta^{-l} + \sum_{l'=L+1}^{\infty} (\sqrt{l'} - \sqrt{l'-1}) \sum_{l=l'}^{\infty} \eta^{-l} \\
 &= \sqrt{L} \frac{\eta^{-L+1}}{\eta-1} + \sum_{l'=L+1}^{\infty} (\sqrt{l'} - \sqrt{l'-1}) \frac{\eta^{-l'+1}}{\eta-1} \\
 &\leq \sqrt{L} \frac{\eta^{-L+1}}{\eta-1} + \sum_{l'=L+1}^{\infty} (\sqrt{L+1} - \sqrt{L}) \frac{\eta^{-l'+1}}{\eta-1} \\
 &= \sqrt{L} \frac{\eta^{-L+1}}{\eta-1} + (\sqrt{L+1} - \sqrt{L}) \frac{\eta^{-L+1}}{(\eta-1)^2} \\
 &\leq \sqrt{L+1} \eta^{-L+1} / (\eta-1),
 \end{aligned}$$

where we swapped the order of summation in the second equation, used the fact that  $\sqrt{l} - \sqrt{l-1}$  is decreasing in  $l$  in the first inequality, and  $\eta - 1 \geq 1$  in the last inequality.

For the second relation, suppose positive real numbers  $a, b, x$  satisfy  $x - a\sqrt{x} \leq b$ . Formulating the left hand side to a squared expression, we get  $(\sqrt{x} - a/2)^2 \leq a^2/4 + b$ , and further  $\sqrt{x} \leq a/2 + \sqrt{a^2/4 + b}$ . This leads to

$$\begin{aligned}
 x &\leq (a/2 + \sqrt{a^2/4 + b})^2 \\
 &\leq 2 \cdot \left( \frac{a^2}{4} + \frac{a^2}{4} + b \right) \\
 &= a^2 + 2b,
 \end{aligned}$$

where we used the inequality  $(x + y)^2 \leq 2(x^2 + y^2)$  in the inequality.

For the third relation, we first substitute  $x_0 := \log(a + b\sqrt{2\log(a+b)/\log 2}) / \log 2$  into  $2^x - a - b\sqrt{x}$ , and obtain

$$\begin{aligned}
 2^{x_0} - a - b\sqrt{x_0} &= a + b\sqrt{2\log(a+b)/\log 2} - a - b\sqrt{x_0} \\
 &= b \frac{\log(a+b)^2 / \log 2 - x_0}{\sqrt{\log(a+b)^2 / \log 2} + \sqrt{x_0}} \\
 &= G(a, b) \log \frac{(a+b)^2}{a + b\sqrt{2\log(a+b)/\log 2}} / \log 2,
 \end{aligned}$$

where we gathered everything besides the numerator in the second line into  $G(a, b)$  for the final equation.  $G(a, b)$  as a function in  $a, b$  is evidently always positive. After this we can see  $a + b\sqrt{2\log(a+b)/\log 2} \leq a + b \cdot \sqrt{2(a+b)} \leq (a+b)^2$  since  $a, b \geq 1$ , so  $2^{x_0} - a - b\sqrt{x_0} \geq 0$ . Now we only need to prove for  $x \geq x_0$ ,  $d(2^x - a - b\sqrt{x})/dx \geq 0$ , which is guaranteed since

$$\begin{aligned}
 \frac{d(2^x - a - b\sqrt{x})}{dx} &= \log 2 \cdot 2^x - \frac{b}{2\sqrt{x}} \\
 &> \log 2 \cdot (2^{x_0} - b\sqrt{x_0}) \\
 &> 0,
 \end{aligned}$$

where we used the fact  $x \geq x_0 \geq \log(1+1)/\log 2 = 1$  and  $2\log 2 > 1$  in the first inequality. This concludes our proof.  $\square$